

Unicode Demystified

*A Practical Programmer's
Guide to the Encoding Standard*

by Richard Gillam

Copyright ©2000–2002 by Richard T. Gillam. All rights reserved.

Pre-publication draft number 0.3.1

Tuesday, January 15, 2002

*To Mark, Kathleen, Laura, Helena, Doug,
John F, John R, Markus, Bertrand, Alan, Eric,
and the rest of the old JTCSV Unicode crew,
without whom this book would not have been possible*

and

*To Ted and Joyce Gillam,
without whom the author would not have been possible*

Table of Contents

Table of Contents v

Preface xv

About this book xvi

How this book was produced xviii

The author's journey xviii

Acknowledgements xix

A personal appeal xx

Unicode in Essence An Architectural Overview of the Unicode Standard 1

CHAPTER 1 Language, Computers, and Unicode 3

What Unicode Is 6

What Unicode Isn't 8

The challenge of representing text in computers 10

What This Book Does 14

How this book is organized 15

Section I: Unicode in Essence 15

Section II: Unicode in Depth 16

Section III: Unicode in Action 17

CHAPTER 2 A Brief History of Character Encoding 19

Prehistory 19

The telegraph and Morse code 20

The teletypewriter and Baudot code 21
Other teletype and telegraphy codes 22
FIELDATA and ASCII 23
Hollerith and EBCDIC 24

Single-byte encoding systems 26

Eight-bit encoding schemes and the ISO 2022 model 27
ISO 8859 28
Other 8-bit encoding schemes 29

Character encoding terminology 30

Multiple-byte encoding systems 32

East Asian coded character sets 32
Character encoding schemes for East Asian coded character sets 33
Other East Asian encoding systems 36

ISO 10646 and Unicode 36

How the Unicode standard is maintained 41

CHAPTER 3 Architecture: Not Just a Pile of Code Charts 43

The Unicode Character-Glyph Model 44

Character positioning 47

The Principle of Unification 50

Alternate-glyph selection 53

Multiple Representations 54

Flavors of Unicode 56

Character Semantics 58

Unicode Versions and Unicode Technical Reports 60

Unicode Standard Annexes 60
Unicode Technical Standards 61
Unicode Technical Reports 61
Draft and Proposed Draft Technical Reports 61
Superseded Technical Reports 62
Unicode Versions 62
Unicode stability policies 63

Arrangement of the encoding space 64

Organization of the planes 64
The Basic Multilingual Plane 66
The Supplementary Planes 69
Non-Character code point values 72

Conforming to the standard 73

General 74
Producing text as output 75
Interpreting text from the outside world 75
Passing text through 76
Drawing text on the screen or other output devices 76
Comparing character strings 77
Summary 77

CHAPTER 4 Combining character sequences and Unicode normalization 79

How Unicode non-spacing marks work 81

Dealing properly with combining character sequences 83

- Canonical decompositions 84**
- Canonical accent ordering 85**
- Double diacritics 87**
- Compatibility decompositions 88**
- Singleton decompositions 90**
- Hangul 91**
- Unicode normalization forms 93**
- Grapheme clusters 94**

CHAPTER 5 Character Properties and the Unicode Character Database 99

Where to get the Unicode Character Database 99

The UNIDATA directory 100

UnicodeData.txt 103

PropList.txt 105

General character properties 107

Standard character names 107

Algorithmically-derived names 108

Control-character names 109

ISO 10646 comment 109

Block and Script 110

General Category 110

Letters 110

Marks 112

Numbers 112

Punctuation 113

Symbols 114

Separators 114

Miscellaneous 114

Other categories 115

Properties of letters 117

SpecialCasing.txt 117

CaseFolding.txt 119

Properties of digits, numerals, and mathematical symbols 119

Layout-related properties 120

Bidirectional layout 120

Mirroring 121

Arabic contextual shaping 122

East Asian width 122

Line breaking property 123

Normalization-related properties 124

Decomposition 124

Decomposition type 124

Combining class 126

Composition exclusion list 127

Normalization test file 127

Derived normalization properties 128

Grapheme-cluster-related properties 128

UniHan.txt 129

CHAPTER 6 Unicode Storage and Serialization Formats 131

- A historical note 132
- UTF-32 133
- UTF-16 and the surrogate mechanism 134
- Endian-ness and the Byte Order Mark 136
- UTF-8 138
- CESU-8 141
- UTF-EBCDIC 141
- UTF-7 143
- Standard Compression Scheme for Unicode 143
- BOCU 146
- Detecting Unicode storage formats 147

**Unicode in Depth A Guided Tour of the
Character Repertoire 149**

CHAPTER 7 Scripts of Europe 151

The Western alphabetic scripts 151

The Latin alphabet 153

- The Latin-1 characters 155*
- The Latin Extended A block 155*
- The Latin Extended B block 157*
- The Latin Extended Additional block 158*
- The International Phonetic Alphabet 159*

Diacritical marks 160

- Isolated combining marks 164*
- Spacing modifier letters 165*

The Greek alphabet 166

- The Greek block 168*
- The Greek Extended block 169*
- The Coptic alphabet 169*

The Cyrillic alphabet 170

- The Cyrillic block 173*
- The Cyrillic Supplementary block 173*

The Armenian alphabet 174

The Georgian alphabet 175

CHAPTER 8 Scripts of The Middle East 177

Bidirectional Text Layout 178

The Unicode Bidirectional Layout Algorithm 181

- Inherent directionality 181*
- Neutrals 184*
- Numbers 185*
- The Left-to-Right and Right-to-Left Marks 186*
- The Explicit Embedding Characters 187*

Mirroring characters 188

Line and Paragraph Boundaries 188

Bidirectional Text in a Text-Editing Environment 189

The Hebrew Alphabet 192

The Hebrew block 194

The Arabic Alphabet 194

The Arabic block 199

Joiners and non-joiners 199

The Arabic Presentation Forms B block 201

The Arabic Presentation Forms A block 202

The Syriac Alphabet 202

The Syriac block 204

The Thaana Script 205

The Thaana block 207

CHAPTER 9 Scripts of India and Southeast Asia 209

Devanagari 212

The Devanagari block 217

Bengali 221

The Bengali block 223

Gurmukhi 223

The Gurmukhi block 225

Gujarati 225

The Gujarati block 226

Oriya 226

The Oriya block 227

Tamil 227

The Tamil block 230

Telugu 230

The Telugu block 232

Kannada 232

The Kannada block 233

Malayalam 234

The Malayalam block 235

Sinhala 235

The Sinhala block 236

Thai 237

The Thai block 238

Lao 239

The Lao block 240

Khmer 241

The Khmer block 243

Myanmar 243

The Myanmar block 244

Tibetan 245

The Tibetan block 247

The Philippine Scripts 247

CHAPTER 10 Scripts of East Asia 251

The Han characters 252

Variant forms of Han characters 261

Han characters in Unicode 263

- The CJK Unified Ideographs area 267*
- The CJK Unified Ideographs Extension A area 267*
- The CJK Unified Ideographs Extension B area 267*
- The CJK Compatibility Ideographs block 268*
- The CJK Compatibility Ideographs Supplement block 268*
- The Kangxi Radicals block 268*
- The CJK Radicals Supplement block 269*

Indeographic description sequences 269

Bopomofo 274

- The Bopomofo block 275*
- The Bopomofo Extended block 275*

Japanese 275

- The Hiragana block 281*
- The Katakana block 281*
- The Katakana Phonetic Extensions block 281*
- The Kanbun block 281*

Korean 282

- The Hangul Jamo block 284*
- The Hangul Compatibility Jamo block 285*
- The Hangul Syllables area 285*

Halfwidth and fullwidth characters 286

- The Halfwidth and Fullwidth Forms block 288*

Vertical text layout 288

Ruby 292

- The Interlinear Annotation characters 293*

Yi 294

- The Yi Syllables block 295*
- The Yi Radicals block 295*

CHAPTER 11 Scripts from Other Parts of the World 297

Mongolian 298

- The Mongolian block 300*

Ethiopic 301

- The Ethiopic block 303*

Cherokee 303

- The Cherokee block 304*

Canadian Aboriginal Syllables 304

- The Unified Canadian Aboriginal Syllabics block 305*

Historical scripts 305

- Runic 306*
- Ogham 307*
- Old Italic 307*
- Gothic 308*
- Deseret 309*

CHAPTER 12 Numbers, Punctuation, Symbols, and Specials 311

Numbers 311

Western positional notation 312
Alphabetic numerals 313
Roman numerals 313
Han characters as numerals 314
Other numeration systems 317
Numeric presentation forms 319
National and nominal digit shapes 319

Punctuation 320

Script-specific punctuation 320
The General Punctuation block 322
The CJK Symbols and Punctuation block 323
Spaces 323
Dashes and hyphens 325
Quotation marks, apostrophes, and similar-looking characters 326
Paired punctuation 331
Dot leaders 332
Bullets and dots 332

Special characters 333

Line and paragraph separators 333
Segment and page separators 335
Control characters 336
Characters that control word wrapping 336
Characters that control glyph selection 339
The grapheme joiner 345
Bidirectional formatting characters 346
Deprecated characters 346
Interlinear annotation 347
The object-replacement character 348
The general substitution character 348
Tagging characters 349
Non-characters 351

Symbols used with numbers 351

Numeric punctuation 351
Currency symbols 352
Unit markers 353
Math symbols 353
Mathematical alphanumeric symbols 356

Other symbols and miscellaneous characters 357

Musical notation 357
Braille 359
Other symbols 359
Presentation forms 360
Miscellaneous characters 361

Unicode in Action Implementing and Using the Unicode Standard 363

CHAPTER 13 Techniques and Data Structures for Handling Unicode Text 365

Useful data structures 366

Testing for membership in a class 366

The inversion list 369

Performing set operations on inversion lists 370

Mapping single characters to other values 374

Inversion maps 375

The compact array 376

Two-level compact arrays 379

Mapping single characters to multiple values 380

Exception tables 381

Mapping multiple characters to other values 382

Exception tables and key closure 382

Tries as exception tables 385

Tries as the main lookup table 388

Single versus multiple tables 390

CHAPTER 14 Conversions and Transformations 393

Converting between Unicode encoding forms 394

Converting between UTF-16 and UTF-32 395

Converting between UTF-8 and UTF-32 397

Converting between UTF-8 and UTF-16 401

Implementing Unicode compression 401

Unicode normalization 408

Canonical decomposition 409

Compatibility decomposition 413

Canonical composition 414

Optimizing Unicode normalization 420

Testing Unicode normalization 420

Converting between Unicode and other standards 421

Getting conversion information 421

Converting between Unicode and single-byte encodings 422

Converting between Unicode and multi-byte encodings 422

Other types of conversion 422

Handling exceptional conditions 423

Dealing with differences in encoding philosophy 424

Choosing a converter 425

Line-break conversion 425

Case mapping and case folding 426

Case mapping on a single character 426

Case mapping on a string 427

Case folding 427

Transliteration 428

CHAPTER 15 Searching and Sorting 433

The basics of language-sensitive string comparison 433

Multi-level comparisons 436

Ignorable characters 438

French accent sorting 439

Contracting character sequences 440

Expanding characters 440

Context-sensitive weighting 441

Putting it all together 441
Other processes and equivalences 442

Language-sensitive comparison on Unicode text 443

Unicode normalization 443
Reordering 444
A general implementation strategy 445
The Unicode Collation Algorithm 447
The default UCA sort order 449
Alternate weighting 451
Optimizations and enhancements 453

Language-insensitive string comparison 455

Sorting 457

Collation strength and secondary keys 457
Exposing sort keys 459
Minimizing sort key length 460

Searching 461

The Boyer-Moore algorithm 462
Using Boyer-Moore with Unicode 465
“Whole words” searches 466

Using Unicode with regular expressions 466

CHAPTER 16 Rendering and Editing 469

Line breaking 470

Line breaking properties 471
Implementing boundary analysis with pair tables 473
Implementing boundary analysis with state machines 474
Performing boundary analysis using a dictionary 476
A couple more thoughts about boundary analysis 477
Performing line breaking 477

Line layout 479

Glyph selection and positioning 483

Font technologies 483
Poor man’s glyph selection 485
Glyph selection and placement in AAT 487
Glyph selection and placement in OpenType 489
Special-purpose rendering technology 491
Compound and virtual fonts 491

Special text-editing considerations 492

Optimizing for editing performance 492
Accepting text input 496
Handling arrow keys 497
Handling discontinuous selection 499
Handling multiple-click selection 500

CHAPTER 17 Unicode and Other Technologies 503

Unicode and the Internet 503

The W3C character model 504
XML 506
HTML and HTTP 507
URLs and domain names 508
Mail and Usenet 509

Unicode and programming languages 512

The Unicode identifier guidelines 512

Java 512

C and C++ 513

Javascript and JScript 513

Visual Basic 513

Perl 514

Unicode and operating systems 514

Microsoft Windows 514

MacOS 515

Varieties of Unix 516

Unicode and databases 517

Conclusion 517

Glossary 519

Bibliography 591

The Unicode Standard 591

Other Standards Documents 593

Books and Magazine Articles 593

Unicode Conference papers 594

Other papers 594

Online resources 595

Preface

As the economies of the world continue to become more connected together, and as the American computer market becomes more and more saturated, computer-related businesses are looking more and more to markets outside the United States to grow their businesses. At the same time, companies in other industries are not only beginning to do the same thing (or, in fact, have been for a long time), but are increasingly turning to computer technology, especially the Internet, to grow their businesses and streamline their operations.

The convergence of these two trends means that it's no longer just an English-only market for computer software. More and more, computer software is being used not only by people outside the United States or by people whose first language isn't English, but by people who don't speak English at all. As a result, interest in software internationalization is growing in the software development community.

A lot of things are involved in software internationalization: displaying text in the user's native language (and in different languages depending on the user), accepting input in the user's native language, altering window layouts to accommodate expansion or contraction of text or differences in writing direction, displaying numeric values according to local customs, indicating events in time according to the local calendar systems, and so on.

This book isn't about any of these things. It's about something more basic, and which underlies most of the issues listed above: representing written language in a computer. There are many different ways to do this; in fact, there are several for just about every language that's been represented in computers. In fact, that's the whole problem. Designing software that's flexible enough to handle data in multiple languages (at least multiple languages that use different writing systems) has traditionally meant not just keeping track of the text, but also keeping track of which encoding scheme is being used to represent it. And if you want to mix text in multiple writing systems, this bookkeeping becomes more and more cumbersome.

The Unicode standard was designed specifically to solve this problem. It aims to be the *universal* character encoding standard, providing unique, unambiguous representations for every character in virtually every writing system and language in the world. The most recent version of Unicode provides representations for over 90,000 characters.

Unicode has been around for twelve years now and is in its third major revision, adding support for more languages with each revision. It has gained widespread support in the software community and is now supported in a wide variety of operating systems, programming languages, and application programs. Each of the semiannual International Unicode Conferences is better-attended than the previous one, and the number of presenters and sessions at the Conferences grows correspondingly.

Representing text isn't as straightforward as it appears at first glance: it's not merely as simple as picking out a bunch of characters and assigning numbers to them. First you have to decide what a "character" is, which isn't as obvious in many writing systems as it is in English. You have to contend with things such as how to represent characters with diacritical marks applied to them, how to represent clusters of marks that represent syllables, when differently-shaped marks on the page are different "characters" and when they're just different ways of writing the same "character," what order to store the characters in when they don't proceed in a straightforward manner from one side of the page to the other (for example, some characters stack on top of each other, or you have two parallel lines of characters, or the reading order of the text on the page zigzags around the line because of differences in natural reading direction), and many similar issues.

The decisions you make on each of these issues for every character affect how various processes, such as comparing strings or analyzing a string for word boundaries, are performed, making them more complicated. In addition, the sheer number of different characters representable using the Unicode standard make many processes on text more complicated.

For all of these reasons, the Unicode standard is a large, complicated affair. Unicode 3.0, the last version published as a book, is 1,040 pages long. Even at this length, many of the explanations are fairly concise and assume the reader already has some degree of familiarity with the problems to be solved. It can be kind of intimidating.

The aim of this book is to provide an easier entrée into the world of Unicode. It arranges things in a more pedagogical manner, takes more time to explain the various issues and how they're solved, fills in various pieces of background information, and adds implementation information and information on what Unicode support is already out there. It is this author's hope that this book will be a worthy companion to the standard itself, and will provide the average programmer and the internationalization specialist alike with all the information they need to effectively handle Unicode in their software.

About this book

There are a few things you should keep in mind as you go through this book:

- This book assumes the reader either is a professional computer programmer or is familiar with most computer-programming concepts and terms. Most general computer-science jargon isn't defined or explained here.

- It's helpful, but not essential, if the reader has some basic understanding of the basic concepts of software internationalization. Many of those concepts are explained here, but if they're not central to one of the book's topics, they're not given a lot of time.
- This book covers a lot of ground, and it isn't intended as a comprehensive and definitive reference for every single topic it discusses. In particular, I'm not repeating the entire text of the Unicode standard here; the idea is to complement the standard, not replace it. In many cases, this book will summarize a topic or attempt to explain it at a high level, leaving it to other documents (typically the Unicode standard or one of its technical reports) to fill in all the details.
- The Unicode standard changes rapidly. New versions come out yearly, and small changes, new technical reports, and other things happen more quickly. In Unicode's history, terminology has changed, and this will probably continue to happen from time to time. In addition, there are a lot of other technologies that use or depend on Unicode, and they are also constantly changing, and I'm certainly not an expert on every single topic I discuss here. (In my darker moments, I'm not sure I'm an expert on any of them!) I have made every effort I could to see to it that this book is complete, accurate, and up to date, but I can't guarantee I've succeeded in every detail. In fact, I can almost guarantee you that there is information in here that is either outdated or just plain wrong. But I have made every effort to make the proportion of such information in this book as small as possible, and I pledge to continue, with each future version, to try to bring it closer to being fully accurate.
- At the time of this writing (January 2002), the newest version of Unicode, Unicode 3.2, was in beta, and thus still in flux. The Unicode 3.2 spec is scheduled to be finalized in March 2002, well before this book actually hits the streets. With a few exceptions, I don't expect major changes between now and March, but they're always possible, and therefore, the Unicode 3.2 information in this book may wind up wrong in some details. I've tried to flag all the Unicode 3.2-specific information here as being from Unicode 3.2, and I've tried to indicate the areas that I think are still in the greatest amount of flux.
- Sample code in this book is almost always in Java. This is partially because Java is the language I personally use in my regular job, and thus the programming language I think in these days. But I also chose Java because of its increasing importance and popularity in the programming world in general and because Java code tends to be somewhat easier to understand than, say, C (or at least no more difficult). Because of Java's syntactic similarity to C and C++, I also hope the examples will be reasonably accessible to C and C++ programmers who don't also program in Java.
- The sample code is provided for illustrative purposes only. I've gone to the trouble, at least with the examples that can stand alone, to make sure the examples all compile, and I've tested them to make sure I didn't make any obvious stupid mistakes, but they haven't been tested comprehensively. They were also written with far more of an eye toward explaining a concept than being directly usable in any particular context. Incorporate them into your code at your own risk!
- I've tried to define all the jargon the first time I use it or to indicate a full explanation is coming later, but there's also a glossary at the back you can refer to if you come across an unfamiliar term that isn't defined.
- Numeric constants, especially numbers representing characters, are pretty much always shown in hexadecimal notation. Hexadecimal numbers in the text are always written using the 0x notation familiar to C and Java programmers.
- Unicode code point values are shown using the standard Unicode notation, U+1234, where "1234" is a hexadecimal number of from four to six digits. In many cases, a character is referred to by both its Unicode code point value and its Unicode name: for example, "U+0041 LATIN CAPITAL LETTER A." Code unit values in one of the Unicode transformation formats are shown using the 0x notation.

How this book was produced

All of the examples of text in non-Latin writing systems posed quite a challenge for the production process. The bulk of this manuscript was written on a Power Macintosh G4/450 using Adobe FrameMaker 6 running on MacOS 9. I did the original versions of the various diagrams in Microsoft PowerPoint 98 on the Macintosh. But I discovered early on that FrameMaker on the Macintosh couldn't handle a lot of the characters I needed to be able to show in the book. I wound up writing the whole thing with little placeholder notes to myself throughout describing what the examples were supposed to look like.

FrameMaker was somewhat compatible with Apple's WorldScript technology, enabling me to do some of the example, but I quickly discovered Acrobat 3, which I was using at the time, wasn't. It crashed when I tried to create PDFs of chapters that included the non-Latin characters. Switching to Windows didn't prove much more fruitful: On both platforms, FrameMaker 6, Adobe Illustrator 9, and Acrobat 3 and 4 were not Unicode compatible. The non-Latin characters would either turn into garbage characters, not show up at all, or show up with very compromised quality.

Late in the process, I decided to switch to the one combination of software and platform I knew would work: Microsoft Office 2000 on Windows 2000, which handles (with varying degrees of difficulty) everything I needed to do. I converted the whole project from FrameMaker to Word and spent a couple of months restoring all the missing examples to the text. (In a few cases where I didn't have suitable fonts at my disposal, or where Word didn't produce quite the results I wanted, I either scanned pictures out of books or just left the placeholders in.) The last rounds of revisions were done in Word on either the Mac or on a Windows machine, depending on where I was physically at the time, and all the example text was done on the Windows machine.

This produced a manuscript of high-enough quality to get reviews from people, but didn't really produce anything we could typeset from. The work of converting the whole mess *back* to FrameMaker, redrawing my diagrams in Illustrator, and coming up with another way to do the non-Latin text examples fell to [name], who [finish the story after we're far enough into production to know how these problems were solved].

[Note to Ross and the AW production people: Is it customary to have a colophon on these things? Because of the difficulty and importance of typesetting this particular material, I'm thinking the normal information in a colophon should be included. Including it in the preface might be difficult, however, as I can't finish this section until at least some of the chapters have been typeset. Should I leave this info here, move it to a real colophon at the back of the book, or toss it?]

The author's journey

Like many in the field, I fell into software internationalization by happenstance. I've always been interested in language—written language in particular—and (of course) in computers. But my job had never really dealt with this directly.

In the spring of 1995, that changed when I went to work for Taligent. Taligent, you may remember, was the ill-fated joint venture between Apple Computer and IBM (later joined by Hewlett-Packard) that was originally formed to create a new operating system for personal computers using state-of-the-art object-oriented technology. The fruit of our labors, CommonPoint, turned out to be too little too late, but it spawned a lot of technologies that found their places in other products.

For a while there, Taligent enjoyed a cachet in the industry as the place where Apple and IBM had sent many of their best and brightest. If you managed to get a job at Taligent, you had “made it.”

I almost didn’t “make it.” I had wanted to work at Taligent for some time and eventually got the chance, but turned in a rather unimpressive interview performance (a couple coworkers kidded me about that for years afterward) and wasn’t offered the job. About that same time, a friend of mine *did* get a job there, and after the person who did get the job I interviewed for turned it down for family reasons, my friend put in a good word for me and I got a second chance.

I probably would have taken almost any job there, but the specific opening was in the text and internationalization group, and thus began my long association with Unicode.

One thing pretty much everybody who ever worked at Taligent will agree on is that working there was a wonderful learning experience: an opportunity, as it were, to “sit at the feet of the masters.” Personally, the Taligent experience made me into the programmer I am today. My C++ and OOD skills improved dramatically, I became proficient in Java, and I went from knowing virtually nothing about written language and software internationalization to... well, I’ll let you be the judge.

My team was eventually absorbed into IBM, and I enjoyed a little over two years as an IBMer before deciding to move on in early 2000. During my time at Taligent/IBM, I worked on four different sets of Unicode-related text handling routines: the text-editing frameworks in CommonPoint, various text-storage and internationalization frameworks in IBM’s Open Class Library, various internationalization facilities in Sun’s Java Class Library (which IBM wrote under contract to Sun), and the libraries that eventually came to be known as the International Components for Unicode.

International Components for Unicode, or ICU, began life as an IBM developer-tools package based on the Java internationalization libraries, but has since morphed into an open-source project and taken on a life of its own. It’s gaining increasing popularity and showing up in more operating systems and software packages, and it’s acquiring a reputation as a great demonstration of how to implement the various features of the Unicode standard. I had the twin privileges of contributing frameworks to Java and ICU and of working alongside those who developed the other frameworks and learning from them. I got to watch the Unicode standard develop, work with some of those who were developing it, occasionally rub shoulders with the others, and occasionally contribute a tidbit or two to the effort myself. It was a fantastic experience, and I hope that at least some of their expertise rubbed off on me.

Acknowledgements

It’s been said that it takes a village to raise a child. Well, I don’t really know about that, but it definitely takes a village to write a book like this. The person whose name is on the cover gets to take the credit, but there’s an army of people who contribute to the content.

Acknowledgements sections have a bad tendency to sound like acceptance speeches at the Oscars as the authors go on forever thanking everyone in sight. This set of acknowledgments will be no different. If you’re bored or annoyed by that sort of thing, I recommend you skip to the next section now. You have been warned.

Acknowledgements

Here goes: First and foremost, I'm indebted to the various wonderful and brilliant people I worked with on the internationalization teams at Taligent and IBM: Mark Davis, Kathleen Wilson, Laura Werner, Doug Felt, Helena Shih, John Fitzpatrick, Alan Liu, John Raley, Markus Scherer, Eric Mader, Bertrand Damiba, Stephen Booth, Steven Loomis, Vladimir Weinstein, Judy Lin, Thomas Scott, Brian Beck, John Jenkins, Deborah Goldsmith, Clayton Lewis, Chen-Lieh Huang, and Norbert Schatz. Whatever I know about either Unicode or software internationalization I learned from these people. I'd also like to thank the crew at Sun that we worked with: Brian Beck, Norbert Lindenberg, Stuart Gill, and John O'Conner.

I'd also like to thank my management and coworkers at Trilogy, particularly Robin Williamson, Chris Hyams, Doug Spencer, Marc Surplus, Dave Griffith, John DeRegnaucourt, Bryon Jacob, and Zach Roadhouse, for their understanding and support as I worked on this book, and especially for letting me continue to participate in various Unicode-related activities, especially the conferences, on company time and with company money.

Numerous people have helped me directly in my efforts to put this book together, by reviewing parts of it and offering advice and corrections, by answering questions and giving advice, by helping me put together example or letting me use examples they'd already put together, or simply by offering an encouraging word or two. I'm tremendously indebted to all who helped out in these ways: Jim Agenbroad, Matitiahu Allouche, Christian Cooke, John Cowan, Simon Cozens, Mark Davis, Roy Daya, Andy Deitsch, Martin Dürst, Tom Emerson, Franklin Friedmann, David Gallardo, Tim Greenwood, Jarkko Hietaniemi, Richard Ishida, John Jenkins, Kent Karlsson, Koji Kodama, Alain LaBonté, Ken Lunde, Rick McGowan, John O'Conner, Chris Pratley, John Raley, Jonathan Rosenne, Yair Sarig, Dave Thomas, and Garret Wilson. **[Be sure to add the names of anyone who sends me feedback between 12/31/01 and RTP.]**

And of course I'd like to acknowledge all the people at Addison-Wesley who've had a hand in putting this thing together: Ross Venables, my current editor; Julie DiNicola, my former editor; John Fuller, the production manager; **[name]**, who copy-edited the manuscript, **[name]**, who had the unenviable task of designing and typesetting the manuscript and cleaning up all my diagrams and examples; **[name]**, who put together the index; **[name]**, who designed the cover; and Mike Hendrickson, who oversaw the whole thing. I greatly appreciate their professionalism and their patience in dealing with this first-time author.

And last but not least, I'd like to thank the family members and friends who've had to sit and listen to me talk about this project for the last couple years: especially my parents, Ted and Joyce Gillam; Leigh Anne and Ken Boynton, Ken Cancelosi, Kelly Bowers, Bruce Rittenhouse, and many of the people listed above.

As always with these things, I couldn't have done it without all these people. If this book is good, they deserve the lion's share of the credit. If it isn't, I deserve the blame and owe them my apologies.

A personal appeal

For a year and a half, I wrote a column on Java for *C++ Report* magazine, and for much of that time, I wondered if anyone was actually reading the thing and what they thought of it. I would occasionally get an email from a reader commenting on something I'd written, and I was always

Preface

grateful, whether the feedback was good or bad, because it meant someone was reading the thing and took it seriously enough to let me know what they thought.

I'm hoping there will be more than one edition of this book, and I really want it to be as good as possible. If you read it and find it less than helpful, I hope you won't just throw it on a shelf somewhere and grumble about the money you threw away on it. Please, if this book fails to adequately answer your questions about Unicode, or if it wastes too much time answering questions you don't care about, I want to know. The more specific you can be about just what isn't doing it for you, the better. Please write me at rtgillam@concentric.net with your comments and criticisms.

For that matter, if you *like* what you see here, I wouldn't mind hearing from you either. God knows, I can use the encouragement.

—*R. T. G.*
Austin, Teaxs
January 2002

SECTION I

Unicode in Essence

*An Architectural Overview of the
Unicode Standard*

CHAPTER 1 *Language, Computers, and Unicode*

Words are amazing, marvelous things. They have the power both to build up and to tear down great civilizations. Words make us who we are. Indeed, many have observed that our use of language is what separates humankind from the rest of the animal kingdom. Humans have the capacity for symbolic thought, the ability to think about and discuss things we cannot immediately see or touch. Language is our chief symbolic system for doing this. Consider even a fairly simple concept such as “There is water over the next hill.” Without language, this would be an extraordinarily difficult thing to convey.

There’s been a lot of talk in recent years about “information.” We face an “information explosion” and live in an “information age.” Maybe this is true, but when it comes right down to it, we are creatures of information. And language is one of our main ways both of sharing information with one another and, often, one of our main ways of processing information.

We often hear that we are in the midst of an “information revolution,” surrounded by new forms of “information technology,” a phrase that didn’t even exist a generation or two ago. These days, the term “information technology” is generally used to refer to technology that helps us perform one or more of three basic processes: the storage and retrieval of information, the extraction of higher levels of meaning from a collection of information, and the transmission of information over large distances. The telegraph, and later the telephone, was a quantum leap in the last of these three processes, and the digital computer in the first two, and these two technologies form the cornerstone of the modern “information age.”

Yet by far the most important advance in information technology occurred many thousands of years ago and can’t be credited with an inventor. That advance (like so many technological revolutions, really a series of smaller advances) was written language. Think about it: before written language, storing and retrieving information over a long period of time relied mostly on human memory, or on

precursors of writing, such as making notches in sticks. Human memory is unreliable, and storage and retrieval (or storage over a time longer than a single person's lifetime) involved direct oral contact between people. Notches in sticks and the like avoids this problem, but doesn't allow for a lot of nuance or depth in the information being stored. Likewise, transmission of information over a long distance either also required memorization, or relied on things like drums or smoke signals that also had limited range and bandwidth.

Writing made both of these processes vastly more powerful and reliable. It enabled storage of information in dramatically greater concentrations and over dramatically longer time spans than was ever thought possible, and made possible transmission of information over dramatically greater distances, with greater nuance and fidelity, than had ever been thought possible.

In fact, most of today's data processing and telecommunications technologies have written language as their basis. Much of what we do with computers today is use them to store, retrieve, transmit, produce, analyze, and print written language.

Information technology didn't begin with the computer, and it didn't begin with the telephone or telegraph. It began with written language.

This is a book about how computers are used to store, retrieve, transmit, manipulate, and analyze written language.

* * *

Language makes us who we are. The words we choose speak volumes about who we are and about how we see the world and the people around us. They tell others who "our people" are: where we're from, what social class we belong to, possibly even who we're related to.

The world is made up of hundreds of ethnic groups, who constantly do business with, create alliances with, commingle with, and go to war with each other. And the whole concept of "ethnic group" is rooted in language. Who "your people" are is rooted not only in which language *you* speak, but in which language or languages your ancestors spoke. As a group's subgroups become separated, the languages the two subgroups speak will begin to diverge, eventually giving rise to multiple languages that share a common heritage (for example, classical Latin diverging into modern Spanish and French), and as different groups come into contact with each other, their respective languages will change under each other's influence (for example, much of modern english vocabulary was borrowed in from French). Much about a group's history is encoded in its language.

We live in a world of languages. There are some 6,000 to 7,000 different languages spoken in the world today, each with countless dialects and regional variations¹. We may be united by language, but we're divided by our languages.

And yet the world of computing is strangely homogeneous. For decades now, the language of computing has been English. Specifically, American English. Thankfully, this is changing. One of the things that information technology is increasingly making possible is contact between people in different parts of the world. In particular, information technology is making it more and more

¹ SIL International's Ethnologue Web site (www.ethnologue.com) lists 6,800 "main" languages and 41,000 variants and dialects.

possible to do business in different parts of the world. And as information technology invades more and more of the world, people are increasingly becoming unwilling to speak the language of information technology—they want information technology to speak *their* language.

This is a book about how all, or at least most, written languages—not just English or Western European languages—can be used with information technology.

* * *

Language makes us who we are. Almost every human activity has its own language, a specialized version of a normal human language adapted to the demands of discussing a particular activity. So the language you speak also says much about what you do.

Every profession has its jargon, which is used both to provide a more precise method of discussing the various aspects of the profession and to help tell “insiders” from “outsiders.” The information technology industry has always had a reputation as one of the most jargon-laden professional groups there is. This probably isn’t true, but it looks that way because of the way computers have come to permeate our lives: now that non-computer people have to deal with computers, non-computer people are having to deal with the language that computer people use.

In fact, it’s interesting to watch as the language of information technology starts to infect the vernacular: “I don’t have the bandwidth to deal with that right now.” “Joe, can you spare some cycles to talk to me?” And my personal favorite, from a TV commercial from a few years back: “No shampoo helps download dandruff better.”

What’s interesting is that subspecialties within information technology each have their own jargon as well that isn’t shared by computer people outside their subspecialty. In the same way that there’s been a bit of culture clash as the language of computers enters the language of everyday life, there’s been a bit of culture clash as the language of software internationalization enters the language of general computing.

This is a good development, because it shows the increasing interest of the computing community in developing computers, software, and other products that can deal with people in their native languages. We’re slowly moving from (apologies to Henry Ford) “your native language, as long as it’s English” to “your native language.” The challenge of writing one piece of software that can deal with users in multiple human languages involves many different problems that need to be solved, and each of those problems has its own terminology.

This is a book about some of that terminology.

* * *

One of the biggest problems to be dealt with in software internationalization is that the ways human language have been traditionally represented inside computers don’t often lend themselves to many human languages, and they lend themselves especially badly to multiple human languages at the same time.

Over time, systems have been developed for representing quite a few different written languages in computers, but each scheme is generally designed for only a single language, or at best a small

collection of related languages, and these systems are mutually incompatible. Interpreting a series of bits encoded with one standard using the rules of another yields gibberish, so software that handles multiple languages has traditionally had to do a lot of extra bookkeeping to keep track of the various different systems used to encode the characters of those languages. This is difficult to do well, and few pieces of software attempt it, leading to a Balkanization of computer software and the data it manipulates.

Unicode solves this problem by providing a *unified* representation for the characters in the various written languages. By providing a unique bit pattern for every single character, you eliminate the problem of having to keep track of which of many different characters this particular instance of a particular bit pattern is supposed to represent.

Of course, each language has its own peculiarities, and presents its own challenges for computerized representation. Dealing with Unicode doesn't necessarily mean dealing with all the peculiarities of the various languages, but it can—it depends on how many languages you actually want to support in your software, and how much you can rely on other software (such as the operating system) to do that work for you.

In addition, because of the sheer number of characters it encodes, there are challenges to dealing with Unicode-encoded text in software that go beyond those of dealing with the various languages it allows you to represent. The aim of this book is to help the average programmer find his way through the jargon and understand what goes into dealing with Unicode.

This is a book about Unicode.

What Unicode Is

Unicode is a standard method for representing written language in computers. So why do we need this? After all, there are probably dozens, if not hundreds, of ways of doing this already. Well, this is exactly the point. Unicode isn't just another in the endless parade of text-encoding standards; it's an attempt to do away with all the others, or at least simplify their use, by creating a *universal* text encoding standard.

Let's back up for a second. The best known and most widely used character encoding standard is the American Standard Code for Information Interchange, or ASCII for short. The first version of ASCII was published in 1964 as a standard way of representing textual data in computer memory and sending it over communication links between computers. ASCII is based on an even-bit byte. Each byte represented a character, and characters were represented by assigning them to individual bit patterns (or, if you prefer, individual numbers). A seven-bit byte can have 128 different bit patterns. 33 of these were set aside for use as control signals of various types (start- and end-of-transmission codes, block and record separators, etc.), leaving 95 free for representing characters.

Perhaps the main deficiency in ASCII comes from the A in its name: American. ASCII is an American standard, and was designed for the storage and transmission of English text. 95 characters are sufficient for representing English text, barely, but that's it. On early teletype machines, ASCII could also be used to represent the accented letters found in many European languages, but this capability disappeared in the transition from teletypes to CRT terminals.

So, as computer use became more and more widespread in different parts of the world, alternative methods of representing characters in computers arose for representing other languages, leading to the situation we have today, where there are generally three or four different encoding schemes for every language and writing system in use today.

Unicode is the latest of several attempts to solve this Tower of Babel problem by creating a universal character encoding. Its main way of doing this is to increase the size of the possible encoding space by increasing the number of bits used to encode each character. Most other character encodings are based upon an eight-bit byte, which provides enough space to encode a maximum of 256 characters (in practice, most encodings reserve some of these values for control signals and encode fewer than 256 characters). For languages, such as Japanese, that have more than 256 characters, most encodings are still based on the eight-bit byte, but use sequences of several bytes to represent most of the characters, using relatively complicated schemes to manage the variable numbers of bytes used to encode the characters.

Unicode uses a 16-bit word to encode characters, allowing up to 65,536 characters to be encoded without resorting to more complicated schemes involving multiple machine words per character. 65,000 characters, with careful management, is enough to allow encoding of the vast majority of characters in the vast majority of written languages in use today. The current version of Unicode, version 3.2, actually encodes 95,156 different characters—it actually does use a scheme to represent the less-common characters using two 16-bit units, but with 50,212 characters actually encoded using only a single unit, you rarely encounter the two-unit characters. In fact, these 50,212 characters include all of the characters representable with all of the other character encoding methods that are in reasonably widespread use.

This provides two main benefits: First, a system that wants to allow textual data (either user data or things like messages and labels that may need to be localized for different user communities) to be in any language would, without Unicode, have to keep track of not just the text itself, but also of which character encoding method was being used. In fact, mixtures of languages might require mixtures of character encodings. This extra bookkeeping means you couldn't look directly at the text and know, for example, that the value 65 was a capital letter A. Depending on the encoding scheme used for that particular piece of text, it might represent some other character, or even be simply part of a character (i.e., it might have to be considered along with an adjacent byte in order to be interpreted as a character). This might also mean you'd need different logic to perform certain processes on text depending on which encoding scheme they happened to use, or convert pieces of text between different encodings.

Unicode does away with this. It allows all of the same languages and characters to be represented using only one encoding scheme. Every character has its own unique, unambiguous value. The value 65, for example, *always* represents the capital letter A. You don't need to rely on extra information about the text in order to interpret it, you don't need different algorithms to perform certain processes on the text depending on the encoding or language, and you don't (with some relatively rare exceptions) need to consider context to correctly interpret any given 16-bit unit of text.

The other thing Unicode gives you is a pivot point for converting between other character encoding schemes. Because it's a superset of all of the other common character encoding systems, you can convert between any other two encodings by converting from one of them to Unicode, and then from Unicode to the other. Thus, if you have to provide a system that can convert text between any arbitrary pair of encodings, the number of converters you have to provide can be dramatically smaller. If you support n different encoding schemes, you only need $2n$ different converters, not n^2 different converters. It also means that when you have to write a system that interacts with the outside world using several different non-Unicode character representations, it can do its internal

processing in Unicode and convert at the boundaries, rather than potentially having to have alternate code to do the same things for text in the different outside encodings.

What Unicode Isn't

It's also important to keep in mind what Unicode isn't. First, Unicode is a standard scheme for representing *plain text* in computers and data communication. It is not a scheme for representing *rich text* (sometimes called “fancy text” or “styled text”). This is an important distinction. Plain text is the words, sentences, numbers, and so forth themselves. Rich text is plain text plus information about the text, especially information on the text's visual presentation (e.g., the fact that a given word is in italics), the structure of a document (e.g., the fact that a piece of text is a section header or footnote), or the language (e.g., the fact that a particular sentence is in Spanish). Rich text may also include non-text items that travel with the text, such as pictures.

It can be somewhat tough to draw a line between what qualifies as plain text, and therefore should be encoded in Unicode, and what's really rich text. In fact, debates on this very subject flare up from time to time in the various Unicode discussion forums. The basic rule is that plain text contains all of the information necessary to carry the semantic meaning of the text—the letters, spaces, digits, punctuation, and so forth. If removing it would make the text unintelligible, then it's plain text.

This is still a slippery definition. After all, italics and boldface carry semantic information, and losing them may lose some of the meaning of a sentence that uses them. On the other hand, it's perfectly possible to write intelligible, grammatical English without using italics and boldface, where it'd be impossible, or at least extremely difficult, to write intelligible, grammatical English without the letter “m”, or the comma, or the digit “3.” Some of it may also come down to user expectation—you can write intelligible English with only capital letters, but it's generally not considered grammatical or acceptable nowadays.

There's also a certain amount of document structure you need to be able to convey in plain text, even though document structure is generally considered the province of rich text. The classic example is the paragraph separator. You can't really get by without a way to represent a paragraph break in plain text without compromising legibility, even though it's technically something that indicates document structure. But many higher-level protocols that deal with document structure have their own ways of marking the beginnings and endings of paragraphs. The paragraph separator, thus, is one of a number of characters in Unicode that are explicitly disallowed (or ignored) in rich-text representations that are based on Unicode. HTML, for example, allows paragraph marks, but they're not recognized by HTML parsers as paragraph marks. Instead, HTML uses the `<P>` and `</P>` tags to mark paragraph boundaries.

When it comes down to it, the distinction between plain text and rich text is a judgment call. It's kind of like Potter Stewart's famous remark on obscenity—“I may not be able to give you a definition, but I know it when I see it.” Still, the principle is that Unicode encodes only plain text. Unicode may be used as the basis of a scheme for representing rich text, but isn't intended as a complete solution to this problem on its own.

Rich text is an example of a “higher-level protocol,” a phrase you'll run across a number of times in the Unicode standard. A higher-level protocol is anything that starts with the Unicode standard as its basis and then adds additional rules or processes to it. XML, for example, is a higher-level protocol that uses Unicode as its base and adds rules that define how plain Unicode text can be used to

represent structured information through the use of various kinds of markup tags. To Unicode, the markup tags are just Unicode text like everything else, but to XML, they delineate the structure of the document. You can, in fact, have multiple layers or protocols: XHTML is a higher-level protocol for representing rich text that uses XML as its base.

Markup languages such as HTML and XML are one example of how a higher-level protocol may be used to represent rich text. The other main class of higher-level protocols involves the use of multiple data structures, one or more of which contain plain Unicode text, and which are supplemented by other data structures that contain the information on the document's structure, the text's visual presentation, and any other non-text items that are included with the text. Most word processing programs use schemes like this.

Another thing Unicode isn't is a complete solution for software internationalization. Software internationalization is a set of design practices that lead to software that can be adapted for various international markets ("localized") without having to modify the executable code. The Unicode standard in all of its details includes a lot of stuff, but doesn't include everything necessary to produce internationalized software. In fact, it's perfectly possible to write internationalized software without using Unicode at all, and also perfectly possible to write completely non-internationalized software that uses Unicode.

Unicode is a solution to one particular problem in writing internationalized software: representing text in different languages without getting tripped up dealing with the multiplicity of encoding standards out there. This is an important problem, but it's not the only problem that needs to be solved when developing internationalized software. Among the other things an internationalized piece of software might have to worry about are:

- Presenting a different user interface to the user depending on what language he speaks. This may involve not only translating any text in the user interface into the user's language, but also altering screen layouts to accommodate the size or writing direction of the translated text, changing icons and other pictorial elements to be meaningful (or not to be offensive) to the target audience, changing color schemes for the same reasons, and so forth.
- Altering the ways in which binary values such as numbers, dates, and times are presented to the user, or the ways in which the user enters these values into the system. This involves not only relatively small things, like being able changing the character that's used for a decimal point (it's a comma, not a period, in most of Europe) or the order of the various pieces of a date (day-month-year is common in Europe), but possibly larger-scale changes (Chinese uses a completely different system for writing numbers, for example, and Israel uses a completely different calendar system).
- Altering various aspects of your program's behavior. For example sorting a list into alphabetical order may produce different orders for the same list depending on language because "alphabetical order" is a language-specific concept. Accounting software might need to work differently in different places because of differences in accounting rules.
- And the list goes on...

To those experienced in software internationalization, this is all obvious, of course, but those who aren't often seem to use the words "Unicode" and "internationalization" interchangeably. If you're in this camp, be careful: if you're writing in C++, storing all your character strings as arrays of `wchar_t` doesn't make your software internationalized. Likewise, if you're writing in Java, the fact that it's in Java and the `String` class uses Unicode doesn't automatically make your software internationalized. If you're unclear on the internationalization issues you might run into that Unicode *doesn't* solve, you can find an excellent introduction to the subject at <http://www.xerox-emea.com/globaldesign/paper/paper1.htm>, along with a wealth of other useful papers and other goodies.

Finally, another thing Unicode isn't is a glyph registry. We'll get into the Unicode character-glyph model in Chapter 3, but it's worth a quick synopsis here. Unicode draws a strong, important distinction between a *character*, which is an abstract linguistic concept such as “the Latin letter A” or “the Chinese character for ‘sun,’” and a *glyph*, which is a concrete visual presentation of a character, such as A or 日. There isn't a one-to-one correspondence between these two concepts: a single glyph may represent more than one character (such a glyph is often called a *ligature*), such as the fi ligature, a single mark that represents the letters *f* and *i* together. Or a single character might be represented by two or more glyphs: The vowel sound *au* in the Tamil language (•) is represented by two marks: one that goes to the left of a consonant character, and another on the right, but it's still thought of as a single character. A character may also be represented using different glyphs in different contexts: The Arabic letter *heh* has one shape when it stands alone (ه) and another when it occurs in the middle of a word (ه).

You'll also see what we might consider typeface distinctions between different languages using the same writing system. For instance, both Arabic and Urdu use the Arabic alphabet, but Urdu is generally written in the more ornate Nastaliq style, while Arabic frequently isn't. Japanese and Chinese are both written using Chinese characters, but some characters have a different shape in Chinese (for example, 骨 in Japanese is 骨 in Chinese).

Unicode, as a rule, doesn't care about any of these distinctions. It encodes underlying semantic concepts, not visual presentations (characters, not glyphs) and relies on intelligent rendering software (or the user's choice of fonts) to draw the correct glyphs in the correct places. Unicode does sometime encode glyphic distinctions, but only when necessary to preserve interoperability with some preexisting standard or to preserve legibility (i.e., if smart rendering software can't pick the right glyph for a particular character in a particular spot without clues in the encoded text itself). Despite these exceptions, Unicode by design does not attempt to catalogue every possible variant shape for a particular character. It encodes the character and leaves the shape to higher-level protocols.

The challenge of representing text in computers

The main body of the Unicode standard is 1,040 pages long, counting indexes and appendices, and there's a bunch of supplemental information—addenda, data tables, related substandards, implementation notes, etc.—on the Unicode Consortium's Web site. That's an awful lot of verbiage. And now here I come with another 500 pages on the subject. Why? After all, Unicode's just a character encoding. Sure, it includes a lot of characters, but how hard can it be?

Let's take a look at this for a few minutes. The basic principle at work here is simple: If you want to be able to represent textual information in a computer, you make a list of all the characters you want to represent and assign each one a number.² Now you can represent a sequence of characters with a

² Actually, you assign each one a bit pattern, but numbers are useful surrogates for bit patterns, since there's a generally-agreed-upon mapping from numbers to bit patterns. In fact, in some character encoding standards, including Unicode, there are several alternate ways to represent each character in bits, all based on the same numbers, and so the numbers you assign to them become useful as an intermediate stage between characters and bits. We'll look at this more closely in Chapters 2 and 6.

sequence of numbers. Consider the simple number code we all learned as kids, where A is 1, B is 2, C is 3, and so on. Using this scheme, the word...

food

...would be represented as 6-15-15-4.

Piece of cake. This also works well for Japanese...

日本語

...although you need a lot more numbers (which introduces its own set of problems, which we'll get to in a minute).

In real life, you may choose the numbers fairly judiciously, to facilitate things like sorting (it's useful, for example, to make the numeric order of the character codes follow the alphabetical order of the letters) or character-type tests (it makes sense, for example, to put all the digits in one contiguous group of codes and the letters in another, or even to position them in the encoding space such that you can check whether a character is a letter or digit with simple bit-masking). But the basic principle is still the same: Just assign each character a number.

It starts to get harder, or at least less clear-cut, as you move to other languages. Consider this phrase:

à bientôt

What do you do with the accented letters? You have two basic choices: You can either just assign a different number to each accented version of a given letter, or you can treat the accent marks as independent characters and given them their own numbers.

If you take the first approach, a process examining the text (comparing two strings, perhaps) can lose sight of the fact that a and à are the same letter, possibly causing it to do the wrong thing, without extra code that knows from extra information that à is just an accented version of a. If you take the second approach, a keeps its identity, but you then have to make decisions about where the code the accent goes in the sequence relative to the code for the a, and what tells a system that the accent belongs on top of the a and not some other letter.

For European languages, though, the first approach (just assigning a new number to the accented version of a letter) is generally considered to be simpler. But there are other situations...

הדרו ליהנה כּי-טוב: כּי לעולם נסדו

...such as this Hebrew example, where that approach breaks down. Here, most of the letters have marks on them, and the same marks can appear on any letter. Assigning a unique code to every

letter-mark combination quickly becomes unwieldy, and you have to go to giving the marks their own codes.

In fact, Unicode prefers the give-the-marks-their-own-codes approach, but in many cases also provides unique codes for the more common letter-mark combinations. This means that many combinations of characters can be represented more than one way. The “à” and “ô” in “à bientôt,” for example, can be represented either with single character codes, or with pairs of character codes, but you want “à bientôt” to be treated the same no matter which set of codes is used to represent it, so this requires a whole bunch of equivalence rules.

The whole idea that you number the characters in a line as they appear from left to right doesn’t just break down when you add accent marks and the like into the picture. Sometimes, it’s not even straightforward when all you’re dealing with are letters. In this sentence...

Avram said **מזל טוב** and smiled.

...which is in English with some Hebrew words embedded, you have an ordering quandary: The Hebrew letters don’t run from left to right; they run from right to left: the first letter in the Hebrew phrase, **מ**, is the one furthest to the right. This poses a problem for representation order. You can’t really store the characters in the order they appear on the line (in effect, storing either the English or the Hebrew “backward,” because it messes up the determination of which characters go on which line when you break text across lines. But if you store the characters in the order they’re read or typed, you need to specify just how they are to be arranged when the text is displayed or printed.

The ordering thing can go to even wilder extremes. This example...



...is in Hindi. The letters in Hindi knot together into clusters representing syllables. The syllables run from left to right across the page like English text does, but the arrangement of the marks within a syllable can be complicated and doesn’t necessarily follow the order in which the sounds they correspond to are actually spoken (or the characters themselves are typed). There are six characters in this word, arranged like this:



Many writing systems have complicated ordering or shaping behavior, and each presents unique challenges in determining how to represent the characters as a linear sequence of bits.

You also run into interesting decisions as to just what you mean by “the same character” or “different characters.” For example, in this Greek word...

Χριστός

...the letter σ and the letter ς are really both the letter sigma (Σ)—it’s written one way when it occurs at the end of a word and a different way when it occurs at the beginning or in the middle. Two distinct shapes representing the same letter. Do they get one character code or two? This issue comes up over and over again, as many writing systems have letters that change shape depending on context. In our Hindi example, the hook in the upper right-hand corner normally looks like this...

र

...but can take some very different forms (including the hook) depending on the characters surrounding it.

You also have the reverse problem of the same shape meaning different things. Chinese characters often have more than one meaning or pronunciation. Does each different meaning get its own character code? The letter Å can either be a letter in some Scandinavian languages or the symbol for the Angstrom unit. Do these two uses get different codes, or is it the same character in both places? What about this character:

3

Is this the number 3 or the Russian letter $з$? Do these share the same character code just because they happen to look a lot like each other?

For all of these reasons and many others like them, Unicode is more than just a collection of marks on paper with numbers assigned to them. Every character has a story, and for every character or group of characters, someone had to sit down and decide whether it was the same as or different from the other characters in Unicode, whether several related marks got assigned a single number or several, just what a series of numbers in computer memory would look like when you draw them on the screen, just how a series of marks on a page would translate into a series of numbers in computer memory when neither of these mappings was straightforward, how a computer performing various type of processes on a series of Unicode character codes would do its job, and so on.

So for every character code in the Unicode standard, there are rules about what it means, how it should look in various situations, how it gets arranged on a line of text with other characters, what other characters are similar but different, how various text-processing operations should treat it, and so on. Multiply all these decisions by 94,140 unique character codes, and you begin both to get an idea of why the standard is so big, and of just how much labor, how much energy, and how much

heartache, on the part of so many people, went into this thing. Unicode is the largest, most comprehensive, and most carefully designed standard of its type, and the toil of hundreds of people made it that way.

What This Book Does

The definitive source on Unicode is, not surprisingly, the Unicode standard itself. The main body of the standard is available in book form as *The Unicode Standard, Version 3.0*, published by Addison-Wesley and available wherever you bought this book. This book is supplemented by various tables of character properties covering the exact semantic details of each individual character, and by various technical reports that clarify, supplement, or extend the standard in various ways. A snapshot of this supplemental material is on the CD glued to the inside back cover of the book. The most current version of this supplemental material, and indeed the definitive source for all the most up-to-date material on the Unicode standard, is the Unicode Web site, at <http://www.unicode.org>.

For a long time, the Unicode standard was not only the definitive source on Unicode, it was the only source. The problem with this is that the Unicode standard is just that: a standard. Standards documents are written with people who will implement the standard as their audience. They assume extensive domain knowledge and are designed to define as precisely as possible every aspect of the thing being standardized. This makes sense: the whole purpose of a standard is to ensure that a diverse group of corporations and institutions all do some particular thing in the same way so that things produced by these different organizations can work together properly. If there are holes in the definition of the standard, or passages that are open to interpretation, you could wind up with implementations that conform to the standard, but still don't work together properly.

Because of this, and because they're generally written by committees whose members have different, and often conflicting, agendas, standards tend by their very nature to be dry, turgid, legalistic, and highly technical documents. They also tend to be organized in a way that presupposes considerable domain knowledge—if you're coming to the topic fresh, you'll often find that to understand any particular chapter of a standard, you have read every other chapter first.

The Unicode standard is better written than most, but it's still good bedtime reading—at least if you don't mind having nightmares about canonical reordering or the bi-di algorithm. That's where this book comes in. It's intended to act as a companion to the Unicode standard and supplement it by doing the following:

- Provide a more approachable, and more pedagogically organized, introduction to the salient features of the Unicode standard.
- Capture in book form changes and additions to the standard since it was last published in book form, and additions and adjuncts to the standard that haven't been published in book form.
- Fill in background information about the various features of the standard that are beyond the scope of the standard itself.
- Provide an introduction to each of the various writing systems Unicode represents and the encoding and implemented challenges presented by each.
- Provide useful information on implementing various aspects of the standard, or using existing implementations.

My hope is to provide a good enough introduction to “the big picture” and the main components of the technology that you can easily make sense of the more detailed descriptions in the standard itself—or know you don’t have to.

This book is for you if you’re a programmer using any technology that depends on the Unicode standard for something. It will give you a good introduction to the main concepts of Unicode, helping you to understand what’s relevant to you and what things to look for in the libraries or APIs you depend on.

This book is also for you if you’re doing programming work that actually involves implementing part of the Unicode standard and you’re still relatively new either to Unicode itself or to software internationalization in general. It will give you most of what you need and enough of a foundation to be able to find complete and definitive answers in the Unicode standard and its technical reports.

How this book is organized

This book is organized into three sections: **Section I, Unicode in Essence**, provides an architectural overview of the Unicode standard, explaining the most important concepts in the standard and the motivations behind them. **Section II, Unicode in Depth**, goes deeper, taking a close look at each of the writing systems representable using Unicode and the unique encoding and implementation problems they pose. **Section III, Unicode in Action**, take an in-depth look at what goes into implementing various aspects of the standard, writing code that manipulates Unicode text, and how Unicode interacts with other standards and technologies.

Section I: Unicode in Essence

Section I provides an introduction to the main structure and most important concepts of the standard, the things you need to know to deal properly with Unicode whatever you’re doing with it.

Chapter 1, the chapter you’re reading, is the book’s introduction. It gives a very high-level account of the problem Unicode is trying to solve, the goals and non-goals behind the standard, and the complexity of the problem. It also sets forth the goals and organization of this book.

Chapter 2 puts Unicode in historical context and relates it to the various other character encoding standards out there. It discusses ISO 10646, Unicode’s sister standard, and its relationship to the Unicode standard.

Chapter 3 provides a more complete architectural overview. It outlines the structure of the standard, Unicode’s guiding design principles, and what it means to conform to the Unicode standard.

Often, it takes two or more Unicode character codes to get a particular effect, and some effects can be achieved with two or more different sequences of codes. **Chapter 4** talks more about this concept, the combining character sequence, and the extra rules that specify how to deal with combining character sequences that are equivalent.

Every character in Unicode has a large set of properties that define its semantics and how it should be treated by various processes. These are all set forth in the Unicode Character Database, and **Chapter 5** introduces the database and all of the various character properties it defines.

The Unicode standard is actually in two layers: A layer that defines a transformation between written text and a series of abstract numeric codes, and a layer that defines a transformation between those abstract numeric codes and patterns of bits in memory or in persistent storage. The lower layer, from abstract numbers to bits, comprises several different mappings, each optimized for different situations. **Chapter 6** introduces and discusses these mappings.

Section II: Unicode in Depth

Unicode doesn't specifically deal in languages; instead it deals in *scripts*, or writing systems. A script is a collection of characters used to represent a group of related languages. Generally, no language uses all the characters in a script. For example, English is written using the Latin alphabet. Unicode encodes 819 Latin letters, but English only uses 52 (26 upper- and lower-case letters). Section II goes through the standard script by script, looking at the features of each script, the languages that are written with it, and how it's represented in Unicode. It groups scripts into families according to their common characteristics.

For example, in **Chapter 7** we look at the scripts used to write various European languages. These scripts generally don't pose any interesting ordering or shaping problems, but are the only scripts that have special upper- and lower-case forms. They're all descended from (or otherwise related to) the ancient Greek alphabet. This group includes the Latin, Greek, Cyrillic, Armenian, and Georgian alphabets, as well as various collections of diacritical marks and the International Phonetic Alphabet.

Chapter 8 looks at the scripts of the Middle East. The biggest feature shared by these scripts is that they're written from right to left rather than left to right. They also tend to use letters only for consonant sounds, using separate marks around the basic letters to represent the vowels. Two scripts in this group are cursively connected, even in printed text, which poses interesting representational problems. These scripts are all descended from the ancient Aramaic alphabet. This group includes the Hebrew, Arabic, Syriac, and Thaana alphabets.

Chapter 9 looks at the scripts of India and Southeast Asia. The letters in these scripts knot together into clusters that represent whole syllables. The scripts in this group all descend from the ancient Brahmi script. This group includes the Devanagari script used to write Hindi and Sanskrit, plus eighteen other scripts, including such things as Thai and Tibetan.

In **Chapter 10**, we look at the scripts of East Asia. The interesting things here are that these scripts comprise tens of thousands of unique, and often complicated, characters (the exact number is impossible to determine, and new characters are coined all the time). These characters are generally all the same size, don't combine with each other, and can be written either from left to right or vertically. This group includes the Chinese characters and various other writing systems that either are used with Chinese characters or arose under their influence.

While most of the written languages of the world are written using a writing system that falls into one of the above groups, not all of them do. **Chapter 11** discusses the other scripts, including Mongolian, Ethiopic, Cherokee, and the Unified Canadian Aboriginal Syllabics, a set of characters used for writing a variety of Native American languages. In addition to the modern scripts,

Unicode also encodes a growing number of scripts that are *not* used anymore but are of scholarly interest. The current version of Unicode includes four of these, which are also discussed in Chapter 11.

But of course, you can't write only with the characters that represent the sounds or words of spoken language. You also need things like punctuation marks, numbers, symbols, and various other non-letter characters. These are covered in **Chapter 12**, along with various special formatting and document-structure characters.

Section III: Unicode in Action

This section goes into depth on various techniques that can be used in code to implement or make use of the Unicode standard.

Chapter 13 provides an introduction to the subject, discussing a group of generic data structures and techniques that are useful for various types of processes that operate on Unicode text.

Chapter 14 goes into detail on how to perform various types of transformations and conversions on Unicode text. This includes converting between the various Unicode serialization formats, performing Unicode compression and decompression, performing Unicode normalization, converting between Unicode and other encoding standards, and performing case mapping and case folding.

Chapter 15 zeros in on two of the most text-analysis processes: searching and sorting. It talks about both language-sensitive and language-insensitive string comparison and how searching and sorting algorithms build on language-sensitive string comparison.

Chapter 16 discusses the most important operations performed on text: drawing it on the screen (or other output devices) and accepting it as input, otherwise known as rendering and editing. It talks about dividing text up into lines, arranging characters on a line, figuring out what shape to use for a particular character or sequence of characters, and various special considerations one must deal with when writing text-editing software.

Finally, in **Chapter 17** we look at the place where Unicode intersects with other computer technologies. It discusses Unicode and the Internet, Unicode and various programming languages, Unicode and various operating systems, and Unicode and database technology.

CHAPTER 2 *A Brief History of Character Encoding*

To understand Unicode fully, it's helpful to have a good sense of where we came from, and what this whole business of character encoding is all about. Unicode didn't just spring fully-grown from the forehead of Zeus; it's the latest step in a history that actually predates the digital computer, having its roots in telecommunications. Unicode is not the first attempt to solve the problem it solves, and Unicode is also in its third major revision. To understand the design decisions that led to Unicode 3.0, it's useful to understand what worked and what didn't work in Unicode's many predecessors.

This chapter is entirely background—if you want to jump right in and start looking at the features and design of Unicode itself, feel free to skip this chapter.

Prehistory

Fortunately, unlike, say, written language itself, the history of electronic (or electrical) representations of written language doesn't go back very far. This is mainly, of course, because the history of the devices using these representations doesn't go back very far.

The modern age of information technology does, however, start earlier than one might think at first—a good century or so before the advent of the modern digital computer. We can usefully date the beginning of modern information technology from Samuel Morse's invention of the telegraph in 1837.³

³ My main source for this section is Tom Jennings, "Annotated History of Character Codes," found at <http://www.wps.com/texts/codes>.

The telegraph, of course, is more than just an interesting historical curiosity. Telegraphic communication has never really gone away, although it's morphed a few times. Even long after the invention and popularization of the telephone, the successors of the telegraph continued to be used to send written communication over a wire. Telegraphic communications was used to send large volumes of text, especially when it needed ultimately to be in written form (news stories, for example), or when human contact wasn't especially important and saving money on bandwidth was very important (especially for routine business communications such as orders and invoices). These days, email and EDI are more or less the logical descendants of the telegraph.

The telegraph and Morse code

So our story starts with the telegraph. Morse's original telegraph code actually worked on numeric codes, not the alphanumeric code that we're familiar with today. The idea was that the operators on either end of the line would have a dictionary that assigned a unique number to each *word*, not each letter, in English (or a useful subset). The sender would look up each word in the dictionary, and send the number corresponding to that word; the receiver would do the opposite. (The idea was probably to automate this process in some way, perhaps with some kind of mechanical device that would point to each word in a list or something.)

This approach had died out by the time of Morse's famous "WHAT HATH GOD WROUGHT" demonstration in 1844. By this time, the device was being used to send the early version of what we now know as "Morse code," which was probably actually devised by Morse's assistant Alfred Vail.

Morse code was in no way digital in the sense we think of the term—you can't easily turn it into a stream of 1 and 0 bits the way you can with many of the succeeding codes. But it was "digital" in the sense that it was based on a circuit that had only two states, on and off. This is really Morse's big innovation; there were telegraph systems prior to Morse, but they were based on sending varying voltages down the line and deflecting a needle on a gauge of some kind.⁴ The beauty of Morse's scheme is a higher level of error tolerance—it's a lot easier to tell "on" from "off" than it is to tell "half on" from "three-fifths" on. This, of course, is also why modern computers are based on binary numbers.

The difference is that Morse code is based not on a succession of "ons" and "offs," but on a succession of "ons" of different lengths, and with some amount of "off" state separating them. You basically had two types of signal, a long "on" state, usually represented with a dash and pronounced "dah," and a short "on" state, usually represented by a dot and pronounced "dit." Individual letters were represented with varying-length sequences of dots and dashes.

The lengths of the codes were designed to correspond roughly to the relatively frequencies of the characters in a transmitted message. Letters were represented with anywhere from one to four dots and dashes. The two one-signal letters were the two most frequent letters in English: E was represented with a single dot and T with a single dash. The four two-signal letters were I (. .), A (. -), N (- .), and M (- -). The least common letters were represented with the longest codes: Z (- - . .), Y (- . - -), J (. - - -), and Q (- - . -). Digits were represented with sequences of five signals, and punctuation, which was used sparingly, was represented with sequences of six signals.

The dots and dashes were separated by just enough space to keep everything from running together, individual characters by longer spaces, and words by even longer spaces.

⁴ This tidbit comes from Steven J. Searle, "A Brief History of Character Codes," found at <http://www.tronweb.super-nova-co-jp/characcodehist.html>.

Prehistory

As an interesting sidelight, the telegraph was designed as a recording instrument—the signal operated a solenoid that caused a stylus to dig shallow grooves in a moving strip of paper. (The whole thing with interpreting Morse code by listening to beeping like we’ve all seen in World War II movies came later, with radio, although experienced telegraph operators could interpret the signal by listening to the clicking of the stylus.) This is a historical antecedent to the punched-tape systems used in teletype machines and early computers.

The teletypewriter and Baudot code

Of course, what you *really* want isn’t grooves on paper, but actual writing on paper, and one of the problems with Morse’s telegraph is that the varying lengths of the signals didn’t lend itself well to driving a mechanical device that could put actual letters on paper. The first big step in this direction was Emile Baudot’s “printing telegraph,” invented in 1874.

Baudot’s system didn’t use a typewriter keyboard; it used a pianolike keyboard with five keys, each of which controlled a separate electrical connection. The operator operated two keys with the left hand and three with the right and sent each character by pressing down some combination of these five keys simultaneously. (So the “chording keyboards” that are somewhat in vogue today as a way of combating RSI aren’t a new idea—they go all the way back to 1874.)

[should I include a picture of the Baudot keyboard?]

The code for each character is thus some combination of the five keys, so you can think of it as a five-bit code. Of course, this only gives you 32 combinations to play with, kind of a meager allotment for encoding characters. You can’t even get all the letters and digits in 32 codes.

The solution to this problem has persisted for many years since: you have two separate sets of characters assigned to the various key combinations, and you steal two key combinations to switch between them. So you end up with a LTRS bank, consisting of twenty-eight letters (the twenty-six you’d expect, plus two French letters), and a FIGS bank, consisting of twenty-eight characters: the ten digits and various punctuation marks and symbols. The three left-hand-only combinations don’t switch functions: two of them switch back and forth between LTRS and FIGS, and one (both left-hand keys together) was used to mean “ignore the last character” (this later evolves into the ASCII DEL character). The thirty-second combination, no keys at all, of course didn’t mean anything—no keys at all was what separated one character code from the next. (The FIGS and LTRS signals doubled as spaces.)

So you’d go along in LTRS mode, sending letters. When you came to a number or punctuation mark, you’d send FIGS, send the number or punctuation, then send LTRS and go back to sending words again. “I have 23 children.” would thus get sent as

I HAVE [FIGS] 23 [LTRS] CHILDREN [FIGS] .

It would have been possible, of course, to get the same effect by just adding a sixth key, but this was considered too complicated mechanically.

Even though Baudot’s code (actually invented by Johann Gauss and Wilhelm Weber) can be thought of as a series of five-bit numbers, it isn’t laid out like you might lay out a similar thing today: If you lay out the code charts according to the binary-number order, it looks jumbled. As with Morse code, characters were assigned to key combinations in such a way as to minimize fatigue, both to the

operator and to the machinery. More-frequent characters, for example, used fewer fingers than less-frequent characters.

Other teletype and telegraphy codes

The logical next step from Baudot's apparatus would, of course, be a system that uses a normal typewriterlike keyboard instead of the pianolike keyboard used by Baudot. Such a device was invented by Donald Murray sometime between 1899 and 1901. Murray kept the five-bit two-state system Baudot had used, but rearranged the characters. Since there was no longer a direct correlation between the operator's hand movements and the bits being sent over the wire, there was no need to worry about arranging the code to minimize operator fatigue; instead, Murray designed his code entirely to minimize wear and tear on the machinery.

A couple interesting developments occur first in the Murray code: You see the debut of what later became known as “format effectors” or “control characters”—the CR and LF codes, which, respectively return the typewriter carriage to the beginning of the line and advance the platen by one line.⁵ Two codes from Baudot also move to the positions where they stayed since (at least until the introduction of Unicode, by which time the positions no longer mattered): the NULL or BLANK all-bits-off code and the DEL all-bits-on code. All bits off, fairly logically, meant that the receiving machine shouldn't do anything; it was essentially used as an idle code for when no messages were being sent. On real equipment, you also often had to pad codes that took a long time to execute with NULLs: If you issued a CR, for example, it'd take a while for the carriage to return to the home position, and any “real” characters sent during this time would be lost, so you'd sent a bunch of extra NULLs after the CR. This would put enough space after the CR so that the real characters wouldn't go down the line until the receiving machine could print them, and not have any effect (other than maybe to waste time) if the carriage got there while they were still being sent.

The DEL character would also be ignored by the receiving equipment. The idea here is that if you're using paper tape as an intermediate storage medium (as we still see today, it became common to compose a message while off line, storing it on paper tape, and then log on and send the message from the paper tape, rather than “live” from the keyboard) and you make a mistake, the only way to blank out the mistake is to punch out all the holes in the line with the mistake. So a row with all the holes punched out (or a character with all the bits set, as we think of it today) was treated as a null character.

Murray's code forms the basis of most of the various telegraphy codes of the next fifty years or so. Western Union picked it up and used it (with a few changes) as its encoding method all the way through the 1950s. The CCITT (Consultative Committee for International Telephone and Telegraph, a European standards body) picked up the Western Union code and, with a few changes, blessed it as an international standard, International Telegraphy Alphabet #2 (“ITA2” for short).

The ITA2 code is often referred to today as “Baudot code,” although it's significantly different from Baudot's code. It does, however, retain many of the most important features of Baudot's code.

Among the interesting differences between Murray's code and ITA2 are the addition of more “control codes”: You see the introduction of an explicit space character, rather than using the all-bits-off signal, or the LTRS and FIGS signals, as spaces. There's a new BEL signal, which rings a bell or

⁵ I'm taking my cue from Jennings here: These code positions were apparently marked “COL” and “LINE PAGE” originally; Jennings extrapolates back from later codes that had CR and LF in the same positions and assumes that “COL” and “LINE PAGE” were alternate names for the same functions.

produces some other audible signal on the receiving end. And you see the first case of a code that exists explicitly to control the communications process itself—the WRU, or “Who are you?” code, which would cause the receiving machine to send some identifying stream of characters back to the sending machine (this enabled the sender to make sure he was connected to the right receiver before sending sensitive information down the wire, for example).

This is where inertia sets in in the industry. By the time ITA2 and its national variants came into use in the 1930s, you had a significant number of teletype machines out there, and you had the weight of an international standard behind one encoding method. The ITA2 code would be the code used by teletype machines right on into the 1960s. When computers started communicating with the outside world in real time using some kind of terminal, the terminal would be a teletype machine, and the computer would communicate with it using the teletype codes of the day. (The other main way computers dealt with alphanumeric data was through the use of punched cards, which had their own encoding schemes we’ll look at in a minute.)

FIELDATA and ASCII

The teletype codes we’ve looked at are all five-bit codes with two separate banks of characters. This means the codes all include the concept of “state”: the interpretation of a particular code depends on the last bank-change signal you got (some operations also included an implicit bank change—a carriage return, for example, often switched the system back to LTRS mode). This makes sense as long as you’re dealing with a completely serial medium where you don’t have random access into the middle of a character stream and as long as you’re basically dealing with a mechanical system (LTRS and FIGS would just shift the printhead or the platen to a different position). In computer memory, where you might want random access to a character and not have to scan backwards an arbitrary distance for a LTRS or FIGS character to find out what you’re looking at, it generally made more sense to just use an extra bit to store the LTRS or FIGS state. This gives you a six-bit code, and for a long time, characters were thought of as six-bit units. (Punched-card codes of the day were also six bits in length, so you’ve got that, too.) This is one reason why many computers of the time had word lengths that were multiples of 6: As late as 1978, Kernighan and Ritchie mention that the `int` and `short` data types were 36 bits wide on the Honeywell 6000, implying that it had a 36-bit word length.⁶

By the late 1950s, the computer and telecommunications industries were both starting to chafe under the limitations of the six-bit teletype and punched-card codes of the day, and a new standards effort was begun that eventually led to the ASCII code. An important predecessor of ASCII was the FIELDATA code used in various pieces of communications equipment designed and used by the U.S. Army starting in 1957. (It bled out into civilian life as well; UNIVAC computers of the day were based on a modified version of the FIELDATA code.)

FIELDATA code was a *seven*-bit code, but was divided into layers in such a way that you could think of it as a four-bit code with either two or three control bits appended, similar to punched-card codes. It’s useful⁷ to think of it as having a five-bit core somewhat on the ITA2 model with two control bits. The most significant, or “tag” bit, is used similarly to the LTRS/FIGS bit discussed before: it switches the other six bits between two banks: the “alphabetic” and “supervisory” banks. The next-most-significant bit shifted the five core bits between two sub-banks: the alphabetic bank was shifted between upper-case and lower-case sub-banks, and the supervisory bank between a supervisory and a numeric/symbols sub-bank.

⁶ See Brian W. Kernighan and Dennis M. Ritchie, *The C Programming Language*, first edition (Prentice-Hall, 1978), p. 34.

⁷ Or at least *I* think it’s useful, looking at the code charts—my sources don’t describe things this way.

In essence, the extra bit made it possible to include both upper- and lower-case letters for the first time, and also made it possible to include a wide range of control, or supervisory, codes, which were used for things like controlling the communication protocol (there were various handshake and termination codes) and separating various units of variable-length structured data.

Also, within each five-bit bank of characters, we finally see the influence of computer technology: the characters are ordered in binary-number order.

FIELDATA had a pronounced effect on ASCII, which was being designed at the same time. Committee X3.4 of the American Standards Association (now the American National Standards Institute, or ANSI) had been convened in the late 1950s, about the time FIELDATA was deployed, and consisted of representatives from AT&T, IBM (which, ironically, didn't actually use ASCII until the IBM PC came out in 1981), and various other companies from the computer and telecommunications industries.

The first result of this committee's efforts was what we now know as ANSI X3.4-1963, the first version of the American Standard Code for Information Interchange, our old friend ASCII (which, for the three readers of this book who don't already know this, is pronounced "ASS-key"). ASCII-1963 kept the overall structure of FIELDATA and regularized some things. For example, it's a pure 7-bit code and is laid out in such a way as to make it possible to reasonably sort a list into alphabetical order by comparing the numeric character codes directly. ASCII-1963 didn't officially have the lower-case letters in it, although there was a big undefined space where they were obviously going to go, and had a couple weird placements, such as a few control codes up in the printing-character range. These were fixed in the next version of ASCII, ASCII-1967, which is more or less the ASCII we all know and love today. It standardized the meanings of some of the control codes that were left open in ASCII-1963, moved all of the control codes (except for DEL, which we talked about) into the control-code area, and added the lower-case letters and some special programming-language symbols. It also discarded a couple of programming-language symbols from ASCII-1963 in a bow to international usage: the upward-pointing arrow used for exponentiation turned into the caret (which did double duty as the circumflex accent), and the left-pointing arrow (used sometimes to represent assignment) turned into the underscore character.

Hollerith and EBCDIC

So ASCII has its roots in telegraphy codes that go all the way back to the nineteenth century, and in fact was designed in part as telegraphy code itself. The other important category of character codes is the punched-card codes. Punched cards were for many years the main way computers dealt with alphanumeric data. In fact, the use of the punched card for data processing predates the modern computer, although not by as long as the telegraph.

Punched cards date back at least as far as 1810, when Joseph-Marie Jacquard used them to control automatic weaving machines, and Charles Babbage had proposed adapting Jacquard's punched cards for use in his "analytical engine." But punched cards weren't actually used for data processing until the 1880s when Herman Hollerith, a U.S. Census Bureau employee, devised a method of using punched cards to collate and tabulate census data. His punched cards were first used on a national scale in the 1890 census, dramatically speeding the tabulation process: The 1880 census figures, calculated entirely by hand, had taken seven years to tabulate. The 1890 census figures took six weeks. Flush with the success of the 1890 census, Hollerith formed the Tabulating Machine company in 1896 to market his punched cards and the machines used to punch, sort, and count them.

Prehistory

This company eventually merged with two others, diversified into actual computers, and grew into what we now know as the world's largest computer company, International Business Machines.⁸

The modern IBM rectangular-hole punched card was first introduced in 1928 and has since become the standard punched-card format. It had 80 columns, each representing a single character (it's no coincidence that most text-based CRT terminals had 80-column displays). Each column had twelve punch positions. This would seem to indicate that Hollerith code, the means of mapping from punch positions to characters, was a 12-bit code, which would give you 4,096 possible combinations.

It wasn't. This is because you didn't actually want to use all the possible combinations of punch positions—doing so would put too many holes in the card and weaken its structural integrity (not a small consideration when these things are flying through the sorting machine at a rate of a few dozen a second). The system worked like this: You had 10 main rows of punch holes, numbered 0 through 9, and two “zone” rows, 11 and 12 (row 0 also did double duty as a zone row). The digits were represented by single punches in rows 0 through 9, and the space was represented with no punches at all.

The letters were represented with two punches: a punch in row 11, 12, or 0 plus a punch in one of the rows from 1 to 9. This divides the alphabet into three nine-letter “zones.” A few special characters were represented with the extra couple one-punch combinations and the one remaining two-punch combination. The others were represented with three-punch combinations: row 8 would be pressed into service as an extra zone row and the characters would be represented with a combination of a punch in row 8, a punch in row 11, 12, or 0, and a punch in one of the rows from 1 to 7 (row 9 wasn't used). This gave you three banks of seven symbols each (four banks if included a bank that used row 8 as a zone row without an additional punch in rows 11, 12, or 0). All together, you got a grand total of 67 unique punch combinations. Early punched card systems didn't use all of these combinations, but later systems filled in until all of the possible combinations were being used.

[would a picture be helpful here?]

In memory, the computers used a six-bit encoding system tied to the punch patterns: The four least-significant bits would specify the main punch row (meaning 10 of the possible 16 combinations would be used), and the two most-significant bits would identify which of the zone rows was punched. (Actual systems weren't always quite so straightforward, but this was the basic idea.) This was sufficient to reproduce all of the one- and two-punch combinations in a straightforward manner, and was known as the Binary Coded Decimal Information Code, BCDIC, or just BCD for short.

IBM added two more bits to BCD to form the *Extended* Binary Coded Decimal Information Code, or EBCDIC (pronounced “EB-suh-dick”), which first appeared with the introduction of the System/360 in 1964. It was backward compatible with the BCD system used in the punched cards, but added the lower-case letters and a bunch of control codes borrowed from ASCII-1963 (you didn't really need control codes in a punched-card system, since the position of the columns on the cards, or the division of information between cards, gave you the structure of the data and you didn't need codes for controlling the communication session. This code was designed both to allow a simple mapping from character codes to punch positions on a punched card and, like ASCII, to produce a reasonable sorting order when the numeric codes were used to sort character data (note that this doesn't mean

⁸ Much of the information in the preceding paragraph is drawn from the Searle article; the source for the rest of this section is Douglas W. Jones, “Punched Cards: An Illustrated Technical History” and “Doug Jones' Punched Card Codes,” both found at <http://www.cs.uiowa.edu/~jones/cards>.

you get the *same* order as ASCII—digits sort after letters instead of before them, and lower case sorts before upper case instead of after).

One consequence of EBCDIC’s lineage is that the three groups of nine letters in the alphabet that you had in the punched-card codes are numerically separated from each other in the EBCDIC encoding: I is represented by 0xC9 and J by 0xD1, leaving an eight-space gap in the numerical sequence. The original version of EBCDIC encoded only 50 characters in an 8-bit encoding space, leaving a large numbers of gaping holes with no assigned characters. Later versions of EBCDIC filled in these holes in various ways, but retained the backward compatibility with the old punched-card system. Although ASCII is finally taking over in IBM’s product line, EBCDIC still survives in the current models of IBM’s System/390, even though punched cards are long obsolete. Backward compatibility is a powerful force.

Single-byte encoding systems

ANSI X3.4-1967 (ASCII-1967) went on to be adopted as an international standard, first by the European computer Manufacturers’ Association as ECMA-6 (which actually came out in 1965, two years before the updated version of ANSI X3.4, and then by the International Organization for Standardization (“ISO” for short) as ISO 646 in 1972. [This is generally the way it works with ISO standards—ISO serves as an umbrella organization for a bunch of national standards bodies and generally creates international standards by taking various national standards, modifying them to be more palatable to an international audience, and republishing them as international standards.]

A couple interesting things happened to ASCII on its way to turning into ISO 646. First, it formalized a system for applying the various accent and other diacritical marks used in European languages to the letters. It did this not by introducing accented variants of all the letters—there was no room for that—but by pressing a bunch of punctuation marks into service as diacritical marks. The apostrophe did double duty as the acute accent, the opening quote mark as the grave accent, the double quotation mark as the diaeresis (or umlaut), the caret as the circumflex accent, the swung dash as the tilde, and the comma as the cedilla. To produce an accented letter, you’d follow the letter in the code sequence with a backspace code and then the appropriate ‘accent mark.’ On a teletype machine, this would cause the letter to be overstruck with the punctuation mark, producing an ugly but serviceable version of the accented letter.

The other thing is that ISO 646 leaves the definitions of twelve characters open (in ASCII, these eleven characters are #, \$, @, [, \,], ^, \, {, |, }, and ~). These are called the “national use” code positions. There’s an International Reference Version of ISO 646 that gives the American meanings to the corresponding code points, but other national bodies were free to assign other characters to these twelve code values. (Some national bodies did put accented letters in these slots.) The various national variants have generally fallen out of use in favor of more-modern standards like ISO 8859 (see below), but vestiges of the old system still remain; for example, in many Japanese systems’ treatment of the code for \ as the code for ¥.⁹

⁹ The preceding information comes partially from the Jennings paper, and partially from Roman Czyborra, “Good Old ASCII,” found at <http://www.czyborra.com/charsets/iso646.html>.

Eight-bit encoding schemes and the ISO 2022 model

ASCII gave us the eight-bit byte. All earlier encoding systems (except for FIELDATA, which can be considered an embryonic version of ASCII) used a six-bit byte. ASCII extended that to seven, and most communication protocols tacked on an eighth bit as a parity-check bit. As the parity bit became less necessary (especially in computer memory) and 7-bit ASCII codes were stored in eight-bit computer bytes, it was only natural that the 128 bit combinations not defined by ASCII (the ones with the high-order bit set) would be pressed into service to represent more characters.

This was anticipated as early as 1971, when the ECMA-35 standard was first published.¹⁰ This standard later became ISO 2022. ISO 2022 sets forth a standard method of organizing the code space for various character encoding methods. An ISO-2022-compliant character encoding can have up to two sets of control characters, designated C0 and C1, and up to four sets of printing (“graphic”) characters, designated G0, G1, G2, and G3.

The encoding space can either be seven or eight bits wide. In a seven-bit encoding, the byte values from 0x00 to 0x1F are reserved for the C0 controls and the byte values from 0x20 to 0xFF are used for the G0, G1, G2, and G3 sets. The range defaults to the G0 set, and escape sequences can be used to switch it to one of the other sets.

In an eight-bit encoding, the range of values from 0x00 to 0x1F (the “CL area”) is the C0 controls and the range from 0x80 to 0x9F (the “CR area”) is the C1 controls. The range from 0x20 to 0x7F (the “GL area”) is always the G0 characters, and the range from 0xA0 to 0xFF (the “GR area”) can be switched between the G1, G2, and G3 characters.

Control functions (i.e., signals to the receiving process, as opposed to printing characters) can be represented either as single control characters or as sequences of characters usually beginning with a control character (usually that control character is ESC and the multi-character control function is an “escape sequence”). ISO 2022 uses escape sequences to represent C1 controls in the 7-bit systems, and to switch the GR and GL areas between the various sets of printing characters in both the 7- and 8-bit versions. It also specifies a method of using escape sequences to associate the various areas of the encoding space (C0, C1, G0, G1, G2, and G3) with actual sets of characters.

ISO 2022 doesn’t actually apply semantics to most of the code positions—the big exceptions are ESC, DEL and the space, which are given the positions they have in ASCII. Other than these, character semantics are taken from other standards, and escape sequences can be used to switch an area’s interpretation from one standard to another (there’s a registry of auxiliary standards and the escape sequences used to switch between them).

An ISO 2022-based standard can impose whatever semantics it wants on the various encoding areas and can choose to use or not use the escape sequences for switching things. As a practical matter, most ISO2022-derived standards put the ISO 646 IRV (i.e., US ASCII) printing characters in the G0 area and the C0 control functions from ISO 6429 (an ISO standard that defines a whole mess of control functions, along with standardized C0 and C1 sets—the C0 set is the same as the ASCII control characters).

¹⁰ The information in this section is taken from Peter K. Edberg, “Survey of Character Encodings,” *Proceedings of the 13th International Unicode Conference*, session TA4, September 9, 1998, and from a quick glance at the ECMA 35 standard itself.

ISO 8859

By far the most important of the ISO 2022-derived encoding schemes is the ISO 8859 family.¹¹ The ISO 8859 standard comprises fourteen separate ISO 2022-compliant encoding standards, each covering a different set of characters for a different set of languages. Each of these counts as a separate standard: ISO 8859-1, for example, is the near-ubiquitous Latin-1 character set you see on the Web.

Work on the ISO 8859 family began in 1982 as a joint project of ANSI and ECMA. The first part was originally published in 1985 as ECMA-94. This was adopted as ISO 8859-1, and a later addition of ECMA-94 became the first four parts of ISO 8859. The other parts of ISO 8859 likewise originated in various ECMA standards.

The ISO 8859 series is oriented, as one might expect, toward European languages and European usage, as well as certain languages around the periphery of Europe that get used a lot there. It aimed to do a few things: 1) Do away with the use of backspace sequences as the way to represent accented characters (the backspace thing is workable on a teletype, but doesn't work at all on a CRT terminal without some pretty fancy rendering hardware), 2) do away with the varying meanings of the "national use" characters from ISO 646, replacing them with a set of code values that would have the same meaning everywhere and still include everyone's characters, and 3) unify various other national and vendor standards that were attempting to do the same thing.

All of the parts of ISO 8859 are based on the ISO 2022 structure, and all have a lot in common. Each of them assigns the ISO 646 printing characters to the G0 range, and the ISO 6429 C0 and C1 control characters to the C0 and C1 ranges. This means that each of them, whatever else they include, includes the basic Latin alphabet and is downward compatible with ASCII. (That is, pure 7-bit ASCII text, when represented with 8-bit bytes, conforms to any of the 8859 standards.) Where they differ is in their treatment of the G1 range (none of them defines anything in the G2 or G3 areas or uses escape sequences to switch interpretations of any of the code points, although you can use the registered ISO 2022 escape sequences to assign the G1 repertoire from any of these standards to the various GR ranges in a generalized ISO 2022 implementation).

The fourteen parts of ISO 8859 are as follows:

ISO 8859-1	Latin-1	Western European languages (French, German, Spanish, Italian, the Scandinavian languages etc.)
ISO 8859-2	Latin-2	Eastern European languages (Czech, Hungarian, Polish, Romanian, etc.)
ISO 8859-3	Latin-3	Southern European languages (Maltese and Turkish, plus Esperanto)
ISO 8859-4	Latin-4	Northern European languages (Latvian, Lithuanian, Estonian, Greenlandic, and Sami)
ISO 8859-5	Cyrillic	Russian, Bulgarian, Ukrainian, Belarusian, Serbian, and Macedonian

¹¹ Much of the information in this section is drawn from Roman Czyborra, "ISO 8859 Alphabet Soup", found at <http://www.czyborra.com/charsets/iso8859.html>, supplemented with info from the ISO Web site, the ECMA 94 standard, and the Edberg article.

Single-byte encoding systems

ISO 8859-6	Arabic	Arabic
ISO 8859-7	Greek	Greek
ISO 8859-8	Hebrew	Hebrew
ISO 8859-9	Latin-5	Turkish (replaces Latin-3)
ISO 8859-10	Latin-6	Northern European languages (unifies Latin-1 and Latin-4)
ISO 8859-11	Thai	Thai
ISO 8859-13	Latin-7	Baltic languages (replaces Latin-4, supplements Latin-6)
ISO 8859-14	Latin-8	Celtic languages
ISO 8859-15	Latin-9	Western European languages (replaces Latin-1, adds the Euro symbol and some forgotten French and Finnish letters)
ISO 8859-16	Latin-10	Eastern European languages (replaces Latin-2, adds the Euro symbol and a few missing Romanian letters)

One of the things that the ISO 8859 family of standards does is show how constraining an 8-bit encoding space can be. You have a whole mess of different versions of the Latin alphabet with considerable overlap, each optimizing for a slightly different collection of alphabets. Likewise, the ones for the non-Latin alphabets tend to make do with fairly constrained collections of characters from their writing systems. But they're a lot better, and a lot more standard, than the mishmash that came before.

Other 8-bit encoding schemes

There are tons of other 8-bit encoding standards. These include a wide variety of vendor-developed encodings (many vendors, including IBM and Microsoft, call these "code pages") that cover the same set of languages as the 8859 series, but generally predate it and put things in a different order. There are also code pages that add other non-letter characters, such as the old DOS box-drawing characters. These include the old Mac and DOS code pages that predate ISO 8859, as well as some code pages that follow and extend ISO 8859. While most of these are downward-compatible with ASCII, they're not always ISO 2022 compliant. Windows code page 1252, the Western European encoding for Windows, for example, is a superset of ISO 8859-1. It includes all the Latin-1 characters, but then goes on to fill the C1 space with additional printing characters. (Sometimes you even see this happen to the C0 space: the old standard ICM PC code page not only had box-drawing characters and what-not up in the C1/G1 range, but also put printing characters down in the C0 range where they could be confused with control characters.) IBM also has a whole mess of EBCDIC-derived code pages devoted to the languages covered by the ISO 8859 family.

There are other eight-bit encoding standards for languages outside the ISO 8859 family. These includes the TIS 620 standard for Thai, an ISO 2022-derived standard that puts the Thai letters in the

G1 range, the ISCII standard (“Indian Script Code for Information Interchange”), which places any of several sets of characters for writing Indian languages in the G1 area. (The various Indian writing systems are all so closely related that the meanings of the code points are common across all of them—you can transliterate from one Indian writing system to another simply by changing fonts [more or less]). The VISCII and TCVN 5712 standards (VISCII [“Vietnamese Standard Code for Information Interchange”] is an Internet RFC; TCVN 5712 is a Vietnamese national standard) for writing Vietnamese are both ASCII-derived, but invade both the C0 and C1 areas with Vietnamese letters (it’s a big set). There’s a version of TCVN 5712 that complies with ISO 2022, but it omits a bunch of letters.

The Japanese JIS X 0201 (“JIS-Roman”) standard is also an ISO 2022-derived 8-bit standard. It places the Japanese Katakana syllabary in the G1 area. Here, though, the G0 area isn’t exactly the same as American ASCII: a few of the “national use” code positions from ISO 646 have Japanese-specific meanings (the one everybody’s most familiar with is that the code position normally occupied by the backslash [\] is instead occupied by the yen sign [¥], so DOS path names on Japanese computers come out with yen signs between their components).

Character encoding terminology

Before we move on to talk about more character encoding standards, we need to take a break and define some terms. It’ll make the discussion in the next section easier and help crystallize some of the fuzzy terminology in the sections we just finished.

It’s useful to think of the mapping of a sequence of written characters to a sequence of bits in a computer memory or storage device or in a communication link as taking places in a series of stages, or *levels*, rather than all at once. The Internet Architecture Board proposed a three-level encoding model. The Unicode standard, in Unicode Technical Report #17, proposes a five-level model, explicitly discussing one level that was merely implied in the IAB model, and interposing an extra level between two of the IAB levels. This discussion borrows the terms from UTR #17, since it was designed specifically to make it easier to talk about Unicode.

The first level is **abstract character repertoire**. This is simply a collection of characters to be encoded. This can usually be expressed simply in words: such phrases as “the Latin alphabet” or “the union of the English, Spanish, French, German, Italian, Swedish, Danish, and Norwegian alphabets” define repertoires. For most Western languages, this kind of thing is pretty straightforward and usually takes the form “The blah-blah alphabet, plus this set of special characters and digits,” but for East Asian languages, it’s more complicated.

One of the other things you start thinking about when defining a character repertoire is just what constitutes a “character.” Here, “character” can simply mean “a semantic unit that has its own identity when encoded.” An encoding standard can decide, for example whether accented letters are first-class characters that get their own identity when encoded, or whether the accents themselves are first-class characters that get their own identity and combine with unaccented letters to form the accented letters. Another example would be whether the different shapes a character can assume in different contexts each get treated as a “character” or whether the underlying and unchanging semantic unit these shapes represent gets treated as a “character.”

Unicode defines “character” as an independent abstract semantic unit, explicitly eschewing glyph variants and presentation forms, and explicitly encodes accents and other combining marks as first-class characters, but other encoding standards made different decisions. Throughout the rest of the

book, we'll use "character" the way Unicode does, distinguishing it from "glyph," rather than simply to mean "something deemed worthy of encoding."

The second level is the **coded character set**, a mapping between the characters in an abstract repertoire and some collection of coordinates. Most coded character sets map the characters to x and y positions in a table, although these coordinates can always be mapped to a set of single numbers.

Allied with the concept of coded character set is the idea of **encoding space**, basically the dimensions of the table containing all the characters. This can be expressed as a pair of numbers ("a 16×16 encoding space") or as a single number ("a 256-character encoding space"). Often (but not always) the encoding space is constrained by the size of the storage unit used to represent a character, and this is used a shorthand way of describing the encoding space ("an 8-bit encoding space"). Sometimes the encoding space is broken down into subsets. Depending on the subset and the standard, these subsets go by various names (e.g., "row," "plane," "column," etc.)

A position in the encoding space is called a **code point**. The position occupied by a particular character is called a **code point value**. A coded character set is a mapping of characters to code points values.

The third level of the encoding model is the **character encoding form**, also called a "storage format." This maps the abstract code point values to sequences of integers of a specific size, called **code units**. For a fixed-size encoding, this is usually a null mapping, but it doesn't have to be. For variable-length encodings, the mapping is more complicated, as some code points will map to single code units and some code points will map to sequences of code units.

The fourth level of the encoding model is the **character encoding scheme**, also referred to as the "serialization format." This is the mapping from code unit values to actual sequences of bits. Two basic things happen here: First, code units that are more than a byte long are converted to sequences of bytes (in other words, which byte goes first when the encoded text is sent over a communications link or written to a storage medium becomes important). Second, a character encoding scheme may tie together two or more character encoding forms, defining extra bytes that get written into the serialized data stream to switch between them, or some additional transformation that happens to the code units from one of the encoding forms to make their numeric ranges disjoint.

The "charset" parameter that shows up in a variety of Internet standards (see Chapter 17) specifically refers to a character encoding scheme.

A character encoding standard usually defines a stack of all four of these transformations: from characters to code points, from code points to code units, and from code units to bytes, although some of the transformations might be implicit. Sometimes, though, you'll see separate definitions of different layers in the model. This often happens with East Asian encodings, for example.

The fifth level is the **transfer encoding syntax**, which is more or less orthogonal to the other four levels, and is usually defined separately from (and independent of) the others. This is an additional transformation that might be applied *ex post facto* to the bytes produced by a character encoding scheme. Transfer encoding syntaxes usually are used for one of two things: They either map a set of values into a more-constrained set of values to meet the constraints of some environment (think Base64 or quoted-printable, both of which we'll talk about in Chapter 17), or they're compression formats, mapping sequences of bytes into shorter sequences of bytes (think LZW or run-length encoding).

Let's go back and apply these terms to the encoding standards we've looked at so far. The various telegraph codes, such as ITA2, divided their repertoires into two coded character sets, FIGS and LTRS, each of which mapped half the repertoire into a five-bit encoding space, usually expressed as a table of two rows of 16 cells each. The character encoding forms were straightforward transformations, mapping the coordinates into five-bit binary numbers in a straightforward manner. The character encoding *schemes* mapped the five-bit code units into bytes in an equally straightforward manner, and added the concept of using the FIGS and LTRS codes for switching back and forth between the two coded character sets. (In fact, if you consider the character encoding scheme the transformation to physical media, this is also where it was decided whether the most significant bit in the byte was the leftmost hole or the rightmost hole on the paper tape.)

ASCII is actually simpler, mapping its repertoire of characters into a 7-bit, 128-character encoding space, usually expressed as an 8×16 table. It comprised a character encoding form that mapped the table positions into integer code units in an obvious manner, and an encoding scheme that mapped these code units to bytes in an equally obvious manner (if you're going to an 8-bit byte, the character encoding scheme can be said to comprise the padding of the value out to eight bits).

ISO 2022 doesn't actually define any coded character sets, instead simply specifying the layout of the 256-character encoding space, but it does specify a character encoding scheme for switching back and forth between, say, the G1 and G2 sets, and for associating actual coded character sets (for example Latin-1 and Latin-2) with the G1 and G2 sets.

Multiple-byte encoding systems

These distinctions between layers become more interesting when you start talking about East Asian languages such as Chinese, Japanese, and Korean. These languages all make use of the Chinese characters. Nobody is really sure how many Chinese characters there are, although the total is probably more than 100,000. Most Chinese speakers have a working written vocabulary of some 5,000 characters or so, somewhat less for Japanese and Korean speakers, who depend more on auxiliary writing systems to augment the Chinese characters.

East Asian coded character sets

When you've got this many characters to mess around with, you start officially defining particular sets of Chinese characters you're interested in. The Japanese government, for example publishes the Gakushu Kanji, a list of 1,006 Chinese characters (the Japanese call them *kanji*) everyone learns in elementary school, the Joyo Kanji, a list of another 1,945 characters used in government documents and newspapers (and, presumably, learned in secondary school), and the Jinmei-yo Kanji, a list of 285 characters officially sanctioned for use in personal names.¹² (Think of these lists as analogous to the vocabulary lists that school boards and other educational institutions in the U.S. put together to standardize which words a kid should know at which grade level.) The other countries that use the Chinese characters publish similar lists. These lists are examples of character-set standards that just standardize an abstract character repertoire—there's no attempt (so far) to assign numbers to them.

The Japanese Industrial Standards Commission (JISC, Japan's counterpart to ANSI) led the way in devising coded character set standards for representing the Chinese characters in computers (there

¹² This section relies heavily on various sections of Ken Lunde, *CJKV Information Processing* (O'Reilly, 1999), as well as the Edberg paper, *op. cit.*

are earlier Japanese telegraph codes, but we don't consider them here). The first of these was JIS X 0201, first published in 1976, which we looked at earlier. It doesn't actually encode the Chinese characters; it's and ISO 2022-derived encoding that places a Japanese variant of ISO 646 in the G0 area and the Japanese Katakana syllabary in the G1 area. It's a fixed-length 8-bit encoding, and it becomes important later.

The first Japanese standard to actually encode the *kanji* characters was JIS C 6226, first published in 1978 and later renamed JIS X 0208 (the work on this standard actually goes back as far as 1969). JIS X 0208 uses a much larger encoding space than we've looked at so far: a 94×94 table. The choice of 94 as the number of rows and columns isn't coincidental: it's the number of assignable code points in the G0 area according to ISO 2022. This gives you an ISO 2022-compliant 7-bit encoding scheme. Each character would be represented with a pair of 7-bit code units (in effect, you have a single 14-bit code point expressed as a coordinate in a 94×94 space being mapped into two 7-bit code units).

JIS X 0208's repertoire includes not only 6,355 *kanji*, divides into two "levels" by frequency and importance, but also the Latin, Greek and Cyrillic alphabets, the Japanese Hiragana and Katakana syllabaries, and a wide variety of symbols and punctuation marks.

JIS X 0208 was supplemented in 1990 by JIS X 0212, a supplemental standard that adds another 5,801 *kanji*, as well as more symbols and some supplementary characters in the Latin and Greek alphabets. It's based on the same 94×94 encoding space used in JIS X 0208. JIS X 0213, published in 2000, adds in another 5,000 Kanji (Levels 3 and 4; I don't know what "level" JIS X 0212 is supposed to be) and follows the same structure as the other two standards.

The other countries that use the Chinese characters followed the JIS model when designing their own national character set standards. These include the People's Republic of China's GB 2312, first published in 1981, Taiwan's CNS 11643, published in 1992, and South Korea's KS X 1001 (aka KS C 5601), first published in 1987. All of these make use of a 94×94 encoding space, laid out in roughly similar ways (although the actual allocation of individual characters, in particular the Chinese characters and the native-language characters in each standard, varies quite a bit).

CNS 11643, by the way is the most ambitious of all of these, encoding a whopping 48,027 characters across 16 "planes" (actually, only 7 are currently used). The "planes" correspond roughly to the different parts of the ISO 8859 standard; they're more or less independent standards that put different characters in the same encoding space. CNS 11643 is set up such that characters at corresponding positions in each plan are variants of each other. In theory, you could take a text consisting of characters from all the places, encode it using only plane 1, and it'd still be legible.

All of these standards define coded character sets; if you take the row and column numbers of the code points and shift them up into ISO 2022's G0 space (by adding 32), you have a simple ISO 2022-derived encoding form. You can then have two encoding schemes based on this—one that encodes the row number before the column number, and another that encodes the column number before the row number. Piece of cake.

Character encoding schemes for East Asian coded character sets

Here's where things start to get interesting. Rather than just treat the East Asian standards as two-byte encodings along the lines above, various packing schemes have been developed, either to save space or to combine characters from two or more of the coded character set standards described above. There are four main categories of these character encoding schemes:

The most familiar of all of these is **Shift-JIS**. It combines the JIS X 0201 and JIS X 0208 coded character sets. The JIS X 0201 code point values are represented as-is as single-byte values. The byte values not used by 201 (the C1 range and the top two rows of the G1 range: 0x80–0x9F and 0xE0–0xEF) are used as leading byte values for two-byte sequences representing the characters from 208. (The second byte would be in the range 0x40–0x7E or 0x80–0xFC.) The transformation from row and column numbers in JIS X 0208 to a two-byte Shift-JIS sequence isn't straightforward, but it's enough to cover all the characters. In fact, Shift-JIS also adds a space for user-defined characters: these are also represented with two-byte sequences, the leading byte of which runs from 0xF0 to 0xFC.

One interesting consequence of Shift-JIS (and the other encodings we'll look at) is that it includes both JIS X 0201 and JIS X 0208, even though all the characters in 201 are also included in 208. They've come to have different meanings: the codes from 201 (the single-byte sequences) are often thought of as “half-width” (taking up half the horizontal space of a typical *kanji* character), which the codes from 208 (the two-byte sequences) are thought of as “full-width” (taking up the full horizontal space of a *kanji* character).

The next most familiar encoding scheme is **Extended UNIX Code**, which, as the name suggests, was originally developed for use with UNIX systems and was standardized in 1991 by a group of UNIX vendors. It makes use of variable-length sequences ranging from one to four bytes and based on the ISO 2022 structure. EUC defines four separate code sets. Code set 0 is always ASCII (or some other national variant of ISO 646), and is represented as-is in the C0 and G0 spaces. Code sets 1, 2, and 3 are represented using the values in the G1 space. Characters in code set 1 are represented unadorned. Characters in code set 2 are preceded with a C1 control character, the SS2 character, 0x8E. Characters in code set 3 are preceded with a different C1 control characters, SS3 or 0x8F. Characters in code set 0 are always represented with one-byte sequences; characters in code sets are always sequences of at least two bytes (and can be as many as four—this is because of the SS2 or SS3 byte appended to the front); characters in code set 1 may be sequences of from one to three bytes.

The assignment of specific characters to the three EUC code sets is locale-specific, so there are variants of EUC for each locale, as follows:

EUC-JP (Japan):	Code set 0:	JIS Roman (i.e., the Japanese variant of ISO 646) [one byte]
	Code set 1:	JIS X 0208, with the row and column numbers shifted up into the G1 range [two bytes]
	Code set 2:	Katakana from JIS X 0201 [SS2 plus one byte]
	Code set 3:	JIS X 0212, with the row and column numbers shifted up into the G1 range [SS3 plus two bytes]
EUC-CN (China):	Code set 0:	Chinese ISO 646 [one byte]
	Code set 1:	GB 2312, with the row and column numbers shifted up to the G1 range [two bytes]

Multiple-byte encoding systems

Code sets 2 and 3: Not used

EUC-TW
(Taiwan):

Code set 0: Taiwanese ISO 646 [one byte]

Code set 1: CNS 11643, plane 1, with the row and column numbers shifted up into the G1 range [two bytes]

Code set 2: All of CNS 11643, the row and column numbers as above, preceded by the plane number, also shifted into the G1 range [SS2 plus three bytes]

Code set 3: Not used.

EUC-KR (Korea):

Code set 0: Korean ISO 646 [one byte]

Code set 1: KS X 1001, with the row and column numbers shifted up into the G1 range [two bytes]

Code sets 2 and 3: Not used.

There there's the family of 7-bit **ISO 2022** encodings. Here you use the standard ISO 2022 machinery and the normal framework of 7-bit ISO 2022 code points to represent the characters. Again, depending on locale, you get different default behavior. The basic idea here is that the system is modal and you can use the ASCII SI and SO (shift in and shift out, 0x0F and 0x0E) to switch back and forth between single-byte mode and double-byte mode. You can also use the SS2 and SS3 codes from ISO 2022 (here in their escape-sequence incarnations, since we're staying within seven bits) to quote characters from alternate two-byte character sets. The various national variants of ISO 2022 then make use of various escape sequences (often running to four bytes) to switch the two-byte mode from one coded character set to another. The various versions give you access to the following character sets:

ISO-2022-JP	JIS-Roman and JIS X 0208
ISO-2022-JP-1 and ISO-2022-JP-2	ISO-2022-JP plus JIS X 0212
ISO-2022-CN	ASCII, GB 2312, CNS 11643 planes 1 and 2
ISO-2022-CN-EXT	ISO-2022-CN plus the rest of CNS 11643
ISO-2022-KR	ASCII and KS X 1001

The actual specifications vary widely from one variant to the next, and they don't all strictly follow the ISO 2022 standard (despite the name). See Ken Lunde's *CJKV Information Processing* for all the gory details.

Finally, there's **HZ**, another 7-bit encoding, which uses sequences of printing characters to shift back and forth between single-byte mode (ASCII) and double-byte mode (GB 2312). The sequence `~{` would switch into double-byte mode, and the sequence `~}` would switch back.

Other East Asian encoding systems

Finally, there are two more important encoding standards that don't fit neatly into the coded-character-set/character-encoding-scheme division we're been considering.

The Taiwanese **Big5** standard is an industry standard (so named because it was backed by five large computer manufacturers) that predates the CNS 11643, the "official" Taiwanese standard, dating back to 1984. (CNS 11643 is downward compatible with Big5, and can be thought of as an extended and improved version of it.) Big5 uses a 94×157 encoding space, providing room for 14,758 characters. It's a variable-length encoding compatible with the ISO 2022 layout: The Taiwanese version of ISO 646 is represented as-is in the C0 and G0 spaces with single-byte sequences. The Chinese characters (and all the usual other stuff) is represented with two-byte sequences: the first byte is in the G1 range; the second byte can be in either the G1 or the G0 range (although not all G0 values are used as trailing bytes).

CCCII ("Chinese Character Code for Information Interchange") is another Taiwanese industry standard that predates Big5 (it was published in 1980) and also influenced CNS 11643. It's a fixed-length encoding where each code point maps to *three* seven-bit code units. It has a $94 \times 94 \times 94$ encoding space: sixteen *layers*, each consisting of six 94×94 planes. The first byte of the code point is the layer and plane numbers, the second the row number, and the third the column number.

Johab is an alternate encoding scheme described in the KS X 1001 standard. KS X 1001 normally encodes only 2,350 Hangul syllables (see Chapter 10 for more on Hangul) actually used in Korean, but it proposes an alternate method that permits encoding of all 11,172 possible modern Hangul syllables (the Korean Hangul script is an alphabetic system, but the letters are commonly grouped into syllabic blocks, which are encoded as units in most Korean encoding systems in a manner analogous to that used for the Chinese characters). It takes advantage of the alphabetic nature of Hangul: instead of just giving each syllable an arbitrary number, the syllable is broken down into its component letters, or *jamo* (there can be up to three in a syllable), each of which can be encoded in five bits. This gives you a fifteen-bit character code, which can be encoded in a two-byte sequence. The G0 space in the leading byte is taken up with the Korean version of ISO 646, as always, but the leading byte of a two-byte sequence invades the C1 range and the second byte can be from the G0, C1, or G1 ranges. Thus, the Johab encoding doesn't fit the ISO 2022 structure.

ISO 10646 and Unicode

It should be fairly clear by now that we have a real mess on our hands. There are literally hundreds of different encoding standards out there, many of which are redundant, encoding the same characters but encoding them differently. Even within the ISO 8859 family, there are ten different encodings of the Latin alphabet, each containing a slightly different set of letters.

This means you have to be very explicit about which encoding scheme you're using, lest the computer interpret your text as characters other than the characters you intend, garbling it in the process (log onto a Japanese Web site with an American computer, for example, and it's likely you'll see gobbledygook rather than Japanese). About the only thing that's safely interchangeable across a wide variety of systems is 7-bit ASCII, and even then you're not safe—you might run into

A Brief History of Character Encoding

a system using one of the national variants of ISO 646, mangling the “national use” characters from ASCII, or your text might pass through an EBCDIC-based system and get completely messed up.

You also can’t easily mix characters from different languages (other than English and something else) without having auxiliary data structures to keep track of which encodings different pieces of the same text are in. Your other choice is to use a code switching scheme such as the one specified in ISO 2022, but this is cumbersome both for the user and for the processing software.

So you have a bunch of problems it’d be really nice to fix:

- A wide variety of mutually-incompatible methods of encoding the various languages, and no way to tell from the data itself which encoding standard is being followed.
- A limited encoding space, leading to an even wider variety of conflicting encoding standards, none of which is terribly complete.
- Variable-length encoding schemes that don’t use numerically disjoint code unit values for the different pieces of a multi-byte character, making it difficult to count characters or locate character boundaries, and which make the data highly susceptible to corruption. (A missing byte can trigger the wrong interpretation of every byte that follows it.)
- Stateful in-band code-switching schemes such as ISO 2022 are complicated, make encoded text in memory harder to process, and require arbitrary-length look-back to understand how to interpret a given byte on random access.

Clearly what was needed was an universal character encoding system, one that would solve these problems and assign an unique code point value to every character in common use. A group formed in 1984 under the auspices of the International Organization for Standardization and the International Electrotechnical Commission to create such a thing. The group goes by the rather unwieldy name of ISO/IEC JTC1/SC2/WG2 (that’s “ISO/IEC Joint Technical Committee #1 [Information Technology], Subcommittee #2 [Coded Character Sets], Working Group #2 [Multi-octet codes]”), or just “WG2” for short. They began working on what became known as ISO 10646 (the connection between 10646 and 646 is not accidental: the number signifies that 10646 is an extension of 646).¹³

In 1988 a similar effort got underway independently of WG2. This was sparked by a paper written by Joe Becker of Xerox setting forth the problems listed above and proposing a list of desirable characteristics for a universal character encoding standard. The Becker paper was the first to use the name “Unicode,” which has stuck. An informal group of internationalization experts from Xerox, Apple, and a few other companies got together to begin work on the Unicode standard (beginning from Xerox’s earlier XCCS encoding standard). This group grew into the Unicode Consortium of today.

Although they had similar goals, they went about them in very different ways. The original version of ISO 10646 was based on 32-bit abstract code point values, arranged into a $256 \times 256 \times 256 \times 256$ encoding space, but not all code point values were actually possible: byte values from 0x00 to

¹³ This section draws on Mike Ksar, “Unicode and ISO 10646: Achievements and Directions,” *Proceedings of the 13th International nicode Conference*, session B11, September 11, 1998, and the Edberg paper, *op. cit.* The information on the early history of ISO 10646 comes from contributions by Mark Davis and Ken Whistler to online discussions by the IETF’s Internationalized Domain Names task force, and were forwarded to me by Mark Davis.

0x1F and 0x80 to 0x9F, corresponding to the C0 and C1 areas in ISO 2022, were illegal in any position in the four-byte code point values, effectively yielding a $192 \times 192 \times 192 \times 192$ encoding space. This was referred to as 192 “groups” of 192 “planes” of 192 “rows” of 192 “cells.”

Recognizing that a scheme with four-byte character codes, involving a quadrupling of the space required for a given piece of text, probably wouldn’t go over well, the early versions of ISO 10646 had a whopping five encoding forms, based on code units of from one to four bytes in size, plus a variable-length scheme that allowed single code points to be represented with sequences of one-byte code units of varying lengths. There were also encoding schemes for the shorter encoding forms that used ISO 2022-style escape sequences to switch the meanings of the code units from one part of the encoding space to another.

The division into “plane” was intended to facilitate the use of shorter encoding forms: If all your characters came from the same plane, you could omit the plane and group numbers from the code unit and just have it consist of the row and cell numbers. There was a “basic multilingual plane” (“BMP” for short) the contained the characters from most of the modern writing systems, the big exception being the Han characters. The Han characters were organized into several planes of their own: for Traditional Chinese, Simplified Chinese, Japanese, and Korean. There was also provision for various special-purpose and private-use planes.

The original version of ISO 10646, with these features, failed to get approval on its first ballot attempt. The prohibition of C0 and C1 byte values not only restricted the available space too much (especially in the basic multilingual plane), but it also caused problems for C implementations because zero-extended ASCII values weren’t legal 10646 code points, making it difficult or impossible to use `wchar_t` to store 10646 code points and maintain backward compatibility. No one liked the multitude of different encoding forms and the escape sequences that were to be used with them. And the Asian countries didn’t like having their writing systems relegated off to their own planes.¹⁴

Meanwhile, the Unicode Consortium started with a 16-bit code point length, equivalent to a single plane of the ISO 10646 encoding space, and declared no byte values off limits. There were no separate encoding forms or schemes, no ambiguous interpretations of variable lengths, just Unicode. (Encoding forms such as UTF-8 were a later development.)

One of the big differences between Unicode and ISO 10646 was the treatment of the Chinese characters. Even though there are certainly significant differences in the use of the Chinese characters in China, Japan, Korean, Vietnam, and Taiwan, they’re still essentially the same characters—there’s great overlap between the five different sets of Chinese characters, and it doesn’t make sense to encode the same character once for each language any more than it makes sense to encode the letter A once each for English, Spanish, French, Italian, and German.

Work had already been going on at Xerox and Apple to devise a unified set of Chinese characters, and this work followed work that had already been going on in China (CCCII, for example

¹⁴ One of Ken Whistler’s messages in the discussion I’m basing this on actually quotes the Chinese national body’s objections. The Chinese not only didn’t like havng their codes in a separate plane, but didn’t like the idea that a single Han character could have one of several code point values depending on the language of the text. It advocates a unified Han area located in the BMP. This seems to put the lie rather effectively to the allegations you hear every so often that Han unification was forced on the Asians by the Americans.

A Brief History of Character Encoding

attempts to unify the various national versions of the Chinese characters, and another effort was underway in China). The Apple/Xerox effort went into early drafts of Unicode and was proposed to ISO as a possible addition to ISO 10646. The Apple/Xerox efforts eventually merged with the various Chinese efforts to form the CJK Joint Research Group, which produced the first version of the unified Chinese character set in ISO 10646 in 1993. This group, now called the Ideographic Rapporteur Group (or “IRG”) is still responsible for this section of ISO 10646 and Unicode and is now a formal subcommittee of WG2.

In addition to unifying the set of Chinese characters, the original idea was to further conserve encoding space by building up things like Korean syllables, rarer Chinese characters, and rarer accented-letter combinations out of pieces: You’d take more than one code point to represent them, but each code point value would have a defined meaning independent of the others. Unicode also had a Private Use Area, which could be used for characters with more specialized applications (writing systems no longer in common use, for example). A higher-level protocol outside of Unicode would be used to specify the meaning of private-use code points.

After the initial ISO 10646 ballot failed, talks began between the two groups to merge their technologies. 10646 retained its 32-bit-ness and its organization into groups, planes, rows, and cells, but it removed the prohibition against C0 and C1 byte values, opening up the entire 32-bit space for encoding. (Actually, it turned into a 31-bit space: code point values with the sign bit turned on were prohibited, making it possible to store an ISO 10646 code point in either a signed or an unsigned data type without worrying about mangling the contents.) The multitude of encoding forms and schemes was simplified: the ISO 2022-style escape sequences were eliminated, as were the 1-byte, 3-byte and variable-byte encodings, leaving just UCS-4, a straightforward representation of the code point values in bytes, and UCS-2, which allowed the Basic Multilingual Plane to be represented with two-byte code unit values.

Both ISO 10646 and Unicode wound up adding and moving characters around: Unicode was forced to give up the idea of exclusively building up Korean syllables and rare Chinese characters out of pieces, and ISO 10646 moved to a unified Chinese-character repertoire and located it in the BMP. From this point, in 1991, on forward, Unicode and ISO 10646 have had their character repertoires and the code points those characters are assigned to synchronized. Unicode 1.1, minus the unified Chinese repertoire, was the first standard published after the merger. The unified Chinese repertoire was finished in 1993, at which point the first version of 10646 was published as ISO/IEC 10646-1:1993. The standard have remained synchronized ever since.

Neither standard went away in the merger; Unicode and ISO 10646 still exist as separate, distinct standards. But great effort is expended to make sure that the two standards remain in sync with each other: there’s a lot of overlapping membership on the two governing committees, and both sides care a lot about staying in sync with each other. From the time of the merger until now, the character repertoires of the two standards have remained effectively identical.

ISO 10646’s organization into planes eventually had to be adopted into Unicode to accommodate the rapidly growing character repertoire. Unicode 2.0 added something called the “surrogate mechanism” (now known as “UTF-16”) to permit representation of characters outside the BMP in a 16-bit framework (a range of code point values was set aside to be used in pairs to represent characters outside the BMP). This permitted access to the first seventeen planes of ISO 10646. WG2 agreed to declare the planes above Plane 16 off limits, and the encoding spaces were also synchronized.

So starting in Unicode 2.0, Unicode's encoding space has expanded from a single 256×256 plane, corresponding to ISO 10646's BMP, to seventeen 256×256 planes, corresponding to Planes 0 to 16 in ISO 10646. This means a Unicode code point value is now a 21-bit value. Unicode and ISO 10646 both also now define a set of three character encoding forms, UTF-32, UTF-16, and UTF-8, which map the 21-bit code points into single 32-bit code units, variable-length sequences of 1 or 2 16-bit code units, and variable-length sequences of from 1 to 4 8-bit code units, respectively. They further go on to define seven encoding schemes (UTF-8 and three flavors each of UTF-16 and UTF-32, arranged for machine architectures with different byte ordering).

While they now define the same coded character set and character encoding forms and schemes, Unicode and ISO 10646 differ in content and in how updates are issued. Among the important differences:

- ISO 10646 defines a formalized method of declaring which characters are supported by an implementation; Unicode leaves this to ad-hoc methods.
- Unicode defines a whole host of semantic information about each of the characters and provides a whole bunch of implementation guidelines and rules; ISO 10646 basically just lists the characters and their code point assignments.
- The Unicode standard can be picked up in any bookstore (and updates between books are published on the Web); ISO 10646 has to be ordered from an ISO member organization.
- The ISO 10646 standard is divided into two parts: ISO 10646-1, which defines the basic architecture and standardizes the code assignments in the Basic Multilingual Plane, and ISO 10646-2, which standardizes the code assignments in the other sixteen planes, which Unicode is published as a single unified standard.

The other important difference is the way the standards are updated. Unicode has a version-numbering scheme; minor updates (where the number to the right of the decimal point increases) are published as technical reports, while major updates result in a whole new edition of the book being published. New editions of ISO 10646 get published far less often, but amendments and corrigenda are issued more frequently.

This means that one standard can periodically get ahead of the other. Usually ISO 10646 runs a little ahead as amendments get adopted—every so often, a bunch of amendments are gathered together to make a new Unicode version. Still, there's always a direct mapping from a particular version of one standard to a particular version of the other. For the versions that have been published to date, that mapping is:

Unicode 1.0	<i>(pre-merger)</i>
Unicode 1.1	ISO 10646-1:1993
Unicode 2.0	ISO 10646-1:1993, plus amendments 1 thru 7
Unicode 2.1	ISO 10646-1:1993, plus amendments 1 thru 7 and 18
Unicode 3.0	ISO 10646-1:2000 (which was formed from ISO 10646-1:1993 and amendments 1 thru 31)
Unicode 3.1	ISO 10646-1:2000 and ISO 10646-2:2001

Unicode 3.2 ISO 10646-1:2000 and amendment 1, plus ISO 10646-2:2001

How the Unicode standard is maintained

Before we go ahead to look at the nuts and bolts of Unicode itself, it's worth it to take a few minutes to look at how the Unicode Consortium actually works and the process by which changes are made to the standard.¹⁵

The Unicode Consortium is a nonprofit corporation located in Silicon Valley. It's a membership organization—members pay dues to belong and have a say in the development of the standard. Membership is pretty expensive, so full voting members are generally corporations, research and educational institutions, and national governments. (There are also various levels of non-voting membership, many of which are individuals.)

The Unicode Consortium basically exists to maintain and promote the Unicode standard. The actual work of maintaining the Unicode standard falls to the Unicode Technical Committee ("UTC" for short). Each member organization gets to appoint one principal and one or two alternate representatives to the Unicode Technical Committee. Representatives of non-voting members can also send people to the UTC who can participate in discussion but not vote. Other people are sometimes permitted to observe but not participate.

The UTC usually meets four or five times a year to discuss changes and additions to the standard, make recommendations to various other standard-setting bodies, issue various resolutions, and occasionally issue new standards. A change to the standard requires a majority vote of the UTC, and a precedent-setting change requires a two-thirds vote.

The UTC doesn't take its job lightly; only the most trivial of actions goes through the whole process in a single UTC meeting. Often proposals are sent back to their submitters for clarification or referred to various outside organizations for comment. Since many other standards are based on the Unicode standard, the organizations responsible for these dependent standards get a prominent place at the table whenever architectural changes that might affect them are discussed. Additions of new writing systems (or significant changes to old ones) require input from the affected user communities. For certain areas, either the UTC or WG2 (or both) maintain standing subcommittees: the Ideographic Rapporteur Group (IRG) is a WG2 subcommittee charged with maintaining the unified Chinese-character repertoire, for example.

A lot of work also goes into keeping Unicode aligned with ISO 10646. Because the Unicode Consortium is an industry body, it has no vote with ISO, but it works closely with the American national body, Technical Committee L2 of the National Committee on Information Technology Standards, an ANSI organization. L2 *does* have a vote with ISO; its meetings are held jointly with the UTC meetings, and most UTC members are also L2 members.

ISO 10646 is maintained by ISO JTC1/SC2/WG2, which is made up of representatives from various national bodies, including L2. Many WG2 members are also UTC members. For things that affect both Unicode and 10646, such as the addition of new characters, both the UTC and WG2 have to agree—effectively, each organization has a veto over the other's actions. Only after both the UTC and WG2 agree on something does it go through the formal ISO voting process. A typical proposal takes two years to wend its way through the whole process, from submission to

¹⁵ The information in this section is drawn from Rick McGowan, "About Unicode Consortium Procedures, Policies, Stability, and Public Access," published as an IETF Internet-Draft, <http://search.ietf.org/internet-drafts/draft-rmcgowan-unicode-procs-00.txt>.

WG2 or the UTC to being adopted as an international standard (or amendment) by an international ballot.

The Unicode Consortium maintains two Internet mailing lists: the “unicore” list, which is limited to Unicode Consortium members, and the “unicode” list, which is open to anybody (go to the Unicode Web site for information on how to join). A lot of proposals start as information discussions on one of these two lists and only turn into formal proposals after being thoroughly hashed out in one of these information discussions. In fact, it’s generally recommended that an idea be floated on the mailing lists before being submitted formally—this cuts down on the back-and-forth after the proposal has been sent to the UTC.

Anyone, not just UTC or Unicode Consortium members, can submit a proposal for a change or addition to the Unicode standard. You can find submission forms and instructions on the Unicode Web site. Don’t make submissions lightly, however: To be taken seriously, you should understand Unicode well enough first to make sure the submission makes sense within the architecture and the constraints Unicode works under, and you’ll need to be prepared to back up the proposal with evidence of why it’s the right thing to do and who needs it. You also need patience: it takes a long time for a proposal to get adopted, and shepherding a proposal through the process can be time consuming.

More background than you really wanted, huh? Well, it’ll be helpful to be armed with it as we march on ahead to explore the Unicode standard in depth. Onward...

CHAPTER 3 *Architecture: Not Just a Pile of Code Charts*

If you're used to working with ASCII or other similar encodings designed for European languages, you'll find Unicode noticeably different from other character encoding standards, and you'll find that when you're dealing with Unicode text, various assumptions you may have made in the past about how you deal with text don't hold. If you've worked with encodings for other languages, at least some characteristics of Unicode will be familiar to you, but even then, some pieces of Unicode will be unfamiliar.

Unicode is more than just a big pile of code charts. To be sure, it *includes* a big pile of code charts, but Unicode goes much further than that. It doesn't just take a bunch of character forms and assign numbers to them; it adds a wealth of information on what those characters mean and how they are used.

The main reason for this is that unlike virtually all other character encoding standards, Unicode isn't designed for the encoding of a single language or a family of closely related languages; Unicode is designed for the encoding of *all* written languages. The current version doesn't give you a way to encode *all* written languages (and in fact, this is such a slippery thing to define that it probably never will), but it gives you a way to encode an extremely wide variety of languages. The languages vary tremendously in how they are written, and so Unicode must be flexible enough to accommodate all of them. This necessitates rules on how Unicode is to be used with each language it works with. Also, because the same encoding standard can be used for all these different languages, there's a higher likelihood that they will be mixed in the same document, requiring rules on how text in the different languages interacts. The sheer number of characters requires special attention, as does the fact that Unicode often provides multiple ways of representing the same thing.

The idea behind all the rules is simple: to ensure that a particular sequence of code points will get drawn and interpreted the same way (or in semantically-equivalent ways) by all systems that handle

Unicode text. In other words, it's not so that there should only be one way to encode "à bientôt," but that a particular sequence of code points that represents "à bientôt" on one system will also represent it on any other system that purports to understand the code points used. This also doesn't mean that every system has to handle that sequence of code points *exactly the same*, but merely that it interpret it as meaning the same thing. In English, for example, there are tons of different typographical conventions you can follow when you draw the word "carburetor," but someone who reads English would still interpret all of them as the word "carburetor." Any Unicode-based system has wide latitude in how it deals with a sequence of code points representing the word "carburetor," as long as it still treats it as the word "carburetor."

All of this means that there's a lot more that goes into supporting Unicode text than supplying a font with the appropriate character forms for all the characters. The whole purpose of this book is to explain all these other things you have to be aware of (or at least *might* have to be aware of). This chapter will highlight all the things which are special about Unicode and attempt to tie them together into a coherent architecture.

The Unicode Character-Glyph Model

The first and most important thing to understand about Unicode is what is known as the *character-glyph model*. Up until the introduction of the Macintosh in 1984, text was usually displayed on computer screens in a fairly simple fashion. The screen would be divided up into a number of equally-sized *display cells*. The most common video mode on the old IBM PCs, for example, had 25 rows of 80 display cells each. There was a video buffer in memory that consisted of 2,000 bytes, one for each display cell. The video hardware would contain a *character generator* chip that contained a bitmap for each possible byte value, and this chip was used to map from the character codes in memory to a particular set of lit pixels on the screen.

Handling text was simple. There were 2,000 possible locations on the screen, and 256 possible characters to put in them. All the characters were the same size, and were laid out regularly from the left to the right across the screen. There was a one-to-one correspondence between character codes stored in memory and visible characters on the screen, and there was a one-to-one correspondence between keystrokes and characters.

We don't live in that world anymore. One reason is the rise of the WYSIWYG ("what you see is what you get" text editor, where you can see on the screen exactly what you want to see on paper, meaning video displays have to be able to handle proportionally-spaced fonts, the mixing of different typefaces, sizes, and styles, and the mixing of text with pictures and other pieces of data. The other reason is that the old world of simple video-display terminals can't handle many languages, which are more complicated to write than the Latin alphabet is.

This means there's been a shift away from translating character codes to pixels in hardware and toward doing it in software. And the software for doing this has become considerably more sophisticated.

The main consequence of this is that on modern computer systems, Unicode or no, there is no longer always a nice simple one-to-one relationship between character codes stored in your computer's memory and actual shapes drawn on your computer's screen. This is important to understand because Unicode *requires* this—a Unicode-compatible system cannot be designed to assume a one-to-one

correspondence between code points in the backing store and marks on the screen (or on paper), or between code points in the backing store and keystrokes in the input.¹⁶

So let's start by defining two concepts: character and glyph. A character is an atomic unit of text with some semantic identity; a glyph is a visual representation of that character.

Let's consider an analogy. Consider the following examples:

13 thirteen 1.3×10^1
dreizehn \$0C Kγ 十三
?? treize 𠄎𠄎𠄎 𐌆𐌆

[The picture got mangled in the pasting process—the ?? is supposed to be native Arabic digits.]

These are eleven different visual representations of the number 13. The underlying semantic is the same in every case: the concept “thirteen.” These are just different ways of depicting the concept of “thirteen.”

Now consider the following:

g g g g

Each of these is a different presentation of the Latin lowercase letter g. To go back to our words, these are all the same *character* (the lowercase letter g), but four different *glyphs*.

Of course, I got these four glyphs by taking the small *g* out of four different typefaces. That's because there's generally only one glyph per character in a Latin typeface. But in other writing systems, that isn't true. The Arabic alphabet, for example, joins cursively even when printed. This isn't an optional feature, as it is with the Latin alphabet; it's the way the Arabic alphabet is *always* written.

٥ ٤ ٣ ٢

¹⁶ Technically, if you restrict the repertoire of characters your system supports enough, you actually *can* make this assumption, but at that point, you're likely back to being a fairly simplistic English-only system, in which case why bother with Unicode in the first place?

These are four different forms of the Arabic letter *heh*. The first is how the letter looks in isolation. The second is how it looks when it joins only to a letter on its right (usually at the end of a word). The third is how it looks when it joins to letters on both sides in the middle of a word. And the last form is how the letter looks when it joins to a letter on its left (usually at the beginning of a word).

Unicode only provides one character code for this letter,¹⁷ and it's up to the code that draws it on the screen (the *text rendering process*) to select the appropriate glyph depending on context. The process of selecting from among a set of glyphs for a character depending on the surrounding characters is called *contextual shaping*, and it's required in order to draw many writing systems correctly.

There's also not always a one-to-one mapping between character and glyphs. Consider the following example:

fi

This, of course, is the letter f followed by the letter i, but it's a single glyph. In many typefaces, if you put a lowercase f next to a lowercase i, the top of the f tends to run into the dot on the i, so often the typeface includes a special glyph called a *ligature* that represents this particular pair of letters. The dot on the i is incorporated into the overhanging arch of the f, and the crossbar of the f connects to the serif on the top of base of the i. Some desktop-publishing software and some high-end fonts will automatically substitute this ligature for the plain f and i.

In fact, some typefaces include additional ligatures. Other forms involving the lowercase f are common, for example. You'll often see ligatures for a-e and o-e pairs (useful for looking erudite when using words like “archæology” or “œnophile”), although software rarely forms these automatically (æ and œ are actually separate letters in some languages, rather than combinations of letters), and some fonts include other ligatures for decorative use.

Again, though, ligature formation isn't just a gimmick. Consider the Arabic letter *lam* (ﻝ) and the Arabic letter *alef* (ﺀ). When they occur next to each other, you'd expect them to appear like this if they followed normal shaping rules:

ﻝﺀ

But they actually don't. Instead of forming a U shape, the vertical strokes of the lam and the alef actually cross, forming a loop at the bottom like this:

ﻝﺀ

¹⁷ Technically, this isn't true—Unicode actually provides separate codes for each glyph in Arabic, although they're only included for backwards compatibility. They're called “presentation forms” and encoded in a separate numeric range to emphasize this. Implementations generally aren't supposed to use them. Almost all other writing systems that have contextual shaping don't have separate presentation forms in Unicode.

Unlike the *f* and *i* in English, these two letters *always* combine together when they occur together. It's not optional. The form that looks like a *U* is just plain wrong. So ligature formation is a required behavior for writing many languages.

A single character may also split into more than one glyph. This happens in some Indian languages, such as Tamil. It's very roughly analogous to the use of the silent *e* in English. The *e* at the end of "bite," for example, doesn't have a sound of its own; it merely changes the way the *i* is pronounced. Since the *i* and the *e* are being used together to represent a single vowel sound, you could think of them as two halves of a *single* vowel character. This is sort of what happens in languages like Tamil. Here's an example of a Tamil split vowel:

௫௪

This looks like three letters, but it's really only two. The middle glyph is a consonant, the letter ழ. The vowel ஫ is shown with a mark on either side of the ழ. This kind of thing is required for the display of a number of languages.

So there's not always a simple straightforward one-to-one mapping between characters and glyphs. Unicode assumes the presence of a character rendering process capable of handling the sometimes complex mapping from characters to glyphs. It doesn't provide separate character codes for different glyphs that represent the same character, or for ligatures representing multiple characters.

Exactly how this all works varies from writing system to writing system (and to a lesser degree from language to language within a writing system). For all the details on just how Unicode deals with the peculiar characteristics of each writing system it encodes, see Section II (Chapters 7 to 12).

Character positioning

Another assumption that has to go is the idea that characters are laid out in a neat linear progression running in lines from left to right. In many languages, this isn't true.

Many languages also make use of various kinds of diacritical marks which are used in combination with other characters to indicate pronunciation. Exactly where the marks get drawn can depend on what they're being attached to. For example, look at these two letters:

ä Ä

Each of these examples is an *a* with an umlaut placed on top of it. But the umlaut needs to be positioned higher when attached to the capital *A* than when attached to the small *a*.

This can get more complicated when multiple marks are attached to the same character. In Thai, for example, a consonant with a tone mark might look like this:

ع

But if the consonant also has a vowel mark attached to it, the tone mark has to move out of the way. It actually moves up and gets smaller when there's a vowel mark:

عَ

Mark positioning can get quite complicated. In Arabic, there's a whole host of dots and marks that can appear along with the actual letters. There are dots that are used to differentiate the consonants from one another when they're written cursively, there are diacritical marks that modify the pronunciation of the consonants, there may be vowel marks (Arabic generally doesn't use letters for vowels—they're either left out or shown as marks attached to consonants), and there may also be reading or chanting marks attached to the letters. In fact, some Arabic calligraphy includes other marks that are purely decorative. There's a hierarchy of how these various marks are placed relative to the letters that can get quite complicated when all the various marks are actually being used.



Unicode expects, again, that a text rendering process will know how to position marks appropriately. It generally doesn't encode mark position at all—it adopts a single convention that marks follow in memory the characters they attach to, but that's it.¹⁸

But it's not just diacritical marks that may have complicated positioning. Sometimes the letters themselves do. For example, many Middle Eastern languages are written from right to left rather than from left to right (in Unicode, the languages that use the Arabic, Hebrew, Syriac, and Thaana alphabets are written from right to left). Unicode stores these characters in the order they'd be spoken or typed by a native speaker of one of the relevant languages. This is known as *logical order*.

Logical order means that the “first” character in character storage is the character that a native user of that character would consider “first.” For a left-to-right writing system, the “first” character is drawn furthest to the left. (For example, the first character in this paragraph is the letter L, which is the character furthest to the left on the first line.) For a right-to-left writing system, the “first” character would be drawn furthest to the right. For a vertically-oriented writing system, such as that used to write Chinese, the “first” character is drawn closest to the top of the page.

This is in contrast with “visual order,” which assumes that all characters are drawn progressing in the same direction (usually left to right). When text is stored in visual order, text that runs counter to the direction assumed (i.e., right-to-left text) is stored in memory in the reverse of the order in which it was typed.

¹⁸ The one exception to this rule has to do with multiple marks attached to a single character. In the absence of language-specific rules governing how multiple marks attach to the same character, Unicode adopts a convention that marks that would otherwise collide radiate outward from the character they're attached to in the order they appear in storage. This is covered in depth in Chapter 4.

Character Positioning

Unicode doesn't assume any bias in layout direction. The characters in a Hebrew document are stored in the order they are typed, and Unicode expects that the text rendering process will know that because they're Hebrew letters, the first one in memory should be positioned the furthest to the right, with the succeeding characters progressing leftward from there.

This gets really interesting when left-to-right text and right-to-left text mix in the same document. Say you have an English sentence with a Hebrew phrase embedded into the middle of it:

Avram said **מזל טוב** and smiled.

Even though the dominant writing direction of the text is from left to right, the first letter in the Hebrew phrase (מ) still goes to the right of the other Hebrew letters—the Hebrew phrase still reads from right to left. The same thing can happen even when you're not mixing languages: in Arabic and Hebrew, even though the dominant writing direction is from right to left, numbers are still written from left to right.

This can get even more fun when you throw in punctuation. Letters have inherent directionality; punctuation doesn't. Instead, punctuation marks take on the directionality of the surrounding text. In fact, some punctuation marks (such as the parentheses) actually *change shape* based on the directionality of the surrounding text (this is called *mirroring*, since the two shapes are usually mirror images of each other). Mirroring is another example of how Unicode encodes meaning rather than appearance—the code point encodes the meaning (“starting parenthesis”) rather than the shape (which can be either ‘(’ or ‘)’) depending on the surrounding text).

Dealing with mixed-directionality text can get quite complicated, not to mention ambiguous, so Unicode includes a set of rules that govern just how text of mixed directionality is to be arranged on a line. The rules are rather involved, but are required for Unicode implementations that claim to support Hebrew and Arabic.¹⁹

The writing systems used for the various languages used on the Indian subcontinent and in Southeast Asia have even more complicated positioning requirements. For example, the Devanagari alphabet used to write Hindi and Sanskrit treats vowels as marks that get attached to consonants (which get treated as “letters”). But a vowel may attach not just to the top or bottom of the consonant, but also to the left or right side. Text generally runs left to right, but when you get a vowel that attaches to the left-hand side of its consonant, you get the effect of a character appearing “before” (i.e., to the left of) the character it logically comes “after.”

क + ि = कि

In some alphabets, a vowel can actually attach to the left-hand side of a group of consonants, meaning this “reordering” may actually involve more than just two characters switching places. Also, in some alphabets, such as the Tamil example we looked at earlier, a vowel might actually appear on *both* the left- and right-hand sides of the consonant it attaches to (this is called a “split vowel”).

¹⁹ This set of rules, the Unicode Bidirectional Text Layout Algorithm, is explained in a lot more detail later in this book: a high-level overview is in Chapter 8, and a look at implementation strategies is in Chapter 16.

Again, Unicode stores the characters in the order they're spoken or typed; it expects the display engine to do this reordering. For more on the complexities of dealing with the Indian scripts and their cousins, see Chapter 9.

Chinese, Japanese, and Korean can be written either horizontally or vertically. Again, Unicode stores them in logical order, and again, the character codes encode the semantics. So many of the punctuation marks used with Chinese characters have a different appearance when used with horizontal text than when used with vertical text (some are positioned differently, some are rotated ninety degrees). Horizontal scripts are sometimes rotated ninety degrees when mixed into vertical text and sometimes not, but this distinction is made by the rendering process and not by Unicode.

Japanese and Chinese text may include annotations (called “ruby” or “furigana” in Japanese) that appear in between the lines of normal text. Unicode includes ways of marking text as ruby and leaves it up to the rendering process to determine how to draw it.

For more on vertical text and ruby, see Chapter 10.

The Principle of Unification

The bottom-line philosophy you should draw from the discussions on the character-glyph model and on character positioning is that *Unicode encodes semantics, not appearances*. In fact, the Unicode standard specifically states that the pictures of the characters in the code charts are there for illustrative purposes only—the pictures of the characters are there to help clarify the meaning of the character code, not to specify the appearance of the character with that code.

The philosophy that Unicode encodes semantics and not appearances also undergirds the principle that Unicode is a plain-text encoding, which we discussed in Chapter 1. The fact that an Arabic letter looks different depending on the letters around it doesn't change what letter it is, and thus doesn't justify different codes for the different shapes. The fact that the letters lam and alef combine into a single mark when written doesn't change the fact that the word still contains the letters lam and alef in succession. The fact that text from some language might be combined on the same line with text from another language whose writing system runs the opposite direction doesn't justify storing either language's text in some order other than the order in which the characters are typed or spoken. In all of these cases, Unicode encodes the underlying meaning and leaves it up to the process that draws it to be smart enough to do so properly.

The philosophy of encoding semantics rather than appearance also leads to another important Unicode principle: the principle of *unification*.

Unlike most character encoding schemes, Unicode aims to be comprehensive. It aims to provide codes for all the characters in all of the world's written languages. It also aims to be a superset of all other character encoding schemes (or at least the vast majority). By being a superset, Unicode can be an acceptable substitute for any of those other encodings (technical limitations aside, anyway), and it can also serve as a pivot point for processes converting text between any of the other encodings.

Other character encoding standards are Unicode's chief source of characters. The designers of Unicode aimed to include all the characters from every computer character encoding standard in reasonably widespread use at the time Unicode was designed, and have continued to incorporate

characters from other standards as it has evolved, either as important new standards emerged, or as the scope of Unicode widened to include new languages. The designers of Unicode drew characters from every international and national standard they could get their hands on, as well as code pages from all the big computer and software manufacturers, telegraphy codes, various other corporate standards, and even popular fonts, in addition to all the non-computer sources they used. As an example of their thoroughness, Unicode includes code-point values for the glyphs the old IBM PC code pages would show for certain ASCII control characters. As another example, Unicode assigns values to the glyphs from the popular Zapf Dingbats typeface.

This wealth of sources led to an amazingly extensive repertoire of characters, but also led to a lot of redundancy. If every character code from every source encoding retained its identity in Unicode (say, Unicode kept the original code values and just padded them all out to the same length and prefixed them with some identifier for the source encoding), they'd never all fit in a 16-bit code space. You'd also wind up with numerous alternate representations for things that anyone with a little common sense would consider to be the same thing.

For starters, almost every language has several different encoding standards. For example, there might be one national standard for each country where the language is spoken, plus one or more corporate standards devised by computer manufacturers selling into that market. Think about ASCII and EBCDIC in American English, for example. The capital letter A encoded by ASCII (as 0x41) is the same capital letter A that is encoded by EBCDIC (as 0xC1), so it makes little sense to have these two different source values map to different codes in Unicode. Then Unicode would have two different values for the letter A. So instead, Unicode *unifies* these two character codes and says that both sources map to the same Unicode value. Thus, the letter A is encoded only once in Unicode (as U+0041), not twice.

In addition to there being multiple encoding standards for most languages, most languages share their writing system with at least one other language. The German alphabet is different from the English alphabet—it adds ß and some other letters, for example—but they're really both just variations on the Latin alphabet. We need to make sure that the letter ß is encoded, but we don't need to create a different letter k for German—the same letter k we use in English will do just fine.

A truly vast number of languages use the Latin alphabet. Most omit some letters from what English speakers know as the alphabet, and most add some special letters of their own. Just the same, there's considerable overlap between their alphabets. The characters that overlap between languages are only encoded once in Unicode, not once for every language that uses them. For example, both Danish and Norwegian add the letter ø to the Latin alphabet, but the letter ø is only encoded once in Unicode.

Generally, characters are not unified across writing-system boundaries. For instance, the Latin letter B, the Cyrillic letter Б and the Greek letter Β are not unified, even though they look the same and have the same historical origins. This is partially because their lowercase forms are all different (b, б, and β, respectively), but mostly because the designers of Unicode didn't want to unify across writing-system boundaries.²⁰ It made more sense to keep each writing system distinct.

²⁰ There are a few exceptions to this base rule: One that comes up often is Kurdish, which when written with the Cyrillic alphabet also uses the letters Q and W from the Latin alphabet. Since the characters are a direct borrowing from the Latin alphabet, they weren't given counterparts in Unicode's version of the Cyrillic alphabet, a decision which still arouses debate.

So the basic principle is that wherever possible, Unicode unifies character codes from its various source encodings, whenever they can be demonstrated beyond reasonable doubt to refer to the same character. But there's a big exception: respect for existing practice. It was important to Unicode's designers (and probably a big factor in Unicode's success) for Unicode to be interoperable with the various encoding systems that came before it. In particular, for a subset of so-called "legacy" encodings, Unicode is specifically designed to *preserve round-trip compatibility*. That is, if you convert from one of the legacy encodings to Unicode and then back to the legacy encoding, you should get the same thing you started with. This means that there are many examples of characters that would have been unified in Unicode that aren't because of the need to preserve round-trip compatibility with a legacy encoding (or sometimes simply to conform to standard practice).

As an example, the Greek lowercase letter sigma has two forms. σ is used in the middle of words, and ς is used at the ends of words. As with the letters of the Arabic alphabet, this is an example of two different glyphs for the same letter. Unicode would normally just have a single code point for the lowercase sigma, but because all the standard encodings for Greek give different character codes to the two different versions of the lowercase sigma, Unicode has to as well. The same thing happens in Hebrew, where the word-ending forms of several letters have their own code point values in Unicode.

If a letter does double duty as a symbol, this generally isn't sufficient grounds for different character codes either. The Greek letter pi (π), for example, is still the Greek letter pi even when it's being used as the symbol of the ratio between a circle's diameter and its circumference, so it's still represented with the same character code. There are exceptions to this, however: the Hebrew letter aleph (\aleph) is used in mathematics to represent the transfinite numbers, and this use of the letter aleph is given a separate character code. The rationale here is that aleph-as-a-mathematical-symbol is a left-to-right character like all the other numerals and mathematical symbols, while aleph-as-a-letter is a right-to-left character. The letter Å is used in physics as the symbol of the Angstrom unit. Å-as-the-Angstrom-sign is given its own character code because some of the variable-length Japanese encodings did.

The business of deciding which characters can be unified can be complicated. Looking different is definitely not sufficient grounds by itself. For instance, the Arabic and Urdu alphabets have a very different look, but the Urdu alphabet is really just a particular calligraphic or typographical variation of the Arabic alphabet. The same set of character codes in Unicode is used to represent both Arabic and Urdu. The same thing happens with Greek and Coptic,²¹ modern and old Cyrillic (the original Cyrillic alphabet had different letter shapes and a bunch of letters that have since disappeared), and Russian and Serbian Cyrillic (in italicized fonts, there are some letters that have a different shape in Serbian from their Russian shape to avoid confusion with italicized Latin letters).

But by far the biggest, most complicated, and most controversial instance of character unification in Unicode is the Han ideographs. The characters originally developed to write the various Chinese languages, often called "Han characters" after the Han Dynasty, were also adopted by various other peoples in East Asia to write their languages, and the Han characters are still used (in combination with other characters) to write Japanese (who call them *kanji*) and Korean (who call them *hanja*).

The problem is that over the centuries, many of the Han characters have developed different forms in the different places where they're used. Even within the same written language—Chinese—you have different forms: In the early 1960s, the Mao regime in the People's Republic of China standardized

²¹ The unification of Greek and Coptic has always been controversial, since they're generally considered to be different alphabets, rather than just different typographical versions of the same alphabet. It's quite possible that Greek and Coptic will be disunified in a future Unicode version.

on simplified versions of many of the more complicated characters, but the traditional forms are still used in Taiwan and Hong Kong.

So the same ideograph can have four different forms: one each for Traditional Chinese, Simplified Chinese, Japanese, and Korean (and when Vietnamese is written with Chinese characters, you might have a fifth form). Worse yet, it's very often not clear what really counts as the "same ideograph" between these languages. Considerable linguistic research went into coming up with a unified set of ideographs for Unicode that can be used for both forms of written Chinese, Japanese, Korean, and Vietnamese. In fact, without this, it would have been impossible to fit Unicode into a 16-bit code space.

[One popular misconception about Unicode, by the way, is that Simplified and Traditional Chinese are unified. This isn't true; in fact, it's impossible, since the same Simplified Chinese character might be used as a stand-in for several different Traditional Chinese characters. Thus, most of the time, Simplified and Traditional Chinese characters get different code point values in Unicode. Only small differences that could be reliably categorized as font-design differences, analogous to the difference between Arabic and Urdu, were unified. For more on Han unification, see Chapter 10.]

In all of these situations where multiple glyphs are given the same character code, it either means the difference in glyph is simply the artistic choice of a type designer (for example, whether the dollar sign has one vertical stroke or two), or it's language dependent and it's expected a user will use an appropriate font for his language (or some mechanism outside Unicode's scope, such as automatic language detection or some kind of tagging scheme, would be used to determine the language and select an appropriate font).

The opposite situation—different character codes being represented by the same glyph—can also happen. One notable example is the apostrophe ('). There's one character code for this glyph when it's used as a punctuation mark and another when it's used as a letter (it's used in some languages to represent a glottal stop, such as in "Hawai'i").

Alternate-glyph selection

One interesting blurring of the line that can happen from time to time is the situation where a character with multiple glyphs needs to be drawn with a particular glyph in a certain situation, the glyph to use can't be algorithmically derived, and the particular choice of glyph needs to be preserved even in plain text.

Unicode has taken different approaches to solving this problem in different situations. Much of the time, the alternate glyphs are just given different code points. For example, there are five Hebrew letters that normally have a different shape when they appear at the end of a word from the shape they normally have. But in foreign words, these letters keep their normal shape even when they appear at the end of a word. Unicode just gives the regular and "final" versions of the letters different code point values. Examples like this can be found throughout Unicode.

Two special characters, U+200C ZERO WIDTH NON-JOINER ("ZWNJ" for short) and U+200D ZERO WIDTH JOINER ("ZWJ" for short), can be used as hints of which glyph shape is preferred in a particular situation. The ZWNJ prevents formation of a cursive connection or ligature in situations where one would normally happen, and the ZWJ causes a ligature or cursive connection where one

would otherwise not occur. These two characters can be used to override the default choice of glyphs.

The Unicode Mongolian block takes yet another approach: Many characters in the Mongolian block have or cause special shaping behavior to happen, but there are still times when the proper shape for a particular letter in a particular word can't be determined algorithmically (except maybe with an especially sophisticated algorithm that recognized certain words). So the Mongolian block includes three "variation selectors," characters that have no appearance of their own, but change the shape of the character that precedes them in some well-defined way.

Beginning in Unicode 3.2, the variation-selector approach has been extended to all of Unicode. Unicode 3.2 introduces sixteen general-purpose variation selectors, which work the same way as the Mongolian variation selectors: They have no visual presentation of their own, but act as "hints" to the rendering process that the preceding character should be drawn with a particular glyph shape. The list of allowable combinations of regular characters and variation selectors is given in a file called `StandardizedVariants.html` in the Unicode Character Database.

For more information on the joiner, non-joiner, and variation selectors, see Chapter 12.

Multiple Representations

Having now talked about the importance of the principle of unification, we have to talk about the opposite property, the fact that Unicode provides alternate representations for many characters.

As we saw earlier, Unicode's designers placed a high premium on respect for existing practice and interoperability with existing character encoding standards, in many cases sacrificing some measure of architectural purity in pursuit of the greater good (i.e., people actually using Unicode).

This means that Unicode includes code-point assignments for a lot of characters that were included solely or primarily to allow for round-trip compatibility with some legacy standard, a broad category of characters known more or less informally as *compatibility characters*. Exactly which characters are compatibility characters is somewhat a matter of opinion, and there isn't necessarily anything special about the compatibility characters that flags them as such. But there's an important subset of compatibility characters that *are* called out as special because they have alternate, preferred, representations in Unicode. Since the preferred representations usually consist of more than one Unicode code point, these characters are said to *decompose* into multiple code points.

There are two broad categories of decomposing characters: those with *canonical decompositions* (these characters are often referred to as "precomposed characters" or "canonical composites") and those with *compatibility decompositions* (the term "compatibility characters" is frequently used to refer specifically to these characters; a more specific term, "compatibility composite" is probably better). A canonical composite can be replaced with its canonical decomposition with no loss of data: the two representations are strictly equivalent, and the canonical decomposition is the character's preferred representation.²²

²² I have to qualify the word "preferred" here slightly. For characters whose decompositions consist of a single other character (so-called "singleton decompositions"), this is true. For multiple-character decompositions, there's nothing that necessarily makes them "better" than the precomposed forms, and you can generally use either representation. Decomposed representations are somewhat easier to deal with in code, though, and

Multiple Representations

Most canonical composites are a combination of a “base character” and one or more diacritical marks. For example, we talked about the character positioning rules and how the rendering engine needs to be smart enough so that when it sees, for example, an `a` followed by an umlaut, it draws the umlaut on top of the `a`, like so: `ä`. The thing is, much of the time, normal users of these characters don’t see them as the combination of a base letter and an accent mark. A German speaker sees `ä` simply as “the letter `ä`” and not as “the letter `a` with an umlaut on top.” So a vast number of letter-mark combinations are encoded using single character codes in the various source encodings, and these are very often more convenient to work with than the combinations of characters would be. The various European character-encoding standards follow this pattern, for example, assigning character codes to letter-accent combinations like `é`, `ä`, `å`, `û`, and so on, and Unicode follows suit.

Because a canonical composite can be mapped to its canonical decomposition without losing data, the original character and its decomposition are freely interchangeable. The Unicode standard enshrines this principle in law: On systems that support both the canonical composites and the combining characters that are included in their decompositions, the two different representations of the same character (composed and decomposed) are required to be treated as identical. That is, there is no difference between `ä` when represented by two code points and `ä` when represented with a single code point. In both cases, it’s still the letter `ä`.

This means that most Unicode implementations have to be smart enough to treat the two representations as equivalent. One way to do this is by *normalizing* a body of text to always prefer one of the alternate representations. The Unicode standard actually provides four different normalized forms for Unicode text.

All of the canonical decompositions involve one or more *combining marks*, a special class of Unicode code points representing marks that combine graphically in some way with the character that precedes them. If a Unicode-compatible system sees an `a` followed by a combining umlaut, it draws the umlaut on top of the `a`. This can be a little more inconvenient than just using a single code point to represent the `a`-umlaut combination, but it does give you an easy way to represent a letter with *more than one* mark attached to it, such as you find in Vietnamese or some other languages: just follow the base character with multiple combining marks.

But this means you can get into trouble with equivalence testing even without having composite characters. There are plenty of cases where the same character can be represented multiple ways by putting the various combining marks in different orders. Sometimes, the difference in ordering can be significant (if two combining marks attach to the base character in the same place, the one that comes first in the backing store is drawn closest to the character and the others are moved out of the way), but in many cases, the ordering isn’t significant—you get the same visual result whatever order the combining marks come in. In cases like this, the different forms are all legal and are required—once again—to be treated as identical. But the Unicode standard provides for a *canonical ordering* of combining marks to aid in testing sequences like these for equivalence.

The other class of decomposing characters is *compatibility composites*, characters with *compatibility decompositions*²³. A character can’t be mapped to its compatibility decomposition without losing data. For example, sometimes alternate glyphs for the same character are given their own character

many processes on Unicode text are based on mapping characters to their canonical decompositions, so they’re “preferred” in that sense. We’ll untangle all of this terminology in Chapter 4.

²³ There are a few characters in Unicode whose canonical decompositions include characters with compatibility decompositions. The Unicode standard considers these characters to be *both* canonical composites *and* compatibility composites.

codes. In these cases, there will be a preferred Unicode code point value representing the character, independent of glyph. Then there will be code point values representing the different glyphs. These are called *presentation forms*. The presentation forms have mapping back to the regular character they represent, but they're not simply interchangeable; the presentation forms refer to specific glyphs, while the preferred character maps to whatever glyph is appropriate for the context. In this way, the presentation forms carry more information than the canonical forms. The most notable set of presentation forms are the Arabic presentation forms, where each standard glyph for each Arabic letter, plus a wide selection of ligatures, has its own Unicode character code. Some rendering engines use presentation forms as an implementation detail, but normal users of Unicode are discouraged from using them and urged to use the non-decomposing characters instead. The same goes for the smaller set of presentation forms for other languages.

Another interesting class of compatibility composites are those that represent a particular stylistic variant of a particular character. These are similar to presentation forms, but instead of representing particular glyphs that are normally contextually selected, they represent particular glyphs that are normally specified through the use of additional styling information (remember, Unicode only represents plain text, not styled text). Examples include superscripted or subscripted numerals, or letters with special styles applied to them. For example the Planck constant is represented using an italicized letter *h*. Unicode includes a compatibility character code for the symbol for the Planck constant, but you could also just use a regular *h* in conjunction with some non-Unicode method of specifying that it's italicized. Characters with adornments such as surrounding circles also fall into this category, as do the abbreviations sometimes used in Japanese typesetting that consist of several characters arranged in a square.

For compatibility composites, the Unicode standard not only specifies the characters they decompose to, but some additional information intended to explain what additional non-text information is needed to express exactly the same thing.

Canonical and compatibility decompositions, combining characters, normalized forms, canonical accent ordering, and various related topics are all dealt with in excruciating detail in Chapter 4.

Flavors of Unicode

Let's take a minute to go back over the character-encoding terms from Chapter 2:

- An **abstract character repertoire** is just a collection of characters.
- A **coded character set** maps the characters in an abstract repertoire to abstract numeric values or positions in a table. These abstract numeric values are called *code points*. (For a while, the Unicode 2.0 standard was referring to code points as a "Unicode scalar values.")
- A **character encoding form** maps code points to series of fixed-length bit patterns known as *code units*. (For a while, the Unicode 2.0 standard was referring to code units as "code points.")
- A **character encoding scheme**, also called a "serialization format," maps code units to bytes in a sequential order. (This may involve specifying a serialization order for code units that are more than one byte long, or specifying a method of mapping code units from more than one encoding form into bytes, or both.)
- A **transfer encoding syntax** is an additional transformation that may be performed on a serialized sequence of bytes to optimize it for some situation (transforming a sequence of eight-bit byte values for transmission through a system that only handles seven-bit values, for example).

For most Western encoding standards, the transforms in the middle (i.e., from code points to code units and from code units to bytes) are so straightforward that they're never thought of as distinct steps. The standards in the ISO 8859 family, for example define coded character sets. Since the code point values are a byte long already, the character encoding forms and character encoding schemes used with these coded character sets are basically null transforms: You use the normal binary representation of the code point values as code units, and you don't have to do anything to convert the code units to bytes. (Although ISO 2022 does define a character encoding scheme that lets you mix characters from different coded character sets in a single serialized data stream.)

The East Asian character standards make these transforms more explicit: JIS X 0208 and JIS X 0212 define coded character sets only: they just map each character to row and column numbers in a table. You then have a choice of character encoding schemes you can use to convert the row and column numbers into serialized bytes: Shift-JIS, EUC-JP, and ISO 2022-JP are all examples of character encoding schemes used with the JIS coded character sets.

The Unicode standard makes each layer in this hierarchy explicit. It comprises:

- An abstract character repertoire that includes characters for an extremely wide variety of writing systems.
- A single coded character set that maps each of the characters in the abstract repertoire to a 21-bit value. (The 21-bit value can also be thought of as a coordinate in a three-dimensional space: a 5-bit plane number, an 8-bit row number, and an 8-bit cell number.)
- Three character encoding forms:
 - UTF-32, which represents each 21-bit code point value as a single 32-bit code unit. UTF-32 is optimized for systems where 32-bit values are easier or faster to process and space isn't at a premium.
 - UTF-16, which represents each 21-bit code point value as a sequence of one or two 16-bit code units. The vast majority of characters are represented with single 16-bit code units, making it a good general-use compromise between UTF-32 and UTF-8. UTF-16 is the oldest Unicode encoding form, and is the form specified by the Java and Javascript programming languages and the XML Document Object Model APIs.
 - UTF-8, which represents each 21-bit code point value as a sequence of from one to four 8-bit code units. The ASCII characters have exactly the same representation in UTF-8 as they do in ASCII, and UTF-8 is optimized for byte-oriented systems or systems where backward compatibility with ASCII is important. For European languages, UTF-8 is also more compact than UTF-16, although for Asian languages UTF-16 is more compact than UTF-8. UTF-8 is the default encoding form for a wide variety of Internet standards.

These encoding forms are often referred to as Unicode Transformation Formats (hence the "UTF" abbreviation).

- Seven character encoding schemes: UTF-8 is a character encoding scheme unto itself because it uses 8-bit code units. UTF-16 and UTF-32 each have three associated encoding schemes:
 - A "big-endian" version that serializes each code unit most-significant-byte first
 - A "little-endian" version that serializes each code unit least-significant-byte first
 - A self-describing version that uses an extra sentinel value at the beginning of the stream, called the "byte order mark," to specify whether the code units are in big-endian or little-endian order.

In addition, there are some allied specifications that aren't officially part of the Unicode standard:

- UTF-EBCDIC, an alternate version of UTF-8 designed for use on EBCDIC-based systems that maps Unicode code points to series of from one to five 8-bit code units.
- UTF-7, a now-mostly-obsolete character encoding scheme for use with 7-bit Internet standards that maps UTF-16 code units to sequences of 7-bit values.
- Standard Compression Scheme for Unicode (SCSU), a character encoding scheme that maps a sequence of UTF-16 code units to a compressed sequence of bytes, providing a serialized Unicode representation that is generally as compact for a given language as that language’s legacy encoding standards and optimizes Unicode text for further compression with byte-oriented compression schemes such as LZW.
- Byte-Order Preserving Compression for Unicode (BOCU), another compression format for Unicode.

We’ll delve into the details of all of these encoding forms and schemes in Chapter 6.

Character Semantics

Since Unicode aims to encode semantics rather than appearances, simple code charts don’t suffice. After all, all they do is show pictures of characters in a grid that maps them to numeric values. The pictures of the characters can certainly help illustrate the semantics of the characters, but they can’t tell the whole story. The Unicode standard actually goes well beyond just pictures of the characters, providing a wealth of information on every character.

Every code chart in the standard is followed by a list of the characters in the code chart. For each character, there’s an entry that gives the following information for the character:

- Its Unicode code point value.
- A representative glyph. For characters, such as accent marks, that combine with other characters, the representative glyph includes a dotted circle that shows where the main character would go—this makes it possible to distinguish COMBINING DOT ABOVE from COMBINING DOT BELOW, for example. For characters that have no visual appearance, such as spaces and control and formatting codes, the representative glyph is a dotted square with some sort of abbreviation of the character name inside.
- The character’s name. The name, and not the representative glyph, is the normative property (the parts of the standard that are declared to be “normative” are the parts you have to follow exactly in order to conform to the standard; parts declared “informative” are there to supplement or clarify the normative parts and don’t have to be followed exactly in order to conform). This reflects the philosophy that Unicode encodes semantics, although there are a few cases where the actual meaning of the character has drifted since the earliest drafts of the standard and no longer matches the name. This is very rare, though.

In addition to the code point value, name, and representative glyph, an entry may also include:

- Alternate names the character might be known by.
- Cross-references to similar characters elsewhere in the standard (this helps to distinguish them from each other).
- The character’s canonical or compatibility decomposition (if it’s a composite character)
- Additional notes on its usage or meaning (for example, the entries for many letters include the languages that use them).

The Unicode standard also includes chapters on each major group of characters in the standard, with information that's common to all of the characters in the group (such as encoding philosophy or information on special processing challenges) and additional narrative explaining the meaning and usage of any characters in the group that have special properties or behavior that needs to be called out.

But the Unicode standard actually consists of more than just *The Unicode Standard*. That is, there's more to the Unicode standard than just the book. The Unicode standard also comprises a comprehensive database of all the characters, a copy of which is included on a CD that's included with the book (the character database changes more frequently than the book, however, so for truly up-to-date information, it's usually a good idea to get the most recent version of the database from the Unicode Consortium's Web site [www.unicode.org]).

Every character in Unicode has a bunch of properties associated with it that define how that character is to be treated by various processes. The Unicode Character Database comprises a group of text files that give the properties for each character in Unicode. Among the properties that each character has are:

- The character's code-point value and name.
- The character's general category. All of the characters in Unicode are grouped into 30 categories, 17 of which are considered normative. The category tells you things like whether the character is a letter, numeral, symbol, whitespace character, control code, etc.
- The character's decomposition, along with whether it's a canonical or compatibility decomposition, and for compatibility composites, a tag that attempts to indicate what data is lost when you convert to the decomposed form.
- The character's case mapping. If the character is a cased letter, the database includes the mapping from the character to its counterpart in the opposite case.
- For characters that are considered numerals, the database includes the character's numeric value. (That is, the numeric value the character represents, not the character's code point value.)
- The character's directionality. (e.g., whether it's left-to-right, right-to-left, or takes on the directionality of the surrounding text). The Unicode Bidirectional Layout Algorithm uses this property to determine how to arrange characters of different directionalities on a single line of text.
- The character's mirroring property. This says whether the character take on a mirror-image glyph shape when surrounded by right-to-left text.
- The character's combining class. This is used to derive the canonical representation of a character with more than one combining mark attached to it (it's used to derive the canonical ordering of combining characters that don't interact with each other).
- The character's line-break properties. This is used by text rendering processes to help figure out where line divisions should go.
- Many more...

For an in-depth look at the various files in the Unicode Character Database, see Chapter 5.

Unicode Versions and Unicode Technical Reports

The Unicode standard also includes a bunch of supplementary documents known as Unicode Technical Reports. A snapshot of these is also included on the CD that comes with the book, but at this point, the CD in the Unicode 3.0 is so out of date as to be nearly useless. You can always find the most up-to-date slate of technical reports on the Unicode Web site (www.unicode.org).

There are three kinds of technical reports:

- **Unicode Standard Annexes** (abbreviated “UAX”). These are actual addenda and amendments to the Unicode standard.
- **Unicode Technical Standards** (abbreviated “UTS”). These are adjunct standards related to Unicode. They’re not normative parts of the standard itself, but carry their own conformance criteria.
- **Unicode Technical Reports** (abbreviated “UTR”). Various other types of adjunct information, such as text clarifying or expanding on parts of the standard, implementation guidelines, and descriptions of procedures for doing things with Unicode text that don’t rise to the level of official Unicode Technical Standards.

There are also:

- **Draft Unicode Technical Reports** (abbreviated “DUTR”). Technical reports are often published while they’re still in draft form. DUTR status indicates that an agreement has been reached in principle to adopt the proposal, but that details are still to be worked out. Draft technical reports don’t have any normative force until the Unicode Technical Committee votes to remove “Draft” from their name, but they’re published early to solicit comment and give implementers a head start.
- **Proposed Draft Unicode Technical Reports** (abbreviated “PDUTR”). Technical reports that are published before an agreement in principle has been reached to adopt them.

The status of the technical reports is constantly changing. Here’s a summary of the slate of technical reports as of this writing (December 2001):

Unicode Standard Annexes

All of the following Unicode Standard Annexes are officially part of the Unicode 3.2 standard.

- **UAX #9: The Unicode Bidirectional Algorithm:** This specifies the algorithm for laying out lines of text that mix left-to-right characters with right-to-left characters. This supersedes the description of the Bidirectional layout algorithm in the Unicode 3.0 book. For an overview of the bidirectional layout algorithm, see Chapter 8. For implementation details, see Chapter 16.
- **UAX #11: East Asian Width.** Specifies a set of character properties that determine how many display cells a character takes up when used in the context of East Asian typography. For more information, see Chapter 10.
- **UAX #13: Unicode Newline Guidelines.** There are some interesting issues regarding how you represent the end of a line or paragraph in Unicode text, and this document clarifies them. This information is covered in Chapter 12.
- **UAX #14: Line Breaking Properties.** This document specifies how a word-wrapping routine should treat the various Unicode characters. Word-wrapping is covered in Chapter 16.

- **UAX #15: Unicode Normalization Forms.** Since Unicode actually has multiple ways of representing a lot of characters, it's often helpful to convert Unicode text to some kind of normalized form that prefers one representation for any given character over all the others. There are four Unicode Normalization Forms, and this document describes them. The Unicode Normalization Forms are covered in Chapter 4.
- **UAX #19: UTF-32.** Specifies the UTF-32 encoding form. UTF-32 is covered in Chapter 6.
- **UAX #27: Unicode 3.1.** The official definition of Unicode 3.1. Includes the changes and additions to Unicode 3.0 that form Unicode 3.1. This is the document that officially incorporates the other Unicode Standard Annexes into the standard and gives them a version number. It also includes code charts and character lists for all the new characters added to the standard in Unicode 3.1.

Unicode Technical Standards

- **UTS #6: A Standard Compression Scheme for Unicode.** Defines SCSU, a character encoding scheme for Unicode that results in serialized Unicode text of comparable size to the same text in legacy encodings. For more information, see Chapter 6.
- **UTS #10: The Unicode Collation Algorithm.** Specifies a method of comparing character strings in Unicode in a language-sensitive manner and a default ordering of all the characters to be used in the absence of a language-specific ordering. For more information, see Chapter 15.

Unicode Technical Reports

- **UTR #16: UTF-EBCDIC.** Specifies a special 8-bit transformation of Unicode for use on EBCDIC-based systems. This is covered in Chapter 6.
- **UTR #17: Character Encoding Model.** Defines a set of useful terms for discussing the various aspects of character encodings. This material was covered in Chapter 2 and reiterated in this chapter.
- **UTR #18: Regular Expression Guidelines.** Provides some guidelines for how a regular-expression facility should behave when operating on Unicode text. This material is covered in Chapter 15.
- **UTR #20: Unicode in XML and Other Markup Languages.** Specifies some guidelines for how Unicode should be used in the context of a markup language such as HTML or XML. This is covered in Chapter 17.
- **UTR #21: Case Mappings.** Gives more detail than the standard itself on the whole business of mapping uppercase to lowercase and vice versa. Case mapping is covered in Chapter 14.
- **UTR #22: Character Mapping Tables.** Specifies an XML-based file format for describing a mapping between Unicode and some other encoding standard. Mapping between Unicode and other encodings is covered in Chapter 14.
- **UTR #24: Script Names.** Defines an informative character property: script name, which would identify which writing system (or “script”) each character belongs to. The script-name property is discussed in Chapter 5.

Draft and Proposed Draft Technical Reports

- **PDUTR #25: Unicode support for mathematics.** Gives a detailed account of all the various considerations involved in using Unicode to represent mathematical expressions, including things like spacing and layout, font design, and how to interpret various Unicode characters in the context of mathematical expressions. It also includes a heuristic for detecting mathematical

expressions in plain Unicode text, and proposes a scheme for representing structured mathematical expressions in plain text. We'll look at math symbols in Chapter 12.

- **DUTR #26: Compatibility Encoding Scheme for UTF-16: 8-bit (CESU-8).** Documents a UTF-8 like Unicode encoding scheme that's being used in some existing systems. We'll look at CESU-8 in chapter 6.
- **PDUTR #28: Unicode 3.2.** The official definition of Unicode 3.2. Together with UAX #27, this gives all the changes and additions to Unicode 3.0 that define Unicode 3.2, including revised and updated code charts with all the new Unicode 3.2 characters. Unicode 3.2 is scheduled to be released in March 2002, so it's likely this will be UAX #28 by the time you read this.

Superseded Technical Reports

The following technical reports have been superseded by (or absorbed into) more recent versions of the standard.

- **UTR #1: Myanmar, Khmer and Ethiopic.** Absorbed into Unicode 3.0.
- **UTR #2: Sinhala, Mongolian, and Tibetan.** Tibetan was in Unicode 2.0; Sinhala and Mongolian was added in Unicode 3.0.
- **UTR #3: Various less-common scripts.** This document includes exploratory proposals for a whole bunch of historical or rare writing systems. It has been superseded by more-recent proposals. Some of the scripts in this proposal have been incorporated into more-recent versions of Unicode: Cherokee (Unicode 3.0), Old Italic (Etruscan; Unicode 3.1), Thaana (Maldivian; 3.0), Ogham (3.0), Runic (3.0), Syriac (3.0), Tagalog (3.2), Buhid (3.2), and Tagbanwa (3.2). Most of the others are in various stages of discussion, with another batch scheduled for inclusion in Unicode 4.0. This document is mostly interesting as a list of writing systems that will probably be in future versions of Unicode. For more information on this subject, check out <http://www.unicode.org/unicode/alloc/Pipeline.html>.
- **UTR #4: Unicode 1.1.** Superseded by later versions.
- **UTR #5: Handling non-spacing marks.** Incorporated into Unicode 2.0.
- **UTR #7: Plane 14 Characters for Language Tags.** Incorporated into Unicode 3.1.
- **UTR #8: Unicode 2.1.** Superseded by later versions.

The missing numbers belong to technical reports that have been withdrawn by their proposers, have been turned down by the Unicode Technical Committee, or haven't been published yet.

Unicode Versions

Many of the technical reports either define certain versions of Unicode or are superseded by certain versions of Unicode. Each version of Unicode comprises a particular set of characters, a particular version of the character-property files, and a certain set of rules for dealing with them. All of these things change over time (although there are certain things that are guaranteed to remain the same—see “Unicode stability policies” later in this chapter).

Which gets us to the question of Unicode version numbers. A Unicode version number consists of three parts: for example “Unicode 2.1.8.” The first number is the *major version number*. This gets bumped every time a new edition of the book is released. That happens when the accumulation of Technical Reports starts to get unwieldy or when a great many significant changes (such as lots of new characters) are incorporated into the standard at once. There's a new major version of Unicode once every several years.

The *minor version number* gets bumped whenever new characters are added to the standard or other significant changes are made. New minor versions of Unicode don't get published as books, but do get published as Unicode Standard Annexes. There has been one minor version of Unicode between each pair of major versions (i.e., there was a Unicode 1.1 [published as UTR #4], a Unicode 2.1 [UTR #8], and a Unicode 3.1 [UAX #27]), but this was broken with the release of Unicode 3.2 (PDUTR #28 at the time of this writing, but most likely UAX #28 by the time you read this).

The *update number* gets bumped when changes are made to the Unicode Character Database. Updates are usually just published as new versions of the database; there is no corresponding technical report.

The current version of Unicode at the time of this writing (January 2002) is Unicode 3.1.1. Unicode 3.2, which includes a whole slate of new characters, is currently in beta and will likely be the current version by the time you read this.

One has to be a bit careful when referring to a particular version of the Unicode standard from another document, particularly another standard. Unicode is changing all the time, and so it's seldom a good idea to nail yourself to one specific version. Most of the time, it's best either to specify only a major version number (or, in some cases, just major and minor version numbers), or to specify an open-ended range of versions (e.g., "Unicode 2.1 or later"). Generally, this is okay, as future versions of Unicode will only add characters—since you're never required to support a particular character, you can pin yourself to an open-ended range of Unicode versions and not sweat the new characters. Sometimes, however, changes are made to a character's properties in the Unicode Character Database, and these could alter program behavior. Usually this is a good thing—changes are made to the database when the database is deemed to have been *wrong* before—but your software may need to deal with this in some way.

Unicode stability policies

Unicode will, of course, continue to evolve. There are, however, a few things you can count on to remain stable:

- Characters that are in the current standard will never be removed from future standards. They may, in unusual circumstances, be deprecated (i.e., their use might be discouraged), but they'll never be taken out altogether and their code point values will never be reused to refer to different characters.
- Characters will never be reassigned from one code point to another. If a character has ambiguous semantics, a new character may be introduced with more specific semantics, but the old one will never be taken away and will continue to have ambiguous semantics.
- Character names will never change. This means that occasionally a character with ambiguous semantics will get out of sync with its name as its semantics evolve, but this is very rare.
- Text in one of the Unicode Normalized Forms will always be in that normalized form. That is, the definition of the Unicode Normalized Forms will not change between versions of the standard in ways that would cause text that is normalized according to one version of the standard not to be normalized in later versions of the standard.
- A character's combining class and canonical and compatibility decompositions will never change, as this would break the normalization guarantee.
- A character's properties may change, but not in a way that would alter the character's fundamental identity. In other words, the representative glyph for "A" won't change to "B", and

“A”’s category won’t change to “lowercase letter.” You’ll see property changes only to correct clear mistakes in previous versions.

- Various structural aspects of the Unicode character properties will remain the same, as implementations depend on some of these things: For example, the standard won’t add any new general categories or any new bi-di categories, it’ll keep characters combining classes in the range from 0 to 255, non-combining characters will always have a combining class of 0, and so on.

You can generally count on characters’ other normative properties not changing, although the Unicode Consortium certainly reserves the right to fix mistakes in these properties.

The Unicode Consortium can change a character’s informative properties pretty much at will, without changing version numbers, since you don’t have to follow them anyway. Again, this shouldn’t actually happen much, except when they feel they need to correct a mistake of some kind.

These stability guarantees are borne out of bitter experience: In particular, characters *did* get removed and reassigned, and characters’ names *did* change in Unicode 1.1 as a result of the merger with ISO 10646, and it caused serious grief. This won’t happen again.

Arrangement of the encoding space

Unicode’s designers tried to assign the characters to numeric values in an orderly manner that would make it easy to tell something about a character just from its code point value. As the encoding space has filled up, this has become harder to do. But the logic still comes through reasonably well.

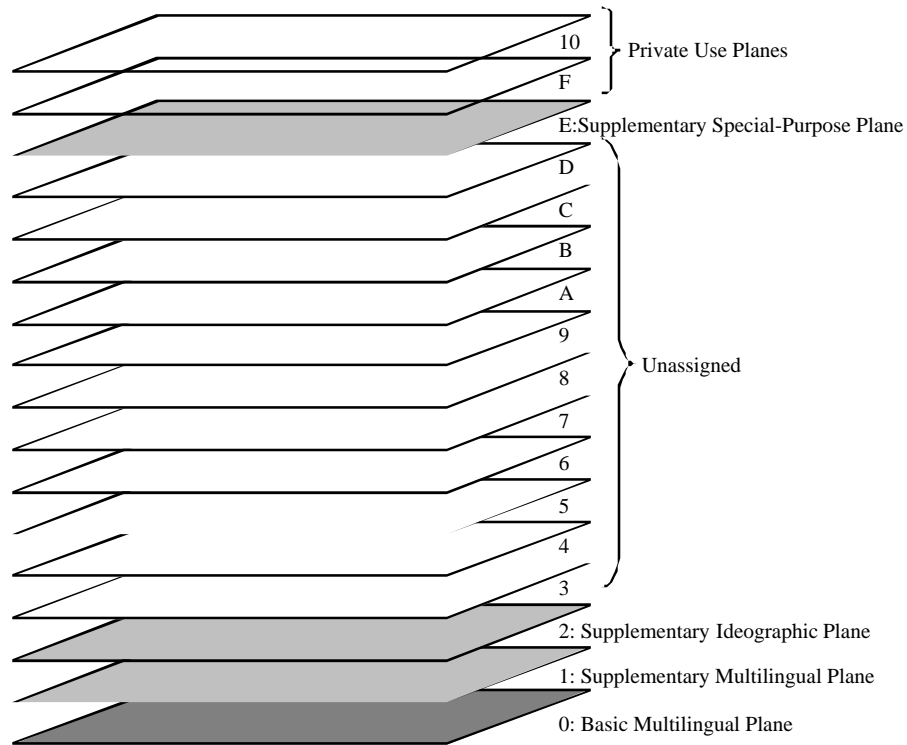
Unicode was originally designed for a 16-bit encoding space, consisting of 256 *rows* of 256 characters each. ISO 10646 was designed for a 32-bit encoding space, consisting of 128 *groups* of 256 *planes* containing 256 rows of 256 characters. So the original Unicode encoding space had room for 65,536 characters. ISO 10646 had room for an unbelievable 2,147,483,648 characters. The ISO encoding space is clearly overkill (experts estimate there’s maybe a million or so characters eligible for encoding), but it was clear by the time Unicode 2.0 came out that the 16-bit Unicode encoding space was too small.

The solution was the *surrogate mechanism*, a scheme whereby special escape sequences known as *surrogate pairs* could be used to represent characters outside the original encoding space. This extended the number of characters that could be encoded to 1,114,112, ample space for the foreseeable future (only 95,156 characters are actually encoded in Unicode 3.2, and Unicode has been in development for twelve years). The surrogate mechanism was introduced in Unicode 2.0 and has since become known as UTF-16. It effectively encodes the first 17 planes of the ISO 10646 encoding space. The Unicode Consortium and WG2 have agreed never to populate the planes above Plane 16, so for all intents and purposes, Unicode and ISO 10646 now share a 21-bit encoding space consisting of 17 planes of 256 rows of 256 characters. Valid Unicode code point values run from U+0000 to U+10FFFF.

Organization of the planes

The Unicode encoding space currently looks like this:

Arrangement of the Encoding Space



Plane 0 is the **Basic Multilingual Plane** (or “BMP” for short). It contains the majority of the encoded characters, including all of the most common. In fact, prior to Unicode 3.1, there were no characters at all encoded in any of the other planes. The characters in the BMP can be represented in UTF-16 with a single 16-bit code unit.

Plane 1 is the **Supplementary Multilingual Plane** (“SMP” for short). It is intended to contain characters from archaic or obsolete writing systems. Why encode these at all? These are mostly here for the use of the scholarly community for papers where they write about these characters. Various specialized collections of symbols will also go into this plane.

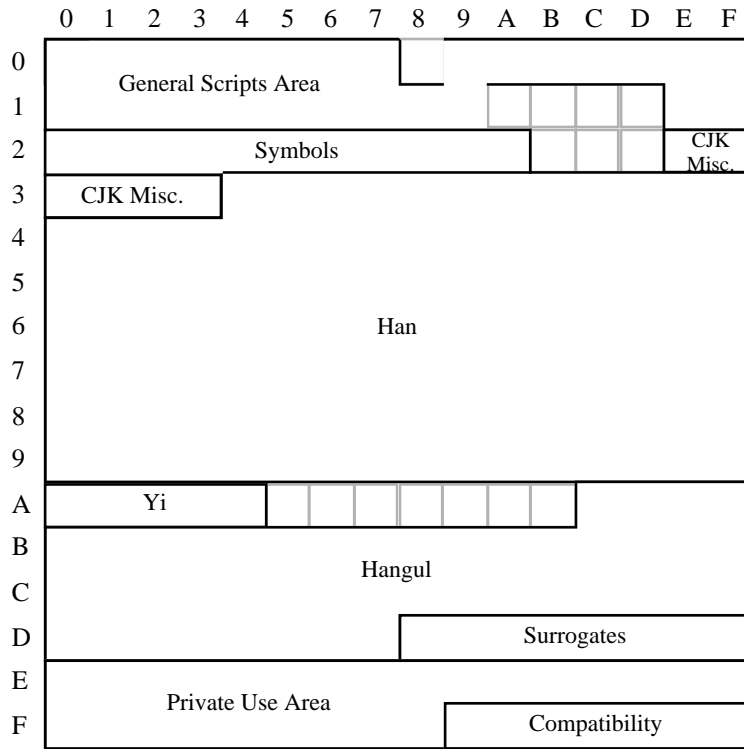
Plane 2 is the **Supplementary Ideographic Plane** (“SIP” for short). It’s an extension of the CJK Ideographs Area from the BMP and contains rare and unusual Chinese characters.

Plane 14 (E) is the **Supplementary Special-Purpose Plane** (“SSP”). It’s reserved for special-purpose characters—generally code points that don’t encode characters as such but are instead used by higher-level protocols or as signals to processes operating on Unicode text.

Planes 15 and 16 (F and 10) are the **Private Use Planes**, an extension of the Private Use Area in the BMP. The other planes are currently unassigned, and will probably remain that way until Planes 1, 2, and 14 start to fill up.

The Basic Multilingual Plane

The heart and soul of Unicode is Plane 0, the Basic Multilingual Plane. It contains the vast majority of characters in common use today, and those that aren't yet encoded will go here as well. Here's a graph of the BMP:



The characters whose code point values begin with 0 and 1 form the **General Scripts Area**. This area essentially contains the characters from all of the alphabetic writing systems, including the Latin, Greek, Cyrillic, Hebrew, Arabic, Devanagari (Hindi), and Thai alphabets, among many others. It also contains a collection of combining marks that are often used in conjunction with the letters in this area. The General Scripts Area breaks down as follows:

Arrangement of the Encoding Space

	00	20	40	60	80	A0	C0	E0	00	20	40	60	80	A0	C0	E0
00/01	ASCII				Latin1				Lat. Ext. A				Lat. Ext. B			
02/03	Lat. Ext. B		IPA		Mod. Ltrs.		Comb. Diac.		Greek							
04/05	Cyrillic						Armenian			Hebrew						
06/07	Arabic						Syriac				Thaana					
08/09							Devanagari				Bengali					
0A/0B	Gurmukhi			Gujarati			Oriya				Tamil					
0C/0D	Telugu			Kannada			Malayalam				Sinhala					
0E/0F	Thai			Lao			Tibetan									
10/11	Myanmar			Georgian			Hangul Jamo									
12/13	Ethiopic								Cherokee							
14/15	Canadian Aboriginal Syllabics															
16/17	Ogham			Runic			Phillipine Scripts				Khmer					
18/19	Mongolian															
1A/1B																
1C/1D																
1E/1F	Latin Extended Additional								Greek Extended							

There are a couple of important things to note about the general scripts area. The first is that the first 128 characters, those from U+0000 to U+7F, are exactly the same as the ASCII characters with the same code point values. This means you can convert from ASCII to Unicode simply by zero-padding the characters out to 16 bits (in fact, in UTF-8, the 8-bit version of Unicode, the ASCII characters have exactly the same representation as they do in ASCII).

Along the same lines, the first 256 characters, those from U+0000 to U+00FF, are exactly the same as the characters with the same code point values from the ISO 8859-1 (ISO Latin-1) standard. (Latin-1 is a superset of ASCII; its lower 128 characters are identical to ASCII.) You can convert Latin-1 to Unicode by zero-padding out to 16 bits. (However, the non-ASCII Latin-1 characters have two-byte representations in UTF-8.)

For those writing systems, such as most of the Indian and Southeast Asian ones, that have only one dominant existing encoding, Unicode keeps the same relative arrangement of the characters as their original encoding had, enabling conversion back and forth by adding or subtracting a constant.

We'll be taking an in-depth look at all of these scripts in Section II. The Latin, Greek, Cyrillic, Armenian, and Georgian blocks, as well as the Combining Diacritical Marks, IPA Extensions, and Spacing Modifier Letters blocks, are covered in Chapter 7. The Hebrew, Arabic, Syriac, and Thaana blocks are covered in Chapter 8. The Devanagari, Bengali, Gurmukhi, Gujarati, Oriya, Tamil, Telugu, Kannada, Malayalam, Sinhala, Thai, Lao, Tibetan, Myanmar, and Khmer blocks are covered in Chapter 9. The Hangul Jamo block is covered in Chapter 10. The Ethiopic, Cherokee, Canadian Aboriginal Syllables, Ogham, Runic, and Mongolian blocks are covered in Chapter 11.

The characters whose code point values begin with 2 (with a few recent exceptions) form the **Symbols Area**. There's all kinds of stuff in this area. It includes a collection of punctuation that can

be used with many different languages (this block actually supplements the punctuation marks in the ASCII and Latin-1 blocks), collections of math, currency, technical, and miscellaneous symbols, arrows, box-drawing characters, and so forth. It breaks down like this:

	00	20	40	60	80	A0	C0	E0	00	20	40	60	80	A0	C0	E0
20/21	Gen. Punc.			Sup/ Sub	Curr- ency	Symb. Diac.	Letterlike			Numbers			Arrows			
22/23	Mathematical Operators						Miscellaneous Technical									
24/25	Ctrl. Pict.	OCR	Enclosed			Box Drawing			Block	Geom. Shapes						
26/27	Misc. Symbols						Dingbats			Misc. Math						
28/29	Braille Patterns						Supp. Arrows			Misc. Math						
2A/2B	Supplemental Mathematical Operators						Bopomofo			Kanbun						
2C/2D																
2E/2F							Supp. CJK Radicals			KangXi Radicals			IDC			
30/31	CJK Punc.	Hiragana		Katakana		Hangul Compat.										
32/33	Enclosed CJK						CJK Compatibility									
34/35	Han															
36/37																
38/39																
3A/3B																
3C/3D																
3E/3F																

The various blocks in the Symbols Area are covered in Chapter 12.

The characters whose code point values begin with 3 (actually, with Unicode 3.0, this group has now slopped over to include some code points values beginning with 2) form the **CJK Miscellaneous Area**. This includes all of the characters used in the East Asian writing systems, except for the three very large areas immediately following. This includes punctuation used in East Asian writing, the phonetic systems used for Japanese and Chinese, various symbols and abbreviations used in Japanese technical material, and a collection of “radicals,” component parts of Han ideographic characters. These blocks are covered in Chapter 10.

The characters whose code point values begin with 4, 5, 6, 7, 8, and 9 (in Unicode 3.0, this area has slopped over to include most of the characters whose code points values begin with 3 as well) constitute the **CJKV Unified Ideographs Area**. This is where the Han ideographs used in Chinese, Japanese, Korean, and (much less frequently) Vietnamese are located.

The characters whose code point values range from U+A000 to U+A4CF form the **Yi Area**, which contains the characters used for writing Yi, a minority Chinese language.

The characters whose code point values range from U+AC00 to U+D7FF form the **Hangul Syllables Area**. Hangul is the alphabetic writing system used (sometimes in conjunction with Han ideographs) to write Korean. Hangul can be represented using the individual letters, or “jamo”, which are encoded in the General Scripts Area, but the jamo are usually arranged into ideograph-like blocks

Arrangement of the Encoding Space

representing whole syllables, and most Koreans look at whole syllables as single characters. This area encodes all possible modern Hangul syllables using a single code point for each syllable.

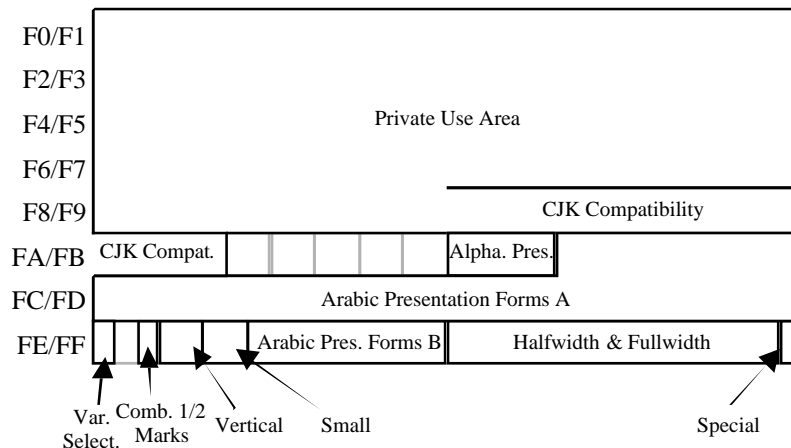
We look at the CJKV Unified Ideographs, Yi, and Hangul Syllables Areas in chapter 10.

The code point values from U+D800 to U+DFFF constitute the **Surrogates Area**. This range of code point values is reserved and will never be used to encode characters. Instead, values from this range are used in pairs as code-unit values by the UTF-16 encoding to represent characters from Planes 1 through 16.

The code point values from U+E000 to U+F8FF form the **Private Use Area**, or “PUA.” This area is reserved for the private use of applications and systems that use Unicode, which may assign any meaning they wish to the code point values in this range. Private-use characters should be used only within closed systems that can apply a consistent meaning to these code points; text that is supposed to be exchanged between systems is prohibited from using these code point values (unless there’s a private agreement otherwise by the sending and receiving parties), since there’s no guarantee that a receiving process would know what meaning to apply to them.

The remaining characters with code point values beginning with F form the **Compatibility Area**. This is a catch-all area for characters that are included in Unicode simply to maintain backward compatibility with the source encodings. This area includes various ideographs that would be unified with ideographs in the CJK Unicode Ideographs except that the source encodings draw a distinction, presentation forms for various writing systems, especially Arabic, halfwidth and fullwidth variants of various Latin and Japanese characters, and various other things. This section isn’t the only area of the encoding space containing compatibility characters; the Symbols Area includes many blocks of compatibility characters, and some are also scattered throughout the rest of the encoding space. This area also contains a number of special-purpose characters and non-character code points.

The Compatibility Area breaks down like this:



The Supplementary Planes

Planes 1 through 16 are collectively known as the Supplementary Planes, and include rarer or more specialized characters.

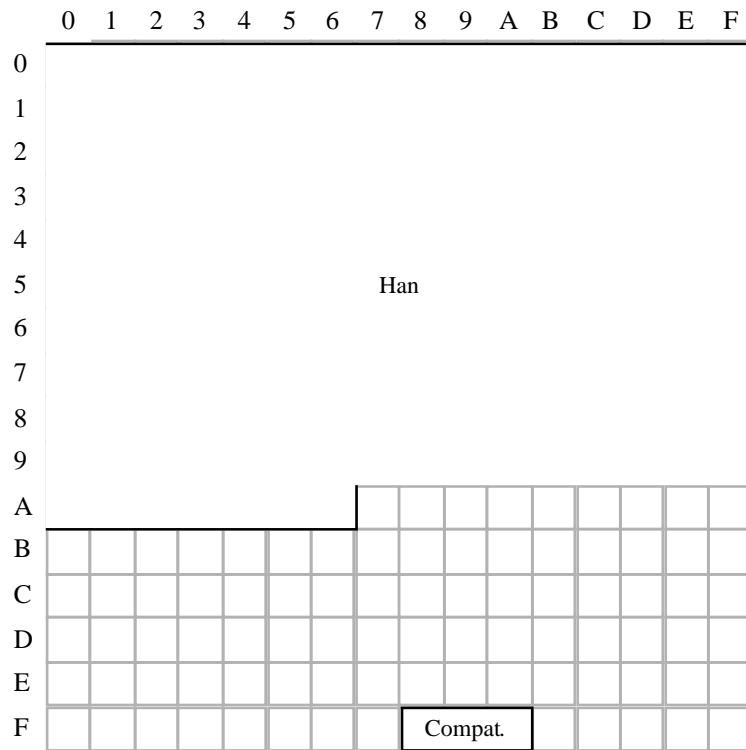
Plane 1 breaks down like this:

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0				Letters												
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D	Music				Math											
E																
F																

The area marked “Letters” includes a number of obsolete writing systems and will expand to include more. The area marked “Music” includes a large collection of musical symbols, and the area marked “Math” includes a special set of alphanumeric characters intended to be used as symbols in mathematical formulas.

Plane 2 breaks down like this:

Arrangement of the Encoding Space



It's given over entirely to Chinese ideographic characters, acting as an extension of the CJKV Unified Ideographs Area in the BMP.

Plane 14 looks like this:

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	Tags															
1																
2																
3																
4																
5																
6																
7																
8																
9																
A																
B																
C																
D																
E																
F																

The only thing in there right now is a small collection of “tag” characters intended to be used for things like language tagging.

Although there aren’t a lot of unassigned code point values left in the BMP, there are thousands upon thousands in the other planes. And except for the private-use areas, Unicode implementations are not permitted to use the unassigned code point values for anything. All of them are reserved for future expansion, and may be assigned to characters in future versions of Unicode. Conforming Unicode implementations can’t use these values for any purpose, and can’t emit text purporting to be Unicode that uses them. This also goes for the planes above Plane 16, even though they may never be used to encode characters. It’s specifically illegal to use the unused bits in a UTF-32 code unit to store other data.

Non-Character code point values

In addition, the code point values U+FFFE and U+FFFF, plus the corresponding code point values from all the other planes, are illegal. They’re not to be used in Unicode text at all. U+FFFE can be used in conjunction with the Unicode byte-order mark (U+FEFF) to detect byte-ordering problems (for example, if a Unicode text file produced on a Wintel PC starts with the byte-order mark, a Macintosh program reading it will read the byte-order mark as the illegal value U+FFFE and know that it has to byte-swap the file in order to read it properly).

U+FFFF is illegal for two main reasons: First, it provided a non-Unicode value that can be used as a sentinel value by Unicode-conformant processes. For example, the `getc()` function in C has to have a return type of `int` even though all of the character values it can return are of type `char`. This is because all `char` values are legal character codes, leaving no values to serve as the end-of-file signal. The `int` value `-1` is the end-of-file signal—you can’t use the `char` value `-1` as end-of-file

because it's the same as `0xFF`, which is a legal character. The Unicode version of `getc()`, on the other hand, could return `unsigned short` (or `wchar_t` on many systems) and still have a non-character value of that type—`U+FFFF`—available to use as the end-of-file signal.

Here's the other reason why `U+FFFF` isn't a legal Unicode code point value: Say you want to iterate over all of the Unicode code point values. You write the following (in C):

```
unsigned short c;
for (c = 0; c <= 0xFFFF; ++c) {
    // etc...
```

The loop will never terminate. This is because the next value after `0xFFFF` is `0`. Designating `U+FFFF` as a non-Unicode value enables you to write loops that iterate over the whole Unicode range in a straightforward manner without having to resort to a larger type (and a lot of casting) for the loop variable or other funny business to make sure the loop terminates.

The corresponding code points in the other planes were reserved for the same reasons, although this is mostly a historical curiosity now. In the original design of ISO 10646, it was expected that each plane would function as a more or less independent encoding space, and that if you only dealt with characters from one plane, you might represent them with 16-bit units (effectively chopping off the plane and group numbers) and have the same problem as described above.

Unicode 3.1 sets aside 32 additional code point values, `U+FDD0` to `U+FDEF`, as non-character code points. Like the others, this basically makes these values available to implementations for their internal use as markers or sentinel values without the implementations having to worry about their being assigned to characters in the future. These are very specifically *not* private use code points and therefore aren't supposed to be used to represent characters. Like the other non-character code points, they're never legal in serialized Unicode text.

Conforming to the standard

So just what does it mean to say that you conform to the Unicode standard? Well, this answer varies depending on what it is that your product does. The answer tends to be both more and less than what most people think.

First of all, one thing that conforming to the Unicode standard does *not* mean is that you have to be able to properly support every single character that the Unicode standard defines. What the Unicode standard requires is that you declare which characters you do support. For the characters you claim to support, *then* you have to follow all the rules in the standard. In other words, if you declare your program to be Unicode conformant (and you're basically doing that if you use the word "Unicode" anywhere in your advertising or docs) and say "Superduperword supports Arabic," then you have to support Arabic the way the Unicode standard says you should. In particular, you've got to be able to automatically select the right glyphs for the various Arabic letters depending on their context, and you've got to support the Unicode bidirectional text layout algorithm. If you don't do these things, then as far as the Unicode standard is concerned, you don't support Arabic.

Here are the rules for conforming to the Unicode standard. These differ somewhat from the rules as set forth in Chapter 3 of the actual Unicode standard, but they produce the same end result. There are certain algorithms that you have to follow (or mimic) in certain cases to be conformant. I haven't included those here, but will go over them in future chapters. There are also some terms used here that haven't been defined yet; all will be defined in future chapters.

General

For most processes, it's not enough to say you support Unicode. By itself, this doesn't mean very much. You'll also need to say:

- **Which version of Unicode you're supporting.** Generally, this is just a shorthand way of saying which characters you support, but in cases where the Unicode versions differ in the semantics they give to characters, or in their algorithms to do different things, you're specifying which versions of those things you're using as well. Typically, if you support a given Unicode version, you also support all previous versions as well.²⁴

Informative character semantics can and do change from version to version. You're not required to conform to the informative parts of the standard, but saying which version you support is also a way of saying which set of informative properties you're using.

It's legal, and in fact often a good idea, to say something like "Unicode 2.1.8 and later" when specifying which version of Unicode you use. This is particularly true when you're writing a standard that uses the Unicode standard as one of its base standards. This permits new versions of the standard (or conforming implementations) to support new characters without going out of compliance. It's also rarely necessary to specify which version of Unicode you're using all the way out to the last version number. Often, it's sufficient to just specify the major revision ("This product supports Unicode 2. x").

- **Which transformation formats you support.** This is only relevant if you exchange Unicode text with the outside world (including writing it to disk or sending it over a network connection), but if you do, you have to specify which of the various character encoding schemes defined by Unicode (the "Unicode Transformation Formats") you support. If you support several, you need to specify your default (i.e., which formats you can read without being told by the user or some other outside source what format the incoming file is in). The Unicode Transformation Formats are discussed in Chapter 6.
- **Which normalization forms you support or expect.** Again, this is basically important if you're exchanging Unicode text with the outside world somehow. This can be thought of as another shorthand way of specifying which characters you support, but is specifically oriented toward telling people what characters can be in an incoming file. The normalization forms are discussed in Chapter 4.
- **Which characters you support.** The Unicode standard doesn't require you to support any particular set of characters, so you need to say which sets of characters you know how to handle properly (of course, if you're relying on an external library, such as the operating system, for part or all of your Unicode support, you support whatever characters it supports).

The ISO 10646 standard has formal ways of specifying which characters you support. Unicode doesn't; it asks that you say, but you can specify them any way you want, and you can specify any characters that you want.

Part of the reason that Unicode doesn't provide a formal way of specifying which characters you support is that this often varies depending on what you're doing with the characters. Which characters you can display, for example, is often governed by the fonts installed on the system

²⁴ Technically, this is only guaranteed as far back as Unicode 2.0—there were some radical changes, including removal and movement of characters, between some of the earlier versions as Unicode and ISO 10646 were brought into sync with each other.

you're running on. You might also be able to sort lists properly only for a subset of languages you can display. Some of this you can specify in advance, but you may be limited by the capabilities of the system you're actually running on.

Producing text as output

If your process produces Unicode text as output, either by writing it to a file or by sending it over some type of communication link, there are certain things you can't do. (Note that this refers to *machine-readable* output; displaying Unicode text on the screen or printing it on a printer follow different rules, outlined below.)

- Your output can't contain any code point values which are unassigned in the version of Unicode you're supporting.
- Your output can't contain the U+FFFE, U+FFFF, or any of the other non-character code point values.
- Your output *is* allowed to include code point values in the private use area, but this is strongly discouraged. Since anyone's allowed to attach any meaning they want to the private-use code points, you can't guarantee someone reading the file will interpret the private-use characters the same way you do (or interpret them at all). [You can, of course, exchange things any way you want within the universe you control, but this doesn't count as exchanging with "the outside world."] You can get around this if there's some kind of private agreement you're expecting the receiving party to uphold, but then you're not technically supporting Unicode anymore; you're supporting a higher-level protocol that uses Unicode as its basis.
- You can't produce a sequence of bytes that's illegal for whatever Unicode transformation format you're using. Among other things, this means you have to obey the shortest-sequence rule: If you're putting out UTF-8, for example, you can't use a three-byte sequence when the character can be represented with a two-byte sequence, and you can't represent characters outside the BMP using two three-byte sequences representing surrogates.

Interpreting text from the outside world

If your program reads Unicode text files or accepts Unicode over a communications link (from an arbitrary source, of course; you can have private agreements with a known source), you're subject to the following restrictions:

- If the input contains unassigned or illegal code point values, you must treat them as errors. Exactly what this means may vary from application to application, but this is basically intended to prevent security holes that could conceivably result from letting an application interpret illegal byte sequences.
- In the input contains malformed byte sequences according to the transformation format it's supposed to be in, you must treat that as an error.
- If the input contains code point values from the private-use area, you can interpret them however you want, but are encouraged to ignore them or treat them as errors. See the caveats above.
- You must interpret every code point value you purport to understand according to the semantics the Unicode standard gives to those values.
- You can handle the code point values you don't claim to support any way that's convenient for you, unless you're passing them through to another process, in which case see below.

Passing text through

If your process accepts text from the outside world and then passes it back out to the outside world (for example, you perform some kind of process on an existing disk file), you have to be sure you don't mess it up. This means that with certain exceptions, your process can't have any side effects on the text—it must do to the text only what you say it's going to do. In particular:

- If the input contains characters that you don't recognize, you can't drop them or modify them in the output. However, you are allowed to drop illegal characters from the output.
- You *are* allowed to change a sequence of code points to a canonically equivalent sequence, but you're *not* allowed to change a sequence to a *compatibility*-equivalent sequence. This will generally occur as part of producing normalized text from potentially unnormalized text. Be aware, however, that you can't claim you produce normalized text unless the process normalizing the text can do so properly on any piece of Unicode text, regardless of which characters you support for other purposes. (In other words, you can't claim you produce text in Normalized Form D if you only know how to decompose the precomposed Latin letters.)
[Note that this guarantees you're producing normalized text according to whatever version of Unicode you support—if someone passes you text that includes characters from *later* Unicode versions, you may still not normalize them properly, but this is okay, as long as you're clear about what version of Unicode you support.]
- You *are* allowed to translate the text to a different Unicode transformation format, or a different byte ordering, as long as you do it correctly.
- You *are* allowed to convert U+FEFF ZERO WIDTH NO-BREAK SPACE to U+2060 WORD JOINER, as long as it doesn't appear at the beginning of a file.

Drawing text on the screen or other output devices

Again, you're not required to be able to display every Unicode character, but for those you purport to display, you've got to do so correctly.

- You can do more or less whatever you want with any characters you encounter that you don't support (including illegal and unassigned code point values). The most common approach is to display some type of “unknown character” glyph. In particular, you're allowed to draw the “unknown character” glyph even for characters that don't have a visual representation, and you're also allowed to treat combining characters as non-combining characters.
[It's better, of course, if you don't do these things—even if you don't handle certain characters, if you know enough to know which ones not to display (such as formatting codes), or can display a “missing” glyph that gives the user some idea of what kind of character it is, that's better.]
- If you claim to support the non-spacing marks, they have to combine with the characters that precede them. In fact, multiple combining marks should combine according to the accent-stacking rules in the Unicode standard (or a more-appropriate language-specific way). Generally, this is governed by the font being used—application software usually can't influence this much.
- If you claim to support the characters in the Hebrew, Arabic, Syriac, or Thaana blocks, you have to support the Unicode bidirectional text layout algorithm.
- If you claim to support the characters in the Arabic block, you have to perform contextual glyph selection correctly.
- If you claim to support the conjoining Hangul jamo, you have to support the conjoining jamo behavior, as set forth in the standard.

- If you claim to support any of the Indic blocks, you have to do whatever glyph reordering, contextual glyph selection, and accent stacking is necessary to properly display that script. (Note that “properly display” gives you some latitude—anything that is legible and correctly conveys the writer’s meaning to the reader is good enough. Different fonts, for example, may include different sets of ligatures or contextual forms.)
- If you support the Mongolian script, you have to draw the characters vertically.
- When you’re word-wrapping lines, you have to follow the mandated semantics of the characters with normative line-breaking properties.
- You’re not allowed to assign semantics to any combination of a regular character and a variation selector that isn’t listed in the `StandardizedVariants.html` file; if the combination isn’t officially standardized, the variation selector has no effect; you can’t define ad-hoc glyph variations with the variation selectors (you can, of course, create your own “variation selectors” in the Private Use Area).

Comparing character strings

When you compare two Unicode character strings for equality, strings that are canonically equivalent should compare equal. This means you’re not supposed to do a straight bitwise comparison without normalizing the two strings first. You can sometimes get around this by declaring that you expect all text coming in from outside to already be normalized, or by not supporting the non-spacing marks.

Summary

In a nutshell, conforming to the Unicode standard boils down to just these rules:

- If you receive text from the outside world and pass it back to the outside world, don’t mess it up, even if it contains characters you don’t understand.
- In order to claim you support a particular character, you have to follow all the rules in the Unicode standard that are relevant to that character and to what you’re doing with it.
- If you produce output that purports to be Unicode text, another Unicode-conformant process should be able to interpret it properly.

CHAPTER 4 *Combining character sequences and Unicode normalization*

One feature of Unicode we've talked a lot about is combining characters. Indeed, this is one of the features of Unicode that gives it its power. However, it may be the single greatest contributor to Unicode's complexity as well. In this chapter, we'll take an in-depth look at combining characters and all of the issues that arise because of them.

Consider the following collection of characters:

á	à	ä	â
é	è	ë	ê
í	ì	ï	î
ó	ò	ö	ô
ú	ù	ü	û

What we have here are the Latin vowels with various diacritical marks added to them. Most of these letters occur fairly often in various European languages. There are twenty letters in this group, but what's interesting to note is that all of the letters in each row have something in common and all of the letters in each column have something in common. Each row consists of characters sharing the same basic letter, and each column consists of characters sharing the same diacritical mark. We can take advantage of this commonality to cut down on the number of characters we have to encode in our encoding. We already have codes for the basic letters:

a e i o u

All we really need to do is add codes for the diacritical marks:

˘ ˘ ˘ ˘

If we use these in combination with the regular letters, we can produce all twenty of the characters in our original chart, and we've only had to allocate four additional character codes, not twenty. Using this system, you'd represent the letter ö using the code for the letter o and the code for the umlaut together.

The ASCII character encoding was originally designed to work that way. On a teletype machine, the backspace control code would actually cause the carriage of the teletype to back up one space, causing the next character to be printed over the top of the preceding character. The carriage-return code originally returned the teletype carriage to the beginning of the line without advancing the platen, allowing you to overwrite a whole line. A number of ASCII characters were designed to be used in conjunction with the backspace and carriage return. For example, the underscore character we're all familiar with today in computers (“_”) was originally intended to be used with the backspace and carriage return characters to underline words. In the same way, you could get ô by sending the letter o, followed by a backspace, followed by the caret or circumflex character “^”. ASCII provided characters for the circumflex (^), tilde (~), and grave accent (`). The apostrophe/single-quote character did double duty (or maybe that's triple duty) as the acute accent, the double-quote character did double duty as the diaeresis or umlaut (the terms “diaeresis” and “umlaut” are more or less interchangeable—they refer to the same mark used for two different things), and the comma did double duty as the cedilla.

Unfortunately, this way of representing accented characters disappeared with the advent of the CRT terminal. Early CRTs weren't sophisticated enough to show two characters in the same display cell, so if they encountered two characters separated by a backspace, the second character would just replace the first one on the screen. The diacritical marks in ASCII lost their meaning as diacritical marks and just turned into symbols, and the backspace fell into disuse (it was still the code transmitted when you hit the backspace key, but you didn't see it used inside stored or transmitted text to glue characters together anymore).

For CRT display, it became necessary to have a separate character code for each combination of letter and diacritical mark, and this is the situation that prevails in most character encodings in use today, at least for European languages.

There are a few problems with this approach. One is that it can be limiting. If you have an unusual combination of base letter and accent that's not encoded, you're out of luck. Every combination of letter and mark requires the allocation of another code. In any given European language, the number of letter-mark combinations is relatively small, but the number of them for *all* languages that use the Latin alphabet is quite large—witness the large number of different “Latin alphabet” encodings in the ISO 8859 standard, for example.

Furthermore, for some other languages, base letters and diacritical marks can legitimately occur in almost any arbitrary combination. Almost any Hebrew letter, for example, can occur with almost any Hebrew vowel point. It can get worse: you can have combinations that include multiple marks

How Unicode non-spacing marks work

applied to the same base letter: In Thai, you can have arbitrary combinations of a consonant (the base letter), a vowel mark, and a tone mark.

So Unicode dispenses with the idea that a single code point always maps to a single “display cell” on the screen. A single code point might map to a single “display cell,” or perhaps multiple code points will, as in the cases we’ve examined above.

Instead, Unicode includes a whole class of characters known as “combining marks” or “non-spacing marks.” (The term “non-spacing mark” also comes from teletype machines—some European teletype machines were designed so that the carriage wouldn’t advance to the next position when they received an accent mark, allowing one to send an accented letter without using a backspace. This practice goes back to the use of “dead keys” on European typewriters for the same purpose.) A non-spacing mark doesn’t display as a self-contained unit; instead, it combines typographically with another character. A sequence of code points consisting of a regular character and one or more non-spacing marks is called a *combining character sequence*.

How Unicode non-spacing marks work

There are three basic rules governing the behavior of Unicode non-spacing marks:

- **A non-spacing mark always combines with the character that precedes it.** So if the backing store contains the following character codes:

U+006F LATIN SMALL LETTER O

U+0302 COMBINING DIAERESIS

U+006F LATIN SMALL LETTER O

that represents the sequence

Öö

and not the sequence

Oö

In other words, the diaeresis attaches to the o that precedes it.

Unicode’s designers could have gone either way with this one. It really doesn’t make much difference whether the mark attaches to the character before it or the one after it. Having the mark attach to the character after it is consistent with keyboards that use “dead key” combinations for entering accented letters, but having the mark attach to the character before it makes certain types of text analysis easier and is a little easier to understand when you’ve got multiple combining marks attached to the same character.

- **If you want to show a non-spacing mark appearing by itself, apply it to a space.** Unicode provides spacing versions of some non-spacing marks, generally for backward compatibility with some legacy encoding, but you can get any non-spacing mark to appear alone by preceding it with a space. In other words, U+0020 SPACE followed by U+0308 COMBINING DIAERESIS gives you a spacing (i.e., non-combining) diaeresis.
- **When multiple non-spacing marks that interact typographically with the base character in the same way are applied to the same character, the marks occurring closest to the base character in the backing store also occur closest to it typographically, unless a different language-specific behavior is more appropriate.** This is a little complicated. Basically, Unicode allows you to apply arbitrarily many combining marks to a character. If they attach to

different parts of the character (say, the top and the bottom), their ordering in the backing store isn't important. In other words, the sequence

```
U+006F LATIN SMALL LETTER O
U+0302 COMBINING CIRCUMFLEX ACCENT
U+0323 COMBINING DOT BELOW
```

and the sequence

```
U+006F LATIN SMALL LETTER O
U+0323 COMBINING DOT BELOW
U+0302 COMBINING CIRCUMFLEX ACCENT
```

are equivalent.

Both look like this:

Ô

But if the two marks attach to the *same* part of the character (say, they both attach to the top), then the order *is* important. In this case, the marks radiate outward from the character in the order in which they appear in the backing store.

In other words, the following sequence:

```
U+0075 LATIN SMALL LETTER U
U+0308 COMBINING DIAERESIS
U+0304 COMBINING MACRON
```

looks like this:

ü

However, the following sequence:

```
U+0075 LATIN SMALL LETTER U
U+0304 COMBINING MACRON
U+0308 COMBINING DIAERESIS
```

looks like this:

ÿ

Note that the marks never collide typographically! Proper display of a Unicode combining character sequence includes positioning all of the marks so that they don't collide.

This is the default behavior, but it isn't necessarily required. In particular, there are cases of language-specific behavior that's different from the default behavior—usually with multiple marks appearing side by side instead of stacked. For example, in Vietnamese, when a vowel has both a circumflex and an acute or grave accent on it, instead of stacking like this...

â

...the accent marks appear next to each other, like this:

ã

In Greek, when a vowel has both an accent and a diaeresis mark on it, the accent is drawn between the dots of the diaeresis...

ÿ̂

...and if a vowel has a breathing mark and an accent on it, they appear side by side:

ÿ̂̈́

Unicode doesn't explicitly specify this kind of behavior (although you can usually tell you're supposed to do it from looking at the representative glyph pictures in the Unicode standard), partially because accent placement on the same letters sometimes varies depending on what language is being displayed. If you're designing a font for a particular language, you need to know what the characters should look like for that language.

Dealing properly with combining character sequences

Systems handling Unicode text need to be careful when dealing with combining character sequences. Generally speaking, they should be thought of as single characters that happen to be stored using multiple code points, just as supplementary-plane characters are representing using two code units in UTF-16 (the so-called "surrogate pair"). This means they should travel as units.

For example, deleting a base character without deleting the non-spacing marks that follow it will cause those marks to attach themselves to the character that precedes the one that was deleted. Inserting characters into the middle of a combining character sequence can have similarly weird consequences. There's nothing in Unicode or in most string-handling libraries that prevents you from doing this kind of thing, just as they don't prevent you from breaking up surrogate pairs; code that manipulates Unicode strings has to be careful to respect combining-character-sequence boundaries.

Unicode also doesn't say anything about what happens when a string or file begins with a non-spacing mark, or when a non-spacing mark follows a paragraph boundary or some kind of control character. These are situations that shouldn't happen if combining character sequences are treated as units, but again nothing prevents them. A good text-display engine might display the non-spacing marks as though they were attached to a space, but this isn't required or guaranteed.

This means, for example, that if you're editing a document and have the insertion point to the right of an accented character, and that character is represented in memory with a combining character sequence, hitting the left arrow should *not* simply back up one code unit in the backing store—this would let the user insert a character between the base letter and the accent or delete the base letter without also deleting the accent.

Search routines similarly have to be careful about this: Searching for "resume" shouldn't find "resumé" just because the é is represented with a combining character sequence that starts with a normal e. You have to be careful to make sure that search hits land on combining-character-sequence boundaries.

Other than the sheer number of characters you have to deal with, combining character sequences may be the single biggest challenge to implementing truly Unicode-friendly text processing code. In effect, even before the introduction of surrogate pairs in Unicode 2.0, Unicode was still effectively a variable-length encoding because of combining character sequences.

The trick is, again, to disabuse yourself of the idea that there's a one-to-one correspondence between "characters" as the user is used to thinking of them and code points (or code units) in the backing store. Unicode uses the term "character" to mean more or less "the entity that's represented by a single Unicode code point," but this doesn't always match the user's definition of "character": a French speaker doesn't see é as two characters stacked on top of each other; he sees it as a single character.

A lot of writing about Unicode uses the term "grapheme" to mean "character as the user understands it." Sometimes this is sort of ambiguous—is an Indic syllable cluster (see Chapter 9) a single "grapheme" or multiple "graphemes"?—but it's important for Unicode-friendly applications to deal with text in their user interfaces as a series of graphemes and not as a series of Unicode code points (or, worse, a series of UTF-16 or UTF-8 code units).

Canonical decompositions

Combining character sequences are great for cutting down on encoding space and allowing for representation of combinations of marks you never thought of, but they have a couple of big disadvantages. They take up more space, and they're harder to process, requiring more sophisticated display technology, among other things.

For these reasons, Unicode also contains a large number of so-called "precomposed characters," code point values representing the combination of a base character and one or more non-spacing marks. Many character encoding standards, including the Latin-1 encoding used in most of Europe, use precomposed characters instead of combining character sequences. Users of these encodings are used to needing only a single code point to represent characters like é and ä, and implementations based on these encodings can adhere to the simple one-to-one relationship between code points and glyphs. Going to Unicode represents a significant step in either complexity or encoding size.

With Latin-1, there's the additional consideration that Latin-1 forms the basis of Unicode's representation of the Latin alphabet. You can convert between Latin-1 and Unicode simply by zero-padding to 16 bits or truncating to 8 bits. This wouldn't be possible if Unicode didn't have precomposed characters.

The rule in Unicode is that all precomposed characters are compatibility characters; that is, everything you can represent using precomposed characters you must also be able to represent without them. Thus, every precomposed character in Unicode has an equivalent combining character sequence. This is known as its *canonical decomposition*, and a character with a canonical decomposition is known as a *canonical composite*.

For instance, consider the letter é. It can be represented using a single code point...

```
U+00E9 LATIN SMALL LETTER E WITH ACUTE
```

How Unicode non-spacing marks work

...but this representation is, technically speaking, only there for compatibility. It *decomposes* into the letter é's *canonical* representation:

```
U+0065 LATIN SMALL LETTER E  
U+0301 COMBINING ACUTE ACCENT
```

This means that there are two perfectly legal representations in Unicode for the same character. One of them is the canonical representation, but the other is perfectly valid (and, in practice, way more likely to be used). And here's the big gotcha: *the two representations are absolutely equivalent*. In other words, a conforming Unicode implementation that supports all three code point values under discussion here is supposed to display the same thing on the screen for both representations (although few implementations actually do right now) and, more importantly, *is supposed to compare the two different sequences as equal*.

There's an important reason for this. An average user doesn't have any direct control over how what he types is stored in memory; that's the province of the programmers who wrote the software he's using. Say I'm writing some document using WhizBangWord, and WhizBangWord stores a combining character sequence in the document when I type the letter é. Say I decide I want to quote something a friend wrote to me in an email and that the email program he's using (SuperDuperMail) stores a canonical composite when you type the letter é. I paste in some text from my friend's email that includes the word "résumé." Now I do a search for the word "résumé" in my document, but WhizBangWord doesn't find "résumé" in the document, even though my friend clearly used it in the section I pasted.

The problem, of course, is that WhizBangWord is only searching for "résumé" using the representation it produces. Since SuperDuperMail uses an alternative representation, WhizBangWord doesn't recognize it. Meanwhile, I, the naive user who doesn't know or care anything about how my writing is represented inside the computer, just think my word processor is broken. This is the behavior the Unicode standard is trying to prevent. We'll talk more about this in Chapter 15, but the basic rule is that code that compares strings has to be smart enough to account for the different possible representations of the same thing in Unicode.

By definition, every canonical composite character in Unicode has a canonical decomposition. The decomposition may be more than two code points long, if, for example, the original character has more than one accent mark. The Unicode Character Database saves space by giving a one- or two-code-point decomposition for every canonical composite. If a character's decomposition is more than two code points long, the first character in its decomposition will be another canonical composite. Implementations following this approach have to repeatedly look up decompositions until they arrive at a sequence consisting exclusively of non-composite characters.

Canonical accent ordering

One of the things that combining characters let you do is attach arbitrarily many combining marks to a base character, leading to arbitrarily long combining character sequences. And since there are many languages where a single base character has at least two diacritical marks attached to it, these combinations do occur in practice. In some cases, there are even precomposed characters that include multiple diacriticals.

One of the crazy things about having both precomposed characters *and* combining character sequences is that *both* can be used together to represent the same character. Consider the letter o with a circumflex on top and a dot beneath, a letter that occurs in Vietnamese. This letter has *five* possible representations in Unicode:

```
U+006F LATIN SMALL LETTER O
U+0302 COMBINING CIRCUMFLEX ACCENT
U+0323 COMBINING DOT BELOW
```

```
U+006E LATIN SMALL LETTER O
U+0323 COMBINING DOT BELOW
U+0302 COMBINING CIRCUMFLEX ACCENT
```

```
U+00F4 LATIN SMALL LETTER O WITH CIRCUMFLEX
U+0323 COMBINING DOT BELOW
```

```
U+1ECD LATIN SMALL LETTER O WITH DOT BELOW
U+0302 COMBINING CIRCUMFLEX ACCENT
```

```
U+1ED9 LATIN SMALL LETTER O WITH CIRCUMFLEX AND DOT BELOW
```

All five of these representations are equivalent. All of them look like this:

Ô

You should get this glyph for any of these five internal representations, and if you compare any two of these sequences to each other, they should compare equal.

Generally you deal with comparing alternative representations of the same thing by converting all of them to their canonical representation and then comparing the canonical representations. This normally involves simply taking any composite characters and decomposing them, but this example shows it isn't always as simple as all that.

As it turns out, the canonical decomposition of U+1ED9 LATIN SMALL LETTER O WITH CIRCUMFLEX AND DOT BELOW is the second representation in the example: U+006F U+0323 U+0302—the letter o, followed by the underdot, followed by the circumflex.

So how do we map the other alternative representations to this one? Well, there's another step one must go through in addition to decomposing any canonical composites. Every Unicode character has a property called its *combining class*. This is a number between 0 and 255. All non-combining characters have a combining class of 0. Some combining characters also have a combining class of 0, but most combining characters have some other number for their combining class.

The combining class specifies how the combining character interacts with its base character. For example, all non-spacing marks that appear above their base character but don't attach to it have a combining class of 230. All marks that attach to the bottom of their base character have a combining class of 202.

So to get the canonical representation for some character, you first decompose the first character in the sequence, and then you take all the combining marks that follow it (which may include additional marks that followed the character you decomposed) and sort them in numerical order by their combining classes. In our example, the underdot has a combining class of 220 and the circumflex has a combining class of 230, so the underdot goes first. (Note that this ordering is arbitrary—there’s no special reason *why* the underdot has to go first; it’s just that *something* has to go first and the committee decided to have the combining classes go from bottom to top.)

Remember the part about picking up combining characters that weren’t in the original decomposition. If you start with o-circumflex followed by the combining underdot (U+00F4 U+0323), you’d begin by decomposing o-circumflex, giving you U+006F U+0302 U+0323. You’d then have to notice that the sequence actually also includes the underdot (U+0323), even though it actually follows the character you decomposed (U+006F U+0302), and *then* sort them into the correct order. The basic rule is that a combining character sequence goes from one character with a combining class of 0 to the next character with a combining class of 0.²⁵

Characters that have the *same* combining class don’t get reordered relative to each other when deriving the canonical version of a combining character sequence. This is because characters with the same combining class interact typographically. Because they interact typographically, their relative order is important: the one that occurs first in the backing store gets drawn closest to the base character.

Double diacritics

One other interesting oddity is that there are a couple of combining characters in Unicode that actually attach to *two* characters. There’s a tilde, for example, that is meant to be drawn over a pair of letters, like so:

ãã

²⁵ This is actually the rule for comparing characters. Text-rendering engines will actually have to consider longer sequences from time to time. This is because there are combining characters with a combining class of 0. Many of these, such as Indic vowel signs, simply get treated as non-combining characters, but a few, such as U+20DD COMBINING ENCLOSING CIRCLE, interact typographically with all other combining marks. The circle must maintain its position in the backing store relative to the other marks so that the rendering engine knows which marks to draw inside the circle and which to draw outside.

These special characters, which are used in Tagalog and in the International Phonetic Alphabet, are treated as normal non-spacing marks—they are stored after the first of the two characters they appear over and are just drawn so that they hang over whatever comes next.

For compatibility with some legacy encodings, the standard also includes pairs of combining characters that when drawn next to each other (i.e., when applied to succeeding characters) produce the same effect. You generally shouldn't use these.

Compatibility decompositions

Canonical composites are just one kind of compatibility character; in fact, they're only one kind of composite character. There are also *compatibility composites*. In fact, Unicode is rife with compatibility composites, which account for 3,165 assigned code point values in Unicode 3.1. All of these characters have assigned code point values in some encoding standard in reasonably wide-spread use. The compatibility characters are those characters from those standards that wouldn't have made it into Unicode on their own merits, but were given their own code point values in Unicode to allow text to be converted from the source encodings to Unicode and back again without losing any of the original information (this is usually referred to as “round-trip compatibility”), and the *compatibility composites* are the compatibility characters with *compatibility decompositions*, mappings to a preferred representation in Unicode.

There are a few important differences between canonical and compatibility decompositions:

- Compatibility decompositions may lose information. Many compatibility composites decompose to some other Unicode character plus some formatting information. The decomposition will just be to the canonical characters; since Unicode doesn't encode formatting information, the additional formatting is lost. (Theoretically, a styled-text system using Unicode as its character encoding could preserve the formatting information, but this is beyond the scope of Unicode.)
- Canonical composites can always be decomposed into pairs of characters, with a full decomposition achievable by repeatedly splitting characters into pairs of characters until no more splitting is possible. This isn't necessarily true with compatibility composites—the intermediate forms necessary to make this possible aren't all encoded as they are with canonical decompositions.

The compatibility decompositions fall into sixteen rough categories, as follows:

Superscripts and subscripts. The encoding block from U+2070 to U+209F contains superscripted and subscripted versions of the Western digits and various math symbols. There are also some scattered in the Latin-1 block and other places. The preferred representation of these things is the normal character codes with additional out-of-band formatting information indicating the subscripting or superscripting.

Font variants. The Letterlike Symbols block from U+2100 to U+214F includes a bunch of symbols that are basically letters with some kind of graphic variation. Many of these are just letters in some particular kind of font. These are marked as compatibility composites. The preferred representation is the regular letter plus out-of-band information indicating the font.²⁶

²⁶ It's worth noting that with many of the letterlike symbols, going to an alternative representation featuring a regular Unicode letter with styling information attached may *still* lose data, even though you get the same

Circled characters. Unicode includes a lot of characters with circles around them. Many are encoded in the Enclosed Alphanumerics block from U+2460 to U+24FF, and many are encoded in the Enclosed CJK Letters and Months block from U+3200 to U+327F. The preferred representation for all of these is their normal representation plus out-of-band information indicating that they're circled. (You could also represent many of these using U+20DD COMBINING ENCLOSING CIRCLE, but the official decompositions don't include this character.)

Halfwidth and fullwidth. The legacy East Asian standards include “fullwidth” variants of the ASCII characters. These normally map to fixed-width glyphs that are sized to fit in the display cells occupied by CJK ideographs, allowing them to be mixed better with the ideographs in text that is laid out vertically. There are also various Japanese and Korean characters that are sized to fill up half a display cell, allowing them to be laid out in pairs in vertical text. These live in the Halfwidth and Fullwidth Forms block from U+FF00 to U+FFEF. Again, the preferred representation is the normal representation for the characters with out-of-band information indicating how they are to be laid out.

Square forms. The legacy Japanese standards also include a large number of abbreviations, some in Japanese characters, some in Latin letters, specially laid out so as to fit in a single display cell (many of these consist of three or four Katakana characters arranged in a little square, hence the name). These are all in the CJK Compatibility block from U+3300 to U+33FF. Again, the preferred representation is the normal characters with additional out-of-band info indicating how they are to be laid out.

Vertical forms. The CJK Compatibility Forms block from U+FE30 to U+FE4F contains a number of Japanese punctuation marks and symbols rotated ninety degrees for use in vertically-laid-out text. The preferred representations for these symbols are the same as for horizontally-laid-out text. The rendering engine should be smart enough to pick the appropriate glyph for the directionality of the text.

Small forms. The Small Form Variants block from U+FE50 to U+FE6F includes a number of ASCII punctuation marks. Taken from the Taiwanese national standard, these are intended to be drawn in a full CJK display cell, but using a much-smaller-than-normal glyph. As with most other compatibility composites, these should be represented using the normal character codes and out-of-band information indicating the desired layout.

Initial, medial, final, and isolated forms. The Arabic alphabet is always written cursorily, with most letters in a word joining to their neighbors and changing shape slightly depending on how they join. In Arabic printing, each letter is generally considered to have four visual forms: initial, medial, final, and isolated, depending on how the letter joins to its neighbors. Unicode just encodes the letters, leaving it up to the rendering engine to select the appropriate glyph for each letter, but it also includes the Arabic Presentation Forms B block (U+FE70 to U+FEFE), which has a separate code point value for each glyph, for backward compatibility with systems whose rendering engines don't do this automatic glyph selection. The preferred representations of these characters are, of course, the glyph-independent versions.

appearance. This is because the Unicode general category changes: The character goes from being a symbol to being a letter. A lot of the time there won't be a practical difference, but it could cause things like search functions to work differently—searching for the word “I” in a mathematical paper, for example, might find the italicized *i* used to represent imaginary numbers.

Non-breaking characters. There are a few characters, such as U+00A0 NON-BREAKING SPACE and U+2011 NON-BREAKING HYPHEN, that have non-breaking semantics. That is, paragraph-layout engines are not supposed to put these characters either right before or right after a line boundary. These are considered compatibility composites, with the preferred representation being the normal version of the character plus out-of-band information preventing the line breaks. (Any character can be given non-breaking semantics by surrounding it with U+2060 WORD JOINER on either side, but these characters' decompositions don't include this.)

Fractions. Unicode includes a number of fractions encoded as single code points. The preferred representation of these fractions uses normal digits and U+2044 FRACTION SLASH.

Miscellaneous. Most of the remaining compatibility composites are, like the fractions, multiple characters encoded using a single code point. Some of these, such as the characters in Arabic Presentation Forms A, are ligatures, special glyphs representing the combination of two or more characters. Many, such as the Roman numerals in the Number Forms block or the parenthesized numbers in the Enclosed Alphanumerics block, are just multiple characters laid out normally. In both cases, the preferred representations just use multiple characters. (In the case of the ligatures, the rendering engine is supposed to be smart enough to form the ligature automatically when it sees the appropriate pair of characters.)

There's also a handful of characters that are *both* canonical composites *and* compatibility composites (that is, they have both canonical and compatibility decompositions, and they're different). This can happen when a character has a canonical decomposition that includes a compatibility composite. The character's canonical decomposition would leave this character intact, while the character's compatibility decomposition would start with the canonical decomposition and then find the compatibility decomposition of the composite character within it.

Singleton decompositions

There's also a small number of compatibility characters with one-character *canonical* decompositions, often referred to as "singleton decompositions." These characters count as "canonical composites," but are fundamentally different from the precomposed characters—the other canonical composites—because they don't actually represent the composition of anything. They're more akin to the compatibility composites, but they have canonical decompositions because replacing one of these characters with its decomposition doesn't lose any data.

You see this kind of thing when a character, for historical reasons, gets assigned two different code point values in Unicode, usually because this had happened in some source standard and Unicode's designers wanted to maintain round-trip compatibility.

By giving one of the two code points a singleton canonical decomposition to the other one, the Unicode standard is effectively discouraging its use by preventing it from appearing in normalized Unicode text (we'll get to normalization in a minute). Generally speaking unless you're specifically trying to maintain round-trip compatibility with some other standard, you should avoid characters with singleton decompositions.

For example, the CJK Compatibility Ideographs block includes a couple hundred Chinese characters that are also encoded in the CJK Unified Ideographs area. The duplicates in the compatibility area are, for all intents and purposes, identical to their counterparts in the main CJK area. Unicode has this second set of code point assignments for these characters because some source standard included duplicate encodings for these characters. For example, the majority of them come from the Korean KS C 5601 standard, which gave a separate code point assignment to each alternate pronunciation of

the character. Unicode rejects this approach to encoding these characters, since they don't differ in appearance or semantics, but retains the duplicate encodings for round-trip compatibility with KS C 5601. But all of the CJK compatibility characters have singleton decompositions to their counterparts in the main CJK Ideographs area, indicating that they really shouldn't be used in Unicode text unless that interoperability is really important.

Hangul

Hangul is the name of the Korean writing system. Some linguists consider it the most perfect writing system in common use. Invented in the fifteenth century by King Sejong, it's basically an alphabetic system, but the letters, called *jamo*, are arranged into blocks representing whole syllables, probably because of the influence of the Chinese characters. Hangul syllables look kind of like Chinese characters and are laid out on the page in a similar fashion (in fact, Korean writing often uses Chinese characters in addition to Hangul syllables).

Hangul jamo come in three categories, reflecting the structure of Korean pronunciation: *choseong*, which are initial consonants, *jungseong*, which are vowels, and *jongseong*, which are trailing consonants. A Korean syllable consists of a sequence of choseong, followed by a sequence of jungseong, optionally followed by a sequence of jongseong, all arranged to fit into a single square display cell.

Unicode provides two alternative methods of encoding Hangul. It includes a complete collection of precomposed Hangul syllable blocks, encoded in the range from U+AC00 to U+D7FF, but also includes codes for the individual jamo in the range from U+1100 to U+11FF. A modern Korean syllable can be represented either with a single code point representing the whole syllable, or with a sequence of code points for the individual jamo.

The jamo code points have conjoining semantics—they're intended to be used together to represent whole syllables (there's another set of jamo, encoded from U+3130 to U+318F, which stand alone, but they're considered compatibility composites). The precomposed syllables are considered to be canonical composites—they all have canonical decompositions into conjoining jamo.

The opposite isn't true, however—it's possible to produce sequences of conjoining jamo that can't be represented using the precomposed syllables. Many of these are mere nonsense, or sequences intended to give special effects (there are invisible filler characters that can be used to represent isolated jamo or syllable blocks with missing jamo), but there are also archaic Korean syllables that can only be represented using the conjoining jamo. For this reason, as well as consistency with the rest of Unicode, the conjoining jamo are the canonical representation of Korean Hangul, even though the precomposed syllables are more commonly used.

A sequence of conjoining Hangul jamo is a combining character sequence, but it's different from other Unicode combining character sequences in that it doesn't consist of a non-combining character followed by one or more combining characters. Instead, all of the jamo sometimes function as combining characters and sometimes as non-combining characters. Basically, all of the jamo are non-combining characters except when they appear together. When they appear together, the Unicode standard gives an algorithm for locating syllable breaks. (The algorithm is actually quite simple: leading consonants are followed by vowels, which are followed by trailing consonants. The appearance of a character out of that sequence marks the beginning of a new syllable. Since every

syllable is supposed to have an initial consonant and a vowel, Unicode provides invisible filler characters to allow the representation to be more regular.)

The canonical decomposition of a precomposed syllable always consists of a single leading consonant, a single vowel, and, optionally, a single trailing consonant. This is possible, even for syllables like *krang*, because all of the compound consonants and all of the diphthongs in Korean are encoded as single code points. The precomposed syllables and conjoining jamo are encoded in such a way as to allow for algorithmic conversion between the two representations. Extra work may still have to be done, however, for some archaic Hangul syllables or nonstandard decomposed representations such as representing the *kr* in *krang* using the separate *k* and *r* jamo. (The handling of nonstandard representations of modern syllables is neither specified nor required by the standard, although very early versions had it—it's an additional language-specific thing an application can handle if it wishes.)

Finally, there's a side effect of the decomposition of precomposed Hangul syllables into conjoining jamo. If you *mix* precomposed syllables and conjoining jamo in the same passage of text, they can combine with *each other*. To see why, consider this character:

뢴

This character can be represented either this way...

U+B8B4 HANGUL SYLLABLE ROEN

...or this way:

U+1105 HANGUL CHOSEONG RIEUL

U+116C HANGUL JUNGSEONG OE

U+11AB HANGUL JONSEONG NIEUN

Likewise, this character...

뢹

...can be represented either this way...

U+B878 HANGUL SYLLABLE ROE

...or this way:

U+1105 HANGUL CHOSEONG RIEUL

U+116C HANGUL JUNGSEONG OE

Notice that the second sequence is a subset of the first. So if you have this sequence...

U+B878 HANGUL SYLLABLE ROE

U+11AB HANGUL JONSEONG NIEUN

...and you decompose U+B878 into conjoining jamo...

U+1105 HANGUL CHOSEONG RIEUL

Compatibility decompositions

U+116C HANGUL JUNGSEONG OE
U+11AB HANGUL JONSEONG NIEUN

...you get the sequence we started with, the decomposed version of `꺠`. This means that this sequence...

U+B878 HANGUL SYLLABLE ROE
U+11AB HANGUL JONSEONG NIEUN

...has to be another alternate representation for `꺠`. So precomposed Hangul syllables have to combine with conjoining jamo.

For more information on Hangul, see Chapter 10.

Unicode normalization forms

Of course, an encoding that provides so many alternative ways of representing characters can give rise to text that is much more difficult than necessary to process. In particular, comparing strings for equality is a big challenge when significantly different sequences of bits are supposed to be treated as equal. One way to deal with this is to require that text be *normalized*, or represented in a uniform manner, or to normalize text at some well-defined point so as to make things like comparing for equality simpler.

Of course, by defining something as the “canonical representation” of a particular idea, you’re nominating that as the form you normalize to. In this way, Unicode 1.x and 2.x could be thought to have either one or two normalized forms. Unicode 3.0 has four normalized forms, adding two more and formalizing the definitions of all four. The four Unicode normalized forms are as follows:

Normalized Form D is the traditional normalized form of Unicode, where all canonical composites (including Hangul syllables) are represented using their canonical decompositions.

Normalized Form KD also eliminates the compatibility composites (the K stands for “kompatibility”), representing everything using both compatibility and canonical decompositions. Form KD is not recommended for general use (e.g., for transmitting things around on the Internet) because it potentially loses data.

Using the decomposed forms of everything, of course, makes everything bigger, which is a common objection in areas where the cost of every additional code point is significant, such as when transmitting things over the Internet, or for people used to using character encodings that have only precomposed characters, so there was a strong call for a form of Unicode that favors the composed forms of everything over the decomposed forms. Two additional normalized forms were created in response to this need.

Normalized Form C favors the canonical-composite forms of the characters in Unicode. Of course, this is more complicated than decomposition because there are things which can *only* be represented using combining character sequences; there is no composite form that can be used as

an alternative. For this reason, you get to Normalized Form C by first converting to Normalized Form D and then composing everything it's possible to compose.

Another additional complication with using composed characters comes if more canonical composites are added to the standard in the future. Text which is in Normalized Form C with version 3.0 of Unicode might not be in Normalized Form C according to a future Unicode version with more canonical composites if the definition of Normalized Form C amounted to “compose everything you can.” To keep the definition of Normalized Form C consistent, the allowable canonical composites are fixed for all time to be just the ones in Unicode 3.0.²⁷ One ironic side effect is that a hypothetical Unicode 4.0 canonical composite would be represented by its *decomposed* form in Normalized Form C.

The Internet environment is an excellent example of an environment where Unicode normalization comes in handy. XML is a Unicode-based standard, and without requiring text to be normalized, matching up XML tags would be very difficult if XML were to follow the Unicode equivalence rules. The World Wide Web Consortium (W3C) deals with this by declaring a rule known as Uniform Early Normalization, which requires that all processes that produce text in the Internet environment produce it in normalized form. Processes that consume or pass text through (which are generally expected to be way more common) can assume that the text they receive is already in normalized form and don't have to do it themselves. This means, among other things, that they can depend on character sequences always being represented with the same sequences of code points, making it possible to compare two strings for equality by doing simple binary comparison on the code points.

Normalized Form KC is the counterpart of Normalized Form KD, taking compatibility decompositions into account as well. It's a little silly to favor the compatibility composites over the regular characters (would you want to convert two capital Is in a row into U+2161 ROMAN NUMERAL TWO everywhere you encountered them, for example?), and Normalized Form KC doesn't do that.

It does the opposite. To get to Normalized Form KC, you first convert to Normalized Form KD (compatibility-decompose everything), then substitute precomposed characters (*canonical compositions*) where possible.

We'll take an in-depth look at implementation strategies for the various pieces of Unicode normalization—canonical decomposition, compatibility decomposition, and canonical reordering—in Chapter 14.

Grapheme clusters

The newest version of Unicode, Unicode 3.2, introduces a new concept called the “grapheme cluster.” Actually, the concept isn't all that new; what Unicode 3.2 essentially does is formalize a concept that was already out there, nailing down a more specific definition and some related character properties and giving the concept a new name.

²⁷ Actually, it isn't even all the precomposed characters in Unicode 3.0. A fair number of Unicode 3.0 characters with canonical decompositions are excluded from Normalized Form C because the decomposed alternatives were actually considered preferable. We talk about this more in Chapter 5.

A grapheme cluster is a sequence of one or more Unicode code points that are to be treated as a single unit by various processes:

- Text-editing software should generally allow placement of the cursor only at grapheme cluster boundaries: Clicking the mouse on a piece of text should place the insertion point at the nearest grapheme cluster boundary, and the arrow keys should move forward and back one grapheme cluster at a time.
- Text-rendering software should never put a line break in the middle of a grapheme cluster (since the individual characters in a grapheme cluster generally interact typographically in ways that make it difficult to separate out the pieces, you generally couldn't put a line break in the middle of a grapheme cluster without deforming it visually in some way).
- Sort orders for different languages generally give a relative ordering for grapheme clusters, not necessarily individual characters. For instance, in Spanish, the sequence *ch* is treated as a separate letter that sorts between *c* and *d*. Therefore, in Spanish, *ch* would generally be considered a grapheme cluster.
- Search algorithms should only count a matching sequence in the text being searched as a hit if it begins and ends on a grapheme cluster boundary.

Exactly what constitutes a grapheme cluster may vary from language to language; what Unicode 3.2 attempts to do is set down a default definition of grapheme cluster that can then be tailored as necessary for specific languages.

The term “grapheme cluster” is new in Unicode 3.2;²⁸ in earlier versions of Unicode, the same concept was generally just called a “grapheme,” although other phrases such as “logical character” or “or “user character” were also thrown around. The problem with “grapheme” is that it has a specific meaning in linguistics, and the Unicode definition didn't agree completely with the common linguistic definition. A grapheme is essentially the smallest meaning-bearing unit of a writing system, as understood by the average user of that writing system.

This definition is necessarily subjective and language-specific: Is *ö* a single grapheme or two? An English speaker would probably say two—an *o* and an umlaut, but a Swede would probably say one: *ö* is a separate letter of the Swedish alphabet. A German speaker could go either way. And a clump of letters that you might want a text editor to treat as a single unit, such as a Hindi syllable cluster, might still be recognized by an average Hindi speaker as several letters—a cluster of graphemes.

Again, exactly what constitutes a grapheme cluster may vary from language to language, user to user, or even process to process. Unicode sets forth a standard default definition of a grapheme cluster that should prevail in the absence of a more application-specific definition. A default grapheme cluster is one of the following things:

- A base character followed by zero or more combining marks, i.e., a normal combining character sequence.
- A Hangul syllable, whether represented using a precomposed-syllable code point, a series of conjoining-jamo code points, or a combination of the two.

²⁸ And may, in fact, change before you read this, although that's somewhat doubtful. There is widespread revulsion to this name in the Unicode community; no one really likes it (a fairly common complaint is that it sounds like some kind of breakfast cereal—“Have you tried new Grapheme Clusters? Now with more fiber!”), but no one can come up with anything better.

- An “orthographic syllable” in one of the Indic scripts: Such a sequence consists of a single consonant or independent vowel, or a series of consonants joined together with *virama* marks, optionally followed by a dependent vowel sign. (For more information on this, see Chapter 9.)
- The CRLF sequence (i.e., U+000D U+000A, the ASCII carriage-return-line-feed combination) is a grapheme cluster.
- A user-defined grapheme cluster, which uses U+034F COMBINING GRAPHEME JOINER to “glue” together what would otherwise be separate grapheme clusters.

Generally speaking, grapheme cluster boundaries don’t affect the way text is drawn, but there are a few situations where it does:

- An enclosing mark surrounds all the characters that precede it, up to the nearest preceding grapheme cluster boundary. U+034F COMBINING GRAPHEME JOINER can thus be used to get an enclosing mark to enclose more than one character:

```
U+0031 DIGIT ONE
U+0032 DIGIT TWO
U+20DD COMBINING ENCLOSING CIRCLE
```

...gets drawn like this:

1②

The 1 and the 2 are separate grapheme clusters, so the circle only surrounds the 2. But if you put a grapheme joiner between them...

```
U+0031 DIGIT ONE
U+034F COMBINING GRAPHEME JOINER
U+0032 DIGIT TWO
U+20DD COMBINING ENCLOSING CIRCLE
```

...you get this:

①2

The grapheme joiner eliminates the grapheme cluster boundary between the 1 and the 2, binding them together as a single grapheme cluster. The circle then surrounds both of them.

- A non-spacing mark applies to all characters up to the nearest preceding grapheme cluster boundary, provided they’re not a conventional combining character sequence (if they are, the non-spacing mark is just part of that combining character sequence). This rule is here mainly to say that non-spacing marks following Hangul syllables apply to the whole syllable, regardless of how it’s represented.

The second definition is rather nebulous: What should you get if you have this, for example?

```
U+0061 LATIN SMALL LETTER A
U+034F COMBINING GRAPHEME JOINER
U+0065 LATIN SMALL LETTER E
U+0303 COMBINING TILDE
```

The answer seems to be that the tilde goes over the e: the grapheme joiner would have no effect in this case. But it’s equally reasonable to assume that the tilde would be centered between the a and the e, or even stretched to extend over both the a and the e: This would be a way to get more double-diacritics without having to encode them separately. Neither seems to be the intent, however; indeed there’s been talk of getting this effect by taking advantage of the enclosing-mark behavior. You’d introduce a new INVISIBLE ENCLOSING MARK character and apply the non-spacing mark to *that*: In the same way that a non-spacing mark after an enclosing mark would go outside the enclosing mark (and effectively be applied to it), the invisible enclosing mark would

cause any non-spacing marks that follow it to be applied to the entire sequence of characters it “encloses.” This isn’t in Unicode 3.2, but might be in a future version of Unicode.

The Unicode 3.2 version of the Unicode Character Database specifies three new character properties to aid in the determination of grapheme cluster boundaries: a grapheme cluster consists of a Grapheme_Base character followed by zero or more Grapheme_Extend characters, optionally followed by a Grapheme_Link character, another Grapheme_Base character, and zero or more Grapheme_Extend characters, and so on.

On top of this basic definition you have to layer some extra rules to get Hangul syllables and the CRLF sequence to be treated as grapheme clusters. In addition, Join_Control characters are defined as transparent to the algorithm: they neither cause nor prevent grapheme cluster boundaries.

As some of the above discussion suggests, at the time of this writing (January 2002), the exact semantics and behavior of grapheme clusters aren’t completely nailed down. This discussion may well deviate some from how it ultimately comes down in the final version of Unicode 3.2.

CHAPTER 5 *Character Properties and the Unicode Character Database*

One of the things that makes Unicode unique is that it goes well beyond just assigning characters to numbers. The Unicode standard also provides a wealth of information on how the characters are to be used, both together and individually. The Unicode standard comprises not just the information in the Unicode book, but also the Unicode Character Database, a collection of files included on the CD that comes with the book.

The Unicode Character Database is the most active part of the standard, and actually changes more frequently than the book does. In fact, changes to the database actually happen more often than Unicode Technical Reports are issued. If there are significant changes to the standard, such as new characters added, a technical report will be issued or a new book will be put out, but error fixes and clarifications can happen to the database without changes to any other documents. Changes to the database do cause Unicode's update version number (the third number) to be bumped. As of this writing (January 2002), the current Unicode version is 3.1.1, with version 3.2 in beta.

The Unicode Character Database may fairly be considered the heart of the standard. Implementations of many processes that operate on Unicode text are based on tables of data that are generated from the Unicode Character Database. This chapter will take an in-depth look at the Unicode Character Database and at the various properties that each Unicode character has.

Where to get the Unicode Character Database

Since the Unicode Character Database can be updated relatively frequently, it's usually a good idea not to rely too heavily on the version on the CD that comes with the Unicode standard. In fact, it's pretty obsolete right now, as the structure of the data files was changed when Unicode 3.1 came out.

It'll be mostly right, but may differ in some particulars from the most current version, and it'll be organized differently.

You can always find the most current version on the Unicode Web and FTP sites. The URL of the Unicode Data page, which includes links to all the files, is

<http://www.unicode.org/unicode/onlinedat/online.html>

The current version of the Unicode Character Database, specifically, is always at

<http://www.unicode.org/Public/UNIDATA/>

The parent directory also includes a number of other useful folders of data. Among these are the following:

<http://www.unicode.org/Public/MAPPINGS/> contains a bunch of files giving mappings of characters between Unicode and other encoding standards. The most interesting part of this section is the VENDORS folder, which contains mappings between Unicode and various vendor-defined encoding schemes. The other sections are useful as well, but haven't been updated in a long time.

<http://www.unicode.org/Public/PROGRAMS/> contains sample code to do various interesting things. It includes demonstration code in C for converting between the Unicode encoding forms, a demonstration program for converting from SHIFT-JIS to Unicode, a sample implementation (in Java) of the Standard Compression Scheme for Unicode, and some other stuff.

<http://www.unicode.org/Public/TEXT/> contains some old Unicode Technical Reports and some other stuff of mainly historical interest.

Finally, <http://www.unicode.org/Public/BETA/> contains early versions of data files proposed for future versions of the Unicode Character Database.

In addition to these folders, there are other folders here with version numbers in their names. These are old updates to the Unicode Character Database. These are mainly of historical interest, but are also useful for figuring out what changed from version to version. For instance, if you have code based on an old version of Unicode, you can compare the files for that version to the ones for the current version to figure out what changes you need to make to your own data tables to make them consistent with the current version of Unicode.

The UNIDATA directory

The UNIDATA folder on the Unicode FTP/Web site (<http://www.unicode.org/Public/UNIDATA/>) is the official repository of the Unicode Character Database. Here's a quick rundown of what's in this directory:

- **ArabicShaping.txt** is a data file that groups Arabic and Syriac letters into categories depending on how they connect to their neighbors. The data in this file can be used to put together a minimally-correct implementation of Arabic and Syriac character shaping.

- **BidiMirroring.txt** is useful for implementing a rudimentary version of mirroring. For the characters whose glyph in right-to-left text is supposed to be the mirror image of their glyph in left-to-right text (the characters with the “mirrored” image), this file identifies those that have another character in Unicode that has the correct glyph in left-to-right text (or a similar glyph). This left you implement mirroring for many (but not all) Unicode characters by simply mapping them to other characters that have the right glyph.
- **Blocks.txt** breaks the Unicode encoding space down into named blocks (e.g., “C0 Controls and Basic Latin” or “Miscellaneous Symbols”) and specifies which block each Unicode character is in.
- **CaseFolding.txt** can be used to implement case-insensitive string comparison . It maps every Unicode code point to a case-insensitive representation—a code point or sequence of code points that represent the same character with case distinctions taken out (this was basically derived by converting each character to uppercase and then converting the result to lowercase).
- **CompositionExclusions.txt** lists characters with canonical decompositions that should not appear in normalized Unicode text.
- **EastAsianWidth.txt** assigns every character to a category that describes how it should be treated in East Asian typography. This is the data file for UAX #11.
- **Index.txt** is a soft copy of the character names index from the Unicode standard book. It gives you a way to quickly look up a character’s code point value if you know its name.
- **Jamo.txt** assigns each of the conjoining Hangul jamo characters a “short name” which is used to derive names for the precomposed Hangul syllables.
- **LineBreak.txt** assigns every Unicode character to a category that determines how it should be treated by a process that breaks text up into lines (i.e., a line-breaking process uses the character categories defined in this file to determine where to put the line breaks). This is the data file for UAX #14.
- **NamesList.txt** is the source file used to produce the code charts and character lists in the Unicode standard. **NamesList.html** explains its structure and use.
- **NormalizationTest.txt** is a test file that can be used to determine whether an implementation of Unicode normalization actually conforms to the standard. It includes a bunch of unnormalized test strings along with what they should turn into when converted to each of the Unicode normalized forms.
- **PropertyAliases.txt**, which is new in Unicode 3.2, gives abbreviated and long names for all the Unicode character properties. The idea is that these names can be used in regular-expression languages and other search facilities to identify groups of characters with a given property.
- **PropertyValueAliases.txt** is also new in Unicode 3.2. For each property listed in PropertyAliases.txt that isn’t a numeric or Boolean property, it gives abbreviated and long names for all the possible values of that property. These names are intended to be used with the names in PropertyAliases.txt to provide shorthand ways of referring to sets of characters with common properties for regular-expression engines and similar facilities.
- **PropList.txt** is one of the two main files defining character properties. It acts as a supplement to the UnicodeData.txt file. **PropList.html** explains what’s in it.
- **ReadMe.txt** gives the version number for the characters in the directory, along with points to other files that explain its contents.
- **Scripts.txt** groups the Unicode characters by “script,” that is, the writing system that uses it.
- **SpecialCasing.txt** gives complex case mappings for various Unicode characters. It lists all the characters that have a non-straightforward case mapping.
- **StandardizedVariants.html**, which is new in Unicode 3.2, lays out just what combinations of regular character and variation selector are legal and what glyphs they represent.

- **UnicodeCharacterDatabase.html** is an overview of all the files in the Unicode Character Database and provides references to the sections of the Unicode standard that deal with each of them.
- **UnicodeData.txt** is the original, and still primary, data file in the Unicode Character Database, defining most of the most important character properties for each character. **UnicodeData.html** explains the contents of UnicodeData.txt.
- **Unihan.txt** is a huge data file containing a wealth of information about the various Han (Chinese) characters.

In addition to the files listed above, you'll find a bunch of files with the word "Derived" in their names. Each of these files contains information that can be derived from information in one or more of the files listed above (most often, UnicodeData.txt). Generally, each of these files isolates a single character property from UnicodeData.txt and lists the characters in groups according to the value they have for that property. This can make it easier to figure out which characters are in certain groups, and makes some of the Unicode data easier to work with. Among the derived-data files:

- **DerivedAge.txt**, which is new in Unicode 3.2, gives, for each character in Unicode, the version of Unicode in which it was introduced.
- **DerivedBidiClass.txt** groups the Unicode characters according to their directionality (i.e., how they're treated by the Unicode bidirectional layout algorithm). The normative versions of these properties are in UnicodeData.txt.
- **DerivedBinaryProperties.txt** lists groups of characters that have a value of "yes" for various properties that have yes/no values. Right now, this file only lists the characters that have the "mirrored" property, the one binary property given in UnicodeData.txt. PropList.txt gives more binary properties, but is already formatted the same as this file.
- **DerivedCombiningClass.txt** groups characters according to their combining class, which is officially specified in UnicodeData.txt.
- **DerivedCoreProperties.txt** lists characters belong to various general groups (e.g., letters, numbers, identifier-starting characters, identifier-body characters, etc.). The information in this file is derived from information in UnicodeData.txt and PropList.txt.
- **DerivedDecompositionType.txt** groups characters with decompositions according to the kind of decomposition they have (canonical versus compatibility). For the compatibility composite characters, it further breaks them down according to the type of their compatibility decomposition (that is, what information it loses). This information is derived from UnicodeData.txt.
- **DerivedJoiningGroup.txt** and **DerivedJoiningType.txt** both present information in the ArabicShaping.txt file organized in a different way. DerivedJoiningGroup.txt groups Arabic and Syriac letters according to their basic shape, and DerivedJoiningType.txt groups them according to the way they join to their neighbors.
- **DerivedLineBreak.txt** presents the same information as LineBreak.txt, but organizes it differently: Instead of listing the characters in code-point order and giving each one's line-break category, it lists the categories and then lists the characters in each category.
- **DerivedNormalizationProperties.txt** gives groups of characters that can be used to optimize implementations of Unicode normalization.
- **DerivedNumericType.txt** and **DerivedNumericValues.txt** group together various characters that represent numeric values. DerivedNumericType.txt groups them together according to whether they are decimal digits or something else, and DerivedNumericValues.txt groups them together by the numeric values they represent. Both are derived from UnicodeData.txt.
- **DerivedProperties.html** is an overview document explaining the contents of the other derived-data files.

You'll often see the filenames of the files in the UNIDATA directory with versions numbers appended (there are also version numbers in the comments in the files). These version numbers are keyed to a particular version of Unicode: UnicodeData-3.1.1.txt, for example, would be the version of UnicodeData.txt that goes with version 3.1.1 of the Unicode standard. You can't always rely on this, however, since not every file changes with every new version of Unicode (there is no UnicodeData-3.1.1.txt, for example—the Unicode 3.1.1 data changes were in other files). You'll also sometimes see one-digit version numbers on some of the files—originally some of the files, especially those that are keyed to Unicode Technical Reports, used a different numbering scheme. The changed for each data file as its corresponding Technical Report got adopted into the standard (i.e., turned into a Unicode Standard Annex).

In the sections below, we'll take a closer look at the most important of the files in the Unicode Character Data, and at the properties they define.

UnicodeData.txt

The nerve center of the Unicode Standard is the UnicodeData.txt file, which contains most of the Unicode Character Database. As the database has grown, and as supplementary information has been added to the database, various pieces of it have been split out into separate files, but the most important parts of the standard are still in UnicodeData.txt.

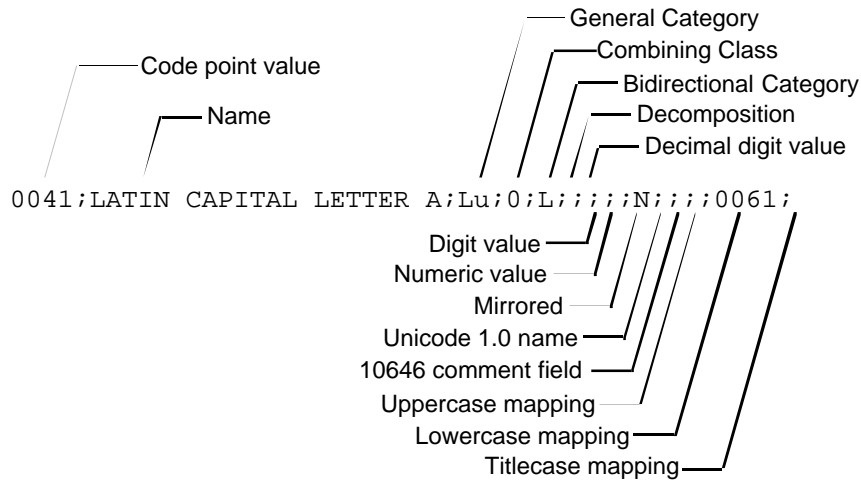
The designers of Unicode wanted the database to be as simple and universal as possible, so it's maintained as a simple ASCII text file (we'll gloss over the irony of having the Unicode Character Database stored in an ASCII text file). Again, for ease of parsing, this file is a simple semicolon-delimited text file. Each record in the database (i.e., the information pertaining to each character) is separated from the next by a newline (ASCII LF, hex 0x0B), and each field in a record (i.e., each property of a single character) is separated from the next by a semicolon.

It's not beautiful, and it's not cutting-edge, but it has the major advantage of being a format that is readable and parsable by just about anything. These days, of course, they could have also used XML, and this is discussed from time to time, but compatibility with existing software that reads this file and easy diff-ability between versions are compelling reasons to keep it in the format it's in. Besides, it's almost trivially easy to write a program to convert between this format and an XML-based format.

Here are a few sample records from the database:

```
0028;LEFT PARENTHESIS;Ps;0;ON;;;Y;OPENING PARENTHESIS;;;
0033;DIGIT THREE;Nd;0;EN;;3;3;3;N;;;;
0041;LATIN CAPITAL LETTER A;Lu;0;L;;;N;;;0061;
022C;LATIN CAPITAL LETTER O WITH TILDE AND MACRON;Lu;0;L;00D5
    0304;;;N;;;022D;
0301;COMBINING ACUTE ACCENT;Mn;230;NSM;;;N;NON-SPACING ACUTE;Oxia;;;
2167;ROMAN NUMERAL EIGHT;Nl;0;L;<compat> 0056 0049 0049 0049;;;8;N;;;2177;
```

A fair amount of the information is immediately accessible (or at least gets that way with practice), but it's easy to get lost in the semicolons. Here's the record for the letter A with the fields labeled:



So, this tells us that the Unicode code point value U+0041 has the name LATIN CAPITAL LETTER A. It has a general category of “uppercase letter” (abbreviated “Lu”), it’s not a combining mark (its combining class is 0), it’s a left-to-right character (abbreviated “L”), it’s not a composite character (no decomposition), it’s not a numeral or digit (the three digit-value fields are all empty), it doesn’t mirror (abbreviated “N”), its name in Unicode 1.0 was the same as its current name (empty field), it maps to itself when being converted to uppercase or titlecase (empty fields), and it maps to U+0061 (LATIN SMALL LETTER A) when being converted to lowercase.

In order from left to right, the properties defined in UnicodeData.txt are:

- The character’s code point value (technically, this isn’t a property; it’s the value we’re assigning properties to), expressed as a four-to-six-digit hexadecimal numeral.
- The character’s name.
- The character’s general category. This tells you whether the character is a letter, a digit, a combining mark, and so on.
- The character’s combining class. This value is used to order combining marks when converting Unicode text to one of the Unicode normalized forms.
- Bidirectional category. This value specifies the character’s directionality and is used by the bidirectional text layout algorithm.
- Decomposition. For decomposing characters, the characters in the decomposition and a tag indicating the decomposition type.
- Decimal digit value. If the character can be used as a decimal digit, the numeric value it represents.
- Digit value. If the character can be used as a digit (decimal or not), the numeric value it represents.
- Numeric value. If the character can be used alone as a numeral, the numeric value it represents.
- Mirrored. Says whether the character adopts a mirror-image glyph when surrounded by right-to-left text.
- Unicode 1.0 name. If the character existed in Unicode 1.0 and had a different name, its old name.
- 10646 comment field. If the character has a comment attached to it in the ISO 10646 standard, it’s given here.
- Uppercase mapping. If the character maps to a different character when converted to uppercase, that character is given here.

- Lowercase mapping. If the character maps to a different character when converted to lowercase, that character is given here.
- Titlecase mapping. If the character maps to a different character when converted to titlecase, that character is given here.

So if we go back to our sampling of entries, we can see interpret it a little more clearly:

- U+0028's current name is LEFT PARENTHESIS, but it used to be STARTING PARENTHESIS. It's a member of the "starting punctuation" (Ps) category, has a combining class of 0, does mirror ("Y"), and has neutral ("ON") directionality.
- U+0033 DIGIT THREE is a decimal digit ("Nd"), has weak left-to-right directionality ("EN"), doesn't mirror ("N"), and has a numeric value of 3.
- U+0041 LATIN CAPITAL LETTER A is an uppercase letter ("Lu"), has strong left-to-right directionality ("L"), and maps to U+0061 LATIN SMALL LETTER A when converted to lower case.
- U+022C LATIN CAPITAL LETTER O WITH TILDE AND MACRON is an uppercase letter with strong left-to-right directionality, maps to U+022D when converted to lower case, and has a canonical decomposition to U+00D0 U+0304.
- U+0301 COMBINING ACUTE ACCENT is a non-spacing mark ("Mn") with a combining class of 230 (which means it appears above the base character), it's transparent to the bi-di algorithm ("NSM"), and is the same as the Greek *oxia* according to ISO 10646.
- U+2167 ROMAN NUMERAL EIGHT is a letter number ("NI") with a numeric value of 8. It maps to a Roman number made out of lower-case letters (U+2177) when mapped to lower case, and it has a compatibility decomposition to "VIII" (U+0056 U+0049 U+0049 U+0049).

We'll actually take a closer look at each of these properties and how they're specified later in this chapter.

PropList.txt

The UnicodeData.txt file is supplemented with a variety of other files that supply additional information about Unicode characters. A bunch of these properties (e.g., East Asian width, Jamo short name, etc) have their own files, but a bunch of them are given in PropList.txt. In particular, PropList.txt lists so-called "binary" properties, categories a character is either in or it isn't. (For example, it's either a digit or it isn't.)

PropList.txt has the format used for most of the other files: Each line represents either a single character (in which case it starts with that character's code point), or a range of characters (in which case it starts with the starting and ending code points in the range, separated with ".."). There's then a semicolon and the property's value. Usually the line ends with a comment that gives the character's name for readability. (Comments are set off with pound signs.)

For example, here are the first two entries in PropList.txt:

```
0009..000D    ; White_space # Cc    [5] <control>..<control>
0020         ; White_space # Zs          SPACE
```

The first line assigns the property "White_space" to the code point values from U+0009 to U+000D. The comment indicates the code points in this range have the "Cc" general category (we'll talk about

that later in this chapter), there are five characters in the range, and the characters at the beginning and end of the range are control characters.

The second line assigns the property “White_space” to the code point value U+0020. The comment indicates that U+0020’s name is SPACE and that it belongs to the “Zs” category.

Files with this format are either ordered by code point value (e.g., LineBreak.txt) or by category value. PropList.txt is sorted by category.

Here’s a rundown of the properties defined in PropList.txt:

- **White_space.** Characters that are considered “white space.” This includes spaces and line and paragraph separators, along with control characters that had these functions in older standards.
- **Bidi_Control.** Special invisible formatting codes that affect the behavior of the bidirectional layout algorithm.
- **Join_Control.** Special invisible formatting codes that affect the glyph selection process (by causing characters to join together or split apart).
- **Dash.** Dashes and hyphens.
- **Hyphen.** Hyphens and characters that behave like hyphens.
- **Quotation_Mark.** Quotation marks.
- **Terminal_Punctuation.** Punctuation marks that occur at the end of a sentence or clause.
- **Other_Math.** Characters that can be used as mathematical operators but weren’t given the “Mathematical symbol” (“Sm”) general category.
- **Hex_Digit.** Characters that are used as digits in hexadecimal numerals.
- **ASCII_Hex_Digit.** ASCII characters in the Hex_Digit category.
- **Other_alphabetic.** Alphabetic characters that weren’t assigned to one of the “letter” general categories.
- **Ideographic.** The Han characters and radicals. All of the characters with the Radical and Unified_Ideograph properties also have this property.
- **Diacritic.** Diacritical marks.
- **Extender.** Characters that elongate, either graphically or phonetically, the letters that precede them.
- **Other_Lowercase.** Characters that weren’t assigned to the “lowercase letter” (“Ll”) general category that are still considered lowercase.
- **Other_Uppercase.** Characters that weren’t assigned to the “uppercase letter” (“Lu”) general category that are still considered uppercase.
- **Noncharacter_Code_Point.** Code point values that are specifically called out by the standard as not representing characters and not being legal in Unicode text.

The Unicode 3.2 version of PropList.txt adds a number of new character properties. At the time of this writing, Unicode 3.2 was still in beta, so the list of properties may change, as may the characters assigned to them. But as of the beta version of Unicode 3.2, they were:

- **Other_Grapheme_Extend.** Characters with the Grapheme_Extend property that don’t automatically get this property by virtue of belonging to some class. Grapheme extenders are characters that always occur in the same grapheme cluster as the character that precedes them.
- **Grapheme_Link.** Characters that cause both the preceding and following characters to be considered part of the same grapheme cluster.

General category

- **IDS_Binary_Operator.** Ideographic description characters that take two operands. This property is here to aid code that interprets ideographic description sequences. See Chapter 10 for more on ideographic description sequences.
- **IDS_Tertiary_Operator.** Ideographic description characters that take three operands. This property is here to aid code that interprets ideographic description sequences. See Chapter 10 for more on ideographic description sequences.
- **Radical.** Characters that are parts of Han characters. This property also exists mostly to aid in interpreting ideographic description sequences.
- **Unified_Ideograph.** Full-blown Han characters.
- **Other_Default_Ignorable_Code_Point.** Invisible formatting characters and other special characters that should be ignored by processes that don't know specifically how to handle them (instead of, say, trying to draw them). This category also included a bunch of unassigned code points that are being set aside now to represent only characters that also should be ignored by processes that don't recognize them.
- **Deprecated.** Characters whose use is strongly discouraged. Because of Unicode's stability policies, this is as close as you can get to simply removing a character from the standard.
- **Soft_Dotted.** Characters whose glyphs include a dot on top that disappears when a top-joining diacritical mark is applied. These characters are all variants of the Latin lower-case i and j.
- **Logical_Order_Exception.** Characters that break Unicode's normal rule about representing things in logical order. This category currently consists of the Thai and Lao left-joining vowel marks, which, unlike all other Indic vowel marks, precede the consonants they modify instead of following them. For more on this, see the Thai and Lao sections of Chapter 9.

General character properties

Each character has a set of properties that serve to identify the character. These include the name, Unicode 1.0 name, Jamo short name, ISO 10646 comment, block and script.

Standard character names

First among these properties, of course, is the character's name, which is given both in the book and in the UnicodeData.txt file. The name is always in English, and the only legal characters for the name are the 26 Latin capital letters, the 10 Western digits, and the hyphen. The name is important, as it's the primary guide to just what character is meant by the code point. The names generally follow some conventions:

- For those characters that belong to a particular script (writing system), the script name is included in the character name.
- Combining marks usually have "COMBINING" in their names. Within certain scripts, combining marks' names may contain the word "MARK" instead. Vowel signs have the phrase "VOWEL SIGN" in their names.
- Letters have names containing the word "LETTER." Syllables have names that usually include either the word "LETTER" or the word "SYLLABLE". Ideographs have names that contain the word "IDEOGRAPH."
- Uppercase letters have the word "CAPITAL" in their names. Lowercase letters have the word "SMALL" in their names.
- Digits have the word "DIGIT" in their names.

Unicode 2.0 was the first version of Unicode that was unified with the ISO 10646 Universal Character Set standard. As part of the unification process, many characters got new names. For instance, U+0028 LEFT PARENTHESIS was called OPENING PARENTHESIS in Unicode 1.0. Sometimes the new names are more European-sounding than the old ones: For example, U+U+002E PERIOD became U+002E FULL STOP. For those characters that changed names in Unicode 2.0, their original names are shown in the “Unicode 1.0 name” field in UnicodeData.txt. If this field is blank, the character either didn’t exist in Unicode 1.0, or it had the same name then that it does now.

The renaming that happened as part of the unification with ISO 10646 was a very unusual occurrence. The names are now carved in stone. No character’s name will ever be changed in future versions of the standard. This occasionally means that a character’s semantics will drift out of sync with its name, but this is very rare.

Algorithmically-derived names

You’ll notice some entries in UnicodeData.txt that look like the following:

```
4E00;<CJK Ideograph, First>;Lo;0;L;;;;;N;;;;;
9FA5;<CJK Ideograph, Last>;Lo;0;L;;;;;N;;;;;
```

The “name” is enclosed in angle brackets. These entries always occur in pairs. The first entry in the pair has a “name” that ends in “First,” and the second entry has a “name” ending in “Last.” These pairs of entries mark the beginnings and endings of ranges of characters that have exactly the same properties. Rather than have nearly-identical entries for every code point from U+4E00 to U+9FA5, for example, the database abbreviates the range by using a pair of entries such as these. The properties shown for these entries apply to every character in the range. The one exception is the name property, which is instead algorithmically generated according to rules in the Unicode standard.

There are ten elided ranges of characters in the Unicode 3.0.1 database:

- The CJK Unified Ideographs Area, running from U+4E00 to U+9FA5.
- The CJK Unified Ideographs Extension A area, running from U+3400 to U+4DB5.
- The Hangul Syllables area, running from U+AC00 to U+D7A3.
- The high-surrogate values for non-private-use characters (i.e., Planes 1 through 14), which run from U+U800 to U+DB7F.
- The high-surrogate values for private-user characters (i.e., Planes 15 and 16), which run from U+DB80 to U+DBFF.
- The low-surrogate values, which run from U+U+DC00 to U+DFF.
- The Private Use Area, running from U+E000 to U+F8FF.
- The CJK Unified Ideographs Extension B area, running from U+20000 to U+2A6D6.
- The Plane 15 Private Use Area, running from U+F0000 to U+FFFFD.
- The Plane 16 Private Use Area, running from U+100000 to U+10FFFD.

The code points in some of the elided ranges don’t have character names (or any other properties, really—their code point values are assigned to a special category, but that’s it). The private-use characters just have a “private use” property. The surrogate code units aren’t really characters at all, but the code point values corresponding to them are put in one of a few special “surrogate” categories for the benefit of UTF-16-based applications that want to treat surrogate pairs as combining character sequences.

General category

The other elided ranges have algorithmically-derived character names. For the three CJK Unified Ideographs ranges, the algorithm for deriving the name of an individual character is simple: it's "CJK UNIFIED IDEOGRAPH-", plus the code point value in hex. For example, the name of U+7E53 is "CJK UNIFIED IDEOGRAPH-7E53". Simple.

The characters in the final elided range, the Hangul syllables, have a more complicated algorithm for deriving their names. It takes advantage of another property of Unicode characters, the "Jamo short name," which is given in the Jamo.txt file.

Most characters' Jamo short names is blank because they're not Hangul jamo (the alphabetic units that combine together to make Hangul syllables). But the characters in the Hangul Jamo block use this property. For them, this property specifies an abbreviated name for that character.

You get the name of a Hangul syllable by decomposing it into jamo and then concatenating the short names together. For example, consider the character U+B8B4 (꺠). This character decomposes into three characters, as follows:

```
U+1105 HANGUL CHOSEONG RIEUL ( ㄹ )
U+116C HANGUL JUNGSEONG OE ( ㅓ )
U+11AB HANGUL JONSEONG NIEUN ( ㄴ )
```

U+1105 has the short name R, U+116C has the short name OE, and U+11AB has the short name N. You concatenate these together and add "HANGUL SYLLABLE" to get the name of the syllable, so the name of the character U+B8B4 (꺠) is HANGUL SYLLABLE ROEN.

Control-character names

The code points corresponding to the C0 and C1 ranges from ISO 2022 (U+0000 to U+001F and U+0080 to U+009F) also get special treatment. They're put in a special "control character" category and aren't given names (the UnicodeData.txt file gives "<control>" as their names. Officially, Unicode doesn't assign semantics to these code points, but preserves them for backward compatibility with the various ISO 2022-compatible standards.

Unofficially, these code points are usually given the semantics they're given in ISO 6429, the same as most one-byte character encodings do. This is reflected in their Unicode 1.0 names, which follow that convention. Over time, various pieces of the Unicode standard have started officially assigning specific semantics to these code points—the Unicode line breaking algorithm treats U+000D and U+000A as carriage return and line feed characters. The Unicode bidirectional layout algorithm treats U+0009 as the tab character, and so on. Despite their lack of official names and their assignment to the special "control character" category, these characters have, over time, picked up other properties that show their true meanings.

ISO 10646 comment

In a few cases, there are comments attached to a character in the ISO 10646 code charts. These are usually alternate names for the character, similar to the informative notes in the character listings in the Unicode standard. For example, U+1095 LATIN SMALL LETTER HV has an ISO 10646 comment of "hwair," the common name for this letter. These comments from ISO 10646 are preserved in the "ISO 10646 comment" field in UnicodeData.txt.

Block and Script

The 21-bit Unicode encoding space is divided up into “blocks,” ranges of code point values with similar characteristics. For example, the range of code point values from U+0000 to U+007F is the “C0 Controls and Basic Latin” block, U+0400 to U+04FF is the Cyrillic block, U+20A0 to U+20CF is the Currency Symbols block, and so on. The official names and boundaries of all the blocks in the Unicode standard is given in the `Blocks.txt` file.

Implementations have, in the past, attempted to identify the more important characteristic—which writing system a character belongs to—by finding out what block it’s in. This is dangerous, because sometimes not everything in a block matches the name of the block. There are currency symbols outside the Currency Symbols block, for example. The Basic Latin block doesn’t just include Latin letters, but a bunch of other characters (digits, punctuation marks, etc.). The Arabic Presentation Forms A block contains two characters (the ornate parentheses) that are “real” characters and not just presentation forms.

What people are usually after when they’re looking at what block a character is in is what writing system, or “script,” the character belongs to. For example, is this code point a Latin letter, a Greek letter, a Chinese ideograph, a Hangul syllable, or what?

This information is given in the `Scripts.txt` file, which assigns every character to a script. Instead of treating the entire Basic Latin (i.e., ASCII) block as belonging to the Latin script, for example, only the letter A through Z and a through z are assigned to the Latin script.

Only letters and diacritical marks are assigned to scripts. The combining and enclosing marks are given a special “script” property—“INHERITED,” which indicates that these code points take on the script of the characters they’re attached to.

`Scripts.txt` doesn’t include every character in Unicode. Symbols, punctuation marks, and digits that are used with more than one writing system aren’t assigned to a script—their “script” property is effectively blank.

General Category

After the code point value and the name, the next most important property a Unicode character has is its general category. There are seven primary categories: letter, number, punctuation, symbol, mark, separator, and miscellaneous. Each of these subdivides into additional categories.

Letters

The Unicode standard uses the term “letter” rather loosely in assigning things to the general category of “letter.” Whatever counts as the basic unit of meaning in a particular writing system, whether it represents a phoneme, a syllable, or a whole word or idea, is given the “letter” category. The one big exception to this rule are marks that combine typographically with other characters, which are categorized as “marks” instead of “letters.” This includes not only diacritical marks and tone marks, but also vowel signs in those consonantal writing systems where the vowels are written as marks applied to the consonants.

General category

Some writing systems, such as the Latin, Greek, and Cyrillic alphabets, also have the concept of “case,” two series of letterforms used together, with one series, the “upper case,” used for the first letter of a sentence or a proper name, or for emphasis, and the other series, the “lower case,” used for most other letters.

Uppercase letter (abbreviated “Lu” in the database). In cased writing systems, the uppercase letters are given this category.

Lowercase letter (“Ll”). In cased writing systems, the lowercase letters are given this category.

Titlecase letter (“Lt”). Titlecase is reserved for a few special characters in Unicode. These characters are basically all examples of compatibility characters— characters that were included for round-trip compatibility with some other standard. Every titlecase letter is actually a glyph representing two letters, the first of which is uppercase and the second of which is lowercase.

For example, the Serbian letter nje (Њ) can be thought of as a ligature of the Cyrillic letter n (Н) and the Cyrillic soft sign (Ь). When Serbian is written using the Latin alphabet (Croatian, which is almost the same language, is), this letter is written using the letter nj. Existing Serbian and Croatian standards were designed so that there was a one-to-one mapping between every Cyrillic character used in Serbian and the corresponding Latin character used in Croatian. This required using a single character code to represent the nj digraph in Croatian, and Unicode carries that character forward. Capital Nje in Cyrillic (Њ) thus can convert to either NJ or Nj in Latin depending on the context. The fully-uppercase form, NJ, is U+01CA LATIN CAPITAL LETTER NJ, and the combined upper-lower form, U+01CB LATIN CAPITAL LETTER N WITH SMALL LETTER J, is considered a “titlecase” letter.

There are three Serbian characters that have a titlecase Latin form: Љ (lje, which converts to lj), Њ (nje, which converts to nj), and Ћ (dzhe, which converts to dž).

These were the only three titlecase letters in Unicode 2.x. Unicode 3.0 added a bunch of Greek letters that also fall into this category. Some early Greek texts represented certain diphthongs by writing a small letter iota underneath the other vowel rather than after it. For example, you’d see “ai” written as αι. If you just capitalized the alpha (“Ai”), you’d get the titlecase version: Αι. In the fully-uppercase version (“AI”), the small iota becomes a regular iota again: ΑΙ. These characters are all in the Extended Greek section of the standard and are only used in writing ancient Greek texts. In modern Greek, these diphthongs are just written using a regular iota: “ai” is just written as αι.

Modifier letter (“Lm”). Just as we have things you might conceptually think of as “letters” (vowel signs in various languages) classified as “marks” in Unicode, we also have the opposite. The modifier letters are independent forms that don’t combine typographically with the characters around them, which is why Unicode doesn’t classify them as “marks” (Unicode marks, by definition, combine typographically with their neighbors). But instead of carrying their own sounds, the modifier letters generally modify the sounds of their neighbors. In other words, conceptually they’re diacritical marks. Since they occur in the middle of words, most text-analysis processes treat them as letters, so they’re classified as letters.

The Unicode modifier letters are generally either International Phonetic Alphabet characters or characters that are used to transliterate certain “real” letters in non-Latin writing systems that don’t seem to correspond to a regular Latin letter. For example, U+02BC MODIFIER LETTER APOSTROPHE is usually used to represent the glottal stop, the sound made by (or, more accurately,

the *absence* of sound represented by) the Arabic letter alef, and so the Arabic letter is often transliterated as this character. Likewise, U+02B2 MODIFIER LETTER SMALL J is used to represent palatalization, and thus is sometimes used in transliteration as the counterpart of the Cyrillic soft sign.

Other letter (“Lo”). This is a catch-all category for everything that’s conceptually a “letter,” but which doesn’t fit into one of the other “letter” categories. Letters from uncased alphabets like the Arabic and Hebrew alphabets fall into this category, but so do syllables from syllabic writing systems like Kana and Hangul, and so do the Han ideographs.

Marks

Like letters, marks are part of words and carry linguistic information. Unlike letters, marks combine typographically with other characters. For example, U+0308 COMBINING DIAERESIS may look like this: “̈” when shown alone, but is usually drawn on top of the letter that precedes it: U+0061 LATIN SMALL LETTER A followed by U+0308 COMBINING DIAERESIS isn’t drawn as “ä”, but as “ä”. All of the Unicode combining marks do this kind of thing.

Non-spacing mark (“Mn”). Most of the Unicode combining marks fall into this category. Non-spacing marks don’t take up any horizontal space along a line of text—they combine completely with the character that precedes them and fit entirely into that character’s space. The various diacritical marks used in European languages, such as the acute and grave accents, the circumflex, the diaeresis, and the cedilla, fall into this category.

Combining spacing mark (“Mc”). Spacing combining marks interact typographically with their neighbors, but still take up horizontal space along a line of text. All of these characters are vowel signs or other diacritical marks in the various Indian and Southeast Asian writing systems. For example, U+093F DEVANAGARI VOWEL SIGN I (•) is a spacing combining mark: U+0915 DEVANAGARI LETTER KA followed by U+093F DEVANAGARI VOWEL SIGN I is drawn as ••—the vowel sign attaches to the left-hand side of the consonant.

Not all spacing combining marks reorder, however: U+0940 DEVANAGARI VOWEL SIGN II (•) is also a combining spacing mark. When it follows U+0915 DEVANAGARI LETTER KA, you get ••—the vowel attaches to the right-hand side of the consonant, but they still combine typographically.

Enclosing mark (“Me”). Enclosing marks completely surround the characters they modify. For example, U+20DD COMBINING ENCLOSING CIRCLE is drawn as a ring around the character that precedes it. There are only ten of these characters: they’re generally used to create symbols.

Numbers

The Unicode characters that represent numeric quantities are given the “number” property (technically, this should be called the “numeral” property, but that’s life). The characters in these categories have additional properties that govern their interpretation as numerals. All of the “number” subcategories are normative. This category subdivides as follows:

General category

Decimal-digit number (“Nd”). The characters in this category can be used as decimal digits. This category includes not only the digits we’re all familiar with (“0123456789”), but similar sets of digits used with other writing systems, such as the Thai digits (“๐๑๒๓๔๕๖๗๘๙”).

Letter number (“NI”). The characters in this category can be either letters or numerals. Many of these characters are compatibility composites whose decompositions consist of letters. The Roman numerals and the Hangzhou numerals are the only characters in this category.

Other number (“No”). All of the characters that belong in the “number” category, but not in one of the other subcategories, fall into this one. This category includes various numeric presentation forms, such as superscripts, subscripts, and circled numbers; various fractions; and numerals used in various numeration systems other than the Arabic positional notation used in the West.

Punctuation

This category attempts to make sense of the various punctuation characters in Unicode. It breaks down as follows:

Opening punctuation (“Ps”). For punctuation marks, such as parentheses and brackets, that occur in opening-closing pairs, the “opening” characters in these pairs are assigned to this category.

Closing punctuation (“Pe”). For punctuation marks, such as parentheses and brackets, that occur in opening-closing pairs, the “closing” characters in these pairs are assigned to this category.

Initial-quote punctuation (“Pi”). Quotation marks occur in opening-closing pairs, just like parentheses do. The problem is that which is which depends on the language. For example, both French and Russian use quotation marks that look like this: «» But they use them differently.

«In French, a quotation is set off like this.»

»But in Russian, a quotation is set off like this.«

This category is equivalent to either Ps or Pe, depending on the language.

Final-quote punctuation (“Pf”). This is the counterpart to the Pi category, and is also used with quotation marks whose usage varies depending on language. It’s equivalent to either Ps or Pe depending on language. It’s always the opposite of Pi.

Dash punctuation (“Pd”). This category is self-explanatory: It encompasses all hyphens and dashes.

Connector punctuation (“Pc”). Characters in this category, such as the middle dot and the underscore, get treated as part of the word they’re in. That is, they “connect” series of letters together into single words. `This_is_all_one_word`. An important example is U+30FB KATAKANA MIDDLE DOT, which is used like a hyphen in Japanese.

Other punctuation (“Po”). Punctuation marks that don’t fit into any of the other subcategories, including obvious things like the period, comma, and question mark, fall into this category.

Symbols

This group of categories contains various symbols.

Currency symbol (“Sc”). Self-explanatory.

Mathematical symbol (“Sm”). Mathematical operators.

Modifier symbol (“Sk”). This category contains two main things: the “spacing” versions of the combining marks and a few other symbols whose purpose is to modify the meaning of the preceding character in some way. Unlike modifier letters, modifier symbols don’t necessarily modify the meanings of letters, and don’t necessarily get counted as parts of words.

Other symbol (“So”). Again, this category contains all symbols that didn’t fit into one of the other categories.

Separators

These are characters that mark the boundaries between units of text.

Space separator (“Zs”). This category includes all of the space characters (yes, there’s more than one space character).

Paragraph separator (“Zp”). There is exactly one character in this category: the Unicode paragraph separator (U+2029), which, as its name suggests, marks the boundary between paragraphs.

Line separator (“Zl”). There’s also only one character in this category: the Unicode line separator, (U+2028) which, as its name suggests, forces a line break without ending a paragraph.

Even though the ASCII carriage-return and line-feed characters are often used as line and paragraph separators, they’re not put in either of these categories. Likewise, the ASCII tab character isn’t considered a Unicode space character, even though it probably should. They’re all put in the “Cc” category.

Miscellaneous

There are a number of special character categories that don’t really fit in with the others. They are as follows:

Control characters (“Cc”). The codes corresponding to the C0 and C1 control characters from the ISO 2022 standard are given this category. The Unicode standard doesn’t officially assign any semantics to these characters (which include the ASCII control characters), but most systems that use Unicode text treat these characters the same as their counterparts in the source standards. For example, most processes treat the ASCII line-feed character as a line or paragraph separator.

General category

The original idea was to leave the definitions of these code points open, as ISO 2022 does, but over time, various Unicode processes and algorithms have attached semantics to these code points, effectively nailing the ISO 6429 semantics to many of them.

Formatting characters (“Cf”). Unicode includes some “control” characters of its own: characters with no visual representation of their own which are used to control how the characters around them are drawn or handled by various processes. These characters are assigned to this category.

Surrogates (“Cs”). The code points in the UTF-16 surrogate range belong to this category. Technically, the code points in the surrogate range are treated as unassigned and reserved, but Unicode implementations based on UTF-16 often treat them as characters, handling surrogate pairs the same way combining character sequences are handled.

Private-use characters (“Co”). The code points in the private-use ranges are assigned to this category.

Unassigned code points (“Cn”). All unassigned and non-character code points, other than those in the surrogate range, are given this category (these code points aren’t listed in the Unicode Character Database—their omission gives them this category—but are listed explicitly in `DerivedGeneralCategory.txt`).

Other categories

Over time, it’s become necessary to draw finer distinctions between characters than the general categories let you do, and it’s also been noticed that there’s some overlap between the general categories. Another set of categories, defined mostly in `PropList.txt`, has been created to capture these distinctions. Among the other categories:

- **Whitespace.** A lot of processes that operate on text treat various characters as “whitespace,” important only insofar as it separates groups of other characters from one another. In Unicode, “whitespace” can be thought of mainly as consisting of the characters in the “Z” (separator) categories. But this is one of those cases of the ISO control characters having meaning—most processes want the code points corresponding to the old ASCII and ISO 8859 whitespace characters (TAB, CR, LF, FF, and NEL) treated as whitespace as well. This property encompasses all of these characters.
- **Bi-di control and join control.** These categories isolate out a few of the characters in the format-control (“Cf”) category that control the behavior on certain specific text processes (in this case, the bi-di algorithm and the glyph selection process).
- **Dash and hyphen.** These categories not only break down the dash-punctuation (“Pd”) category into dashes and hyphens, but also encompass a few characters from the connector-punctuation (“Pc”) category that get used as hyphens (specifically, the Katakana middle dot we looked at earlier).
- **Quotation marks** brings together all the characters from the Ps, Pe, Pi, Pf, and Po categories that represent quotation marks.
- **Terminal punctuation** isolates out those characters in the other-punctuation (“Po”) category that represent sentence- or clause-ending punctuation.
- **Math** brings together all the characters that represent mathematical symbols. This includes not only those characters assigned to the math-symbol (“Sm”) category, but characters from the other

categories that have other uses beyond their use as math symbols. The non-Sm characters with the Math property are assigned to the “Other_Math” property in PropList.txt.

- **Hex digit** and **ASCII hex digit** single out the characters that can be used as digits in hexadecimal numerals.
- **Alphabetic** singles out alphabetic and syllabic characters. It encompasses everything in the “letter” (“L”) categories other than the characters with the “ideographic” property, plus all the combining marks and vowel signs used with them. PropList.txt lists all the characters with the “alphabetic” property that aren’t in one of the L categories and assigns them to the “Other_Alphabetic” category.
- **Ideographic** isolates out those characters in the “other letter” (“Lo”) category that represent whole words or ideas (basically, the Chinese characters).
- **Diacritic** brings together those characters from the modifier-letter (“Lm”), modifier-symbol (“Sk”), non-spacing-mark (“Mn”), and combining spacing-mark (“Mc”) categories that are used as diacritical marks.
- **Extender** isolates those characters from the modifier-letter (“Lm”) category (plus one character—the middle dot—from the other-punctuation (“Po”) category) that elongate, either phonetically or graphically, the letter that precedes them.
- **Uppercase** and **Lowercase** bring together all the upper- and lower-case letters. PropList.txt supplements the characters assigned to these general categories (“Lu” and “Ll”) with characters from other categories whose forms are the forms of upper- or lower-case letters (it assigns these extra characters to the “Other_Uppercase” and “Other_Lowercase” categories).
- **Noncharacter code points** isolates out the code point values from the “unassigned” (“Cn”) category that are illegal code point values in Unicode text (as opposed to merely not having a character assigned to them yet).

Again, Unicode 3.2 adds a few new properties to the mix. At the time Unicode 3.2 entered beta, these included:

- **Default-ignorable code points** are code points or characters that should simply be ignored by processes that don’t know about them. For example, the default behavior of a rendering process, when confronted with a code point it doesn’t recognize, is to draw some kind of “missing” glyph. If the character has this property, it should instead just not draw anything. This category basically consists of characters in the “formatting character” (“Cf”) or “control character” (“Cc”) general categories, plus a bunch of characters called out specifically in PropList.txt as “Other_Default_Ignorable_Code_Point.” This is one of the few properties that is actually given to *unassigned* code points; this is so that current implementations can “do the right thing” even with future default-ignorable characters that may be added to the standard.
- **Soft-dotted** calls out characters whose glyphs include a dot that disappears when certain diacritical marks are applied. This basically means the lowercase Latin letters i and j: If you apply an acute accent to the letter i, for example, the dot disappears: í. This property exists to make case-mapping algorithms easier.
- **Logical-order exceptions** are the few characters that break Unicode’s logical-order rule: This basically means the left-joining Thai and Lao vowel marks, which precede in memory the consonants they attach to, rather than following them, as every other Indic vowel mark does. This property exists to make life easier for the Unicode Collation Algorithm, some rendering algorithms, and possibly other processes that care about backing-store order. For more on Thai and Lao, see Chapter 9.
- **Deprecated** characters are characters that Unicode’s designers would remove from the standard if they could. They’re retained for backward compatibility, but their use and interpretation is strongly discouraged.

- **Radical, CJK_Unified_Ideograph, IDS_Binary_Operator, and IDS_Tertiary_Operator** are used to help interpret ideograph description sequences. See Chapter 10 for more information on ideographic description sequences.

The `DerivedCoreProperties.txt` file explicitly lists the characters in the Math, Alphabetic, Uppercase, and Lowercase categories, along with the characters that are legal in various positions in a programmatic identifier according to the Unicode identifier guidelines (see Chapter 17).

Properties of letters

We looked at a bunch of letter-related properties in the preceding sections: First, there are the various subdivisions of the “letter” category (Lu, Ll, Lt, Lm, and Lo). Then there are the Alphabetic and Ideographic categories. Finally, there’s the case: defined by the Lu, Ll, and Lt general categories and also by the separate Uppercase and Lowercase categories.

For cased letters (i.e., characters in the Lu, Ll, and Lt) categories, Unicode provides additional information about how the characters in the different cases relate to each other. In particular, it provides information on how to map a character to its counterpart in the opposite case (or in titlecase), and a set of mappings that allow you to normalize away case differences when comparing strings.

Each character has uppercase, lowercase, and titlecase mappings in `UnicodeData.txt`. If the character should be replaced with another character when converted to uppercase, lowercase, or titlecase, the code point value for the character it should be replaced with appears in the appropriate place in `UnicodeData.txt`.

So the entry for U+0061 LATIN SMALL LETTER A (a) gives U+0041 (A) as its uppercase and titlecase mappings. (The titlecase mapping is the same as the uppercase mappings for all but a handful of characters.) The entry for U+0393 GREEK CAPITAL LETTER GAMMA (Γ) gives U+03B3 (γ) as its lowercase mapping.

Converting a string to uppercase or lowercase is simple—just replace every character in the string with its uppercase or lowercase mapping (if it has one). Converting to titlecase is almost as simple—replace the first letter of every word with its titlecase mapping.

SpecialCasing.txt

But the mappings aren’t always as simple as they appear at first glance. The mappings in the `UnicodeData.txt` file are always single-character mappings. Sometimes a character expands into two characters when its case is changed. Sometimes a mapping is conditional—the mapping is different depending on the context or language.

These special mappings are relegated to their own file—`SpecialCasing.txt`. Like `UnicodeData.txt`, `SpecialCasing.txt` is structured as a series of one-line entries for the characters, each consisting of a semicolon-delimited series of fields. The first field is the code point value the mappings are being specified for. The next three fields are the lowercase, titlecase, and uppercase mappings for the characters, each of which is a space-delimited series of code point values. There’s an optional fifth field—a comma-delimited series of tags indicating the context in which that particular set of

mappings happens. As always, the line may end with a comment, which is set off from the actual data with a pound sign (#).

Here's an example entry:

```
00DF; 00DF; 0053 0073; 0053 0053; # LATIN SMALL LETTER SHARP S
```

This is the entry for the German ess-zed character (ß), which generally functions like a double s (in classical German spelling, the ss form was used after short vowels and the ß form after long vowels, although the ss form is gradually being used for both, while the ß form is gradually falling into disuse). First we have the character itself, U+00DF. Next, we have its lowercase mapping. Since it's already a lowercase letter, its lowercase mapping is also 00DF. The titlecase mapping (which should never occur in practice, since ß never occurs at the beginning of a word) is Ss, U+0053 U+0073. The uppercase mapping is SS, U+0053 U+0053.

There are several tags used for context-sensitive mappings²⁹:

- `Final_Sigma` is used for the Greek letter sigma and indicates that the mapping is effective when the letter is the last letter in the word but not the *only* letter in the word.
- `More_Above` indicates the mapping is effective only when the character has one or more top-joining combining marks after it.
- `After_Soft_Dotted` indicates the mapping is effective only when the character follows a `Soft_Dotted` character (basically, an i or j).
- `Before_Dot` indicates the mapping is effective only when the character precedes U+0307 COMBINING DOT ABOVE.

In addition, the condition list may include one or more two-letter abbreviations for languages (the abbreviations come from the ISO 639 standard), which state that a particular mapping only happens if the text is in that language.

The last three of these conditions are used in Lithuanian-specific case mappings: in Lithuanian, the letter i keeps its dot when an accent is applied to it, an effect that can be achieved by adding U+0307 COMBINING DOT ABOVE. These conditions are used to control when this extra character should be added and removed.

Here's a classic example of the use of the special conditions:

```
03A3; 03C2; 03A3; 03A3; Final_Sigma; # GREEK CAPITAL LETTER SIGMA
# 03A3; 03C3; 03A3; 03A3; # GREEK CAPITAL LETTER SIGMA
```

These two entries are for the Greek capital letter sigma (Σ, U+03A3). This letter has two lower-case forms: one (ς, U+03C2) that appears only at the ends of words, and another (σ, U+03C3) that appears at the beginning or in the middle of a word. The `FINAL_SIGMA` at the end of the first entry indicates that this mapping only applies at the end of a word (and when the letter isn't also the only letter in the word). The second entry, which maps to the initial/medial form of the lowercase sigma, is commented out (begins with a #) because this mapping is in the `UnicodeData.txt` file.

²⁹ These are the tags in the Unicode 3.2 version of `SpecialCasing.txt`, which simplifies the approach used in the Unicode 3.1 version (but produces the same effect). It's possible that these tags will change before Unicode 3.2 is final.

And here's the other classic example:

```
0049; 0131; 0049; 0049; tr Not_Before_Dot; # LATIN CAPITAL LETTER I
0069; 0069; 0130; 0130; tr; # LATIN SMALL LETTER I
# 0131; 0131; 0049; 0049; tr; # LATIN SMALL LETTER DOTLESS I
# 0130; 0069; 0130; 0130; tr; # LATIN CAPITAL LETTER I WITH DOT ABOVE
```

The “tr” in the conditions column tells us that these are language-specific mappings for Turkish. In Turkish, there's an extra letter (ı, U+0131) that looks like a lower-case *i* with the dot missing. This letter's upper-case counterpart is what we'd normally think of as the capital I: (U+0049). The regular lowercase *i* (U+0069) keeps the dot when converted to uppercase so that you can tell the two letters apart (İ, U+0130). Again, the last two mappings in the example are commented out because they're covered in UnicodeData.txt.

CaseFolding.txt

CaseFolding.txt is used for case-insensitive string comparisons. It prescribes a language-independent mapping for each character that erases case distinctions. Generally, it maps each character to its lowercase equivalent, but it includes some extra mappings that map some lowercase letters to variant forms (for example, ß maps to “ss”). Generally, the mappings were derived by converting everything to uppercase using the mappings in the UnicodeData.txt and SpecialCasing.txt files and then back to lowercase using just the UnicodeData.txt file. Each entry consists of the original code point value, a code (“C,” “F,” “S,” or “I”) indicating the type of mapping, and the mapping (one or more code point values). The code tells you whether the mapping causes the text to get bigger (“F”) or not (“C”). In cases where the standard mapping causes the text to expand, a second mapping that doesn't change the text length is provided (and tagged with “S”). (The Turkish I mappings are called out with an “I” tag, since they can cause funny things to happen.)

For more information on case mapping, see Chapter 14. For more information on string comparisons, including case folding, see Chapter 15.

Properties of digits, numerals, and mathematical symbols

We looked at a bunch of number-related properties earlier in this chapter: First, there are the various numeric general categories (“Nd,” “NI,” and “No”). There are also the hex-digit properties and the Math property.

For characters in the “N” categories (characters that are used to represent numeric values), there are three additional properties that indicate the numeric value they represent. These properties are given in the UnicodeData.txt file.

Characters that can be used as decimal digits (generally, the characters in the “Nd” category, although some “No” characters also) have a **decimal digit value** property indicating the value the character has when used as a decimal digit. This is shown in the UnicodeData file as an ASCII decimal digit. For example, the Western digit “2”, the Arabic digit “٢”, and the Thai digit “๒” all represent the value 2 in numeration systems that use Arabic positional notation, and thus all have the decimal-digit-value property of 2. All characters with a decimal-digit-value property are identified as “decimal” in the DerivedNumericProperties.txt file.

Characters that can be used as digits in any numeration system, including decimal, have a **digit value** property indicating their numeric value when used as digits. This is shown in the UnicodeData file as a decimal numeral using ASCII digits. The decimal-digit characters all have a digit value that's the same as their decimal-digit value. The only characters that have a digit value but not a decimal-digit value are certain presentation forms of the decimal digits (circled, parenthesized, followed by a period, etc.). The idea here is that these characters represent whole numerals, not just digits, but are technically still decimal digits. The DerivedNumericProperties.txt file identifies these characters as "digit."

Finally, all other characters representing numeric values have a **numeric value** property that gives the numeric value they represent. This is given in the UnicodeData.txt file as a decimal numeral using ASCII digits. These characters also show up in DerivedNumericProperties.txt as "numeric." Characters in this category include digits from numeration systems that don't use Arabic positional notation, such as the Ethiopic 10 symbol (፩) or the Tamil 100 symbol (௮௮); presentation forms of numbers 10 and above (such as the circled 13); Roman numerals; and fractions. If the numeric value is a fraction, it's shown as a fraction in the UnicodeData file, using the ASCII slash as the fraction bar. For example, the character ½ has a numeric value of "1/2", as shown in the UnicodeData file.

The DerivedNumericValues.txt file lists all characters with the numeric property according to their numeric values.

Sometimes, characters in the "letter" categories are used as digits or numerals. For instance, the Latin letters are often used as digits in numerals with radices above 10 (such as the letters A through F being used to represent digit values from 11 to 15 in hexadecimal notation). The Latin letters are also used in Roman numerals. Ancient Greek and Hebrew texts used the letters in their respective alphabets to form numerals. And various Han ideographs are used to represent numeric values (and sometimes even used as decimal digits). All of these applications are considered specialized uses by Unicode, and the Unicode standard gives none of these characters numeric properties.

We'll look more closely at the various number characters in Unicode in Chapter 12.

Layout-related properties

The Unicode characters have a bunch of properties associated with them that specify how various processes related to text rendering (i.e., drawing text on the screen or on paper) should handle them.

Bidirectional layout

The Unicode standard includes a very detailed specification of how characters from right-to-left and left-to-right writing systems are to be arranged when intermixed when they appear together in a single line of text. Central to this specification are the characters bi-di categories, which are given in the UnicodeData.txt file using these codes:

L	Strong left-to-right characters.
R or AL	Strong right-to-left characters.
EN and AN	Digits, which have weak left-to-right directionality.
ES, ET, and CS	Punctuation marks used with numbers. Treated like digits when they occur

as part of a number and like punctuation the rest of the time.

WS and ON	Neutrals. Whitespace and punctuation, which take on the directionality of the surrounding text.
NSM	Non-spacing marks, which take on the directionality of the characters they attach to.
B and S	Block and segment separators, which break text up into units that are formatted independently of each other by the bi-di algorithm.
BN	Invisible formatting characters that are invisible to the bi-di algorithm.
LRE, RLE, LRO, RLO, and PDF	Special formatting characters that override the default behavior of the bi-di algorithm for the characters they delimit.

Bi-di category is officially specified in the UnicodeData.txt file, but is also listed in the DerivedBidiClass.txt file, which groups characters according to their bi-di category.

The characters in the LRE, RLE, LRO, RLO, and PDF categories—the invisible control characters that affect the bi-di algorithm’s treatment of other characters—have the “Bidi_Control” property in PropList.txt.

For a conceptual overview of the bi-di algorithm see chapter 8. For in-depth information on implementing the bi-di algorithm, including more detailed treatment of these categories, see Chapter 16.

Mirroring

Another important property related to the bi-di algorithm is the “mirrored” property. This is also specified in the UnicodeData.txt file as a “Y” or “N” in the “mirrored” column. (The DerivedBinaryProperties.txt file also lists all the characters with the “mirrored” property.)

A “Y” means that the character has the “mirrored” property. Characters with the “mirrored” property have a different glyph shape when they appear in right-to-left text than they have when they appear in left-to-right text. Usually the two glyphs are mirror images of each other, which is where the word “mirrored” comes from.

Consider the parentheses. The same code, U+0028, is used to represent the mark at the beginning of a parenthetical expression (the character is called “LEFT PARENTHESIS” rather than “STARTING PARENTHESIS”, its old name, an unfortunate side effect of the ISO 10646 merger) regardless of the direction of the surrounding text. The concave side always faces the parenthetical text (the text that comes after it). In left-to-right text that means it looks like this: (

In right-to-left text, it has to turn around, so it looks like this:)

Not all of the characters that mirror come in pairs. For those that do, the BidiMirroring.txt file gives the code point value of the character with the mirror-image glyph (for U+0028 LEFT PARENTHESIS, this, of course, would be U+0029 RIGHT PARENTHESIS, and vice versa). Some

of the mappings in `BidiMirroring.txt` are tagged as “[BEST FIT]”, meaning the match is only approximate. The file also lists all the characters with the “mirrored” property that can’t be mirrored simply by mapping them to another character.

Arabic contextual shaping

The letters of the Arabic and Syriac alphabets join together cursively, even in printed text. Because of this, they adopt a different shape depending on the surrounding characters. The `ArabicShaping.txt` file provides information that can be used to implement a simple glyph-selection algorithm for Arabic and Syriac.

Each line gives information for one character and consists of four fields separated with semicolons:

- The code point value.
- An abbreviated version of the character’s name.
- A code indicating the character’s joining category. R means the character can connect to the character on its right, D means the character can connect to characters on both sides, and U means the character doesn’t connect to characters on either side. L could conceivably be used for characters that only connect to neighbors on the left, but there are no characters in this category. C is also used for characters that cause surrounding characters to join to them but don’t change shape themselves. (There’s also a T category for characters that are transparent to the shaping algorithm, but they’re not listed in this file.)
- The name of the character’s joining group. A lot of Arabic letters share the same basic shape and are differentiated from one another by adding dots in various places. This column identifies the character’s basic shape (this is usually the name of one of the letters with that shape).

Characters not listed in `ArabicShaping.txt` either have a joining type of T (transparent, basically all the non-spacing marks and invisible formatting characters) or U (non-joining, basically everything else).

`DerivedJoiningType.txt` repeats information from `ArabicShaping.txt`, but organizes it by joining type rather than by code point value. `DerivedJoiningGroup.txt` does the same thing for joining groups.

It’s important to note that the information from `ArabicShaping.txt` is *informative*, not normative. It basically represents a minimal set of information necessary to produce legible Arabic. Implementations can, and often do, go to more trouble and produce better-looking text. For more information on Arabic and Syriac, see Chapter 8.

Two invisible formatting characters that affect the joining behavior of the surrounding text are given the `Join_Control` property in `PropList.txt`.

East Asian width

Just as Middle Eastern typography introduces interesting rendering challenges, so does East Asian typography, where lines of text often run vertically from the top of the page to the bottom.

Traditionally, Chinese characters (and, often, the characters used with them) are effectively monospaced: all the characters fit inside the same-size box (which is usually, but not always, square). When Western text is mixed with Asian text, decisions have to be made about how to lay out the mixture. Western letters can be blown up to fill an entire display cell or can be left as their normal

proportionally-spaced selves. Similarly, sometimes multiple Asian characters are placed together in a single cell to represent an abbreviation.

Originally, no special distinction between the two behaviors was enshrined in the Asian encoding standards, but one evolved: The JIS X 0201 and 0208 standards both include the Latin alphabet and the Katakana characters, so encoding schemes based on both, such as SHIFT-JIS, wound up encoding these characters twice: once as single-byte values and again as double-byte values (i.e., pairs of 8-bit code units). Over time, the single-byte codes came to represent half-width (or “hankaku”) glyphs, which are either proportionally spaced or sized to fill half a standard display cell, and the double-byte codes came to represent full-width (“zenkaku”) glyphs, which fill a whole display cell.

Some characters in Unicode are specifically categorized as half-width or full-width, but most are either implicitly one or the other or have ambiguous semantics. The `EastAsianWidth.txt` file nails down all these semantics by giving each character one of the following codes:

F	Explicitly full-width (has “FULLWIDTH” in its name)
H	Explicitly half-width (has “HALFWIDTH” in its name)
W	Implicitly full-width
Na	Implicitly half-width
A	Ambiguous. Generally half-width, but full-width in an East Asian context.
N	Neutral. Never occurs in East Asian typography. Treated same as Na.

This information is also in `DerivedEastAsianWidth.txt`, which lists everything in order by width category rather than by code point value. For more information on the East Asian Width property, see Chapter 10.

Line breaking property

Another integral part of the process of rendering text is breaking up long passages of text into lines. The Unicode standard also specifies how a line breaking process should treat the various characters. This information is in `LineBreak.txt`, which includes both normative and informative line-breaking categories. The normative categories are as follows:

BK	Forces a line break.
CR	Carriage return. Same as BK , except when followed by LF .
LF	Line feed. Same as BK , but causes CR to behave differently. (The CR-LF combination we all know and love from DOS thus gets treated as a single line terminator, not two.)
CM	Combining mark. Never a break before; otherwise transparent.
SG	High surrogate. No break before a low surrogate.
SP	Space. Generally, a break is allowed after a series of spaces.
ZW	Zero-width space. Generally, the same as SP .
GL	Non-breaking (“glue”). Never a break before or after.

CB Contingent break opportunity. Breaks depend on info outside the text stream.

The informative properties give a good default line-break behavior, but implementations are free to do a better (or more language-specific) job if they want to go to more trouble.

As with the other properties, there's a `DerivedLineBreak.txt` file that gives the same information as `LineBreak.txt`, but organized by category rather than code point value.

For more information on line breaking, including information on the other categories, see Chapter 16.

Normalization-related properties

We already looked at the Unicode normalization forms and at basically how Unicode normalization works. A lot of the Unicode Character Database is given over to this important topic, so we'll take a closer look at the normalization-related properties here. For more information on decomposition and normalization refer back to Chapter 4. For in-depth information on implementing Unicode normalization, see Chapter 14.

Decomposition

As we saw before, many Unicode characters are said to “decompose”: that is, they're considered equivalent to (and, generally speaking, less preferable than) other Unicode characters or sequences of other Unicode characters. For those characters that decompose, the `UnicodeData.txt` file gives their decomposition.

For canonical composites, the decomposition is either a single Unicode code point value (a singleton decomposition) or two Unicode code point values, the first of which may itself have a canonical decomposition (this is this way to save space and to help implementations that also want to save space). To get the full canonical decomposition for a character thus means performing the decomposition recursively until you arrive at a sequence of non-decomposing characters.

Compatibility decompositions, on the other hand, can be any number of characters long. This is because the intermediate forms necessary to have every decomposition be two characters long, don't always exist.

The same field in `UnicodeData.txt` is used for both canonical and compatibility decompositions, since a single character can only have one or the other, not both. The decomposition is preceded with a tag to specify what kind of decomposition it is.

Decomposition type

The decomposition in `UnicodeData.txt` for each character is optionally preceded by a tag that specified the type of the decomposition. If the tag is missing, it's a canonical decomposition. If the tag is present, it's a compatibility decomposition and the tag is there to give you an idea of what type of information would be lost if you converted the character to its compatibility decomposition. The tags are as follows:

**** identifies a compatibility composite that's equal to its compatibility decomposition plus some additional font or styling information (for example, a script or black-letter variant of some letter used as a symbol).

<noBreak> identifies a variant of a space or hyphen that prevents line breaks. A no-break variant is the same as its compatibility decomposition with U+2060 WORD JOINER both before and after it. (U+2060 isn't included in the decomposition, however.)

<initial>, **<medial>**, **<final>**, and **<isolated>** are used with the presentation forms of the Arabic letters to identify which particular glyph for the letter the presentation form is intended to represent.

<circle> identifies characters that are equivalent to their compatibility decomposition with a circle drawn around it. The same effect can be achieved (for single characters) by following the character with U+20DD COMBINING ENCLOSING CIRCLE, but this character isn't included in the decomposition.

<super> and **<sub>** mark characters that are the same as their compatibility decompositions drawn as superscripts or subscripts.

<vertical> is used to identify presentation forms for various punctuation marks used with East Asian text. These presentation forms are equivalent to the glyphs that are to be used for their compatibility mappings when used in vertical text (this generally means they're roughly equivalent to the horizontal forms rotated ninety degrees).

<wide> marks fullwidth (zenkaku) variants of various Latin characters used in East Asian text.

<narrow> identifies halfwidth (hankaku) variants of various Japanese characters used in abbreviations.

<small> identifies small variants of various ASCII punctuation marks and symbols sometimes used in Chinese text.

<square> indicates that the character is equivalent to the characters in its compatibility decomposition arranged in a single square East Asian display cell. It's used for various abbreviations used in East Asian text.

<fraction> identifies the "vulgar fractions," single code points representing fractions, which can all be represented using U+2044 FRACTION SLASH (which *is* included in these characters' decompositions).

<compat> precedes all compatibility decompositions that don't fit into one of the other categories.

The decomposition types are informative and somewhat approximate. There is a `DerivedDecompositionType.txt` file that organizes the Unicode characters according to their decomposition type (if they decompose).

Combining class

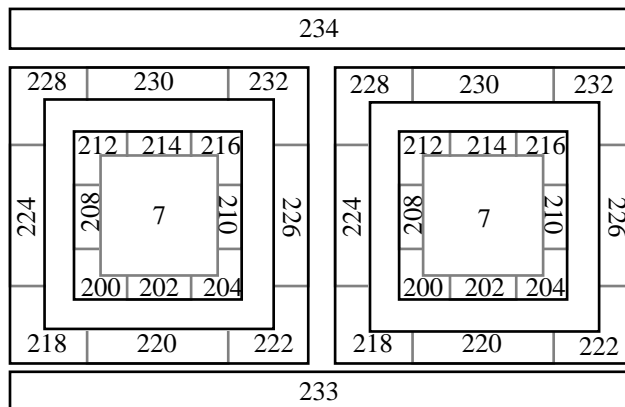
The combining class is a number that says how a combining character interacts typographically with the character it attaches to. This is used primarily to arrange a series of combining characters into canonical order.

All non-combining characters have a combining class of 0. A fair number of combining marks also have a combining class of 0—this indicates that the mark combines typographically with its base character in a way that should prevent reordering of any marks attached to the same character. For example, the enclosing marks have a combining class of 0. Any combining marks that follow the enclosing mark shouldn't be reordered to appear before the enclosing mark in storage, since whether they appear before or after the enclosing mark in storage controls whether or not they're drawn inside the enclosing mark.

Here's a quick breakdown of the combining classes:

- 0 Non-combining characters, or characters that surround or become part of the base character
- 1 Marks that attach to the interior of their base characters, or overlay them
- 7-199 Various marks with language- or script-specific joining behavior
- 200-216 Marks that attach to their base characters in various positions
- 218-232 Marks that are drawn in a specific spot relative to their base characters but don't touch them
- 233-234 Marks that appear above or below two base characters in a row
- 240 The iota subscript used in ancient Greek

For the more general combining classes, they are arranged like this:



Again, there's a `DerivedCombiningClass.txt` file that organizes the characters by combining class rather than by code point value.

Composition exclusion list

This CompositionExclusions.txt file contains the composition exclusion list for Unicode Normalized Form C. Form C favors composed forms over decomposed forms for characters that can be represented either way. This file lists canonical composites that can't appear in text normalized for Form C.

Why would you exclude certain canonical composites from Normalized Form C? There are four main groups of excluded characters:

Script-specific exclusions. For some writing systems, the more common way of representing the character is with the decomposed form, and converting to the composed form doesn't actually help anyone.

Singleton decompositions. Some characters have a single-character canonical decomposition. In these cases, the "composite" form is just an equivalent character that was given a separate code point value for compatibility. The "decomposed" form is the form that should almost always be used.

Non-starter decompositions. The composite character has a decomposition whose first character doesn't have a combining class of 0—that is, it's a composite combining character. What we really want to do with these characters is have them be incorporated into a composite form that includes the base character. A composite combining mark is still a combining mark, and since the idea behind Form C (in general) is to get rid of combining marks when possible, composing things into composite combining marks doesn't help.

Post-composition version characters. The most important category consists of canonical composites added to Unicode after Unicode 3.0. Normalized Form C is defined only to use canonical composites from Unicode 3.0. This is for backward compatibility. If you didn't exclude new canonical composites, you could take a file in Unicode 3.0 (or an earlier version), convert it to Normalized Form C, and change the file's Unicode version. If the file was produced by a program that only understood Unicode 3.0, converting it to Form C could make the file (or at least certain characters in the file) unreadable to the program that produced it. To get around this problem, Normalized Form C can't convert to canonical composites that are added to future versions of Unicode. The CompositionExclusions file refers to these characters as "post-composition version characters."

There are thirteen post-composition version characters in the Unicode 3.1 version of CompositionExclusions.txt. These are all precomposed versions of the musical symbols in the new Musical Symbols block. (All of the characters in these characters' decompositions are also new for Unicode 3.1, so it probably would have been safe not to put them in the composition exclusion list, but it's important to uphold the rules anyway.)

Normalization test file

NormalizationTest.txt is a fairly exhaustive set of test cases for implementations of Unicode normalization. It gives a bunch of sequences of Unicode characters, along with what you should get when you convert each of them to each of the Unicode normalized forms. It can be used (in fact, it *should* be used) to test normalization implementations for conformance.

We look at this file in a lot more detail in Chapter 14.

Derived normalization properties

Finally, there's the `DerivedNormalizationProperties.txt` file, which contains a bunch of categories that can be used to optimize normalization implementations:

- The **Full_Composition_Exclusion** category (or “Comp_Ex” in versions of Unicode before 3.2) includes characters that can't occur in any of the normalized Unicode forms.
- The **NFD_NO** category lists characters that can't occur in Normalized Form D. If you're normalizing to Form D, you can leave characters alone until you hit a `Comp_Ex` or `NFD_NO` character.
- The **NFC_NO** category lists characters that can't occur in Normalized Form C. If you're normalizing to Form C, you can leave characters alone until you hit a `Comp_Ex` or `NFC_NO` character.
- The **NFC_MAYBE** category lists characters that usually don't occur in Normalized Form C. If you're normalizing to Form C and hit an `NFC_MAYBE` character, you may have to check context to determine if you have to do anything.
- The **NFKD_NO** category lists characters that can't occur in Normalized Form KD.
- The **NFKC_NO** and **NFKC_MAYBE** categories are analogous to `NFC_NO` and `NFC_MAYBE` for Normalized Form KC.
- The **Expands_On_NFD** category lists characters that turn into more than one character when converted to Normalized Form D.
- The **Expands_On_NFC** category lists characters that turn into more than one character when converted to Normalized Form C.
- The **Expands_On_NFKD** category lists characters that turn into more than one character when converted to Normalized Form KD.
- The **Expands_On_NFKC** category lists characters that turn into more than one character when converted to Normalized Form KC.
- The **FNC** category lists characters where you get a different result if you convert them to Normalized Form KC than you get if you first case-fold them (using the mappings in `CaseFolding.txt`) and *then* map them to Normalized Form KC. For each character, the mapping straight from the original character to the case-folded Normalization Form KC version is also included, allowing you to fast-path these two conversions.

All of these categories can be (and, in fact, were) derived from the actual normalizations in `UnicodeData.txt`. You can use them in a normalization implementation to help avoid doing unnecessary work.

Grapheme-cluster-related properties

Unicode 3.2 introduces a few new properties intended to nail down how text is parsed into grapheme clusters:

- **Grapheme_Base** includes all the characters that aren't in one of the other categories. `Grapheme_Base` characters are grapheme clusters unto themselves unless connected to other characters with `Grapheme_Extend` or `Grapheme_Link` characters.
- **Grapheme_Extend** includes all combining marks and a few other characters that go in the same grapheme cluster as the character that precedes them.
- **Grapheme_Link** includes Indic viramas and the combining grapheme joiner, which cause the characters that precede and follow them to belong to the same grapheme cluster.

Unihan.txt

Finally, there's the Unihan.txt file. One of the most important scripts in Unicode is the Chinese characters (the "Han characters" or "CJK Ideographs"). Unicode 3.1 includes over 70,000 Han characters, and these characters have additional properties beyond those assigned to the other Unicode characters.

Chief among these properties are mappings to various source standards. This is the way Unicode defines the meanings of the various Han characters: it specifies exactly where they came from. This also lets you see just which characters from which source standards got unified together in Unicode. All of these mappings, plus a lot of other useful data, are in Unihan.txt.

For each Han character, the Unihan.txt gives at least some of the following pieces of information:

- Mappings to codes for the same character in a huge number of national and vendor character encoding standards from various Asian countries.
- Mappings to entries for the character in a wide variety of dictionaries.
- Mappings between equivalent Simplified Chinese and Traditional Chinese characters.
- Mappings between other variant forms of the same character.
- Pronunciations of the character in various languages.
- An English definition of the character.
- The character's stroke count.

For more information on the Han characters and Han unification, see Chapter 10.

CHAPTER 6 *Unicode Storage and Serialization Formats*

As we saw earlier on, the Unicode standard comprises a single, unified coded character set containing characters from most of the world’s writing systems. The past few chapters—indeed, most of this book—focus on the coded character set. But Unicode also comprises three encoding forms and seven encoding schemes, and there are a number of other encoding forms and schemes out there that aren’t actually part of the Unicode standard but are frequently used with it.

The encoding forms and schemes are where the rubber meets the road with Unicode. The coded character set takes each character and places it in a three-dimensional encoding space consisting of seventeen 256×256 planes. The position of a character in this space—its *code point value*—is given as a four- to six-digit hex value ranging from zero to $0x10FFFF$, a 21-bit value.

The Unicode character encoding forms (or “storage formats”) take the abstract code point values and map them to bit patterns in memory (“code units”). The Unicode character encoding schemes (or “serialization formats”) take the code units in memory and map them to serialized bytes in a serial, byte-oriented medium such as a disk file or a network connection. Collectively, these different encoding forms and schemes are called the Unicode Transformation Formats (or “UTFs” for short).

The mappings from abstract code point values to serialized bytes is sometimes straightforward and sometimes not. The variety of different storage and serialization formats exist because depending on the application, you may want to make different tradeoffs between things like compactness, backward compatibility, and encoding complexity. In this chapter, we’ll look at all of them, and in Chapter 14 we’ll go back and look at way of actually implementing them.

A historical note

It's important to note that much of the terminology surrounding the presentation of Unicode in bits is relatively new, going back only a few years in Unicode's history, and some of the terms we discuss here were originally called by different names.

When it was first designed, Unicode was a fixed-length 16-bit standard. The abstract encoding space was sixteen bits wide (a single 256×256 plane), and there was basically one character encoding form—a straightforward mapping of 16-bit code points to 16-bit unsigned integers in memory. There was one official encoding scheme, which prefixed a special sentinel value to the front of the Unicode file to allow systems to auto-detect the underlying byte order of the system that created the file.

Fairly early on in the life of the standard, it became clear that a special encoding system based on Unicode but tailored for backward compatibility with ASCII was important to have. This encoding went through several names and incarnations before coming to be called UTF-8. The term "UTF" ("Unicode Transformation Format") came from the fact that the encoding was an algorithmic transformation you could perform on the 16-bit Unicode values.

By the time Unicode 2.0 came out, it was clear that an encoding space with only 65,536 cells wasn't going to be big enough, and the "surrogate mechanism" was invented: It reserved 2,048 positions in the encoding space and designated that they were to be used in pairs: one from the 1,024 "high surrogate" values followed by one from the 1,024 "low surrogate" values. The surrogates were meaningless in themselves, but *pairs* of surrogates would be used to represent more characters.

The ISO 10646 standard had started with a 31-bit encoding space—32,768 256×256 planes. It had two encoding forms: UCS-4, which mapped the 31-bit abstract code point values to 32-bit unsigned integers in memory (the extra bit—the sign bit—was always off, allowing a UCS-4 value to be stored in either signed or unsigned integer types with no change in numeric semantics), and UCS-2, which allowed Plane 0 to be represented with a 16-bit serialized value simply by lopping off the top two zero bytes of the abstract code point value.

Unicode 2.0 with the surrogate mechanism made it possible for Unicode to cover the first seventeen planes of the ISO 10646 encoding space: Plane 0 with a simple Unicode code point, and Planes 1 through 16 with surrogate pairs. Originally, this wasn't thought of as extending Unicode to a 21-bit encoding space: some characters were just represented with two code points. But the mapping from surrogate pairs to Planes 1 through 16 of the ISO 10646 encoding space, and the abstract code point values from ISO 10646 took hold in Unicode parlance, where they were called "Unicode scalar values."

The surrogate mechanism was adopted into the ISO 10646 standard as "UTF-16," and this name began to take hold in the popular parlance. By the time Unicode 3.0 came out, Unicode was being described as having a 21-bit encoding space, and UTF-16 was put on an equal footing with UTF-8 as alternative ways to represent the 21-bit abstract Unicode values in bits. The "Unicode scalar value" became the "code point value" (a code point was now 21 bits wide instead of 16) and the term "code unit" was coined to refer to the 16-bit units from UTF-16 (or the bytes of UTF-8). A third form, UTF-32, was added to Unicode later as a counterpart to UCS-4 from the ISO 10646 standard, and several new serialization formats were introduced. Meanwhile, WG2 agreed to limit ISO 10646 to the same 21-bit encoding space Unicode was using, and we get to today's situation.

Throughout the discussion that follows, I'll try to refer back to the history as it's important, or when older terms may still be in common use.

UTF-32

The simplest (and newest) of the Unicode encoding forms is UTF-32, which was first defined in Unicode Standard Annex #19 (now officially part of Unicode 3.1). To go from the 21-bit abstract code point value to UTF-32, you simply zero-pad the value out to 32 bits.

UTF-32 exists for three basic reasons: 1) It's the Unicode standard's counterpart to UCS-4, the four-byte format from ISO 10646, 2) it provides a way to represent every Unicode code point value with a single code unit, which can make for simpler implementations, and 3) it can be useful as an in-memory format on systems with a 32-bit word length. Some systems either don't give you a way to access individual bytes of a 32-bit word or impose a performance penalty for doing so. If memory is cheap, you can gain performance by storing text internally as UTF-32.

UTF-32 and UCS-4 didn't historically refer to the same thing: Originally, UCS-4 mapped ISO 10646's 31-bit encoding space to a 32-bit value, while UTF-32 mapped Unicode's 21-bit encoding space to a 32-bit value. In other words, the values from 0x110000 to 0x7FFFFFFF were legal UCS-4 code unit values but not legal UTF-32 code unit values (values of 0x80000000 and above were illegal in both encodings).

ISO 10646 still conceptually has a 31-bit encoding space, but WG2, the group behind ISO 10646, has agreed to restrict actual allocation of code point values to characters to just the first seventeen planes, making UCS-4 and UTF-32 functionally identical from here on out.

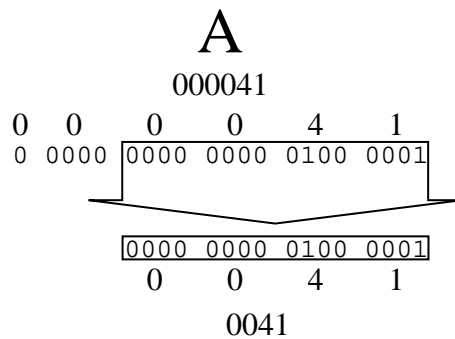
The big downside to UTF-32, of course, is that it wastes a lot of space. Eleven bits in every code unit are always zero and effectively wasted. In fact, since the vast majority of actual text makes use only of characters from the BMP, *sixteen* bits of the code unit are effectively wasted almost all the time.³⁰ And you don't generally gain a lot in implementation efficiency, since a single character as the user sees it might still be represented internally with a combining character sequence. That is, even with UTF-32 you may have to deal with single "characters" as the user sees it being represented with multiple code units (the fact that Unicode considers each of these units to be a "character" isn't generally all that helpful).

³⁰ Theoretically, of course, you could save some of that space by using 24-bit code units, but this has never been seriously suggested because no machine architecture can comfortably handle 24-bit values except by padding them out to 32 bits. A 24-bit encoding would thus effectively be limited to use as a serialization format, and there are more compact ways to represent Unicode in serialized data.

UTF-16 and the surrogate mechanism

So that brings us to UTF-16. UTF-16 is the oldest Unicode encoding form, although the name “UTF-16” only goes back a couple years.

UTF-16 maps the 21-bit abstract code point values to sequences of 16-bit code units. For code point values in the BMP, which represent the vast majority of characters in any typical written document, this is a straightforward mapping: you just lop the five zero bits off the top:

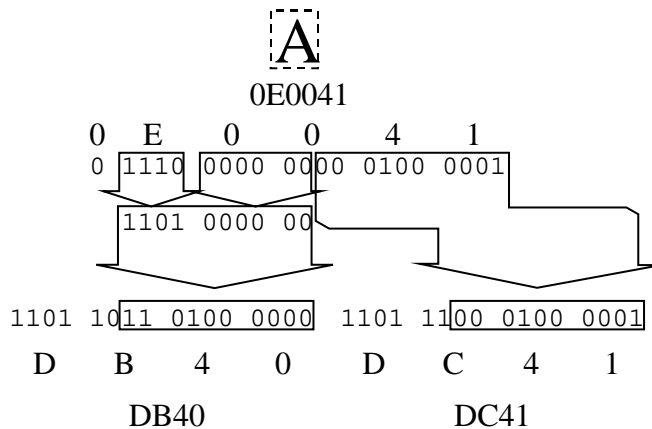


For characters from the supplementary planes, the transformation is more complicated. To represent supplementary-plane characters, Unicode sets aside 2,048 code point values in the BMP and reserves them. These 2,048 code point values will never be assigned to actual characters, which frees up their corresponding 16-bit code unit values to be used for the supplementary-plane characters.

The reserved code point values run from U+D800 to U+DFFF. This range is called the “surrogate range,” and code unit values in this range are called “surrogates.” The name is somewhat historical. Originally the surrogates were thought of as full-fledged code point values that “stood in” or “acted as surrogates” for supplementary-plane characters. Now they’re just reserved code point values.

The surrogate range is divided in half. The range from U+D800 to U+DBFF contains the “high surrogates” and the range from U+DC00 to U+DFFF contains the “low surrogates.” A supplementary-plane character is represented with a single code unit value from the high-surrogate range followed by a single code unit value from the low-surrogate range (a so-called “surrogate pair”). Unpaired or out-of-order surrogates are illegal and don’t represent anything.

The mapping from a supplementary-plane character to a surrogate pair is rather complicated. First, you subtract 0x10000 from the original code point value to get a 20-bit value. You then split the result into two ten-bit sequences. The first ten-bit sequence becomes the lower ten bits of the high-surrogate value, and the second ten-bit sequence becomes the lower ten bits of the low-surrogate value. For U+E0041 TAG LATIN CAPITAL LETTER A, that mapping looks like this:



So U+E0041 turns into 0xDB40 0xDC41 when converted to UTF-16.

People have objected to the surrogate mechanism as a corruption of the original Unicode ideal, which was to have a *fixed-length* 16-bit encoding. They're right, of course, that a compromise had to be made here, but they're wrong in thinking this opens up the whole can of worms that usually comes with variable-length encoding schemes.

In a typical variable-length encoding, you can't tell anything by looking at a single code unit (which is usually a byte). That byte may be a character unto itself, the second byte of a two-byte character, or in some cases the first byte of a two-byte character. You could only tell whether you were actually at the beginning of a character or not by looking at context. And if you lost a byte of the text in transmission, it could throw off the reading process, making it see the rest of the text as random garbage.

UTF-16 gets around this problem by using distinct ranges of values to represent one-unit characters, first units of two-unit characters, and second units of two-unit characters. A high-surrogate value can *only* be the first unit of a two-unit character, a low-surrogate value can *only* be the second unit of a two-unit character, and a non-surrogate value can *only* be a single-unit character. If you ever see an unpaired surrogate, it's an error, but the only character that gets corrupted is the character the surrogate represented half of: the surrounding text remains intact.

ISO 10646 specifies an encoding form called UCS-2, which also uses 16-bit code units. It's important to note that UCS-2 is *not* the same thing as UTF-16. The difference is the surrogate mechanism. UCS-2 can only represent BMP characters, every code point is represented with a single code unit, and the code units from U+D800 to UDFFF just aren't used. Supplementary-plane characters simply can't be represented using UCS-2. In other words, UTF-16 is UCS-2 plus the surrogate mechanism. (Since prior to Unicode 3.1 there were no characters in the supplementary planes, a lot of Unicode implementations didn't actually support the surrogate mechanism; you'll sometimes see these characterized as "supporting UCS-2" as a euphemistic way of saying they don't handle surrogate pairs.)

Some developers object to the additional code complexity involved in supporting UTF-16 as a representation format, but this shouldn't generally be all that big a deal. The reason for this is that

code points outside the Basic Multilingual Plane are only used for relatively rare characters. Mostly, the non-BMP range will contain characters for historical (i.e., dead) languages, specialized sets of symbols, and very specialized control signals. Except for scholarly documents in certain fields, non-BMP characters will never occur, or will occur only extremely rarely. Even in such specialized documents, they'll generally only show up in examples and other places where they'll be the minority of the text.

The one big exception to this is the CJK ideographs, which are living characters used to write living (and, in fact, thriving) languages. Even here, however, the most common characters (tens of thousands of them) are encoded in the BMP. Most ideographic text should only make occasional forays outside the BMP to get the characters they need. The biggest thing that'll require frequent trips outside the BMP is certain people's names, as the naming of children is the single biggest source of newly-coined ideographs. But even here, the number of times someone's name appears in a document about him or her is still way less than the number of times other words occur.

Another thing to keep in mind is that most of the processing required to handle surrogate pairs in memory is also required to support combining character sequences. There are certain operations that require a little more support to handle surrogate pairs, but for the most part, systems that correctly handle combining character sequences should also work right with surrogate pairs.

Endian-ness and the Byte Order Mark

Most experienced programmers are familiar with the concept of byte order, or “endian-ness.” Most personal-computer architectures of today evolved from old machine architectures that had an 8-bit word length. While their modern counterparts have longer word lengths, they still allow access to individual 8-bit fragments of a word. In other words, successive memory-address values refer to successive 8-bit fragments of words, not successive whole words. A 32-bit word on most modern machines thus takes up four memory locations.

How a 32-bit value is distributed across four eight-bit memory locations (or how a 16-bit value is distributed across two eight-bit memory locations) varies from machine architecture to machine architecture. That is, if you write 0xDEADBEEF into memory location 0, you're actually writing values into memory locations 0, 1, 2, and 3. Some machines, such as the PowerPC and SPARC families of processors, store the most significant byte of the word in the lowest-numbered memory location—that is, memory location 0 would contain 0xDE, and memory location 3 would contain 0xEF. These architectures are said to be *big-endian*. Other machines, such as the Intel x86/Pentium family of processors, store the most-significant byte of the word in the *highest*-numbered memory location—that is, memory location 0 would contain 0xEF and memory location 3 would contain 0xDE. These architectures are said to be *little-endian*.

When a program is operating in memory, endian-ness isn't important. In our example, if you write a 32-bit value (say, 0xDEADBEEF) to memory location 0, you just do it, and it's of no concern to you which memory location contains which byte. Later you read the value from memory location 0 back and it comes back as 0xDEADBEEF, just as you wrote it, again without your needing to be aware of which bytes were stored in which specific memory locations. This is because those memory locations will always be treated in a consistent manner. They're either always going to be accessed by the same processor, or at least they're shared by multiple processors that all share the same architecture. You only get in trouble in the rare instances when you sometimes treat memory locations 0 to 3 as a single 32-bit value and sometimes treat them as something else (e.g., four eight-bit values or two 16-bit values).

Where endian-ness starts to be an issue is when you start transferring data from a single computer's memory to some other location, such as another computer or a storage device. Writing data to the disk or sending it over a communications link of some kind requires *serializing* it: writing it out as a series of values of the same size. Again, since the smallest "word" size on most modern machines is eight bits, this means writing and reading the data as a sequence of bytes.

This, of course, requires taking all of the units that are longer than 8 bits and breaking them up into series of bytes. This requires you to make a decision as to which order you're going to write the bytes in when you break up a larger value.

This is pretty much always done in order from lowest-numbered memory location to highest-numbered memory location. This again means that if the same computer both writes something to the disk and reads it from the disk, or if two computers of the same processor architecture share a file over a network, endian-ness again isn't an issue, but if processors of *different* architectures share data over a network, or if you share data on a removable medium between machines of different architectures, you can run into serious trouble...

...unless, of course, you've secured some kind of an agreement beforehand as to what the serialized format's byte order is, or provided some way for the reading process to detect it.

Unicode does both of these things. It defines two encoding schemes based UTF-16: UTF-16BE and UTF-16LE, which are, respectively, the big-endian and little-endian versions of UTF-16. Since this distinction in meaningless in memory, the in-memory representation is just called UTF-16; UTF-16BE and UTF-16LE are only used to refer to UTF-16 text in some kind of serialized format, such as a disk file.

"UTF-16" can also be used to refer to serialized UTF-16, giving you a third encoding scheme. Unicode provides a mechanism for automatically detecting the endian-ness of serialized UTF-16 text whose endian-ness isn't announced by the enclosing protocol. This mechanism is called the *byte-order mark*, or "BOM" for short. ("BOM" is pronounced "bomb," not "bee-oh-em," so [WARNING: lame humor follows] be careful discussing Unicode serialization while passing through airport security checkpoints.)

The code point U+FEFF is designated as the Unicode byte-order mark. Its byte-swapped counterpart, 0xFFFE, on the other hand, is not a legal Unicode value. Thus, if a file you're reading contains any 0xFFFEs, you know it's either not Unicode text, or it's the wrong endian-ness and you need to byte-swap all the incoming character codes to get real Unicode. The writing process would include BOMs to allow this detection to happen on the receiving end.

The convention, thus, is to start off a Unicode text file (or a transmission including Unicode text) with a BOM anytime there's no external mechanism available to specify its endian-ness. So a file or transmission tagged as "UTF-16" rather than "UTF-16BE" or "UTF-16LE" would be expected to begin with a BOM. (The other formats could also begin with a BOM, but it'd be redundant.) If it doesn't, the convention is to assume it's UTF-16BE.

Since the BOM is useful as a way to announce the endian-ness of a Unicode text file, it can also be useful for telling whether a text file is a Unicode file in the first place, or for telling the difference between different variant versions of Unicode. You can find more information on how this works under "Detecting Unicode storage formats" below.

You need to be somewhat careful when using the BOM, however, because it has a double meaning. In Unicode 1.0, the code point U+FEFF only had the single meaning. It was a byte-order mark and only a byte-order mark. In normal Unicode text processing, it was simply a no-op.

But beginning in Unicode 2.0, the code point U+FEFF also acquired the semantic “zero-width non-breaking space” (often abbreviated “ZWNBSP”). A zero-width non-breaking space is used to “glue” two characters together. It has no visual representation, but will prevent a text processor that supports it from breaking a line between the characters on either side of the ZWNBSP. You might use it, for example, in the word “e-commerce” to keep “e-” from getting stranded on a line by itself, or to keep a Chinese phrase together on one line (Chinese text can usually have line breaks in the middle of words and phrases). When ZWNBSP occurs between two characters that naturally stick together on a single line (e.g., the letters in any English word), it has no effect.

Of course, if the ZWNBSP occurs at the beginning of a document, it’s basically a no-op since there’s nothing to “glue” the first character to. *But*, if you concatenate two Unicode text files together using something like the Unix utility `cat`, you’re now inadvertently gluing the last character of the first file to the first character of the second file.

Because of this, the double meaning of U+FEFF quickly came to be seen as a mistake. So starting in Unicode 3.2, the “zero-width non-breaking space” meaning of U+FEFF is deprecated. Despite the name (which unfortunately can’t change), U+FEFF goes back to being only a byte-order mark. A new character, U+2060 WORD JOINER, has been introduced in Unicode 3.2 to take over the task of “gluing” two words together on a single line. Conforming Unicode implementations are now allowed to convert U+FEFF ZERO WIDTH NO-BREAK SPACE to U+2060 WORD JOINER.

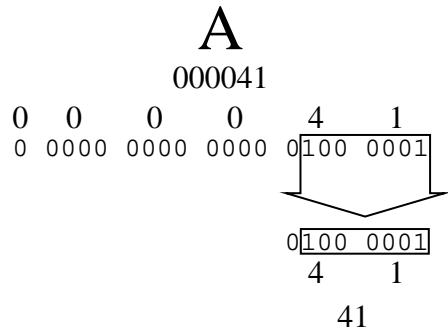
Of course, you also run into the endian-ness problem in UTF-32, so there are also three character encoding schemes based on UTF-32: UTF-32BE, UTF-32LE, and tagged UTF-32, which work the same way as their UTF-16-based counterparts.

UTF-8

Then there’s UTF-8, the 8-bit Unicode encoding form. UTF-8 was designed to allow Unicode to be used in places that only support 8-bit character encodings. A Unicode code point is represented using a sequence of anywhere from 1 to 4 8-bit code units.

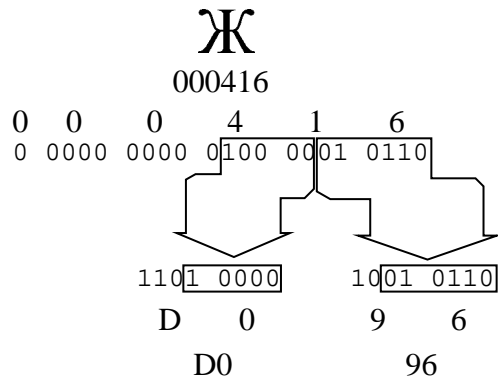
One vitally important property of UTF-8 is that it’s 100% backward compatible with ASCII. That is, valid 7-bit ASCII text is also valid UTF-8 text. This means that UTF-8 can be used in any environment that supports 8-bit ASCII-derived encodings and that environment will still be able to correctly interpret and display the 7-bit ASCII characters. (The characters represented by byte values where the most significant bit is set, of course, aren’t backward compatible—they have a different representation in UTF-8 than they do in the legacy encodings.)

So Unicode code point values from U+0000 to U+007F are mapped to UTF-8 in a very straightforward manner:



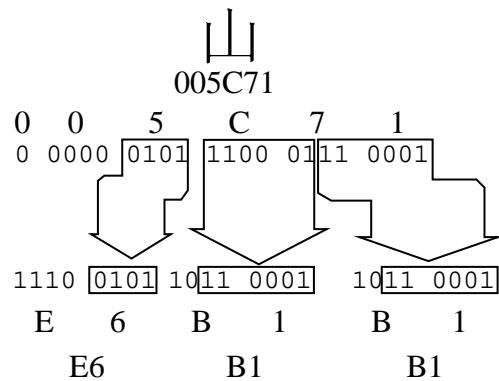
You simply isolate the last seven bits and zero-pad it out to eight bits.

Code point values from U+0080 to U+07FF take two bytes, as follows:



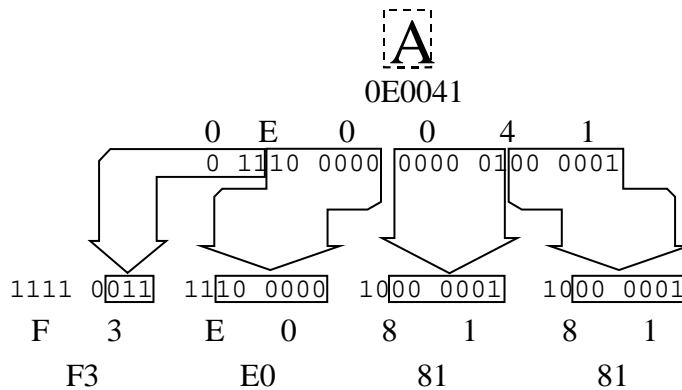
The last six bits, plus 0x80, become the second byte, and the next five bits, plus 0xC0, become the first byte.

Code point values from U+0800 to U+FFFF (the rest of the BMP) turn into three-byte sequences, as follows:



The first four bits, plus 0xE0, become the first byte, the next six bits, plus 0x80, become the second byte, and the last six bits, plus 0x80, become the last byte.

Finally, the code points in Planes 1 through 16 (the non-BMP characters) turn into four-byte sequences, as follows:



The first three bits, plus 0xF0, become the first byte, and the other three bytes each consists of the next six bits from the original code point, plus 0x80.

This actually can extend all the way out to cover the entire 31-bit ISO 10646 range, using 5- and 6-byte sequences. But this is mostly of historical interest (and to explain old UTF-8 implementations) now that WG2 has committed to never put any characters into Planes 18 and above.

As with UTF-16, UTF-8 avoids many of the problems with other multi-byte encodings by ensuring that different ranges of values are used in different positions in the character:

- 00–7F Single-byte character
- 80–BF Trailing byte
- C0–DF Leading byte of two-byte character
- E0–EF Leading byte of three-byte character
- F0–F7 Leading byte of four-byte character
- F8–FB Illegal (formerly leading byte of five-byte character)
- FC–FD Illegal (formerly leading byte of six-byte character)
- FE–FF Illegal

The only ambiguous thing here is the trailing byte. You can tell a byte is a trailing byte, but you can't tell which byte of the character it is, or how many bytes the character is. But because of UTF-8's design, you know you'll never have to scan forward or back more than three bytes to find out.

One side effect of UTF-8's design is that many code points can more than one potential representation in UTF-8. For example, U+0041 could be represented not only as 0x41, but also as 0xC1 0x81 or 0xE0 0x81 0x81. The standard stipulates that the shortest possible representation for any code point is the only legal one. Prior to Unicode 3.1, it was legal for UTF-8 implementations to *interpret* these so-called "non-shortest form" byte sequences. This created a potential security hole, so Unicode 3.1 tightened the definition to disallow *both* interpretation and generation of non-shortest-form sequences.

There's another version of "non-shortest-form" UTF-8 that also needs to be dealt with: representation of supplementary-plane characters using six-byte sequences. For example, instead of representing U+E0041 in UTF-8 as 0xF3 0xE0 0x81 0x81, as in the example above, you could conceivably represent it as 0xED 0xAD 0x80 0xED 0xB1 0x81, which is what you get if you convert it first to UTF-16, producing 0xDB40 0xDC41, and then convert the two surrogate code units to UTF-8.

Even Unicode 3.1's tightening of the UTF-8 definition still premitted this: as with other non-shortest-form sequences, it was illegal to produce byte sequences like this, but legal to interpret them. Unicode 3.2 closes this loophole: UTF-8 sequences representing code point values in the surrogate range (U+D800 to U+DFFF) are now completely illegal. That is, any three-byte UTF-8 sequence whose first byte is 0xED and whose second byte is anything from 0xA0 to 0xBF is illegal.

Since conversion of code points to UTF-8 results in a sequence of bytes, UTF-8 is effectively a character encoding scheme unto itself, and not just a character encoding form.

CESU-8

Draft Unicode Technical Report #26 proposes an encoding scheme with the rather unwieldy name of "Compatibility Encoding Scheme for UTF-16: 8-bit," or "CESU-8" for short.

CESU-8 is a variant of UTF-8. It treats the BMP characters the same, but deals with the supplementary-plane characters differently. In CESU-8, supplementary-plane characters are represented with six-byte sequences instead of four-byte sequences. In other words, the six-byte representation of supplementary-plane characters that is now illegal in UTF-8 is the *preferred* representation of these characters in CESU-8.

CESU-8 is what you get if you convert a sequence of Unicode code point values to UTF-16 and then convert the UTF-16 code units to UTF-8 code units as if they were code points. Supplementary-plane characters thus get turned into pairs of three-byte sequences, where each three-byte sequence represents a surrogate.

U+E0041 TAG LATIN CAPITAL LETTER A, which is 0xDB40 0xDC41 in UTF-16 and 0xF3 0xE0 0x81 0x81 in UTF-8, is 0xED 0xAD 0x80 0xED 0xB1 0x81.

CESU-8 is what you get if your program uses UTF-16 as its internal representation format and its conversion to UTF-8 doesn't know about surrogate pairs (either because it was written before the surrogate mechanism was devised or because it was written by someone who assumed there'd never be character allocations in the supplementary planes). It's safe to say it'll never get accepted as an official part of the standard, but it might get blessed as a technical report as a way of documenting the internal behavior of a lot of existing systems.

UTF-EBCDIC

One of the main uses of UTF-8 is to allow the use of Unicode text in systems that were designed for ASCII text. This is why regular ASCII text is also legal UTF-8. UTF-8 isn't the first international

character encoding to be in some way backwards-compatible with ASCII—quite a few encodings for various languages are backwards-compatible with ASCII in the same way and for the same reason.

Well, the other 8-bit encoding for Latin text that has huge bodies of text encoded in it is EBCDIC, IBM's Extended Binary-Coded Decimal Information Code. There's a lot of legacy data out there in EBCDIC, and a lot of legacy systems that are designed to assume that textual data is encoded in EBCDIC. UTF-EBCDIC (sometimes called "UTF-8-EBCDIC"), much like UTF-8, is designed to allow the use of Unicode in environments originally designed for EBCDIC.

EBCDIC isn't just ASCII with the characters listed in a different order; it includes characters that aren't in ASCII and leaves out characters that ASCII includes. Classical EBCDIC used an 8-bit byte, but encoded only 109 characters, leaving "holes" in the encoding space. Later versions of EBCDIC filled in the "holes," extending EBCDIC to a full 256-character encoding. As there are variant versions of 8-bit ASCII that put different characters into the space above 0x7F, so too are there different variants of EBCDIC that fill in the "holes" with different characters.

There's a set of 65 control signals and 82 printing characters that are common to all current variants of EBCDIC, distributed across the 8-bit encoding space. UTF-EBCDIC aims to preserve these values, using the remaining values to encode the rest of the Unicode characters as multi-byte combinations. It does this through the use of a two-step process that makes use of an intermediate format called "UTF-8-Mod" or "I8".

I8 is similar to UTF-8, except that the transformation to I8 not only preserves all of the values from 0x00 to 0x7F as single-byte sequences, but also the values from 0x80 to 0x9F (the "C1 controls" from the Latin-1 character set). This range covers all of the EBCDIC control characters, as well as all of the printing characters that EBCDIC and ASCII have in common. It then proceeds in UTF-8 style to encode the other Unicode code points as multi-byte sequences. Like UTF-8, I8 is designed so that the leading byte of a multi-byte sequence indicates how many bytes are in the sequence, and so that completely distinct ranges of byte values are used for single-byte characters, leading bytes of multi-byte sequences, and trailing bytes of multi-byte sequences. Because I8 preserves a greater range of values as single-byte sequences, the remaining Unicode characters may map to sequences of up to five bytes. (In UTF-8, characters in the BMP mapped to up to three bytes, with the non-BMP range taking four bytes; in UTF-EBCDIC, characters in the BMP may map to up to four bytes, and non-BMP characters may map to either four- or five-byte sequences.)

The mapping from I8 to actual UTF-EBCDIC is a simple and reversible table lookup: The single-byte values are mapped to the values the corresponding characters have in EBCDIC, and the bytes of the multi-byte characters are mapped into the remaining spaces. The resulting text can be handled properly by most processes expecting to see regular EBCDIC text.

Unlike UTF-8, you can't tell by doing simple bit masking whether an individual byte of a UTF-EBCDIC sequence is a character unto itself or part of a multi-byte sequence, but since you *can* tell this from looking at an I8 sequence, and since you can get from I8 to UTF-EBCDIC with a simple table lookup, you can also use a simple table lookup to identify the leading, trailing, and standalone bytes in a UTF-EBCDIC sequence.

I8, by the way, is an intermediate representation only, and not a character encoding form unto itself. Also, UTF-EBCDIC is suitable for use only in homogeneous EBCDIC-based systems and not for general interchange. In particular, you can't use UTF-EBCDIC in an ASCII or UTF-8-based environment.

Because UTF-EBCDIC is designed only for internal use on EBCDIC-based systems and not for interchange, it's not an official part of the Unicode standard. It's published as Technical Report #16 as a way of officially documenting the practice for people writing code for EBCDIC-based systems.

UTF-7

There are other, less common, Unicode transformation formats you might run across. The most common of these is UTF-7, which was designed for use in environments that were designed specifically for 7-bit ASCII and can't handle 8-bit characters. In particular, the original version of the Simple Mail Transfer Protocol (SMTP), which is still in common use, wouldn't work with 8-bit character values and required this approach.

UTF-7 does this by using a scheme similar to the numeric-character reference scheme used in HTML and XML: most of the ASCII characters are just themselves, but a few are used to signal sequences of characters that specify otherwise unrepresentable Unicode values. Thus, unlike the other Unicode transformation formats, UTF-7 is a stateful encoding, with values whose interpretation depends on context. (By way of analogy, consider the "<" sequence in HTML: to tell whether the "l" is a real letter *l* or not, you have to scan both directions to see if it occurs between a & and a ;—trailing bytes in UTF-8 could be thought of as "stateful" because you can't tell specifically which position they occupy in the sequence, but you know they're trailing bytes, and the number of characters that you have to examine in order to find the beginning of the sequence is bounded.)

UTF-7 represents the letters A through Z and a through z, the digits 0 through 9, certain punctuation marks (in particular, ",()'-./:?", plus, optionally, a number of other characters), and the space, tab, carriage return, and line feed the same way they're represented in ASCII. The other Unicode characters are represented using a variant of the Base64 encoding: the + character signals the beginning of a sequence of Base64 text. A sequence of Base64 ends when a non-Base64 character is encountered. The - character can be used to force a switch out of Base64.

UTF-7 uses a modified version of Base64 that avoids using the = character, since the = character has a special meaning in some of the places where UTF-7 was intended to be used.

UTF-7 was never part of the Unicode standard and was never published as a Unicode Technical Report, but was published by the Internet Engineering Task Force as RFC 2152 (and earlier as RFC 1642). These days, it's pretty much been supplanted by regular Base64, MIME, and various other protocols. For more on Unicode in Internet mail, see Chapter 17.

Standard Compression Scheme for Unicode

One of the major reasons for resistance to Unicode when it first came out was the idea of text files taking up twice as much room as before to store the same amount of actual information. For languages like Chinese and Japanese that were already using two bytes per character, this wasn't a problem, but the idea of using two bytes per character for the Latin alphabet was anathema to a lot of people.

The concern is certainly legitimate: the same document takes up twice as much space on a disk and twice as long to send over a communications link. A database column containing text would now take

twice as much disk space. In an era of slow file downloads, for example, the idea of waiting twice as long to download an email message is pretty unpalatable when it already takes too long. The idea that processing speed, transmission speed across networks, and disk space all increase so rapidly that “it won’t be a problem in a couple years” always rang hollow to me. Just because I have more disk space doesn’t mean I want to waste it.

Of course, these days processing speed, disk space, and even modem speeds are such that sending, receiving, and storing text isn’t all that big a deal. Even long text documents don’t take up much room. The bandwidth problems we’re running up against today don’t have anything to do with sending text back and forth—they have to do with sending graphics, sound, and video back and forth. Today the size premium that Unicode imposes doesn’t seem like a big deal, although it certainly did ten years ago when Unicode was getting its start.

However, it’s still a big deal when you have really large concentrations of text in one place. If you’re General Motors and you’ve got millions of parts to keep track of, doubling the size of the “Description” field in your part database is going to have a big impact. If you’re a big wire service and you’re sending an endless stream of news stories out over the wire to your subscribers, doubling the size of all of your news stories is going to have a huge impact.

The Unicode standard has always been silent on this issue, saying that compression was a “higher-level protocol” and that there were plenty of effective compression algorithms out there that could take care of this problem. True enough, although using something like Lempel-Ziv-Welch compression to compress simple text files seems a little like overkill and introduces noticeable hassle. If all you’re trying to do is get back to about the same level of efficiency you had when you were using a single-byte encoding and do it in a way that makes it simple to compress and decompress, the standard compression algorithms don’t offer a lot of relief.

Many people started using UTF-8 for this purpose. It’s easy enough to compress and decompress, and it represents ASCII text with the same efficiency as ASCII. Besides, it’s less work to retrofit things to use UTF-8 than to go whole-hog with UTF-16 or UTF-32.

The last reason is certainly true enough, and is one reason why UTF-8 will always be with us, and is likely to always be the most popular way of exchanging Unicode data between two entities in a heterogeneous environment. But UTF-8 does impose a small penalty on Latin-1 text. In languages like Spanish and French that use a lot of accented letters, the fact that an accented letter uses two bytes in UTF-8 and only one in Latin-1 begins to make a noticeable difference.

And it doesn’t help speakers of non-Latin languages. Greek, for example, takes two bytes per character in UTF-8 but only one in a native Greek encoding. Worse yet, Japanese, Chinese, and Korean, which generally take two bytes per character in their native legacy encodings, take *three* bytes per character to represent in UTF-8! Finally, there are some languages, such as Thai, that take only one byte per character in their native legacy encodings that also take three bytes in UTF-8.

The bottom line is that using UTF-8 for compression is a very Eurocentric (even English-centric) thing to do. For English, UTF-8 has the same storage efficiency as the legacy encoding, and twice the efficiency of UTF-16, and for most other European languages, UTF-8 imposes maybe a 15% penalty over the legacy encoding, but that’s still better than the 100% penalty UTF-16 imposes. But for Greek, Cyrillic, Hebrew, Arabic, and a couple other languages, UTF-8 and UTF-16 impose the same penalty over the native encoding. For Japanese and Chinese, UTF-16 imposes little or no storage penalty, but UTF-8 imposes a 50% penalty. And for Thai, UTF-16 imposes a 100% penalty over the legacy encoding, but UTF-8 imposes a 200% penalty over the native encoding!

So unless you know you'll always be restricted to working in just the languages where UTF-8 gives you good efficiency, UTF-8 isn't effective as a means of saving space.

Earlier I mentioned how going from ASCII to UTF-16 would really hurt you if you were a wire service sending out thousands of new stories and other information every day. Well, it was a wire service, Reuters, that first put forth a solution to the problem, a simple and efficient compression scheme for Unicode that was the forerunner of what is now known as the Standard Compression Scheme for Unicode, or SCSU for short.

SCSU is a relatively simple compression scheme designed specifically for Unicode text. It takes advantage of the properties of Unicode text to achieve its efficiencies, and produces results that are generally comparable in size to the legacy encodings. Text in SCSU can then be compressed with a general-purpose compression scheme such as LZW to achieve even greater compression ratios (in fact, LZW will generally do a better job on SCSU-encoded Unicode than on regular UTF-16 or UTF-8).

SCSU is a file format, not an algorithm. It lends itself to a multitude of different compression algorithms that can trade off compression speed versus compactness in different ways. The algorithm for decompressing SCSU, on the other hand, remains the same and is simple and straightforward (JPEG and MPEG compression are designed around the same principle: one simple way to decompress, lots of different ways to compress).

SCSU is a stateful encoding, meaning it has "modes" and you can't jump into the middle of an SCSU-compressed file and tell what you're looking at without potentially seeking all the way back to the beginning. This means it's generally suitable as an interchange or storage format, but not terribly well-suited to internal processing.

SCSU has two basic modes: single-byte mode and Unicode mode. In single-byte mode, the ASCII printing characters and the most common ASCII control characters are represented as themselves. This means that most ASCII text files (those that only use the most common control characters) are interpretable as SCSU.

The byte values from 0x80 to 0xFF, as usual, are where most of the magic happens. Unlike in Latin-1 and most other encodings, where they have a fixed interpretation, and unlike in UTF-8, where they have no meaning by themselves but are used in variable-length combinations to form characters, the values from 0x80 to 0xFF have a variable interpretation in SCSU. Their default interpretation is the same as in Latin-1, so most Latin-1 text files are interpretable as SCSU (remember that Latin-1 is *not* interpretable as UTF-8).

Remember, however, that only the most common ASCII control characters are represented unchanged. The others are used by the compression algorithm for various things. One set of control characters is used to control the interpretation of the values from 0x80 to 0xFF. In all cases, these values are used to represent a contiguous range of 128 characters from somewhere in the Unicode encoding range. This range is called a "window" and the default window is the Latin-1 range. There are a bunch of predefined windows that cover other interesting ranges, such as the Cyrillic window or the Hebrew window. You can represent text that all falls into one 128-character range of Unicode characters simply by having one control character in the stream to set the appropriate window and then using the values from 0x80 to 0xFF to represent the characters in that window.

If the Unicode text being encoded spans more than one window (the ASCII window doesn't count, since it's always accessible), you intersperse more control codes to shift the window back and forth. There are also "quote" control characters that shift the window only for the very next character instead of for all following characters.

SCSU also lets you define windows, so that if the characters you need to encode aren't in any of the predefined windows, you can define one that does cover them and have characters that take you to it.

Of course, there are scripts, the Han characters being the most notable example, that span a much larger range than 128 characters and for which the window-shifting scheme would be unwieldy. This is why SCSU has a Unicode mode. You include a control character that shifts into Unicode mode, and the bytes that follow are interpreted as big-endian UTF-16 (there's a range of byte values that are used as control codes in Unicode mode to allow you to shift back into byte mode). There's also a code that shifts into Unicode mode for a single Unicode character, rather than for all following text.

This scheme is relatively easy to implement a decoder for, and allows for considerable discretion in the design of encoders. A simple encoder could simply shift to Unicode mode at the beginning of the text stream and pass the rest of the text through unchanged, taking care to insert appropriate escape characters if the input text includes the tag values that SCSU uses in Unicode mode. Or a simple encoder could pass through Latin-1 characters unchanged and insert the quote-single-Unicode-character tag before each non-Latin-1 character in the input stream. A sophisticated encoder, on the other hand, can make extensive use of all of the features of the encoding, combined with lots of lookahead when examining the input stream, to achieve very high levels of compression.

A good SCSU encoder should be able to produce results for most text that are similar in size to the same text encoded in the common legacy formats. SCSU-encoded text will almost always offer a significant savings over regular UTF-16 or UTF-8 text, and should almost never (with a well-designed encoder) impose a penalty of more than a few bytes.

For more on implementing SCSU, see Chapter 14.

BOCU

An interesting alternative to SCSU for Unicode compression was proposed in a recent paper on IBM's developerWorks Web site.³¹ It's called "Binary Ordered Compression for Unicode," or "BOCU" for short. It provides compression ratios comparable for SCSU, but with a number of interesting advantages. Chief among them: a set of Unicode strings encoded in BOCU sorts in the same order as the unencoded strings would, making it useful for encoded short snippets of text in environments where the sort order is still important (think database applications). The algorithm is also a lot simpler than SCSU.

The basic idea is that instead of independently transforming each code point to a sequence of bytes using some fixed mapping, you transform each code point to a sequence of bytes by subtracting it from the preceding code point in the text and encoding the *difference* as a series of bytes. BOCU arranges things so that differences of plus or minus 0x3F or less are represented with single bytes, differences of plus or minus 0x40 to 0x2F3F are represented with two-byte sequences, and larger

³¹ Mark Davis and Markus Scherer, "Binary-Ordered Compression for Unicode," IBM DeveloperWorks, <http://www-106.ibm.com/developerworks/unicode/library/u-binary.html>.

differences are represented with three-byte sequences. The lead bytes of these sequences are then arranged in such a way as to preserve the relative sorting order or unencoded strings.

If, instead of just subtracting the current character from the preceding character, you adjust the value you subtract from to be in the middle of its Unicode block and consider a little more context than just the immediately preceding character, you can arrive at an encoding that's almost as efficient as SCSU (and in some cases, more so).

Variants of BOCU is used internally in parts of the International Components for Unicode, the popular Unicode support library. It's unclear whether it'll catch on to a greater extent, but it's an interesting idea worth thinking about if you need the things it can give you.

Detecting Unicode storage formats

The Unicode 2.0 standard, in its discussion of the byte-order mark, talked about how it could be used not just to tell whether a Unicode file was the proper endianness, but whether it was a Unicode file at all. The idea is that the sequence 0xFE 0xFF (in Latin-1, a lowercase y with a diaeresis followed by the lowercase Icelandic letter “thorn”) would pretty much never be the first two characters of a normal ASCII/Latin-1 document. Therefore, you could look at something you knew was a text file and tell what it was: If the first two bytes were 0xFE 0xFF, it was Unicode; if 0xFF 0xFE, it was byte-swapped Unicode, and if anything else, it was whatever the default encoding for the system was.

As Unicode transformation formats have proliferated, so too has the idea of using the byte-order mark at the beginning of the file as a way of identifying them. The current chart looks like this:

If the file starts with...	...the file contains...
0xFE 0xFF	UTF-16
0xFF 0xFE	byte-swapped UTF-16
0x00 0x00 0xFE 0xFF	UTF-32
0xFF 0xFE 0x00 0x00	byte-swapped UTF-32
0xEF 0xBB 0xBF	UTF-8
0xDD 0x73 0x73 0x73	UTF-EBCDIC
0x0E 0xFE 0xFF	SCSU (recommended; others are possible)
anything else	non-Unicode or untagged Unicode

The basic rule behind the use of these things is simple: If there's any other way to specify (or tell) which format some piece of text is in, don't rely on the byte-order mark. If the byte-order mark isn't there, then you're stuck. You don't know what the file is. If the byte-order mark *is* there, there's still the remote possibility the file's in some other format and just happens to start with those bytes (for example, is a file that starts with 0x0000 0xFEFF a UTF-32 file starting with a byte-order mark or a UTF-16 file starting with a null character and a zero-width non-breaking space?).

In other words, use some means other than looking at the text itself to identify what the encoding is. Designing a system, for example, that uses “.txt” as the filename extension for all Unicode text files and then looks for the BOM to tell whether it's UTF-16, UTF-8, or ASCII isn't terribly bulletproof. Instead, use different extensions (or different types, if you're using a typed file system), allow only one or two types, rely on the user to tell you, or something else like that.

If you adopt some kind of protocol that *requires* that the byte-order mark always be there, you're effectively specifying a higher-level protocol for specifying the character encoding. In other words, the file isn't really a Unicode text file anymore; it's in a format that contains text in some Unicode format, preceded by a tag telling what the format is (to handle the BOM correctly, you've pretty much got to treat it this way to begin with). There are various higher-level protocols available already (XML, for example) that do this just as well.

The basic philosophy to follow in most situations is to treat the different Unicode transformation formats as entirely distinct character encodings or code pages, no more related than Latin-1 and Latin-2. The byte-order mark only works for telling apart different forms of Unicode; something like XML allows you to also tell these forms of Unicode apart from whole other encoding schemes. Process text internally in one format (for example, UTF-16 in the native byte order) and treat incoming data in all other formats, be they different flavors of Unicode or different encodings altogether, as a code-conversion problem.

SECTION II

Unicode in Depth

A Guided Tour of the Character Repertoire

CHAPTER 7 *Scripts of Europe*

Okay, having taken a look at the overall architecture of Unicode and the various aspects of the Unicode standard that are relevant no matter what you're doing with it, it's time to take a closer look. The next six chapters will be a guided tour through the Unicode character repertoire, highlighting characters with special properties and examining the specifics of how Unicode is used to represent various languages.

Unicode is organized into blocks, with each block representing characters from a specific writing system, or *script*. In each of the following chapters, we'll look at collections of scripts with similar properties. First we'll look at the issues they have in common with one another, and then we'll look at the unique characteristics of each individual script. Along the way, we'll stop to discuss any characters with special properties you may have to know about.

This section will concentrate mostly on the characters that make up “words” in the various languages (i.e., letters, syllables, ideographs, vowel points, diacritical marks, tone marks, and so forth), but will also include a chapter that deals with digits, punctuation, symbols, and so forth. We'll also highlight interesting digits, punctuation marks, and symbols that belong to specific scripts as we go through the individual scripts.

The Western alphabetic scripts

It's entirely possible that one of the reasons computer technology flourished so much in the United States (and, later, Europe) before it caught on in other places has to do with our system of writing (this is not, of course, to discount all of the other reasons, such as economic factors, World War II,

etc.). Compared to many other writing systems, the Latin alphabet is simple and lends itself well both to mechanical means of reproduction and to various types of automated analysis.

The classic twenty-six-letter Latin alphabet, as used to write English, is pretty simple. There are only twenty-six letters, a very manageable number, and very few typographic complications. The letters represent both consonant and vowel sounds, unlike most other writing systems, and they lay out very simply, marching from left to right across the page, one after another. The letters don't change shape depending on the letters around them, they don't change order depending on context, they don't combine together into ligatures, and they don't carry complex systems of marks with them to indicate meaning. Words are always separated with spaces, making it easy to detect word boundaries.

Of course, not all of these things are true all the time, and there was a time in history when each of them was false. All of the writing systems used in the West and Middle East evolved out of the Phoenician alphabet, which was written from right to left and didn't have vowels. The Greek alphabet was the first one to assign real letters to vowel sounds, and the Latin alphabet picked it up from there. Originally, the Greek alphabet was also written from right to left, but a form of writing called *boustrophedon* developed. In *boustrophedon*, lines of text alternated direction: one line would run from right to left, then the next line would run from left to right, then the next from right to left again, and so on in zigzag fashion down the page. Gradually, this died out and all lines were written from left to right.

Boustrophedon was generally seen just in Greek writing, but there are early examples of Etruscan and Latin writing that use it, so it's not 100% unheard-of for Latin-alphabet text to be written right to left. Spaces between words also weren't always used. This doesn't appear to have become common practice until around the 8th century or so. You see spaces in some earlier writing, and you also see dots between words, but you very frequently see earlier writing where the words just all run together and you have to know the language to know the word boundaries.

There are also plenty of cases in which the shapes of the letters can change depending on context. In handwritten English text, for example, the letters don't all join the same way, giving rise to variations in the shapes of the letters. For example, at least in the way I learned to write, the letters *b* and *o* join to the succeeding letter above the baseline, rather than at the baseline. This changes the shape of the next letter slightly to account for the higher joiner:

The image shows the word "Cannon" written in a cursive, handwritten style. The letters are connected in a way that demonstrates how the shape of a letter can change depending on the letter it follows. Specifically, the 'a' and 'n' in "Cannon" are joined to the following letter above the baseline, which affects the shape of the 'n' that follows.

[It'd be really nice if we could replace this drawing with either an example from a font that does this, or a better-looking example done by someone with decent handwriting.]

Note how the shape of the letter "n" changes: The "a" joins to the "n" along the baseline (and the first "n" joins to the second "n" along the baseline). But the "o" joins to the "n" above the baseline.

There are other cases where the letters change shape. In *boustrophedon*, the shape of a letter on a right-to-left line was the mirror image of its shape on a left-to-right line, enabling the reader to tell easily the direction of a given line of text. As recently as two hundred years ago, there were two forms of the letter *s*, one that was used at the end of a word and one that was used at the beginning or in the middle of a word (in fact, this is still true in the Greek alphabet). In modern ornamental typefaces, you'll still see special versions of certain characters that are intended to be used at the ends of words.

And of course, in good typography, letters still sometimes combine into ligatures. We've talked about the *fi* and *fl* ligatures already. In older writing, *ae* and *oe* ligatures were also common, and other forms, such as *ct* and *st*, still happen in ornamental typefaces.

As for systems of marks, you do occasionally see diacritic marks on letters in English words (such as in the word “naïve”), although they’re almost never required to understand the text (we still recognize “naïve”). This isn’t true in most other languages that use the Latin alphabet, but in almost all of them (Vietnamese being the big exception), the number of letter-mark combinations is pretty small.

The point is not that the various more complex typographical effects *never* happen in English or other languages using the Latin alphabet, but that they’re not *required*. Modern English text is perfectly legible with absolutely no typographic interaction between characters, no diacritical marks, and strict left-to-right ordering. In fact, you can go further: you can eliminate the one wrinkle unique to the Greek and Latin family of scripts—upper- and lower-case forms—and you can even make all the characters the same width so that they line up as if they were written on graph paper, and the text is still perfectly legible.

THE QUICK BROWN FOX JUMPS OVER THE LAZY DOG

...may not look as nice as...

The quick brown fox jumps over the lazy dog

[It’d be better to find an expert-set font with some real ligatures and swashes, but I don’t have access to any right now (if we have to change the text of the example to show off the ligatures and swashes better, that’s fine).]

...but no one who reads English would look at it and have trouble figuring out what it says. With many other writing systems, this isn’t true.

These characteristics lend themselves well to mechanical processing, making it possible, for example, to build a typewriter for the Latin alphabet (the Arabic, Devanagari, and Chinese writing systems, in contrast, don’t lend themselves at all well to being written with a typewriter).

This chapter covers the main writing systems used to write the languages of Europe: the Latin, Greek, Cyrillic, Armenian, and Georgian alphabets, along with the International Phonetic Alphabet. These writing systems share the following characteristics: simple left-to-right line layout, minimal (or no) typographic interaction between characters, minimal (or no) use of diacritical marks, and spaces between words. Most of these writing systems are descended in some way from the Greek alphabet. In addition, this chapter covers the collections of diacritical marks used with these (and some other) writing systems.

The Latin alphabet

Since this is an English-language book directed primarily toward an English-speaking audience, and since the bulk of text stored in computer systems today is still in English or other European languages, the logical place to begin our tour of the Unicode character repertoire is with the Latin alphabet.

The Latin alphabet is called that, of course, because it was originally used to write Latin (it’s also frequently called the Roman alphabet because it was developed by the ancient Romans). The Latin alphabet evolved out of the Etruscan alphabet, which in turn evolved from the Greek alphabet. Later, some letters that weren’t used in Etruscan were borrowed straight into the Latin alphabet from the

Greek alphabet. The Latin alphabet acquired the form we're familiar with sometime around the 1st century BC. It originally had twenty-three letters:

A B C D E F G H I K L M N O P Q R S T V X Y Z

In fact, Y and Z were late additions that were only used for words borrowed in from Greek. The letters I and V were used as both consonants and vowels, and J and U were just variant forms of I and V. W was just a special ornamental ligature that was used when two Vs appeared next to each other. It wasn't until the Renaissance that I and U became consistently used as vowels and J, V, and W as consonants. That gave us the twenty-six-letter Latin alphabet we're so familiar with today:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

The lower-case versions of the letters date back to about the third or fourth century AD, but were originally just variant forms of the capital letters. Capital and small letters were being used simultaneously the way they are now somewhere around the eighth century.

The Latin alphabet spread throughout Europe first with the spread of the Roman Empire and later with the spread of the Roman Catholic Church. In more recent times, it has been applied to a wide variety of non-European languages as well. Today, the Latin alphabet is used to write more different spoken languages than any other single writing system.

Since ASCII is the most common character encoding scheme used in computers today, and forms the basis for so many other encoding systems, the designers of Unicode decided to adopt the ASCII character repertoire directly into Unicode unchanged. Not only do the characters appear in Unicode in the same order, but they even have the same numeric code-point values (zero-extended out to 21 bits, of course). In this way, converting between ASCII and UTF-16 or UTF-32 is a simple matter of zero-padding out to the appropriate length or truncating to 8 bits, and converting between ASCII and UTF-8 is a no-op.

ASCII includes fifty-two letters: the upper- and lower-case versions of the twenty-six Latin letters. Unfortunately, only two modern languages use the original Latin alphabet with no additional letters: English and Dutch. Even English historically had other letters in its alphabet.

With a couple of notable exceptions, no single language that uses the Latin alphabet adds more than a handful of extra letters. However, they don't all add the *same* extra letters (many languages don't use all twenty-six letters of the basic Latin alphabet either).

For instance, Spanish, Italian, Portuguese, French, and many other languages all use accent marks on some vowels to clarify pronunciation or distinguish between two similar words (or forms of the same word). So in addition to the regular twenty-six letters, you also need the five vowels with acute and grave accents: á, é, í, ó, ú, à, è, ì, ò, ù.

French and Portuguese also use the circumflex accent, giving us five more letters: â, ê, î, ô, û. German, French, Portuguese, and a number of other languages also use the diaeresis, adding six more letters: ä, ë, ï, ö, ü, ÿ. (You may think of the " symbol as an "umlaut." The terms "umlaut" and "diaeresis" both usually refer to the same mark—a double dot—used to mean different things. The Unicode standard chose to use "diaeresis" as a generic name for the double dot.)

Most of the forms we've just discussed may not technically rise to the definition of "different letters." In most languages that use these marks, they're just considered marks that get applied to the letters: "é" is seen by French speakers as "e with an acute accent," for example. Unicode supports this type of thinking through the use of combining character sequences: there's a character code that represents the acute accent and can be used in conjunction with the other letters. However, because it's more

convenient for most kinds of processing, most character encoding standards just dedicate a new code point value to the accented forms of the various letters. Unicode allows this too in most cases, for backward compatibility.

Then there are additional characters that *really do* count as different letters in various languages. The letters ä and ö, for example, are thought of as marked forms of a and o by German speakers, but as entirely different letters by Swedish speakers. In addition, German adds the letter ß (“ess-zed” or “sharp s,” originally a ligature of “ss” or “sz”) to its alphabet. Danish, Swedish, and Norwegian add the letter å. Danish and Norwegian also add the letter ø. French adds the letters ç and œ. Icelandic adds æ (which also appears in Norwegian), ý, þ, and ð. The last two characters, called “thorn” and “edh,” were also used in Old English. Portuguese adds ã and õ, and Spanish adds ñ.

As you go further beyond the Western European languages, or back into historical writing system, or look at phonetic notations, the number of new letters and variants begins to increase much more quickly. The Unicode standard includes 819 Latin letters. We’ll take a quick look at all of them.

The Latin-1 characters

Unicode didn’t actually adopt the ASCII characters directly in unchanged, as I said above. They adopted in the characters from the ISO 8859-1 standard unchanged. ISO 8859-1 incorporated the 128 ASCII characters unchanged and added 128 more. Unicode preserves all 256 of these characters in their original order and with their original numeric code-point values (zero-padded to 21 bits). You can convert between ISO 8859-1 and UTF-16 or UTF-32 by zero-padding to the appropriate number of bits or truncating to 8 bits. (Conversion between Latin-1 and UTF-8 is, unfortunately, more complicated.)

ISO 8859-1 is the most well-known of the ISO 8859 family of character encoding standards first put together by the European Computer Manufacturers’ Association (ECMA) in the mid-1980s. ISO 8859 is a family of fourteen encoding standards that together provide a way to encode the characters of every (or at least almost every) language used in Europe. Each ISO 8859 standard is a superset of the ASCII standard: the 7-bit values from 0x00 to 0x7F represent the same characters they do in ASCII. The standards vary in their treatment of the other 128 values—the ones from 0x80 to 0xFF. Out of the 14 standards in the ISO 8859 family, ten of them are variations of the Latin alphabet (the other four are Greek, Cyrillic, Arabic, and Hebrew). These are commonly referred to as “Latin-1,” “Latin-2,” “Latin-3,” and so on through “Latin-9.” ISO 8859-1, also known as Latin-1, covers the languages of Western Europe. It’s probably the most widespread 8-bit ASCII-based encoding standard, and was the original base encoding for the HTML standard and various other Internet standards.

The Latin-1 standard was designed to cover the main Western European languages: English, French, Italian, Spanish, Portuguese, German, Dutch, Swedish, Danish, Norwegian, Icelandic, and Finnish. It includes all the extra letters discussed in the previous section—á, à, â, ã, ä, å, æ, ç, ð, é, è, ê, ë, í, ì, î, ï, ñ, ó, ò, ô, õ, ø, þ, ù, û, ü, ý, and ÿ—except for œ. Except for ß, which doesn’t have an uppercase form, and ÿ, Latin-1 provides both upper- and lower-case version of all these letters.

The Latin Extended A block

There’s considerable overlap between the character repertoires of the other ISO 8859 encodings, and none of them seemed to be used often enough to justify keeping their arrangements intact in Unicode the way ISO 8859-1’s arrangement was. The letters from the ISO Latin-2, Latin-3, Latin-4, and Latin-5 (i.e., 8859-2, 8859-3, 8859-4, and 8859-9 respectively) that didn’t also appear in Latin-1 were combined together into Unicode’s Latin Extended A block, which extends from U+0100 to U+017F. Duplicates and non-letters (in addition to characters that were already in the Latin-1 block)

were removed, and the remaining letters were arranged in alphabetical order. Unlike in the ASCII and Latin-1 blocks, case pairs are adjacent, rather than being arranged into two separate series, one for lowercase letters and one for uppercase letters.

In addition to the letters from the Latin-2, Latin-3, Latin-4, and Latin-5 standards, the designers of Unicode added other related letters that are used in the various European languages, including some overlooked characters from French (œ) and Finnish (š and ž), and a few others.

The Latin Extended A block brings to the party all the additional letters necessary to write Croatian (ć, č, đ, š, ž), Czech (č, ď, ě, ě, ř, š, ť, ů, ž), Esperanto (ĉ, ĝ, ĥ, ĵ, ŝ, ŭ), Estonian (š, ž), French (œ, Ÿ), Finnish (š, ž), Greenlandic (ĩ, κ, ũ), Hungarian (ő, ű), Latin (ā, ā, ē, ē, ī, ī, ō, ō, ū, ū), Latvian (ā, ē, ģ, ī, ķ, ļ, ņ, ō, r, ū), Lithuanian (ą, ę, ę, į, ū, ū), Maltese (ċ, ġ, ħ, ż), Polish (ą, ć, ę, ł, ń, ś, ź, ż), Romanian (ă, ș, ț), Sami (đ, ŋ, t), Slovak (č, ď, ľ, ň, ř, š, ť, ž), Slovenian (č, đ, ž), Turkish (ğ, ı, İ, ş), Welsh (ŵ, ŷ), and many other languages.

A couple other characters which aren't strictly necessary are also included: There is a single code (U+0149) representing n preceded by an apostrophe (ñ) that's intended for Afrikaans, but this is better represented by two separate code points: one for the apostrophe and one for the n (U+02BC U+006E). In Catalan, a centered dot is often placed between two *l*s (i.e., *ll*) to indicate they should be given the *l* sound instead of the *y* or *ly* sound that a double *l* normally has in Spanish and Catalan. The Latin Extended A block includes a single code point (U+013F and U+0140) representing an *l* followed by a centered dot (ł): this, too, is better represented using two code points (for example, U+006C U+00B7). Finally, a single code point representing *i* followed by *j* (*ij*; U+0132 and U+0133) is provided for use in Dutch, because these two letters often form a ligature in Dutch, but it's better to just use the regular *i* and *j* with a font that automatically forms the ligature. These characters are drawn from the ISO 6937 standard, a character encoding standard for videotext applications that never caught on in the computer world.

Also included is a code point (U+017F) representing the “long s” (ſ), an alternate glyph for the lowercase *s*. Up until fairly recently, this glyph was used at the beginning or the middle of a word, and what we now know as the lowercase *s* was used only at the end of a word. In handwritten text, this glyph just looks like an elongated *s* (ſ), and this form is still used as a character in the International Phonetic Alphabet (the integral sign in mathematics is basically this same glyph). But in printing, the bottom on the *s* was cut off, with a serif added, making it look like an *f* without the crossbar (or, in some fonts, an *f* with only part of the crossbar). This style was in effect during the American Revolution, which is why you see things like “Prefident” and “Congrefs” in the Constitution and Declaration of Independence. The long *s* is still used in the Fraktur typefaces often used to print German. Technically, this is a compatibility character, and one should represent the letter *s* using its normal code-point value and rely on the rendering engine or font to select the right glyph. Occasionally (such as in writing this paragraph), this code-point value is useful when you specifically want this glyph.

It's worth taking a second to look at one character in the Latin Extended A block. Consider for a moment U+0110 LATIN CAPITAL LETTER D WITH STROKE, Đ. It looks exactly like U+00D0 LATIN CAPITAL LETTER ETH, and, for that matter, like U+0189 LATIN CAPITAL LETTER AFRICAN D. Even though these three characters look exactly the same and never occur in the same language, they're not unified into a single code point value. This is because they have very different uses, and in particular because their lower-case forms are very different: The lower-case form of U+0110 is U+0111, which looks like this: đ. The lower-case form of U+00D0 is U+00F0, which looks like this: ð. And the lower-case form of U+0189 is U+0256, which looks like this: ɖ.

Another character in the Latin Extended A block that deserves a close look is the letter `ı`, U+0131 LATIN SMALL LETTER DOTLESS I, which is used in Turkish and Azeri (when Azeri is written with the Latin alphabet). Despite its name, this letter isn't really an `i` at all, but a completely different vowel. The fact that the new vowel's glyph was basically created by dropping the dot from the lowercase `i` poses a bit of a problem when it comes to deciding what the upper-case forms of these two letters should be. Logically, what we normally think of as the uppercase `I`, which doesn't have a dot, becomes the uppercase `ı` in Turkish or Azeri. To distinguish it, the uppercase form of `i` has a dot in these two languages: `İ`. This form is represented in Unicode as U+0130 LATIN CAPITAL LETTER I WITH DOT ABOVE.

This poses a problem when mapping between uppercase and lowercase. Now you need to know the language. In every language but Turkish and Azeri, `i` maps to `I` when going from lowercase to uppercase, but in Turkish and Azeri, `i` maps to `İ`. In every language but these two, `I` maps to `i` when going from uppercase to lowercase, but in these two languages, `I` maps to `ı`. It could be argued that it would make more sense to just have U+0049 LATIN CAPITAL LETTER I have a different glyph (i.e., `Ī`) in Turkish and Azeri and introduce two new codes that always represent the upper- and lower-case forms of the letter `i`, but the ISO 8859 encodings opted for the alternative case-mapping behavior instead, and Unicode does it the same way out of respect for existing practice.

The Latin Extended B block

Things start to get more exotic, and arguably somewhat more disorganized, with the Latin Extended B block. Latin Extended B, which extends from U+0180 to U+024F, contains a large variety of miscellaneous Latin letters that come from a wide variety of sources. They include not only more letter-diacritic combinations, but several letters with multiple diacritics, letters with interesting typographic changes, such as upside-down or mirror-image letters or letters with elongated strokes, and whole new letters, some borrowed in or adapted from other alphabets, some old letters that have fallen into disuse, and some totally new coinages.

Among the characters in this block are letters used to write various African and Native American languages (such as U+0181, `Ɓ`; U+0188, `Ɔ`; U+0194, `Ƴ`; U+01B7, `Ʒ`), various less-common European languages, including the Celtic languages covered by the ISO 8859-14 ("Latin-8") standard, archaic characters formerly used to write various languages, including such Old English characters as wynn (`ƿ`) and yogh (`ƿ`), letters used in transcribing text in languages that don't normally use the Latin alphabet, and various phonetic symbols that aren't part of the current International Phonetic Alphabet.

Unicode's designers have tried to keep this section in alphabetical order and to keep case pairs together, but this block has ended up arranged into several alphabetical series, reflecting groups of letters added at different times, and sometimes case pairs have wound up widely separated, as sometimes a character added at one time for one purpose where it doesn't have a case partner turns out to also be used for another purpose where it does have a case partner. For example, there are a lot of characters that were originally added as part of the International Phonetic Alphabet, which is always lower-case, that are also used in the normal writing of some languages, where they have uppercase counterparts. In these situations, the lowercase version of the letter is in the IPA block (see below) and the uppercase version is in this block.

Among the characters in this section are the various characters used in the Pinyin system of transcribing Chinese into Latin letters, various characters for Zhuang (the largest non-Chinese language spoken in China—it's actually related to Thai), less-common Slovenian, Croatian, and Romanian letters, letters for Livonian (a rapidly dying minority language spoken in Latvia and related to Estonian), and some of the letters necessary for Vietnamese (this block rounds out the basic

Vietnamese alphabet—with the characters in this block, plus the ones in the previous blocks, you can write Vietnamese, needing the combining marks only to add tone marks to the vowels).

There are a few characters in this block that deserve a closer look. From U+01C4 to U+01CC and U+01F1 to U+01F3 are the Croatian digraphs. Serbo-Croatian is written using both the Latin and Cyrillic alphabets: The Serbs and Montenegrins use the Cyrillic alphabet, and the Croats and Bosnian Muslims use the Latin alphabet, although most Serbo-Croatian speakers can read both alphabets. Serbo-Croatian spelling is standardized so that there's a simple one-to-one mapping between the letters in the Serbian (i.e., Cyrillic) alphabet and the Croatian (i.e., Latin) alphabet. Because of this, converting from the Serbian spelling of a word to the Croatian spelling is a simple affair.

One thing that complicates this is that three of the letters in the Serbian alphabet map to pairs of letters in the Croatian alphabet: the letter Љ becomes dž, the letter Ђ becomes lj, and the letter Ћ becomes nj. It's, of course, simpler for a computer to switch between the Serbian and Croatian alphabets if the letter pairs in the Croatian alphabet are represented using a single code point, which was done some vendor standards. For compatibility, single code points representing these letter pairs (or “digraphs”) are also included in the Latin Extended B block, although it's generally better to use the separate letters instead.

Each of these letter pairs comes in three flavors: one with two capital letters (the “uppercase” version), which is intended to be used when all of the letters in the word are capitalized, one with two small letters, and one with the first letter capitalized and the second letter small (the “titlecase” version), intended to be used when only the first letter of the word is capitalized. Both the full-uppercase and titlecase versions of the letter map to the uppercase version in the Serbian alphabet.

(It's interesting to note that the ij digraph in the Latin Extended A block, used in Dutch, doesn't have a titlecase version. This is because when it appears at the beginning of a word, both the I and the J are actually capitalized, as in “IJssel Meer.”)

Another interesting letter to take note of is U+01C3, the letter !. No, this isn't an exclamation point. It's a letter. It represents the retroflex click, the sound you make by pushing your tongue against the back of your upper gums and pulling it away quickly. It's a letter (and a sound) in several African languages. This is an example of how the principle of unification is applied: This character has exactly the same glyph as the exclamation point, but completely different semantics leading to its being treated differently by various processes operating on text (as a letter rather than as a punctuation mark). Because Unicode encodes semantics rather than appearance, these two uses aren't unified into a single code-point value.

This is also an example of one of the pitfalls you can run into using Unicode: because they look the same, it's very likely that some users, depending on the software they're using, will type the exclamation point where they really mean to type the retroflex-click character, and so processes operating on the text may actually need to detect and account for this erroneous use of the exclamation point. There are numerous groups of Unicode characters that potentially have this problem.

The Latin Extended Additional block

The Latin Extended Additional block can be considered a block of compatibility characters, and so was positioned at the end of the General Scripts Area, rather than with the other Latin blocks at the beginning. It contains a bunch of extra accented forms of the Latin letters that are used in various languages. These characters resulted from the merger with ISO 10646, and, except for the ones that obviously come from ISO8859-14, no records were kept on their original sources.

There aren't any characters in this block that can't be represented using a combining character sequence—new Latin characters that *can't* be represented using a combining character sequence are added to the end of the Latin Extended B block instead. Thus, none of the characters in the Latin Extended Additional block are strictly necessary. They're here for compatibility with existing standards or to make life easier for implementations that can't support combining character sequences.

The Latin Extended Additional block is arranged into two alphabetical series with case pairs adjacent to each other. The first series contains a miscellany of accented Latin letters. The second series is dedicated specifically to Vietnamese. All of the basic letters of the Vietnamese alphabet (a ă â b c d đ e ê g h i k l m n o ô õ p q r s t u v x y) are encoded in the other Latin blocks, but to put tone marks on the vowels, you have to use combining character sequences. (Vietnamese is a tonal language—the way in which a word is said, what English speakers think of an intonation, is part of the pronunciation in Vietnamese and many other Asian languages. Different intonations on the same syllable actually make different words. Tone marks are used to distinguish between intonations. There are five tone marks in Vietnamese, and they're all attached to the vowel: ắ ằ ẳ ẵ ặ) The Latin Extended Additional block includes all of the Vietnamese vowels with each of the tone marks applied to them.

The International Phonetic Alphabet

The International Phonetic Alphabet, or “IPA,” is an international standard system of symbols for notating the sound of human speech. Real alphabets are basically phonetic in character, but tend over time to acquire non-phonetic properties—the pronunciation of a word changes over time, but its spelling doesn't. The phonetic dimension is lost or obscured, but the spelling may still tell you something about the meaning or etymology of the word in question. IPA doesn't mess with any of this; it's simply a system for linguists, speech pathologists, dialogue coaches, and others interested in the mechanics or sound of speech to accurately notate how somebody says something.

IPA was first developed in 1886 and has undergone many revisions since then. It's based on the Latin alphabet, but includes lots of modified forms of letters. Sometimes, glyphs that would normally be thought of as simply typographic variants of the same letter, such as a and α, are different symbols in IPA, representing different sounds. There generally isn't a systematic connection between the shapes of the characters in IPA and the sounds they represent, other than that the symbol for a particular sound is usually based on the shape of some Latin letter that has that sound or a similar sound. In a lot of cases, this means that lots of symbols based on a particular letter arise: sometimes these are variant type styles, sometimes the letter is reversed or turned upside-down, and sometimes the basic shape of the letter is altered: for example, a stroke is elongated or a tail or hook of some kind is added. Some Greek letters are also borrowed into the IPA.

IPA also isn't a completely exact method of notating speech sounds. One of the guiding principles behind it was that if you had a collection of similar sounds, they'd only be given different symbols if some language treated them as different sounds. For example, consider the aspirated and unaspirated versions of the *t* sound: both sounds are formed by pressing your tongue against the back of your upper jaw, applying some air pressure from your diaphragm, and then letting the air escape by pulling the tongue away. The aspirated form allows some extra air to escape, producing a more explosive sound. The *t* at the beginning of *toy* is an aspirated *t*, while the *t* in *stamp* is an unaspirated *t*. English doesn't make a distinction between these sounds; we use one in certain places, such as at the beginnings of words, and the other in other places, but we consider both to be the same sound. Quite a few other languages have both of these sounds and consider them to be different sounds, represented with different letters, so the IPA gives them different symbols. Because the IPA covers so many languages, this then gives us symbols that allow us to make a useful distinction between the

initial *t* in most English words, which is aspirated, and the initial *t* in most Italian words, which isn't, even though both languages consider themselves to have just one *t* sound.³²

These types of distinctions aren't always available in IPA, leading to the use of all kinds of auxiliary marks to help clarify various distinctions. Specialized symbols for things like disordered speech are also used by specialized users. The bottom line is that there are dozens of ways of transcribing the same utterance using IPA, and which one is used depends on who's doing the transcribing and why they're doing the transcribing.

In Unicode, a good many of the IPA characters are encoded in the various Latin blocks we've already discussed, or in the Greek blocks. Many are simply unadorned Latin letters, and many others are variants that were already being used as letters in some language. Unicode doesn't make a distinction between regular letters and IPA symbols, although it doesn't make a distinction between IPA symbols and other kinds of symbols, such as the distinction between the “esh” character—the elongated *s* used to represent the *sh* sound in *sheep* (ʃ)—and the integral sign, which looks just like it (this character is different from the long *s* we discussed earlier, which has the same shape in handwritten writing, but a different shape when printed (I)). This means there's no IPA block containing all the IPA characters. Instead, there's an IPA Extensions block, running from U+0250 to U+02AF, which contains just the IPA symbols that hadn't already been encoded somewhere else.

IPA is a caseless script—effectively, all IPA characters are lowercase. Even in cases where the capital form of a letter is used to make some kind of phonemic distinction with the lowercase form, a small-cap form (a version of the capital letter altered to fit the normal dimensions of a lowercase letter) is always used. At various times, IPA characters used in transcription become part of the normal alphabet for writing some language (this has happened with a number of African languages, for example). In these cases, the IPA character picks up an uppercase form when it passes into common usage. Since the uppercase forms aren't really IPA characters, they're encoded in the Latin Extended B block rather than with their “mates” in the IPA Extensions block. The Unicode standard includes cross-references for these case pairs.

Because the IPA characters are unified with the normal Latin and Greek letters, it doesn't make a lot of sense to give the IPA Extensions block an ordering based on the character's pronunciations. Instead, they're arranged in rough alphabetical order according to the regular Latin letters they most resemble. As with the Latin Extended B block, the IPA Extensions block includes a couple of alphabetic series for historical reasons.

One essential feature of IPA is the use of various diacritical marks and other markings to clarify the meanings of the primary “letter” symbols. These extra marks aren't encoded into the IPA Extensions block either. Instead, they're kept with the other diacritical marks in the Combining Diacritical Marks and Spacing Modifier Letters blocks, which are discussed in the next section.

Diacritical marks

One of the principles Unicode is built on is the idea of dynamic composition—you can represent the marked form of a letter using two code points, one representing the letter followed by another one representing the mark.

³² This is perhaps not exactly the best of examples: The IPA doesn't really have two different symbols for these two sounds. It's got one symbol—not too surprisingly, the letter *t*—for the “unvoiced alveolar plosive.” It then adds a second symbol—a superscript *h*—to show aspiration. So the *t* at the beginning of *toy* is actually represented in IPA with two symbols: *t^h*.

Quite a few of the letters in the various Latin blocks are marked forms—base letters with some kind of mark applied. All of these characters can be represented using two code points. To make this possible, Unicode includes a whole block of characters—the Combining Diacritical Marks block, which runs from U+0300 to U+036F. The characters in this block are special because they specifically have combining semantics. They always modify the character that precedes them in storage. This means that these characters are generally never to be considered alone, except for the purpose of figuring out that they have combining semantics. Instead, the combination of a combining mark and the character it follows are to be treated as a *single character* by any code processing the text. This single character has the combined semantics of its two constituent characters and is called a *combining character sequence*. More than one combining mark can be applied to the same character—a sequence of combining marks are all considered to combine with the character that precedes the first combining mark. If several marks in the sequence are supposed to attach to the same part of the base character, the one that occurs first in storage is the one drawn closest to the base character, with the others radiating outward.

Technically, combining marks aren't “characters” at all. Instead, they're code points that modify the semantics of other characters in some prescribed way. But it's simpler to just call them “characters” like everything else.

For an exhaustive treatment of combining character sequences, see Chapter 4.

There are combining marks scattered all throughout the Unicode standard. Generally speaking, combining marks used with only one particular script are included in that script's encoding block. What's encoded in the Combining Diacritical Marks block are combining marks that are used primarily with the Latin alphabet or with IPA, and combining marks that are commonly used with multiple scripts.

The great thing about combining character sequences is that they allow you to encode characters that hadn't been previously foreseen by Unicode's designers, and that they can cut down on the number of code point values needed to usefully encode some script. A huge variety of marked Latin letters were given their own code point values, but there are still others that are used in real languages that have to be represented using combining character sequences (and there will be more over time, since the need to maintain the stability of Normalized Form C makes new precomposed characters a lot less useful). The combining marks also allow a certain amount of “future-proofing” against new letter-mark combinations that might be coined in transcriptions of languages not previously transcribed.

The downside of combining character sequences is that they can complicate handling of Unicode text. They fix it so that a single character can be represented using multiple code points, which complicates things like counting the number of characters in a field or document. Also, since every letter-mark combination that gets its own code point value is considered simply an alternative representation of the multiple-code-point representation, Unicode implementations generally have to be prepared to see a particular letter-mark combination encoded either way and to treat them the same.

Of course, one way of dealing with this is to simply decline to support the combining marks and combining character sequences. This is perfectly legal, and at least with the Latin alphabet, can be done with little loss of functionality, since most of the letter-mark combinations in common use are also given single-code-point representations. The ISO 10646 standard, which is code-point-for-code-point compatible with the Unicode standard, actually defines three implementation levels, indicating an implementation's support for combining character sequences. A Level 1 implementation supports only precomposed characters; a Level 2 implementation only supports the Indic vowel signs (see Chapter 9), and a Level 3 implementation supports both the precomposed characters and all of the combining characters. So one way of dealing with the combining-character problem is to declare

yourself a Level 1 implementation of ISO 10646. Unfortunately, this really isn't feasible with a lot of scripts other than the Latin alphabet, as they often provide no way of representing certain important letter-mark combinations other than combining character sequences (the Arabic alphabet is one example).

For the Latin letters, at least, an implementation can also support the combining marks without requiring special font technology by mapping the text to Normalized Form C before passing characters through to the font engine. Thus, if you have a font that has a glyph for the single-code-point representation of é, you can still draw é even when it's represented using two code points.

Of course, the “right” thing to do is to have a font engine that understands the combining character sequences and does the right thing automatically. It's important to note that this isn't always as straightforward as it might first appear. This is because rendering a combining character sequence isn't always a simple matter of taking a glyph representing the letter and a glyph representing the mark and drawing them positioned appropriately next to each other. Sometimes it's more complicated.

The Unicode standard specifically calls out a few examples. For instance, several Eastern European languages use a caron or hacek (basically, a little v) over some of their consonants (č š ž), but this doesn't work well when applied to a character with an ascender, such as a d or l. In these cases, the caron can extend up outside of its own line and collide with the previous line of text. Or you can get inconsistent and ugly line spacing trying to avoid this. One solution is to move the caron to one side to get it out of the way, but the solution you usually see is to turn it into an apostrophe and move it to the side. Thus, when you combine U+0064 LATIN SMALL LETTER D (d) with U+030C COMBINING CARON (ˇ), you get this: d'. It's important to note that this isn't the same thing as d followed by an apostrophe, and shouldn't compare equal. This could potentially be a problem if the user types a d followed by an apostrophe, although this probably wouldn't happen with a user using an appropriate keyboard for the language.

Interesting things also happen with the marks that attach to the underside of a letter, such as the cedilla (,) and ogonek (˙). These marks change shape and morph into various types of hooks and things depending on the font or the language. They may also attach to various parts of the letter depending on where they're a good place to attach.

In particular, the cedilla will often be drawn as a comma underneath the letter instead of a cedilla. This is especially common with letters such as k or n where there's no good spot for the cedilla to attach, but also occurs with letters like s and t in some languages. Whether the form with the cedilla or the form with the comma is used depends heavily on the language and on the design of a particular font. Unicode provides a couple of explicit characters with a comma below the letter instead of a cedilla, but their use is generally discouraged. In the same way, using combining character sequences with U+0326 COMBINING COMMA BELOW instead of U+0327 COMBINING CEDILLA in languages where the cedilla can be drawn either way might cause problems. You could end up with the same glyph represented two different and incompatible ways internally without extra work to account for the difference. However, properly-designed input methods should generally prevent problems with this.

An especially interesting phenomenon happens when the lowercase g is combined with a cedilla, as happens in Latvian. Because of the descender on the g, there's no room to draw the cedilla under the g without it crashing into the next line of text, so most fonts actually draw it *on top of* the g as an upside-down comma (ǧ). In this case, a mark that normally attaches to the bottom of the letter moves to the top. Again, there's a U+0312 COMBINING TURNED COMMA ABOVE combining mark that can be used to produce this glyph. Without extra Latvian-specific code that treats a combining

character sequence using this mark the same as a combining character sequence using the cedilla, you can have different and incompatible representations for the same glyph.

The lowercase *i* also has interesting behavior. Generally, the dot on the *i* disappears when a diacritical mark is drawn on top of the *i*: *í î ï*. If you were to apply a diacritical mark to the Turkish dotless *i* (*i*), you'd get the same appearance but a different semantic (and unequal representations). In particular, if you take U+0131 LATIN SMALL LETTER DOTLESS I and follow it with U+0307 COMBINING DOT ABOVE, you don't get the regular lowercase *i*. You get something that looks just like it, but it's actually the Turkish *ı* with a "dot" diacritical mark on top. (Fortunately, you never see diacritical marks used with this letter in practice, so this isn't a big problem.)

In some languages, the dot on the *i* actually sticks around when an accent or other mark is applied. Fonts for languages where this happens, such as Lithuanian, might just do this, but the representation can be forced by explicitly using U+0307 COMBINING DOT ABOVE. (This gives rise, by the way, to a set of special case-conversion rules for Lithuanian that take care of adding and removing U+0307 as needed.)

Some of the combining marks actually overlay the letter they're applied to. Examples include U+0335 COMBINING SHORT STROKE OVERLAY and U+0337 COMBINING SHORT SOLIDUS OVERLAY. Exactly where these marks get drawn depends on the letter. A stroke might draw across the middle of an *l*, but across the vertical stroke of a *d* above its bowl.

The precomposed forms of letters with overlay marks, such as *ø*, *ł*, or *đ*, however, don't decompose to representations using the combining marks. Thus we have another case where you can get similar-looking glyphs with incompatible representations. Again, in practice, this shouldn't be a problem, since (for example) a Danish keyboard implementation would produce the proper representation of the letter *ø* and you wouldn't have to worry about the other representation that might look the same.

Interesting things can also happen sometimes when multiple combining marks are applied to the same letter. In Vietnamese, for example, there are a few vowels with circumflexes. When a tone mark is also applied to one of these letters, it appears *next to*, rather than on top of, the circumflex. The grave accent appears to the left of the circumflex (*à*), and the acute accent and hook appear to the right of the circumflex (*á ă*). The grave and acute accents actually touch the circumflex in many fonts. This breaks the general rule that multiple accents stack. These examples are Vietnamese-specific, and other examples where multiple accent marks combine typographically in interesting ways are also language-specific, but they're cases, again, where more work than simple accent stacking might be necessary to properly display multiple combining marks (the decomposed representation for the Vietnamese vowels, by the way, always has the tone mark coming last).

Another interesting note concerns the double-acute combinations in Hungarian (*ő, ű*). While each of these letters looks like it has two acute accents above it, you *don't* get these glyphs by combining the base letter with two successive combining acute accents. The double acute is a separate combining mark and isn't the same as two single acutes.

One more issue is what happens to accent marks on capital letters. Again, you have the problem of the accent sticking up above the character stream and crashing into the previous line of text. Different typefaces handle this in different ways: by making the base letter shorter, by making sure the line spacing is wide enough to accommodate any accents, by altering the accent in some way (in some German fonts, for example, the two dots of the umlaut may appear on either side of an *A* or *O*) or by omitting the accent on capital letters. The last solution is fairly common in French typography, and may lead to situations where accents in the character storage don't actually appear in the rendered text. This is okay.

All of the above examples generally qualify as language-specific, and should be dealt with properly by software or fonts designed specifically for text in that language. There's no guarantee, however, that generalized Unicode software (in particular, fonts designed to display any Unicode character) will do everything right for every possible language Unicode can be used to represent.

Just because a Unicode implementation supports combining character sequences doesn't mean that it will properly draw *every possible* combining character sequence, or even every possible sequence involving characters it says it supports. Some are nonsensical, or at least extremely unusual, and there's no guarantee (and no requirement) that some of the language-specific tricks discussed above will always work. The designers of Unicode have thought of all this stuff, but when it comes to actual implementations, your mileage may vary.

Isolated combining marks

So what if you *want* to show the combining diacritical marks in isolation? What if you're writing a book like this and you want to talk about the marks without attaching them to letters? Well, the Unicode convention is to attach them to a space. In other words, if you want to see `ˆ` all by itself, you can represent it as U+0020 SPACE followed by U+0301 COMBINING ACUTE ACCENT.

(Actually, the space is the “official” way of representing an isolated combining mark, but certain other characters will work too. Control characters and paragraph separators will also work, since you can't attach combining marks to them.)

Now none of this accounts for the various diacritical marks that are encoded in the ASCII and Latin-1 blocks. These include U+005E CIRCUMFLEX ACCENT (^), U+005F LOW LINE (_), U+0060 GRAVE ACCENT (`), U+007E TILDE (~), U+00A8 DIAERESIS (¨), U+00AF MACRON (-), U+00B4 ACUTE ACCENT (´), U+00B8 CEDILLA (¸), as well as a bunch of other characters that do double duty as regular punctuation and diacritical marks. Whether these characters should have combining semantics is ambiguous: there are older teletype systems where they did: some of them functioned as dead-key characters—after imprinting one of these characters, the carriage wouldn't advance. This meant you could get, for example, ä by sending the umlaut followed by the a. And on all teletype machines, you could make any character a combining character by following it with a backspace. The backspace would back the carriage up to the position of the character you wanted to add the mark to (or back to the mark you wanted to add a character to). However, this usage died out with the advent of the VDT and these characters lost their combining semantics.

The designers of Unicode opted to explicitly give all these characters *non*-combining semantics: the characters from Latin-1 have a compatibility decomposition to a space and the combining form of the same mark. (The ones from ASCII don't, so that ASCII text isn't affected by Unicode normalization.)

Some of the characters in the ASCII and Latin-1 block did double duty as diacritical marks and other things. For example, the tilde (~), originally intended as a diacritical mark, has also gotten used as a swung dash and as a mathematical operator (it's a logical-not operator in a lot of programming languages, for example), and the glyph that's used for this code point value is often centered in the line rather than raised. The caret (^) has similarly done double duty as a circumflex accent and various other things, including a mathematical operator (in various programming languages, it's either a bitwise-exclusive-or operator or an exponentiation operator). And the degree sign at position U+00B0 (°) has done double duty as the ring in å. For all of these cases and a few others, the Spacing Modifier Letters block includes clones that have unambiguous semantics—they're only diacritical marks and they always have non-combining semantics.

Spacing modifier letters

There is a class of diacritical marks that appear to the right of the characters they modify instead of above, below, or through them. Instead of treating these as combining marks, Unicode treats them as ordinary characters—they have non-combining semantics. Since they modify the pronunciation or other properties of letters and are therefore used in the middle of words, Unicode classifies them as “letters” and assigns them to the “modifier letter” (“Lm”) general character category. These characters are grouped together in the Spacing Modifier Letters block, which runs from U+02B0 to U+02FF.

Most of the characters in this block come from the IPA or from some other phonetic-transcription system. Some are just superscript letters, such as the superscript *h* for aspiration or the superscript *j* for palatalization, but there are a lot of other types of marks. These include various marks for noting accented syllables, long vowels, glottal stops, tones, and so on. As mentioned in the previous section, this block also includes clones of various characters from the other blocks that did double duty as diacritical marks. These clones are always diacritical marks and always have non-combining semantics.

There are two characters in this block that deserve special mention: One is the rhotic hook (U+02DE, r^{h}). This character basically represents the “r” sound in words like “car,” “fear,” and “purr.” The “r” sound in these words isn’t really a separate sound, but a modification to the vowel. This is called rhotacization, and in the International Phonetic Alphabet, it’s represented by attaching this little pigtail thing (the rhotic hook) to the right-hand side of the symbol representing the vowel (the American pronunciation of “pear” in IPA is $\text{p}\text{ɚ}^{\text{h}}$).

Unlike the other characters in the Spacing Modifier Letters block, the rhotic hook doesn’t merely appear after the character it modifies; it *attaches to it*. The shape of the character changes to include the pigtail. Because it attaches to the right-hand side of the letter (i.e., the trailing edge in left-to-right layout), it’s not considered to have combining semantics. A character is only considered a combining character if it attaches to some side of the base character other than the trailing edge in its primary layout direction. If it attaches to the trailing edge in the primary layout direction (the right-hand side in left-to-right text), it’s merely considered to form a ligature with the character it follows. This is a subtle (and, in my humble opinion, rather meaningless) distinction, but it’s important and we’ll see lots of examples of it when we look at the Indic scripts in a few chapters.

The other character that bears a closer look is U+02BC MODIFIER LETTER APOSTROPHE. Yes, this is our old friend the apostrophe. There are a whole mess of characters in Unicode that look like the apostrophe, but this one is the real thing. It gets used for lots of things. In IPA, for example, it’s used to represent the ejective consonants, which are formed by closing your throat and using just the air in your mouth. Most of the rest of the time, it’s used to represent the glottal stop. (The glottal stop is a break in the sound caused by momentarily closing your glottis [the back of your throat]. It’s best thought of as the “sound” between the “uh”s in “uh-uh”.) This is how it’s used, for example, in “Hawai’i.”

One thing this character isn’t is a single quotation mark. Unicode has a whole host of other characters for quotation marks, and we’ll get to those in Chapter 12. It also isn’t the apostrophe used in contractions like “ma’am” or “isn’t.” This character is generally considered a punctuation mark, rather than a letter. You use the closing-single-quote character to represent this apostrophe as well.

It’s also worth calling attention to the “tone letters” in the Spacing Modifier Letters block. The IPA has five special characters for representing five distinguishable levels of tone in tonal languages (̂ ̃ ̄ ̅ ̆); these five characters are represented in the Spacing Modifier letters block. The IPA also specifies a bunch of forms for representing tonal contours (e.g., rising, falling, falling then rising,

etc.); these aren't given separate codes in Unicode. Instead, these characters are treated as ligatures: If you follow a low-tone character (◡) with a high-tone character (◑) in storage, for example, the two code points are supposed to be drawn as a single rising-tone character (◡◑ **[this is sort of right—see picture in TUS, p. 178]**). Exactly how the various level-tone characters are to be combined to create various contour-tone characters isn't rigorously defined in the Unicode standard, but these forms aren't all laid out in the IPA charts, either; they're apparently generally ad-hoc forms whose generation is more or less self-explanatory.

The Greek alphabet

The Greek alphabet is the ancestor of both the Latin and Cyrillic alphabets. It first developed somewhere around 700 BC, and had largely acquired its present form by about 400 BC. Prior to 400 BC, there were a number of regional variants of the Greek alphabet; the Latin alphabet evolved from one of these, which is why some of the Latin letters look different from their Greek counterparts. By 400 BC, all Greek-speaking peoples had standardized on a single version of the alphabet, which has come down to us more or less intact.³³

The Greek alphabet was originally written right to left or *boustrophedon*, but by 500 BC was always written left to right. As with the Latin alphabet, the lower-case letters were more recent developments, having appeared somewhere around AD 800. Diacritical marks and spaces developed around the same time, but took a long time to become standardized.

As anyone who's ever pledged a fraternity knows, the Greek alphabet has twenty-four letters, in both upper and lower case:

Α	Β	Γ	Δ	Ε	Ζ	Η	Θ	Ι	Κ	Λ	Μ
Ν	Ξ	Ο	Π	Ρ	Σ	Τ	Υ	Φ	Χ	Ψ	Ω
α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ
ν	ξ	ο	π	ρ	σ	τ	υ	φ	χ	ψ	ω

As with the Latin alphabet, the letters of the Greek alphabet don't generally form ligatures and generally don't change their shape depending on context. The one big exception to this is the lower-case sigma, which has a different shape (σ) at the beginning of a word, or in the middle of the word, than it does at the end of a word (ς).³⁴

The Greek alphabet used to have twenty-eight letters; four (stigma (Ϛ), digamma (Ϝ), koppa (Ϟ), and sampi (Ϸ)) are now obsolete. But for a long time after these letters fell into disuse, the Greeks were still using the letters of their alphabet as digits: the first nine letters of the alphabet represented values from 1 to 9, the next nine letters represented the multiples of 10 from 10 to 90, and the third nine letters represented multiples of 100 from 100 to 900. Since you needed twenty-seven different letters to represent all the numbers up to 999, digamma, koppa, and sampi continued to be used in numerals

³³ My sources for most of the historical information on the Greek alphabet, and on how it works, is Leslie Threatte, "The Greek Alphabet," in *The World's Writing Systems*, Peter T. Daniels and William Bright ed., Oxford: Oxford University Press, 1995, pp. 271-280, and Robert K. Rittner, "The Coptic Alphabet," *op. cit.*, pp. 287-290.

³⁴ Actually, in some Greek writing styles, the small letter sigma only has one form, and it doesn't look like either of these.

long after they ceased to be used as letters. (A prime mark (´) appeared after a sequence of letters used as a numeral to distinguish it from a word.) Today the Greeks use the same digits we use, but you still see the old alphabetic numbers used in numbered lists and outlines.

A number of diacritical marks are used with the Greek alphabet. There were three accent marks: the *oxia*, or acute accent (´), the *varia*, or grave accent (`), and the *perispomeni*, or circumflex accent (̂— as this example shows, the *perispomeni* often looks like a tilde rather than a circumflex). Researchers think these marks originally indicated variations of tone, but this usage died out a very long time ago, and they have generally persisted just as historical quirks of spelling. Modern scholars generally treat the acute and circumflex accents as stress accents and ignore the grave accent. Graphically these symbols are somewhat different from their counterparts used with the Latin letters—the acute and grave accents are generally drawn with a much steeper angle than is used with Latin letters, and the *perispomeni* appears, depending on font design, as a circumflex, tilde, or inverted breve.

Two other marks, the so-called “breathing marks,” evolved to represent the “h” sound. In some locations, the letter eta (H) was a consonant and represented the “h” sound, but in other places, this letter was used as a long E, and this usage eventually won out, leaving no symbol to represent the “h” sound.

It appears the “h” sound only appeared at the beginning of a word, and only at the beginnings of words that began with vowels. A mark that looks like a comma (´), the *psili* or “smooth-breathing mark,” was used to indicate the absence of an “h” sound, and another mark that looks like a reversed comma, the *dasia* or “rough-breathing mark” (͂), was used to indicate the presence of the “h” sound. In transliterating ancient texts, the rough-breathing mark is transliterated as an h, and the smooth-breathing mark is ignored. As with the accent marks, the breathings persist as historical quirks of spelling—the initial “h” sound has long since disappeared from modern Greek pronunciation.

In the early 1970s, a system of spelling known as “monotonic Greek” developed that dispensed with all the unused diacritical marks (the old system is now known as “polytonic Greek”). The acute accent, called a *tonos* (which just means “accent”) is now used to mark stressed syllables, and the other accents and breathing marks are no longer used. This system has been widespread in Greek typography since the late 1970s, and was officially adopted by the Greek government for teaching in the schools in 1982.

Another mark, the *koronis*, which looks just like the smooth-breathing mark (i.e., like an apostrophe), was used in ancient Greek to indicate the elision of two vowels across a word boundary. This mark also isn’t used anymore.

The diaeresis (also called the *dialytika* or *trema*) is used in Greek the same way it’s sometimes used in English: to indicate that two vowels in a row are to be treated as separate vowels (i.e., with a syllable division between them) rather than as a diphthong. To use an example from English, you’ll sometimes see “cooperate” spelled with a diaeresis on the second o: coöperate. The diaeresis indicates that the first four letters are pronounced as two syllables, rather as a single syllable rhyming with “loop.” This kind of thing is optional (and rather rare) in English, but mandatory in Greek.

With all these different marks, you’d often see Greek vowels with multiple marks. If a vowel has both a breathing mark and an acute or a grave, the breathing mark appears to the left of the accent (͂´). An acute or grave with a diaeresis appears above the diaeresis or between the dots of the diaeresis (̂̄). The circumflex accent always appears above the other marks (̂̄). When the marks occur on a diphthong, the convention is for the marks to appear on the second vowel rather than the first. When applied to capital letters, the marks appear to the letter’s left rather than over it.

Finally, there's the curious case of the iota subscript, or *ypogegrammeni*. Beginning in the thirteenth century, scribes began writing the three diphthongs αι (ai), ηι (ei), and ωι (oi) with the iota *under*, rather than following, the other vowel: α̣ η̣ ω̣. This practice has stuck for writing ancient texts, but isn't used in modern Greek. The iota subscript was only used when both letters would be lower-case. When the first letter is capitalized, the iota moves to its customary place to the *right* of the other letter: the "iota-adscript" or *prosgegrammeni* (Αι Ηι Ωι). **[In a good Greek font, I think the iota subscript is smaller than a normal iota, rather than looking just like it, as in the preceding example.]** When all the letters in the word are capitalized, the iota turns into a regular iota (ΑΙ ΗΙ ΩΙ).

Punctuation in modern Greek is basically the same as for the languages using the Latin alphabet, with two main exceptions: The Greek question mark (or *erotimatiko*) looks like a Latin semicolon (;), and the Greek colon or semicolon (the *ano teleia*) is a single centered dot (•).

The Greek block

In Unicode, most of the characters needed to write modern monotonic Greek are encoded in the Greek block, which runs from U+0370 to U+03FF. This block is basically based on the ISO Greek standard (ISO 8859-7), which is basically the same as the Greek national standard. The characters that occur in ISO 8859-7 occur in Unicode in the same relative positions, facilitating conversion between the two standards. The Greek block in Unicode has been augmented with some additional characters, and a few characters from the ISO standard have been unified with characters in other Unicode blocks.

Modern Greek uses punctuation and digits from the ASCII and Latin-1 blocks, and also uses combining diacritical marks from the Combining Diacritical Marks block. As with the Latin blocks, the Greek block includes a few diacritical marks, and the Unicode standard declares the diacritical marks in the Greek block to have non-combining semantics. The combining versions of these marks are in the Combining Diacritical Marks block: the *oxia* and *varia* are represented using the regular acute- and grave-accent characters, the *dialytika* with the regular diaeresis character, and the *psili* and *dasia* are represented using the regular combining-comma-above and combining-reversed-comma-above characters, which are also used for other uses.

The *perispomeni* isn't unified with any other existing combining marks; since its form can be either a circumflex accent or a tilde, it didn't make sense to unify it with either. The *koronis* and the *ypogegrammeni* have their own code-point values in the Combining Diacritical Marks block as well. Finally, the Combining Diacritical Marks block includes a combined *dialytika-tonos* character whose use is now discouraged. Today, use the separate acute-accent and diaeresis combining characters.

The Greek block also includes numeral signs and the Greek question mark and semicolon (the question-mark and semicolon characters are also here for historical reasons—they're now considered equivalent to the regular semicolon character [U+003B] and the middle-dot character[U+00B7]).

The Greek block contains not only the basic upper- and lower-case Greek Letters, but precomposed forms of the vowels with the *tonos* and *dialytika* marks. Keeping with existing practice, the initial/medial and final forms of the lowercase sigma are given separate code-point values, even though this violates the general rule about unifying glyphic variants. It'd be reasonable to have a system always use just one of these codes in storage and draw the correct glyph based on context, but the forms are encoded separately for backwards compatibility with existing encodings and out of a desire for consistency with existing practice. Because of this, code that maps from upper case to lower case has to be careful to take into account the different forms of the lowercase sigma.

The Greek alphabet

The Unicode Greek block also adds the obsolete Greek letters and some alternate forms of the modern Greek letters. Normally these alternate forms would just be considered glyphic variants and the form used in any font would be up to the font designer, but some of these alternate forms are used specifically as symbols rather than letters and are encoded separately for that purpose (these characters all have “SYMBOL” in their names). Interestingly, these symbols are encoded in the Greek block rather than one of the symbol blocks. Finally, the Greek block includes letters from the Coptic alphabet, about which more later.

Ancient polytonic Greek can also be represented using combining character sequences, although algorithms that map characters from lower case to upper case have to be careful with text that uses the combining form of the *ypogegrammeni*. (The combining *ypogegrammeni* should be drawn as a *prosgegrammeni* (iota adscript) when combined with a capital letter, and text that maps to all uppercase characters needs to convert the combining *ypogegrammeni* into a regular capital letter *iota*.)

When representing polytonic Greek using combining character sequences, the *dialytika* should go first in the sequence if there is one, then the breathing mark, then any accent marks, and finally the *ypogegrammeni* if there is one. (This is another set of cases where there’s a language-specific variation of the normal Unicode accent-stacking behaviors; this is the order of marks that produces the correct Greek-specific visual arrangement of the marks.)

The Greek Extended block

Like the Latin Extended Additional block, the Greek Extended block, which runs from U+1F00 to U+1FFF, includes precomposed forms for all of the vowel-mark combinations needed to write ancient polytonic Greek. It also includes forms of the vowels with macrons and breves, which are sometimes used in Greek poetry.

It includes both the *iota-subscript* and *iota-adscript* versions of the ancient diphthongs that use the *iota subscript*. The *iota-adscript* forms are considered to be titlecase letters; the full-uppercase version is two separate letters: the base letter and the *iota*.

The Greek Extended block also includes a bunch of isolated diacritical marks and combinations of diacritical marks with no letter under them; like the ones in the main Greek block, these all have non-combining semantics. After all, the whole purpose of the characters in the Greek Extended block is to avoid using combining character sequences.

Like the Latin Extended Additional block, every character in the Greek Extended block can be represented using a combining character sequence; the characters in this block are here only for compatibility. If new Greek characters are added that can’t be represented using combining character sequences, they’ll be added to the main Greek block, not to this one.

The Coptic alphabet

The Coptic alphabet was the last writing system used to write ancient Egyptian, and it survives as the script for liturgical materials in the Coptic Orthodox Church. The Egyptians basically adopted the Greek alphabet for their language, supplementing it with letters from the older Demotic script when necessary for sounds that the Greek alphabet didn’t have letters for.

Although the Coptic alphabet is based on the Greek alphabet, the letterforms are in many cases different from their forms in the modern Greek alphabet: for example, the sigma (or “semma”) in the Coptic alphabet looks like the Latin letter C (or the Cyrillic letter s) rather than the Greek letter sigma.

Coptic was written left to right, and didn't use diacritical marks, with the exception of a horizontal bar over syllables that don't have a vowel.

The designers of Unicode unified the Coptic and Greek alphabets; the Greek block includes code point values for the Coptic letters that aren't also in the Greek alphabet, and you use the same codes for the letters they have in common as you use to represent Greek. Because the Coptic letters look different, however, you have to use a font specifically designed for Coptic rather than a normal Greek font or a full-blown Unicode font that includes glyphs for all the Unicode characters.

This means that to write a document that uses both Coptic and Greek, you pretty much have to use a styled-text format based on Unicode; with Unicode plain text, something (or someone) would have to understand the text to know which font to use for each character. If you use a styled-text format, you can indicate the font changes or attach language markings that a rendering process can use to choose proper fonts. (Although this situation is highly specialized, it's been used as an argument for adding a language-tagging facility to Unicode plain text, but more important situations involve the use of the Han ideographs. More on this later.)

There's a growing consensus that unifying Greek and Coptic was a mistake—there's a good chance they'll be disunified in a future version of Unicode.

The Cyrillic alphabet

Tradition says that the Cyrillic alphabet was invented by St. Cyril and St. Methodius in the early 860s in preparation for a mission expedition from Constantinople to the Slavic people in Moravia. The idea was to create an alphabet for writing the liturgical texts in the native Slavic language.³⁵ There's no evidence of any written Slavic prior to 860 or so.

One problem with this account is that there were actually *two* alphabets created around that time for writing the Slavic languages: the Glagolitic and the Cyrillic. No one seems to know for sure the origin of either alphabet; it seems clear that what we now know as the Cyrillic alphabet is derived from the Greek alphabet, with extra letterforms for the sounds found in the Slavic languages that aren't also found in Greek. But it's unclear what the source of these additional letterforms are. Various theories have been advanced, none completely convincing.

The Glagolitic alphabet, on the other hand, while it has the same letters as the Old Cyrillic alphabet, has completely different letterforms for them, letterforms that don't seem to have an obvious origin. Paul Cubberley, in his article in *The World's Writing Systems*, theorizes that the Glagolitic alphabet is actually loosely based on early cursive Greek and was probably developed by the Slavs in the decades leading up to 860; St. Cyril formalized this alphabet, coined new letters for the non-Greek sounds, and spread it throughout the Slavic world. What we now know as the Cyrillic alphabet (both names came into use long after the alphabets themselves) seems to have arisen after the Glagolitic alphabet. Cubberley theorizes that it was created by Cyril's followers in Bulgaria in the 890s from the "more dignified" uncial (printed) Greek alphabet, with the non-Greek letters adapted from the Glagolitic alphabet. Some of the letterforms in what we now know as the Glagolitic alphabet appear to be later back-formations from the Cyrillic alphabet that happened during the period where both scripts were in use in the same place.

³⁵ My sources for information on the Cyrillic alphabet are Paul Cubberley, "The Slavic Alphabets," in *The World's Writing Systems*, pp. 346-355, and Bernard Comrie, "Adaptations of the Cyrillic Alphabet," *op. cit.*, pp. 700-726.

Over time, the Cyrillic alphabet seems to have won out in popularity over the Glagolitic alphabet, with the Glagolitic alphabet persisting as late as the nineteenth century in the Slavic areas under the influence of the Roman Catholic church before finally giving way to the Latin alphabet.

The Cyrillic alphabet gradually predominated in the Slavic regions under the influence of the Byzantine church, spreading with the spread of the Russian Empire and, later, the USSR.

Unlike the Greek alphabet, which was standardized early among the various Greek-speaking peoples, the Cyrillic alphabet was never really standardized. It had various redundant letters and competing letterforms. A first attempt at spelling reform happened under Peter the Great in 1708 with the creation of a “civil script” intended specifically for the writing of secular texts. Some later changes were proposed by the Academy of Sciences in 1735 and 1738.

What we now know as the Cyrillic alphabet, the modern Russian alphabet, was standardized (and further reformed) by the Bolshevik regime in 1918. This alphabet spread to the various countries under the control of the Soviet Union, with various changes to accommodate the different languages.

The modern Russian alphabet has thirty-three letters. Like the Latin and Greek alphabets, the letters come in upper-/lower-case pairs, there are both roman and italic versions of the script, spaces are used between words, and text is written from left to right. Modern languages using the Cyrillic alphabet use the same punctuation marks that are normally used with the Latin alphabet. Although, like Greek, the letters of the Cyrillic alphabet were once used as numerals, today the same “Arabic” numerals used with the Latin alphabet are used with the Cyrillic alphabet.

These are the letters of the modern Russian alphabet:

А Б В Г Д Е Ж З И Й К Л М Н О П Р С Т У Ф Х Ц Ч Ш Щ Ъ Ы Ь Э Ю Я
а б в г д е ж з и й к л м н о п р с т у ф х ц ч ш щ ъ ы ь э ю я

There are two interesting letters in the Russian alphabet: the “hard sign” (ъ) and the “soft sign” (ь). These two letters really don’t function like most “letters”; instead, they function more like diacritical marks. Instead of representing distinct phonemes, they modify the pronunciation of the letters around them.

To understand how the hard and soft signs work, it helps to understand *palatalization*, one feature of Russian pronunciation that surfaces in spelling. To an American English speaker, the easiest way to understand palatalization is to think of it as the addition of a “y” sound. A palatalized consonant ends with the tongue moving like you’re going to follow it with the “y” sound. Consider the way many people pronounce the word “Tuesday”: it’s best approximated as “tyoozday.” But the “y” isn’t really given its full value: it isn’t pronounced “tee-oozday.” Instead, the “y” is kind of swallowed up into the “t”. Many English speakers turn this sound into a “ch” sound, pronouncing the word as “choozday,” but if you don’t do that, the sound at the beginning of “Tuesday” is a palatalized t. Likewise, the scraping sound at the beginning of “Houston” the way most people pronounce it is a palatalized h.

Russian has a whole bunch of palatalized consonants, and they’re generally represented using the same letters as their non-palatalized counterparts. The distinction between the two is generally made in one of two ways.

The first way is through the use of the vowels. A palatalized vowel is a vowel with the “y” sound attached to the front of it; for example, the palatalized version of “ah” would be “yah.” The Russian alphabet has ten vowels in five pairs: а/я (“a”/”ya”), э/е (“e”/”ye”), ы/и (roughly “i”/”yi”, but really

и has what we'd call the “ee” sound and ы has a darker “i” sound not found in English), о/ё (“o”/“yo”), and у/ю (“u”/“yu”). (Actually, ё and э are relatively rare in Russian spelling—the “e” and “yo” sounds are usually also represented by the letter е.)

The palatalized versions of the vowels are used to indicate an initial “y” sound (“Yeltsin” is spelled “Елцин”) or a “y” sound between two vowels (“Dostoyevsky” is spelled “Достоевский”). When they follow a consonant, however, the palatalized version of the vowel actually indicates that the *consonant* is palatalized. The soft sign (ь) is used to indicate that a consonant is palatalized when the consonant is followed by another consonant or occurs at the end of the word.

Normally the hard sign (ъ) isn't needed; the absence of the soft sign or the use of a non-palatalized vowel indicates the absence of palatalization. The hard sign is only used now before a palatalized vowel to indicate that the vowel, not the consonant that precedes it, is palatalized (i.e., the hard sign indicates the presence of the full “y” sound between the consonant and the vowel, rather than a change in the pronunciation of the consonant).

Diacritical marks aren't generally used with the Russian alphabet. The two big exceptions to this are the use of a breve over the letter і (Й) to indicate it has the “y” sound in a diphthong (for example, “Sergei” is spelled “СергеЙ”). This is called the “short i” in Russian and is generally considered a separate letter of the alphabet. A diaeresis over the letter е is sometimes used to indicate it has the “yo” sound (ё), but this generally only happens in pedagogical texts and in rare cases where two different words with two different pronunciations would otherwise look the same. Most of the time, for example, “Pyotr” is spelled “Петр” rather than Пётр.

Like the Latin alphabet, the Cyrillic alphabet is used for writing a whole host of languages, mostly Slavic languages and the languages of various minority peoples from the former Soviet Union. In addition, Mongolian can be written either in its native script or using the Cyrillic alphabet (for more on Mongolian, see Chapter 11). Like the Latin alphabet, the alphabet is a little different for every language that uses the Cyrillic alphabet. However, most languages that use the Latin alphabet, at least among the European languages, generally just add diacritical marks to the existing letters to cover their own sounds. The different languages that use the Cyrillic alphabet are much more likely to add whole new letterforms to the alphabet, and different languages with the same sounds are much more likely to represent them in different ways.

The variation is noticeable even among the Slavic languages that use the Cyrillic alphabet. Belarusian basically uses the same alphabet that Russian does, with the addition of a “short u” (ў) to represent the “w” sound, but the letter I looks different: instead of using the Russian I shape (И), it looks like the Latin letter I. (Interestingly, the short i still looks like the Russian short i.) The Ukrainian alphabet actually uses *both* versions of the letter I, and adds a rounded E (Є). The Bulgarian alphabet uses the same letters as the Russian alphabet, but treats the Russian hard sign (ъ) as a regular vowel (“Bulgaria,” for example, is spelled “България”).

The Serbian alphabet deviates further. It drops the hard and soft signs and the palatalized vowels from the Russian alphabet and adds a couple of new letters to represent palatalized consonants: Љ, Њ, Џ, and Ћ. It also adds the letter Ј, which is used for the “y” sound in front of a vowel, rather than using a whole different vowel letter, as is done in Russian (for example, the Russians would spell “Yugoslavia” “Югославия,” but the Yugoslavs spell it “Југославија”). The Macedonian alphabet take the same basic approach as the Serbian alphabet, but uses different letter shapes for some of the letters: Ѓ is rendered as Ѓ́ in Macedonian, and Ћ is rendered as Ћ́. The Macedonian alphabet also adds another new letter: S, which doesn't represent the “s” sound, but the “dz” sound.

Once you leave the confines of the Slavic languages, all bets are off. The Cyrillic alphabet is used to write more than fifty non-Slavic languages, and each one adds a few new letter shapes and/or a few new accented versions of the basic Russian letters. As with the Slavic languages, even languages from the same language family having the same basic repertoires of sounds will often use completely different letterforms to represent them.

The Cyrillic block

The Unicode Cyrillic block runs from U+0400 to U+04FF. The first part of this block, the range from U+0400 to U+045F, is based on the ISO Cyrillic standard (ISO 8859-5), and the characters in this range have the same relative positions as they do in the ISO standard, facilitating conversion between ISO 8859-5 and Unicode. This range includes all the characters necessary to write Russian, Bulgarian, Belarusian, Macedonian, Serbian, and (except for one missing character) Ukrainian.

The characters from the ISO standard are supplemented with a whole host of additional characters. The range from U+0460 to U+0489 contains the characters from the Old Cyrillic alphabet that have dropped out of the modern alphabet. This allows the writing of historical texts and liturgical texts in Old Church Slavonic. (Note that the shapes of some of the letters in the modern Cyrillic alphabet are different from their shapes in the Old Cyrillic alphabet, but that they've been unified anyway. This means you need to use a special Old Cyrillic font to get the right shapes for all the letters.) This range also includes some diacritical marks that were used with Old Cyrillic but have also fallen into disuse (and can't be unified with the marks in the Combining Diacritical Marks block).

The rest of the Cyrillic block contains letters needed for the various non-Slavic languages that use the Cyrillic alphabet. It's broken roughly into two ranges: a range of unique letter shapes, and a range of precomposed letter-mark combinations that can also be represented using combining character sequences. The Unicode standard doesn't provide an official list of the languages this range is intended to cover, but they include the one missing Ukrainian letter and the following additional languages: Abkhassian, Altay, Avar, Azerbaijani, Bashkir, Chukchi, Chuvash, Kabardian, Kazakh, Khakassian, Khanty, Khirgiz, Kildin Sami, Komi, Mari, Moldavian, Mongolian, Nenets, Tatar, Tajik, Turkmen, Uzbek, and Yakut. (A number of these languages, especially in the wake of the fall of the Soviet Union, are also written with other alphabets.)

One interesting character worth looking at is the palochka, which also looks like the letter I. It's different from the Belarusian I in two regards: 1) It doesn't occur in upper- and lower-case flavors; it always looks like an uppercase I, and 2) it's not really a letter, but a diacritical mark. It's used in a number of languages to represent the glottal stop (the way an apostrophe is in the Latin alphabet) or to turn a preceding consonant into an ejective consonant (the way an apostrophe does in the IPA).

A couple of additional notes about the Cyrillic alphabet: A few languages borrow additional letterforms from the Latin alphabet; the Unicode standard makes specific mention of the letters Q and W being used in Kurdish. The Unicode Technical Committee has resisted adding these letters to the Cyrillic block; they say the right way to handle this situation and others like it is to use the code point values in the Latin blocks to represent these letters. Also, there have been many alphabet changes in some of the languages that use the Cyrillic alphabet; the Unicode Cyrillic block only encodes the modern letters used in these languages, not the forms which have become obsolete.

The Cyrillic Supplementary block

Unicode 3.2 adds a new Cyrillic Supplementary block, running from U+0500 to U+052F, which adds a bunch of letters for writing the Komi language.

One character, U+058A ARMENIAN HYPHEN (֊) deserves special mention. This character works the same way as U+00AD SOFT HYPHEN: It's used in Armenian text to mark legal positions where a word may be hyphenated at the end of a line. It's invisible unless it occurs at the end of a line. (For more information on the soft-hyphen characters, see Chapter 12.)

The Georgian alphabet

The Georgian alphabet, like the Armenian alphabet, is essentially a one-language alphabet, being used basically just to write Georgian (although it once was also used to write Abkhastian). There are a number of different traditions for the origin of the Georgian alphabet (including an Armenian tradition attributing it to St. Mesrop, who is said to have invented the Armenian alphabet), but none seems to have much currency.³⁷ The first example of Georgian writing dates from 430. The Georgian alphabet seems to be based, at least in its ordering and general characteristics, on the Greek alphabet, although the letterforms appear to have been independent inventions.

There have been a number of different versions of the Georgian alphabet. The oldest is called *asomtavruli* ("capital letters"). Another form, called *nuskhuri*, developed in the ninth century. The two forms were used in religious texts, originally freely intermixed, but later with *nuskhuri* predominating and *asomtavruli* taking on the role of capital letters. This style was known as *khutsuri*, or "ecclesiastical writing," and was used in liturgical texts. Today, it's mostly died out, although it's still used occasionally in liturgical settings.

Modern Georgian writing, or *mkhedruli* ("soldier's writing"), developed in the tenth century from *nuskhuri*. Unlike *khutsuri*, it's caseless: there are no upper- or lower-case forms. Effectively, all the letters are lowercase. Like lowercase letters in the Latin, Greek, Cyrillic, and Armenian alphabets, the letters generally stay within what English printers call the "x-height" (the vertical space defined by the height of the lowercase "x"), with various letters having ascenders, descenders, or both. There are display fonts for Georgian that get rid of the ascenders and descenders and make all the letters the height of a letter with an ascender, making them look like capital letters; these fonts are used primarily for titles and headlines. They're not true capital letters, however. A Georgian linguist named Akaki Shanidze made an attempt to reintroduce the *asomtavruli* into modern Georgian writing as capital letters, but it didn't catch on.

The Georgian alphabet has stayed relatively unchanged since the tenth century, although there have been a few reform movements. There was a reform of the alphabet in the 1860s under which a couple characters were added to the alphabet and five letters corresponding to obsolete sounds were removed from the alphabet.

The modern *mkhedruli* alphabet consists of thirty-eight letters, as follows:

ა ბ გ დ ე ვ ზ ლ თ ი კ ლ მ ნ ო პ ჟ რ ს ტ ჯ უ ფ ქ ღ ყ შ ჩ ც ძ წ კ ხ ჯ ჳ ჴ ჵ

As in the other alphabets examined in this chapter, Georgian is written from left to right with spaces between the words. Diacritical marks are not used, and the letters don't interact typographically. For the most part, Georgian uses the same punctuation as is used with the Latin alphabet, with the exception of the period, which generally looks like the Armenian period. As with the other alphabets,

³⁷ My sources for this section are Dee Ann Holisky, "The Georgian Alphabet," in *The World's Writing Systems*, pp. 364-369, supplemented by Nakanishi, pp. 22-23.

at one time the letters of the Georgian alphabet were used as digits, but this has died out in favor of the Arabic numerals.

The Georgian block in Unicode extends from U+10A0 to U+10FF. It contains two series of letters: a series of *asomtavruli* and a series of *mkhedruli*. The *nuskhuri* have been unified with the *mkhedruli*; this means you need to use a special font to write in Old Georgian that provides the correct letterforms for the unified letters. The *asomtavruli* letters are categorized as uppercase letters by the standard, and the *mkhedruli* are categorized as lowercase letters, but algorithms for converting between cases should generally leave the *mkhedruli* alone instead of converting them to the *asomtavruli*. (Unicode doesn't define case mappings between the *asomtavruli* and the *mkhedruli*; you want to do that only if you're operating on Old Georgian rather than modern Georgian.)

There's only one punctuation mark in the Georgian block: the Georgian paragraph mark, which appears at the ends of paragraphs. The Georgian period has been unified with the Armenian period; use the Armenian period (U+0589) in Georgian text.

CHAPTER 8 *Scripts of The Middle East*

Most of the alphabetic writing systems descend from the ancient Phoenician alphabet, but they do so along several different lines of descent. In the last chapter, we looked at a collection of scripts that all descended from the ancient Greek alphabet (although the relationship between the Greek alphabet and the Armenian and Georgian alphabets seems to be more distant and speculative, these two alphabets do share many common characteristics with the ancient Greek alphabet).

In this chapter, we'll look at a group of scripts that descended down a different line: they descend from the ancient Aramaic alphabet. These scripts have several important things in common, many of which distinguish them from the scripts we looked at in the last chapter:

- They're written from right to left, instead of from left to right like most other scripts.
- They're uncased. Instead of there being two sets of letterforms, one used most of the time and the other used at the beginning of sentences and proper names and for emphasis, as the European scripts do, each of these scripts has only one set of letterforms (in fact, only the European scripts are cased).
- The letters of these alphabets generally only represent consonant sounds. Most vowel sounds are either not represented at all and are filled in mentally by the reader (th qvln t f wrtng Nglsh lk ths), or are represented by applying various kinds of marks to the letters representing the consonants.
- Like the European scripts, the Middle Eastern scripts use spaces between words.

The four scripts in this group that have encodings in the Unicode standard are Hebrew, Arabic, Syriac, and Thaana.

In addition to the common characteristics listed above, the Arabic and Syriac scripts have an additional important characteristic in common: the letters are cursively connected, the way handwritten English often is, even when they're printed. This means that the shapes of the letters can

differ radically (much more radically than in handwritten English, actually) depending on the letters around them.

Bidirectional Text Layout

By far the most important thing the scripts of the Middle East have in common (at least from the point of view of computer programmers trying to figure out how to draw them on a computer screen and represent them in memory) is that they're written from right to left rather than from left to right. By itself, this isn't so hard to do on a computer, but if you mix text written using one of these alphabets with text written using one of the alphabets in the previous chapter, figuring out how to position all the letters correctly relative to each other can get to be a rather complicated affair. Because you can end up with a line of text with individual snippets of text on the line running in two different directions, these four scripts are often referred to as "bidirectional" scripts. Computer programmers who deal with these scripts often usually abbreviate this to "bi-di."³⁸

You actually don't even have to be mixing text in different languages to get into bidirectional text layout. In all of these languages, even though letters and words are written from right to left, *numerals* are still written from left to right. This means that text in these languages that contains *numbers* is bidirectional, not just text in these languages containing foreign words or expressions. Because of this, the right-to-left scripts can correctly be thought of as bidirectional in and of themselves, and not just when mixed with other scripts.

Let's take a look at just how this works. Consider the Hebrew expression *mazel tov*, which means "congratulations." (It actually translates literally as "good luck," but it used in the same way English speakers use the word "congratulations," not the way we'd use "good luck.") If you wrote "congratulations" in English, it'd look like this:

congratulations

In Hebrew, the same thing would look like this³⁹:

מזל טוב

Notice that it's justified against the right-hand margin. This is because the first letter is the one furthest to the right. Just as we start each line of text flush against the left-hand margin, Hebrew text starts each line flush against the right-hand margin.

Notice I said the *first* letter is against the right-hand margin. These letters transliterate as *mazel tov*. The letter Mem (מ), representing the *m* sound, is furthest to the right, moving to the left, by Zayin (ז, *z*), Lamed (ל, *l*), a space, Tet (ט, *t*), Vav (ו, *o*), and finally Bet (ב, *v*).

³⁸ I've seen numerous different spellings for this term, although it's always pronounced "bye-dye." "BIDI" seems to be fairly common, but to me it always looks like it's pronounced "bidy," so I'll standardize on "bi-di" here.

³⁹ Many thanks to John Raley and Yair Sarig for their help translating the Hebrew examples in this chapter.

So how do you represent this in memory? Well, the straightforward way to do it is for the first letter to go first in the backing store:

```
U+05DE HEBREW LETTER MEM (מ)
U+05D6 HEBREW LETTER ZAYIN (ז)
U+05DC HEBREW LETTER LAMED (ל)
U+0020 SPACE
U+05D8 HEBREW LETTER TET (ט)
U+05D5 HEBREW LETTER VAV (ו)
U+05D1 HEBREW LETTER BET (ב)
```

By now, you're probably getting rather impatient. After all, this is all so simple as to almost go without saying. It starts to get interesting, however, when you mix Hebrew text with, say, English text:

Avram said מזל טוב and smiled.

Notice that the Hebrew expression looks exactly the same when it's mixed with the English: the מ still goes furthest to the right, with the other letters extending to the left after it. The same thing happens when English is embedded into Hebrew:

אברהם אמר congratulations וחייך.

The same thing happens when you have a number in a Hebrew sentence:

יש לי 23 ילדים.

The number in this sentence is twenty-three; it's written the same way it would be in English.

You can even have these things embedded multiple levels deep:

Avram said "יש לי 23 ילדים." and smiled.

Because this is basically an English sentence, the prevailing writing direction is left to right (1) and the sentence is justified against the left-hand margin. But the Hebrew sentence reads from right to left (2). However, even within the Hebrew sentence, the number reads from left to right (3):

① Avram said "יש לי 23 ילדים." and smiled.
② ←
③ →

This means that the eye traces a sort of zigzag motion as it reads this sentence:

Avram said "יש לי 23 ילדים." and smiled.
→

Now things get interesting. When you mix directions on the same line of text, there's no longer an obvious way of ordering the characters in storage. There are two basic approaches to how a sentence such as...

Avram said **מזל טוב** and smiled.

...should be represented in memory. The first approach is to store the characters in *logical order*, the order in which the characters are typed, spoken, and read:

```
U+0073 LATIN SMALL LETTER S
U+0061 LATIN SMALL LETTER A
U+0069 LATIN SMALL LETTER I
U+0064 LATIN SMALL LETTER D
U+0020 SPACE
U+05DE HEBREW LETTER MEM
U+05D6 HEBREW LETTER ZAYIN
U+05DC HEBREW LETTER LAMED
U+0020 SPACE
U+05D8 HEBREW LETTER TET
U+05D5 HEBREW LETTER VAV
U+05D1 HEBREW LETTER BET
U+0020 SPACE
U+0061 LATIN SMALL LETTER A
U+006E LATIN SMALL LETTER N
U+0064 LATIN SMALL LETTER D
```

If you store the characters this way, the onus is now on the text-rendering process (the code that draws characters on the screen) to switch the letters around on display so that they show up arranged correctly. The code that accepts the text as input, on the other hand, doesn't need to do anything special: each character, as it's typed, is simply appended to the end of the stored characters in memory.

Then there's the opposite approach, somewhat paradoxically called *visual order*: the characters are stored in memory in the order in which the eye encounters them if it simply scans in a straight line from one end of the line to the other. This can be done either with a right-to-left bias or a left-to-right bias (both exist in the real world). That leads to this ordering (with a left-to-right bias):

```
U+0073 LATIN SMALL LETTER S
U+0061 LATIN SMALL LETTER A
U+0069 LATIN SMALL LETTER I
U+0064 LATIN SMALL LETTER D
U+0020 SPACE
U+05D1 HEBREW LETTER BET
U+05D5 HEBREW LETTER VAV
U+05D8 HEBREW LETTER TET
U+0020 SPACE
U+05DC HEBREW LETTER LAMED
U+05D6 HEBREW LETTER ZAYIN
U+05DE HEBREW LETTER MEM
U+0020 SPACE
U+0061 LATIN SMALL LETTER A
U+006E LATIN SMALL LETTER N
U+0064 LATIN SMALL LETTER D
```

In this example, the Latin letters occur in the same relative order that they did in the previous example, but the Hebrew letters are reversed. If you store the characters this way, a text-rendering process can be fairly straightforward: the characters are just drawn in succession from left to right. The onus is now on the code that accepts keyboard input, which has to rearrange the characters into the correct order as they're typed. (The simplest system of all just treats everything as unidirectional text, which puts the onus on the user to type the Hebrew text backwards. This is suitable, perhaps, for putting together examples in books like this one, but not for doing actual Hebrew text editing.)

Visual order is easier for the rendering engine, but poses serious difficulties for most other processes. You can really only store individual lines of text in visual order—if you store a whole paragraph that way, it becomes incredibly hard to divide into lines. And simple operations such as putting a paragraph break into the middle of an existing paragraph (in the middle of the backwards text) suddenly become ambiguous. There are still visual-order systems out there, but they're generally older systems based on older technology. All modern Hebrew encodings, including Unicode, use logical order. (The same goes for Arabic, although visual order was always less prevalent on Arabic systems.)

Unicode goes the other logical-order encodings one better by very rigorously and precisely specifying exactly how that logical ordering is to be done. This is the subject of the next section.

The Unicode Bidirectional Layout Algorithm

The Unicode bidirectional text layout algorithm (or “bi-di algorithm,” as it is most commonly called) is possibly the most complicated and difficult-to-understand aspect of the Unicode standard, but it's also pretty vital. The idea is that with Unicode giving one the ability to represent text in multiple languages as simply as one can represent text in a single language, you'll see more mixed-language text, and you have to make sure that any sequence of characters is interpreted the same (and looks the same, ignoring differences in things like font design) by all systems that claim to implement the standard. If the Unicode standard didn't specify exactly what order the characters should be drawn in when you have mixed-directionality text in a Unicode document, then the same document could look different on different systems, and the differences wouldn't be merely cosmetic—they'd actually affect the sense of the document.

Unfortunately, the business of laying out bidirectional text can get rather wild. Most of the complicated cases won't occur very often in real-world documents, but again, everything has to be specified rigorously for the standard to work right. (As it turns out, the first version of the Unicode bi-di algorithm, which appeared in Unicode 2.0, wasn't specified rigorously enough—holes were discovered as people began trying to write actual implementations of it. The Unicode 3.0 standard includes a more rigorous spec, and updated versions have also been published as Unicode Standard Annex #9. UAX #9 also includes two reference implementations of the algorithm, one in Java and one in C.)

I'm not going to try to cover every last detail of the Unicode bi-di algorithm here; you can refer to the Unicode Standard and UAX #9 for that. Instead, I'll try to capture the effect the Unicode bi-di algorithm is supposed to produce and not worry about the details of how you actually produce it.

Inherent directionality

Let's start with the basics. The Unicode bi-di algorithm is based on the concept of *inherent*

directionality, that is, each character in the standard is tagged with a value that says how it behaves when laid out on a horizontal line with other characters. Those characters that Unicode classifies as “letters” (a general category that basically includes all characters that make up “words,” including not only actual letters, but also syllables, ideographs, and non-combining diacritical marks) are classified as either left-to-right characters or right-to-left characters.

If you have two or more left-to-right characters in a row in storage, they are always to be displayed, no matter what, so that each character is drawn to the left of all the characters on the same line that follow it in storage and to the right of all characters on the same line that precede it in storage. In other words,

```
U+0041 LATIN CAPITAL LETTER A
U+0042 LATIN CAPITAL LETTER B
U+0043 LATIN CAPITAL LETTER C
```

will always be drawn as

ABC

regardless of what’s going on around them (provided they’re being drawn on the same line, of course). Right-to-left characters are the reverse: If you have two or more right-to-left characters in a row in storage, they are always to be displayed, no matter what, so that each character is drawn to the right of all the characters that follow it in storage and to the left of all characters that precede it (again, assuming they come out on the same line). In other words,

```
U+05D0 HEBREW LETTER ALEF (א)
U+05D1 HEBREW LETTER BET (ב)
U+05D2 HEBREW LETTER GIMEL (ג)
```

will always be drawn as

גבא

regardless of what’s going on around them.

What happens when you have a left-to-right character and a left-to-right character next to each other in memory depends on context. First, the characters are organized into *directional runs*, sequences of characters having the same directionality. Within the runs, the characters are displayed as shown above, which means that the characters on either side of the run boundary in memory will very likely not be drawn adjacent to each other. In other words, let’s say you concatenate the two sequences we just looked at together into a single sequence:

```
U+0041 LATIN CAPITAL LETTER A
U+0042 LATIN CAPITAL LETTER B
U+0043 LATIN CAPITAL LETTER C
U+05D0 HEBREW LETTER ALEF
U+05D1 HEBREW LETTER BET
U+05D2 HEBREW LETTER GIMEL
```

This can be drawn either as

ABCאבג

or as

אבגABC

Which of these two arrangements is correct depends on context. If the paragraph consists of just these six characters, which arrangement is correct depends on whether the paragraph is considered a Hebrew paragraph or an English paragraph, and this is usually specified using a higher-level protocol such as out-of-band style information. In the absence of a higher-level protocol, the first character of the paragraph determines it, in which case the first arrangement would win out, since the paragraph starts with a left-to-right letter (there are ways or overriding this behavior even in plain text—we'll get to them later).

The interesting thing to note about both examples, though, is that the characters on either side of the run boundary in storage, C and alef, aren't drawn next to each other in either arrangement. Drawing these two characters next to each other would violate the rule about how the characters within a directional run are to be arranged.

Combining marks are basically unaffected by any of this ordering. A combining mark doesn't have any directionality of its own; it always combines with the character that precedes it in storage, and it takes on the directionality of the character it attaches to. So if we throw a couple of combining marks into the previous mix...

```
U+0041 LATIN CAPITAL LETTER A
U+030A COMBINING RING ABOVE
U+0042 LATIN CAPITAL LETTER B
U+0043 LATIN CAPITAL LETTER C
U+05D0 HEBREW LETTER ALEF
U+05D1 HEBREW LETTER BET
U+05BC HEBREW POINT DAGESH OR MAPIQ
U+05D2 HEBREW LETTER GIMEL
```

...they don't break up the directional runs, and the directionality doesn't affect which character they attach to, or how they attach to that character. The ring attaches to the A, and the dagesh (a dot which appears in the middle of the character) attaches to the bet, and you get

ÀBCאבג

In other words, combining character sequences are treated as units for the purposes of laying out a line of bidirectional text. Control characters and other characters that have no visual presentation at all (except for a few we'll look at soon) are similarly transparent to the algorithm (i.e., it functions as though they're not there).

Neutrals

So much for the easy part. Now consider what happens if you have spaces. Consider the example sentence we used earlier:

Avram said מוזל טוב and smiled.

The spaces in “Avram said” and “and smiled” are the same code-point value (U+0020) as the space in “מוזל טוב” and the spaces at the boundaries between the English and the Hebrew. The space is said to have *neutral* directionality, as do a bunch of other punctuation marks. The basic idea here is that the designers of Unicode didn’t want to have two different sets of spaces and punctuation marks that were the same except that one set was supposed to be used for left-to-right characters and the other set for right-to-left characters. Instead you have one set of spaces and punctuation and they have neutral directionality.

So what does that mean? Very simple. If a neutral character is flanked on both sides by characters of the same directionality, it takes on the directionality of the surrounding characters. In fact, if a sequence of neutral characters is flanked on both sides by characters of the same directionality, they all take on that directionality. For example, the following sequence...

```
U+0041 LATIN CAPITAL LETTER A
U+002E FULL STOP
U+0020 SPACE
U+0042 LATIN CAPITAL LETTER B
U+05D0 HEBREW LETTER ALEF
U+002E FULL STOP
U+0020 SPACE
U+05D1 HEBREW LETTER BET
```

...would get drawn like this:

A. B א.ב

If there’s a conflict—that is, if a neutral is flanked by characters of opposite directionality, then the overall directionality of the paragraph wins. So if you have the following sequence...

```
U+0041 LATIN CAPITAL LETTER A
U+0042 LATIN CAPITAL LETTER B
U+0020 SPACE
U+05D0 HEBREW LETTER ALEF
U+05D1 HEBREW LETTER BET
```

...the space takes on the directionality of the A and B (because the A is the first thing in the paragraph). This means it goes to the right of the B, between the English letters and the Hebrew letters:

AB אב

Interestingly, if the overall paragraph direction had been specified as right-to-left, the space would still appear in the same place. This is because it now takes on the directionality of the Hebrew

letters, causing it to appear to the right of the א as part of the run of Hebrew letters:

אב AB

Numbers

Numbers complicate things a bit. Consider the following example:

Avram said “יש לי 23 ילדים.”

The number is twenty-three, written just as it would be in English. Notice that it’s flanked by two snippets of Hebrew text: ילדים and יש לי. The sequence to the right, “יש לי”, occurs in storage (and was typed) before the sequence to the left (“ילדים”). That is, the two snippets, together with the “23,” constitute a single Hebrew sentence (“I have 23 children”).

יש לי 23 ילדים.

...and are ordered accordingly.

Contrast that with this example:

The first two books of the Bible are שמות and בראשית.

It has the same superficial structure as the other example: three pieces of English text interspersed with two pieces of Hebrew text. But here the *leftmost* Hebrew word, “בראשית” (*bereshith*, the Hebrew name for Genesis) occurs first in memory and the rightmost Hebrew word (שמות, or *shemoth*, the Hebrew name for Exodus).

The difference is that “and” and “23” are treated differently by the bi-di algorithm. The word “and” gets treated as part of the outer English sentence, breaking the Hebrew up into two independent units. The “23,” on the other hand, because it’s a number, gets treated as part of the inner Hebrew sentence, even though its individual characters run from left to right. Both Hebrew snippets get treated as part of the same Hebrew sentence, rather than as independent Hebrew words.

Because of this effect, digits are said to have *weak* directionality, while letters are said to have *strong* directionality.

Certain characters that often appear with, or as part of, numbers, such as commas, periods, dashes, and currency and math symbols, also get special treatment. If they occur within a number (or, for some of them, before or after a number), they get treated as part of the number; otherwise, they get treated as neutrals.

Also, numbers are only special within right-to-left text. If a number occurs inside a sequence of normal left-to-right text, it just gets treated as part of that sequence.

The Left-to-Right and Right-to-Left Marks

The basic bi-di algorithm won't give you perfect results with every possible sentence you can contrive. For example, consider this sentence we looked at earlier:

Avram said “יש לי 23 ילדים.” and smiled.

Notice the period at the end of the Hebrew sentence. The period is a neutral character, like the quotation mark and space that follow it. According to the Unicode bi-di algorithm, since these three characters have a left-to-right character on one side and a right-to-left character on the other, they should take on the overall directionality of the paragraph. So the sentence should look like this:

Avram said “. יש לי 23 ילדים” and smiled.

In other words, with the default behavior, the period would show up at the *beginning* of the Hebrew sentence, instead of at the end where it belongs. In styled text, you could use styling information to control where the period shows up, but we need a way to make it show up in the right place even in plain text. The solution is two special Unicode characters:

```
U+200E LEFT-TO-RIGHT MARK
U+200F RIGHT-TO-LEFT MARK
```

These two characters have one purpose and one purpose only: to control the positioning of characters of neutral directionality in bidirectional text. They don't have any visual presentation of their own, and don't affect other text processes, such as comparing strings or splitting a run of text up into words. But the Unicode bi-di algorithm sees them as strong left-to-right and right-to-left characters.

This makes them useful for two things. One of them is controlling the directionality of neutral characters, which is what we need here. So if you insert a right-to-left mark into our sentence between the period and the closing quotation mark...

Avram said “RLM יש לי 23 ילדים” and smiled.

...the period shows up in the right place. The right-to-left mark (shown as “RLM” above) is a strong right-to-left character just like the Hebrew letters. Since the period is now flanked by right-to-left characters, it gets treated as a right-to-left character and appears in its proper place at the end of the sentence. (The “RLM” is just there to show where the right-to-left mark goes; it's actually invisible.)

The other thing that the left-to-right and right-to-left marks are good for is controlling the overall directionality of a paragraph. In the absence of a higher-level protocol, the directionality of a paragraph is the directionality of the first character (of strong directionality) in the paragraph. This is usually the right answer, but not always. Consider this sentence:

מזל טוב is Hebrew for “congratulations.”

The first character in this paragraph is the Hebrew letter mem (מ), which is a strong right-to-left character. This means that the computer will give the paragraph right-to-left directionality. In effect, it'll treat this sentence as a Hebrew sentence with some embedded English, even though it's clearly an English sentence with some embedded Hebrew. What you'll actually see is this:

..is Hebrew for “congratulations מזל טוב

But if you stick a left-to-right mark at the beginning of the sentence, *it* is now the first character in the paragraph with strong directionality. It's a left-to-right character, so the paragraph gets treated as a left-to-right paragraph, even though it begins with a Hebrew letter, producing the arrangement we

want. (Again, since the left-to-right mark is invisible, its presence doesn't do anything but affect the arrangement of the other characters.)

The Explicit Override Characters

Unicode provides a few other invisible formatting characters that can be used to control bidirectional layout. There are occasional situations where you want to explicitly override the inherent directionality of a Unicode character. For example, you might have part numbers that include both English and Hebrew letters, along with digits, and you want to suppress the reordering that would normally happen. The left-to-right override (U+202D) and right-to-left override (U+202E) characters force all of the characters that follow them to be treated, respectively, as strong left-to-right and strong right-to-left characters. The pop-directional-formatting character (U+202C) cancels the effect of the other two, returning the bi-di algorithm to its normal operation. For example, we saw earlier that the following sequence of characters...

U+0041 LATIN CAPITAL LETTER A

U+0042 LATIN CAPITAL LETTER B

U+0043 LATIN CAPITAL LETTER C

U+05D0 HEBREW LETTER ALEF

U+05D1 HEBREW LETTER BET

U+05D2 HEBREW LETTER GIMEL

...is drawn like this:

אבגABC

But if we insert the proper override characters...

U+0041 LATIN CAPITAL LETTER A

U+0042 LATIN CAPITAL LETTER B

U+0043 LATIN CAPITAL LETTER C

U+202D LEFT-TO-RIGHT OVERRIDE

U+05D0 HEBREW LETTER ALEF

U+05D1 HEBREW LETTER BET

U+05D2 HEBREW LETTER GIMEL

U+202C POP DIRECTIONAL FORMATTING

...you get this instead:

גבאABC

The Explicit Embedding Characters

Sometimes you wind up with multiple levels of embedding when mixing languages. For example, consider the following sentence:

“The first two books of the Bible are בראשית and שמות.” אברהם אמר

Here, you've got a Hebrew sentence with an English sentence embedded in it as a quote, and then you've got a couple Hebrew words embedded in the English sentence. You can get this kind of effect with numbers in right-to-left text without doing anything special, but you can't do it with this kind of thing. With the basic bi-di algorithm, you get this:

אברהם אמר “The first two books of the Bible are” בראשית and שמות.

The two Hebrew words in the English sentence get treated as part of the Hebrew sentence instead, breaking the English sentence up into two separate English snippets. Unicode provides two more invisible formatting characters to allow you to get the desired effect here, the left-to-right embedding (U+202A) and the right-to-left embedding (U+202B) characters. The pop-directional-format character (U+202C) is also used to terminate the effects of these two characters.

If you precede the English sentence with U+202A LEFT-TO-RIGHT EMBEDDING and follow it with U+202C POP DIRECTIONAL FORMAT, they force the Hebrew words to be treated as part of the English sentence, producing the desired effect.

Mirroring characters

There are a bunch of neutral characters, parentheses and brackets chief among them, that not only take on the directionality of the characters around them, but have a different glyph shape depending on their directionality. For example, in both of these examples...

Avram (the neighbor) asked for my help.

אברהם (השכן) ביקש את עזרתי.

...the parenthesis at the beginning of the parenthetical expression is represented by U+0028 LEFT PARENTHESIS, even though it looks one way [(] in the English sentence and a different way [)] in the Hebrew sentence. The same goes for the parenthesis at the other end of the parenthetical expression. (The name “LEFT PARENTHESIS” is an unfortunate side effect of the unification of Unicode with ISO 10646, which used that name.) In other words, the code point U+0028 encodes the semantic “parenthesis at the beginning of a parenthetical expression,” not the glyph shape (.

The Unicode standard includes several dozen characters with the mirroring property, mostly either various types of parentheses and brackets, or math symbols. Many of these, such as the parentheses and brackets, occur in mirror-image pairs that can simply exchange glyph shapes when they occur in right-to-left text, but many don’t occur in mirror-image pairs. Fonts containing these characters have to include alternate glyphs that can be used in a right-to-left context.

Line and Paragraph Boundaries

It’s important to note that bidirectional reordering takes place on a line-by-line basis. Consider our original example:

Avram said מזל טוב and smiled.

If it’s split across two lines with the line break in the middle of the Hebrew expression such that only one of the Hebrew words can appear on the first line, it’ll be the *first* Hebrew word, the one to the right, even though this makes the words appear to come in a different order than when the whole sentence is on one line. That is, when the sentence is split across two lines, you get...

Avram said מִזְלֵ
טוֹב and smiled.

...and not...

Avram said טוֹב
מִזְלֵ and smiled.

In other words, bidirectional reordering takes place *after* the text has been broken up into lines.

Also, the effects of bidirectional reordering are limited in scope to a single paragraph. The Unicode bi-di algorithm never has to consider more than a single paragraph at a time. Paragraph separators also automatically cancel the effect of any directional-override or directional-embedding characters.

Bidirectional Text in a Text-Editing Environment

Bidirectional text poses a few challenging problems for text-editing applications. Consider our earlier example:

ABCאבג

This is represented in memory as follows:

```
U+0041 LATIN CAPITAL LETTER A
U+0042 LATIN CAPITAL LETTER B
U+0043 LATIN CAPITAL LETTER C
U+05D0 HEBREW LETTER ALEF
U+05D1 HEBREW LETTER BET
U+05D2 HEBREW LETTER GIMEL
```

Let's say you click to place the insertion point here:

ABCאבג|

This position is ambiguous. It can be between the letter C and the letter א, or it can be at the end of the string, after the letter ג. The same is true if you put the insertion point here:

ABCאבג|

This spot can be either of the same two positions. So if you put the insertion point in either of these spots and start typing, where should the new characters go?

This can be tricky. One solution is to put the newly-typed character at whatever position in storage will cause it to appear at the visual position where the user clicked. For instance, if you click here...

ABC|אבג

...and type an English letter, it goes at the end of the English letters, before the Hebrew. But if you type a Hebrew letter, it goes at the end, after the other Hebrew letters. On the other hand, if you click here...

ABCאבג|

...the opposite happens. If you type an English letter, it goes at the end, after the Hebrew, but if you type a Hebrew letter, it goes at the beginning of the Hebrew letters.

Most text editors actually handle this the opposite way, which is easier to implement. They identify each visual position with a particular position in the character storage, and draw *two* insertion points on the screen to show the two different places in the text where newly-typed characters might appear:

ABC|אבג|

This example shows the position in the storage between the English and Hebrew letters: The leftward tick on the first insertion point shows that a new English letter would go there, and the rightward tick on the second insertion point shows that a new Hebrew letter would go there.

Selection highlighting can pose similar challenges. Consider this example:

Avram said מזל טוב and smiled.

If you drag from here...

Avram|said מזל טוב and smiled.

...to here...

Avram said מזל טוב| and smiled.

...this can mean two different things. In a lot of applications, this gesture will select a contiguous range of characters in storage, starting where you clicked and ending where you released. This produces a visually discontinuous selection:

Avram said טוב מזל and smiled.

This is called *logical selection*. This is a little simpler to handle in code because we're always dealing with logically continuous ranges of text, but it's more complicated to draw and can be jarring for users who aren't used to it. The opposite approach, *visual selection*, selects a visually contiguous piece of text...

Avram said טוב מזל and smiled.

...even though this represents two separate and distinct ranges of text in the character storage. This looks more natural, but requires a lot more behind-the-scenes juggling on the part of the text editor.

To see what I'm talking about, see what effect the two different selection types have on a copy-and-paste operation. Let's say you make that selection gesture on an editor that uses logical selection:

Avram said טוב מזל and smiled.

You've selected a logically continuous range of text. If you pick "Cut" from the Edit menu, you now have this:

Avram טוב and smiled.

If you turn around and immediately choose "Paste" from the Edit menu, you should get back your original text again. This is simple to do: the insertion point marks the spot in the storage where you deleted the text, and it gets inserted in one logical lump, with the bi-di algorithm reordering things to look like they did before. Now let's say you did the same selection gesture on an editor that supported visual selection:

Avram said טוב מזל and smiled.

You pick "Cut" from the Edit menu and get this:

Avram מזל and smiled.

Again, if you turn around and immediately choose "Paste," you should get the original text back again. In the same way, if you were to paste the text into any other spot in the document, the result you get should keep the appearance it had when you originally selected it. To put the text back the way it was in one paste operation, thus, requires the program to keep the two pieces of text you originally selected separate and to paste them into *two separate places* in the storage such that they'll be drawn next to each other on the screen at the specified position. This type of thing can require some interesting behind-the-scenes juggling. For on the whole issue of dealing with bidirectional text in a text editing environment, and on implementing the bi-di algorithm, see Chapter 16.

Okay, with this primer on bidirectional text layout under our belts, we can go on to look at the individual alphabets in Unicode which require it...

The Hebrew Alphabet

What we now know as the Hebrew alphabet dates back to about the second century BC, and has come down to use more or less unchanged.⁴⁰ It evolved out of the Aramaic alphabet, which in turn grew out of the Phoenician alphabet. There was an earlier Hebrew alphabet dating back to about the ninth century BC, but it gradually died out as the ancient Israelites were dominated by a succession of Aramaic-speaking groups, including the Assyrians, Babylonians, and Persians. When the Holy Land was conquered by the Greeks, the Hebrew language made a comeback and a new alphabet for it grew out of the Aramaic alphabet. The modern Hebrew alphabet is often called “square Hebrew” to distinguish it from the more rounded form of the older Hebrew alphabet, which continued to be used by the Samaritans up to relatively recent times.

Like the Phoenician alphabet from which it ultimately derived, the Hebrew alphabet consists of twenty-two letters, all consonants:

א ב ג ד ה ו ז ח ט י כ ל מ נ ס ע פ צ ק ר ש ת

Since the Hebrew alphabet, like all the alphabets in this chapter, is written from right to left, it’s shown here justified against the right-hand margin. The first letter of the alphabet, alef (א), is furthest to the right.

The letters of the Hebrew alphabet don’t generally interact typographically, and there are no upper- and lower-case forms, but five of the letters have a different shape when they appear at the end of a word: כ becomes ך, מ becomes ם, נ becomes ן, פ becomes ף, and צ becomes ץ.

The fact that the Hebrew alphabet has no vowels came to be a problem as, over time, pronunciation diverged from spelling and variant pronunciations developed, so some of the consonants came to do double duty as vowels in certain contexts and certain positions within words. But as detail-perfect interpretation of Biblical texts became more and more important, a more precise system was needed. Because the original consonants were considered so important in Biblical texts, the solution that developed didn’t involve the creation of new letters; instead, a system of *points*, various marks drawn above, below, and sometimes inside the consonants, outside the normal line of text, evolved. The first pointing systems evolved around AD 600, and the modern system of points was standardized by about 1200.

There are basically three categories of points. The first is diacritical marks that change the sound of a consonant. The most important of these is the *dagesh*, a dot drawn inside a Hebrew letterform. Exactly what the *dagesh* means depends on the letter it’s being used with, but it most often gives a letter a stopped sound where it’d normally have a more liquid sound: the letter bet (ב) for example, is pronounced like the letter *b* when it’s written with a *dagesh* (בּ), but closer to a *v* without it. The *rafe*, a line above the letter, is used less often, and has the opposite effect. בֿ is pronounced like *v*.

⁴⁰ My source for information on the Hebrew alphabet is Richard L. Goerwitz, “The Jewish Scripts,” in *The World’s Writing Systems*, pp. 487-498.

The Hebrew letter shin (ש) can be pronounced either as *s* or as *sh*. Sometimes a dot is drawn over the letter to indicate which pronunciation it has: If the dot is drawn over the rightmost tine of the shin (שׁ), it's pronounced *sh*; if it's drawn over the leftmost tine (שׂ), it's pronounced *s*.

Then there are the vowel points: There are eleven points that are used to indicate the vowel sounds that follow a consonant. With one exception, they're all drawn under the letter:

שׁ שׂ שׁ שׂ שׁ שׂ שׁ שׂ שׁ שׂ שׁ שׂ

The *sheva*, a pair of dots under a letter, is used to indicate the absence of a vowel sound following that letter:

שְּׁ

The vowel points are generally not used in modern Hebrew, except when their absence is ambiguous, in learning materials, and in Biblical and other liturgical texts. In other words, Hebrew normally looks something like this...

הָדוּ לַיהוָה כִּי-טוֹב: כִּי לְעוֹלָם חֶסֶד

...but since this is a Bible verse (Psalm 107:1, to be exact), it's usually written like this:

הָדוּ לַיהוָה כִּי-טוֹב: כִּי לְעוֹלָם חֶסֶד

Finally, there are the *cantillation marks*, an elaborate system of marks originally used to indicate accent and stress, and therefore also as punctuation, which now generally only appear in Biblical and liturgical texts, where they are used to indicate accent and intonation when a text is chanted.

Diacritical marks (such as the *dagesh*) and vowel points generally are drawn on different parts of a letter, but either may be drawn on the same side of a letter as certain cantillation marks. The rules for how the various marks are positioned relative to each other when multiple marks are attached to the same letter can get rather complicated.

The Hebrew alphabet is not only used for Hebrew, but for a selection of other languages spoken by the Jews of the Diaspora, especially Yiddish and Ladino. Certain letter-point combinations are more common in these languages, and sometimes these languages add extra diacritical marks beyond those used for Hebrew.

In addition, certain pairs of Hebrew letters are treated as distinct letters with their own phonetic values, in a manner similar to the way the double *l* (ll) is treated as a separate letter having the “y” sound in Spanish. In particular, there are three double-letter combinations—װ, ןן, and ןי—are used in Yiddish and some other languages as additional letters. That is, in these languages, ןן isn't just two ןs in a row, but a whole different letter.

Like many of the other alphabets we've looked at, the letters of the Hebrew alphabet were once used to write numbers as well (the system is a little more complicated than other alphabetical number systems); today, the regular "Arabic" numerals are used, although the old alphabetic numbering system is still often used for item numbers in numbered lists. Modern Hebrew also uses European punctuation now—the old native Hebrew punctuation's use is restricted to Biblical and liturgical texts.

The Hebrew block

The Unicode Hebrew block, which runs from U+0590 to U+05FF, contains all the marks necessary to write Hebrew. In addition to the twenty-two basic letters, the word-final forms of the five letters that have word-final forms are given their own separate code points. This seems to violate the normal Unicode rule about not encoding glyphic variants, but is in keeping with the practice in other Hebrew encodings. In fact, sometimes the word-final forms of these letters aren't actually used at the ends of words—this happens with many borrowed foreign words, for example. Because of this, the choice of which form of the letter to use becomes a matter of spelling; you *can't* have the rendering software automatically pick the proper glyph.

The letters have the same relative ordering in Unicode as they do in the ISO Hebrew standard (ISO 8859-8), but the ISO standard doesn't include the points.

The Unicode Hebrew block also includes all of the various points and cantillation marks, encoded as combining marks. The dots that distinguish the two versions of װ are encoded in this block as well; what they look like when attached to other letters isn't defined by the standard; the standard specifically calls this "an error."

Unlike the Latin or Greek alphabets, where accented forms of letters are usually represented with their own code point values, there are too many letter-mark combinations in Hebrew for this to be feasible. The three letter pairs that are treated as separate letters in Yiddish are also given single code points in the Hebrew block.

A few pointed Hebrew letters, mainly used in Yiddish, are encoded in the Alphabetic Presentation Forms block (see Chapter 12), but pointed Hebrew letters are almost always represented using combining character sequences, even in non-Unicode encodings, and these presentation forms are rarely supported in software and fonts.

The Arabic Alphabet

The Arabic alphabet, like the Hebrew alphabet, evolved from the ancient Aramaic alphabet, but developed quite a bit later: the first inscriptions using the Arabic alphabet date from the fourth century AD.⁴¹ Unlike the Hebrew alphabet, the Arabic alphabet developed from a cursive form of the Aramaic alphabet, where the letters are connected and words are written without lifting the pen. As the letters became more connected, they also started to look more the same, eventually evolving into fourteen families of similarly-shaped letters in a twenty-eight-letter alphabet. Marks used to differentiate between letters with the same shape began to evolve somewhere around the seventh

41 My sources for this section are Thomas Bauer, "Arabic Writing," in *The World's Writing Systems*, pp. 559-568, and Thomas Milo, "Creating Solutions for Arabic: A Case Study," in *Proceedings of the Seventeenth International Unicode Conference*, September 6, 2000, session TB3.

century; the patterns of dots used for this purpose today date to the ninth century. Even though they look like diacritical marks, they're not: they're integral parts of the letters.

The Arabic alphabet has twenty-eight letters, as follows:

ا ب ت ث ج ح خ د ذ ر ز س ش ص ض ط ظ ع غ ف ق ك ل م ن ه و ي

Like the Hebrew alphabet, the letters in the Arabic alphabet are all consonants, although three of them do double duty as the long vowels. The short vowels (and occasionally the long vowels) aren't normally written, but when they are, they take the form of various marks that are written above and below the normal sequence of consonantal letters. You generally only see the vowels in the Koran or other liturgical texts, in teaching materials, in cases where the word might otherwise be misread, and for decorative purposes in signs, titles, and calligraphy. There is also an elaborate system of cantillation marks and other markings that are used exclusively in Koranic texts.

In addition to the vowel signs, there are a couple of other marks: the *hamza* (ء) is used to indicate a glottal stop, usually in combination with one of the long-vowel letters (for example, $\text{أ} + \text{ء}$ produces أ). Depending on the letters involved, the *hamza* can appear either above or below the letter it's applied to. The long vowels with the *hamza* are sometimes treated as different letters, since the presence or absence of the *hamza* is often a matter of spelling.

More rarely, you'll see the *madda*, a curved horizontal stroke, used with the letter alef, one of the long vowels, to indicate a glottal stop followed by the vowel sound (أ, which is shorthand for أ). Another mark, the *shadda* (ّ), is used to mark something called "consonant gemination," basically the doubling of a consonant sound. (To understand consonant gemination, think about the way most Americans say the word "fourteen": most of us say "fort-teen," pronouncing a distinct *t* at the end of "fort" and another distinct *t* at the beginning of "teen," rather than "for-teen." The *t* is a geminated consonant.)

In addition, two Arabic letters take special shapes under certain circumstances: sometimes these are considered separate letters. The letter alef (أ) is sometimes written as آ at the end of a word. This form is called the *alef maksura*. It looks like the letter yeh (ي), except that it doesn't have the dots (the letter yeh also loses its dots, but keeps its identity, when it's combined with the *hamza*: أ). The letter teh (ت) takes a special form (ة) at the end of a word when it represents the feminine ending. This is called the *teh marbuta*.

The Arabic alphabet is used for many languages besides Arabic: because of the importance of the Koran in Islam, and the doctrine that the Koran (written in Arabic using Arabic letters) is eternal and uncreated, the Arabic alphabet is used almost everywhere throughout the Muslim world. In some of these places, they've also adopted the Arabic language, but in many, they've just adopted the Arabic alphabet to write their own language. (This phenomenon, of course, isn't unusual: the Latin alphabet basically followed the influence of the Roman Catholic Church, the Greek and Cyrillic alphabets the influence of the Greek Orthodox church, and the Hebrew alphabet followed the spread of Judaism.) At one time, in fact, the Arabic alphabet was used everywhere in the Muslim world. The few places where it isn't now are the results of political forces (the Latin alphabet being imposed by the Ataturk regime in Turkey, for example, or the Soviet Union imposing the Cyrillic alphabet on the Muslim people within its borders).

As with the other alphabets used for many different languages, many letters have been added to the Arabic language to adapt it for other languages. These new letters tend to take the shapes of existing Arabic letters. They wind up being differentiated in a number of ways: the use of dots or other marks, the use of miniaturized versions of other letters as diacritics, differences in the way the letters join to their neighbors, or variant shapes of the same letter in Arabic turning into distinct letters in some other languages.

Perhaps the biggest thing that sets the Arabic alphabet apart from the alphabets we've looked at so far is its cursive nature: generally speaking, the letters in a word all connect to their neighbors, turning a single word into one continuous pen stroke (often with dots and other marks around it), much like in standard handwritten English. Much more than in handwritten English, however, the letters change shape greatly depending on the letters around them.

The flowing nature and high elasticity of the letters of the Arabic alphabet mean that it lends itself particularly well to calligraphy. Because depictions of people and animals aren't allowed inside Muslim mosques, calligraphy has long been the primary form of Muslim religious art, and Arabic calligraphy can be extraordinarily beautiful. Many different calligraphic styles have been developed for writing Arabic.

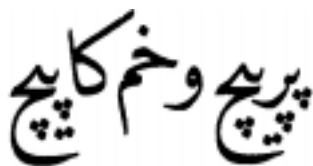
One thing they all have in common is that a letter tends to join to the letter that follows it either on the left-hand side or on the bottom. In handwritten Arabic, individual words tend to slope down and to the left, frequently overlapping the words on either side⁴²:



This means that in handwritten Arabic, spaces between words aren't really necessary. It also makes printing Arabic quite difficult and involved.

It also means that when Arabic text is justified against both margins, extra space on the line isn't taken up by widening the spaces between words, as it is in the other languages we've looked at so far. Instead, words can be elongated or shortened by using different forms of the various letters, or by inserting extenders called *kashidas* or *tatweels*, such as the long one in the example above.

The various ways that the letters can combine can also make it a fairly complicated affair to determine where the dots, vowels and other marks that may accompany the letters go⁴³:



In modern times, simplified versions of the Arabic alphabet that are easier to print have developed. Most modern printed Arabic, especially that generated using computers, uses a simplified version of

⁴² This example is taken from the Milo paper cited above, section 2, p. 10.

⁴³ *Ibid.*

the Arabic alphabet that gives each letter up to four shapes. The four shapes are designed so that all the letters connect along a common baseline. Because the baseline is horizontal rather than slanted, spaces *are* used between words.

For example, the letter heh look like this by itself:

ه

But three of them in a row look like this:

ههه

Likewise, the letter sheen looks like this when it stands alone:

ش

But the long bowl on the left goes away when it connects with a neighboring letter. Three *sheens* in a row look like this:

ششش

(This is a good example, by the way, of how some sequences of Arabic letters can just turn into wiggly lines without the dots there to help discern the letters—in handwritten or better printed Arabic, the shapes of the various humps can vary to help differentiate the letters, but this distinction is lost in simplified printed Arabic.)

The four letter shapes are usually called *initial*, *medial*, *final*, and *independent*, but these names are somewhat misleading. They tend to suggest that the initial and final forms are used only at the beginning and ends of words, with the medial form used elsewhere and the independent form used only in one-letter words and when talking about a letter by itself. However, some letters don't have all four forms, leading to breaks in the middle of a word. The letter alef (ا), for example, doesn't connect to the letter that follows it, so the initial form, rather than the medial form, of a letter is used when it follows alef. For example, despite the breaks in the line, this is all one word:

الأفغاني

Even in simplified Arabic, there are a couple of special cases. The letter lam (ل) and the letter alef (ا) would look like this if their normal initial and final forms were used:

لا

But they don't actually join that way. Instead of forming a U shape, the two vertical strokes actually cross each other, forming a loop at the bottom like this:

لا

If the lam is joined to another character on its right, however, you don't get the loop, however. Still, instead of forming a U shape, you get a somewhat more zigzaggy form like this:

لا

These two combinations are often called the *lam-alef ligatures*, and are required, even in simplified Arabic printing. The U shape never happens.

In simplified Arabic printing, the *kashida*, normally a curved line, becomes a straight line, in order to harmonize better with the letters that surround it.

This style of Arabic printing loses some of the grace and legibility of the more traditional forms⁴⁴...



...but it's still completely legible to Arabic speakers, and a lot simpler to deal with in computerized or mechanized systems. Even in languages such as Urdu and Farsi that have traditionally used the more ornate Nastaliq form of the Arabic alphabet, the simplified letterforms are often used today.

At least in some places, unlike the other alphabets we've looked at so far, Arabic doesn't use the so-called Arabic numerals ("1234567890"). Instead, *real* Arabic numerals (the Arabs call these "Hindi numerals," since they originated in India, even though the Indian forms are different yet again) look like this:

١ ٢ ٣ ٤ ٥ ٦ ٧ ٨ ٩ ٠

Interestingly, even though Arabic writing runs from right to left, numbers, even when written with the native Arabic numerals, run from left to right. 1,234.56 looks like this:

⁴⁴ This example is also from the Milo paper, section 2, p. 5.

١٢٣٤٥٦

There are actually two forms of native Arabic digits. The following forms are used in Urdu and Farsi:

۱ ۲ ۳ ۴ ۵ ۶ ۷ ۸ ۹ ۰

The Arabic block

The Unicode Arabic block, which runs from U+0600 to U+06FF, is based on the international standard Arabic encoding, ISO 8859-6, which is in turn based on the an encoding developed by the Arab Organization for Standardization and Metrology (or “ASMO” for short). The characters they have in common have the same relative positions in both encodings. A couple of the long vowels with the *hamza* and *madda* are encoded as precomposed forms, but the *hamza* and *madda*, as well as the vowel marks, are encoded as combining marks. The normal Arabic usage, like the normal Hebrew usage, is to use combining character sequences to represent the marked letters. Not counting the marked forms, each letter has a single code point—it’s up to the text rendering process to decide what glyph to draw for each letter. The *teh marbuta* and *alef maksura*, whose use is a matter of spelling and are therefore not mere contextual forms, get their own code points. The *kashida* also gets a code point, although you should generally rely on the rendering software to insert *kashidas*, rather than putting them into the encoded text.

Although Arabic punctuation is generally unified with the normal Latin punctuation marks, some of them have such different shapes in Arabic that they’re given their own code points in this block. This block also contains two complete series of Arabic digits: one set for the versions used with Arabic and one set for the versions used with Urdu and Farsi. Even though some of the digits have the same shape, Unicode provides two complete series of digits to simplify number formatting algorithms (the algorithms that convert numeric values to their textual representations).

Unicode rounds out the Arabic block with a selection of Koranic annotation marks and a wide selection of additional letters and punctuation marks used to write various non-Arabic languages using the Arabic alphabet. The Unicode standard doesn’t give a complete list of supported languages, but they include Adighe, Baluchi, Berber, Dargwa, Hausa, Kashmiri, Kazakh, Kirghiz, Kurdish, Lahnda, Malay, Pashto, Persian (Farsi), Sindhi, Uighur, and Urdu.

Joiners and non-joiners

There are two special Unicode characters that deserve special mention here: U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER. Like the directionality-control characters we looked at earlier in this chapter, these characters are invisible and they’re transparent to most processes that operate on Unicode text. They’re used for one and only one thing: they control cursive joining. The zero-width non-joiner breaks the cursive connection between two letters that would otherwise connect. In other words...

```
U+0644 ARABIC LETTER LAM
U+0647 ARABIC LETTER HEH
U+062C ARABIC LETTER JEEM
```

...look like this...

لهج

...but if you add in the non-joiner...

```
U+0644 ARABIC LETTER LAM
U+200C ZERO WIDTH NON-JOINER
U+0647 ARABIC LETTER HEH
U+200C ZERO WIDTH NON-JOINER
U+062C ARABIC LETTER JEEM
```

...the letters break apart, giving you this:

لهج

This is useful not only for putting together examples in books like this, but has real uses in some languages. In Farsi, for example, the plural suffix isn't connected to the word it modifies.

The zero-width joiner does the opposite thing: it forces a cursive connection where one wouldn't normally occur. In other words, if you take...

```
U+0647 ARABIC LETTER HEH
```

...which normally looks like this by itself...

ه

...and put a zero-width joiner on each side, you get the medial form:

ه

The zero-width joiner and non-joiner can be used together to break up a ligature without breaking up a cursive connection. For example,

```
U+0644 ARABIC LETTER LAM
U+0627 ARABIC LETTER ALEF
```

normally looks like this:

لا

If you stick the zero-width non-joiner between them...

```
U+0644 ARABIC LETTER LAM
U+200C ZERO WIDTH NON-JOINER
U+0627 ARABIC LETTER ALEF
```

...it breaks up both the ligature and the cursive connection:

ل ا

If you actually wanted the letters to connect but not form the ligature (say, to produce the example above of how these two letters *don't* connect), you'd surround the non-joiner with joiners:

```
U+0644 ARABIC LETTER LAM
U+200D ZERO WIDTH JOINER
U+200C ZERO WIDTH NON-JOINER
U+200D ZERO WIDTH JOINER
U+0627 ARABIC LETTER ALEF
```

The joiner next to each letter causes the letter to take its cursive connecting form, but since the letters aren't next to each other anymore, they don't form the ligature. That gives you this:

ل ا

The Arabic Presentation Forms B block

The Arabic Presentation Forms B block, which runs from U+FE70 to U+FEFE, includes codes for the contextual forms of the Arabic letters used in simplified Arabic printing. That is, there are separate code-point values for the initial, medial, final, and independent forms of the letter beh, for example. There are also code points representing the various forms of the lam-alef ligature.

With a font that includes the glyphs for all of these code points, a system can implement a rudimentary form of Arabic cursive joining: the system can map the canonical versions of the Arabic letters from the main Arabic block to the glyphs using these code-point values.

The algorithm is a little more complicated than just mapping to the initial-form glyph when the letter's at the beginning of the word, the final-form glyph when the letter's at the end, and so on. This is because of the letters that never join on certain sides, such as alef, which never connects to the letter that follows it. The Unicode standard classifies all the Arabic letters into categories based on whether they only connect with letters on the left (left-joining), whether they only with letters on the right (right-joining), whether they connect with letters on both sides (dual-joining), or don't connect with anything (non-joining). The glyph-selection algorithm can use these categories to help figure out which glyphs to use: For example, if a dual-joining letter is preceded by a left-joining or dual-joining letter and followed by a right-joining or dual-joining letter, you'd use the medial-form glyph to represent it. But if a dual-joining letter is preceded by a left-joining or dual-joining letter, but followed by a *left*-joining letter, you'd use the final form. For example, consider our earlier example:

ل ه ج

All three letters are dual-joining letters. Since the letter heh in the middle, a dual-joining letter, is flanked by two dual-joining letters, we use the medial form to represent it. (Lam is a dual-joining letter, but we use the initial form because it's only got a letter on one side. We use the final form of jeem for the same reason.)

But if you were to substitute alef for lam, you get this:

ا هـ ج

Alef is a right-joining letter, which means we use the initial form of heh instead of the medial form. And since alef doesn't have a letter on its right-hand side, we use the independent form of alef.

The ArabicShaping.txt file in the Unicode Character Database contains the classifications of the letters that enable this algorithm to work right. Simple Arabic cursive shaping is discussed in more detail in Chapter 16.

The Arabic Presentation Forms A block

The Arabic Presentation Forms A block, which runs from U+FC00 to U+FDFF, contains a bunch of Arabic ligatures, that is, single code points representing various combinations of two or three Arabic letters, as well as a bunch of contextual forms of letter-diacritic combinations and a couple phrase ligatures (entire phrases encoded as single glyphs with single code points). These are here for compatibility with some existing systems, but are very rarely used. The basic idea is to produce more authentic-looking Arabic rendering by providing a wider selection of glyphs than you get with the normal four contextual forms of each letter. Because of the fluidity of the Arabic alphabet, however, a ligature-based approach doesn't work all that well: the number of ligatures necessary to draw things right can get huge really quickly, as each letter's shape potentially depends on the shapes of the letters on each side, something you can't really do at the ligature boundaries. Depending on the design of a font, some of the ligatures might be relevant, but there's probably no single font for which they'd all be relevant. Generally, fonts that do something more than the simple Arabic shaping described in the preceding section don't use the code points in the Arabic Presentation Forms A block—they use internal glyph codes that don't correspond to Unicode code points.

Two characters in this block, however—U+FD3E ORNATE LEFT PARENTHESIS and U+FD3F ORNATE RIGHT PARENTHESIS—are not compatibility characters. As their names suggest, they represent more ornate parentheses that are sometimes used in Arabic as alternatives to the regular parentheses.

The Syriac Alphabet

The Syriac language is a member of the Aramaic family of languages that is spoken by a variety of minority communities—mostly Christian—throughout the Middle East.⁴⁵ It's the principal liturgical language, and also serves as the vernacular for many members of these communities. Its alphabet is also used to write modern Aramaic, and many Syriac-speaking communities also use the Syriac alphabet for writing other languages, particularly Arabic (Arabic written using the Syriac alphabet is called *Garshuni*.) A theological schism among the Syrian Christians in the fifth century led to a bifurcation of the language: there are two dialects of Syriac, Eastern and Western, and the Eastern and Western communities have different, though related, writing styles.

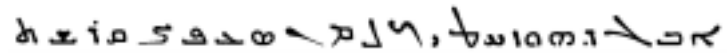
The origins of the Syriac alphabet are unclear, but it clearly belongs to the same family of scripts that include the Arabic and Hebrew alphabets. The alphabet has the same twenty-two letters as the

45 My sources for the section on Syriac are Peter T. Daniels and Robert D. Hoberman, "Aramaic Scripts for Aramaic Languages," in *The World's Writing Systems*, pp. 499-510, and the very detailed linguistic and historical information in the Syriac block description of the Unicode standard itself, pp. 199-205.

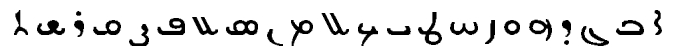
The Syriac Alphabet

Hebrew alphabet, and many of these letters have very similar forms to the forms they have in Hebrew, although Syriac letters are connected like Arabic letters. Like Hebrew and Arabic, Syriac is written from right to left. The oldest dated manuscript using the Syriac alphabet dates from AD 411.

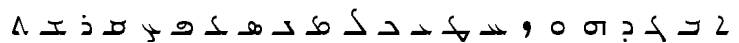
Again, owing to the split between the Syrian Christian communities, there are several different styles of the Syriac alphabet, and the shapes of many of the letters vary significantly between the styles. The oldest style is the Estrangelo. The twenty-two letters of the alphabet look like this in the Estrangelo style:



After the split, the styles diverged into the Serto...



...and the Nestorian:



Today all three type styles are used in both the Eastern and Western communities, although certain type styles continue to be favored by various communities.

As with the Arabic and Hebrew alphabets, the letters are all consonants, although several do double duty as vowels in certain contexts. As with Arabic, some of the letterforms are similar enough that the practice of using dots in various places to differentiate them developed early on.

Like Arabic and Hebrew, fully-vocalized text is represented by supplementing the basic consonantal writing with various types of vowel points. Unlike Arabic and Hebrew, however, it's more common, especially in modern Aramaic, for the points to be included than omitted.

Again because of the split in the Syrian Christian communities, there are two different systems of vowel points. The older of the two uses dots in various places to indicate the vowel sounds, and this system is still used exclusively in Eastern Syriac texts. The more recent system, attributed to Jacob of Edessa, uses marks based on the shapes of the Greek vowels. Western Syriac texts use a combination of this pointing system and the earlier version used by the Eastern communities. In addition, *Garshuni* uses the Arabic vowel signs instead of the Syriac ones. There is also a very complicated system of additional diacritical marks that are used with Syriac letters to mark various pronunciation and grammatical differences.

Historical Syriac scripts have a unique system of punctuation, but in modern Syriac and Aramaic texts, the punctuation marks used with Arabic are generally used. As with most other alphabets, numbers were historically written using the letters of the alphabet—when sequences of letters were used as numerals, they were marked with a barbell-shaped line above them⁴⁶:

⁴⁶ This example is ripped off from the Unicode standard, p. 200.



This mark is still sometimes used to mark abbreviations and contractions. In modern texts, the European digits are now used for numbers, just as they are with the Hebrew alphabet.

Syriac writing, like Arabic writing, is written cursively, with the letters connected together and spaces between the words. Properly rendering cursive Syriac shares many of the problems of properly rendering cursive Arabic, but systems can do a credible job using a system much like that frequently used to render Arabic: each letter has up to four glyphs, and the selection of a glyph depends on whether the letter joins to the letter preceding it, the letter following it, both, or neither. As with the Arabic alphabet, the Syriac alphabet contains a number of letters that only join to the letter that precedes them, meaning that a glyph-selection algorithm has to make use of the same types of combining classes that are used with simple Arabic rendering schemes. The Syriac letter alaph actually has *five* contextual glyph forms, so the rules for Syriac are a little more complicated. There is also a set of obligatory ligatures similar to the lam-alef ligatures in Arabic, but the actual set depends on the type style.

The Syriac block

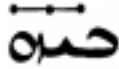
The Unicode Syriac block runs from U+0700 to U+074F. The type styles all use the same basic letters, so the letters are encoded once and you get one type style or another by selecting different fonts. Since both the dot and the Greek-letter forms of the vowel points are used together in Western Syriac texts, they're encoded separately in this block, rather than being unified.

There are a few extra letters beyond the basic twenty-two. These include a few letters used in *Garshuni*, two variant forms of the letter semkath (ܫ and ܫ) that are used in the same texts in different spelling situations, a ligature that's treated like a separate letter in some languages, and a special letter with ambiguous semantics. The letters dalath and rish (ܕ and ܕ) differ only in the placement of the dot. Some early manuscripts don't use the dot and thus don't differentiate between these two letters, so Unicode includes a code point value that looks like this (ܕ) that can be used for either letter.

The Syriac block also includes a wide selection of additional diacritical marks used with Syriac. A bunch of Syriac diacritical marks are unified with the general-purpose marks in the Combining Diacritical marks block; the Unicode standard gives a complete table of the marks from this block that are used with Syriac.

The Syriac block also includes the historical punctuation marks used with Syriac. The modern punctuation marks are unified with those in the Arabic and Latin blocks.

The Syriac block also includes a special invisible formatting character: the Syriac Abbreviation Mark. This character is used to represent the barbell-shaped mark that's used to mark numbers and abbreviations. It has no visual presentation of its own and is transparent to all text-processing algorithms, but marks the starting point of the abbreviation bar. The Syriac Abbreviation Mark precedes the first character that goes under the abbreviation bar, and the abbreviation bar extends from there to the end of the word. For example,



...which means “on the 15th,” (the bar is over the two letters being used as the digits in “15”)⁴⁷, is represented like this:

```
U+0712 SYRIAC LETTER BETH
U+070F SYRIAC ABBREVIATION MARK
U+071D SYRIAC LETTER YUDH
U+0717 SYRIAC LETTER HE
```

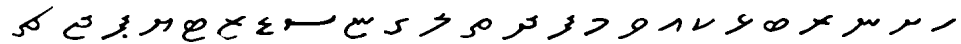
The basic cursive-joining algorithm used with Syriac is essentially the same as the one used with Arabic, with a few modifications to deal with various Syriac peculiarities (for example, the letter *alaph* (ܐ) has five contextual forms instead of four). The `ArabicShaping.txt` file in the Unicode Character Database includes joining-class designations for all the Syriac letters. The Unicode standard also gives a table of ligatures used with the various Syriac type styles.

Unicode doesn’t include compatibility code-point values representing the various glyph shapes of the Syriac letters as it does with Arabic. This means a simple Syriac shaping mechanism (one that doesn’t use the facilities in the advanced font technologies) would have to assign the various contextual forms to Private Use code points.

The Thaana Script

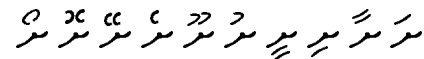
⁴⁷ This example also comes from the Unicode standard, p. 200.

Thaana is the name of the alphabet used to write Dhivehi, the language of the Republic of the Maldives, a group of islands in the Indian Ocean south of India.⁴⁸ Dhivehi is related to Sinhala, and the oldest examples of written Dhivehi, dating from around 1200, are written in a script called Evela, based on the Sinhala script of the time (for more on Sinhala, see Chapter 9). The current script, the Gabuli Tana, or simply “Thaana,” developed in the seventeenth century under the influence of the Arabic alphabet. The alphabet has twenty-four letters:



Interestingly, the first nine letters of the alphabet derive not from the Arabic letters, but from the Arabic *numerals*, with the next nine deriving from the Telugu numerals (see Chapter 9). The last six letters evolved out of Arabic letters.

The Thaana script is written from right to left, like Arabic, Hebrew, and Syriac. Like Hebrew, the letters aren’t cursively connected. Spaces are used between words. Like Arabic and Hebrew, the letters represent consonants, with vowels (or the absence thereof) represented through the use of diacritical marks written above or below the letters. There are ten vowel signs, or *fili*, here shown on the second letter of the alphabet, *shaviani* (سر):



In addition, there’s an eleventh “vowel sign,” the *sukun*, that represents the *absence* of a vowel:



One of the letters, *alifu* (ا), similarly represents the *absence* of a consonant sound, and is used as the base letter for word-initial vowels and vowels that follow other vowels. (Interestingly, this letter isn’t the first letter in the alphabet, like the corresponding letter in the other three alphabets we’ve looked at in this chapter). *Alifu* with *sukun* (ا̇)—that is, the null consonant with the null vowel attached to it—is used to represent the glottal stop.

Unlike in Arabic and Hebrew, the vowel marks are compulsory: you never see Thaana writing without the vowel signs.

In addition to the basic alphabet, there are fourteen extra letters, all based on the regular letters with extra dots attached, that are used for writing Arabic loanwords. There’s a history of writing Arabic loanwords using the Arabic alphabet, but this practice is beginning to die out.

Numbers in Thaana are written using either European or native Arabic numerals, with European numerals being somewhat more common. The Arabic decimal-point and thousands-separator characters are used in numerals regardless of whether European or Arabic digits are used.

48 My source for the section on Thaana is James W. Gair and Bruce D. Cain, “Dhivehi Writing,” in *The World’s Writing Systems*, pp. 564-568, supplemented with a couple tidbits from Nakanishi, pp. 38-39.

The Thaana Script

The same punctuation is used with Thaana now as is used with Arabic, although there's an older practice of using a period where we'd normally expect a comma (between phrases) and two periods where we'd use one period (at the end of a sentence).

The Thaana block

The Unicode Thaana block extends from U+0780 to U+07BF. It includes the basic twenty-four letters, the additional fourteen letters used for Arabic, the ten *fili*, and the *sukun*. Thaana punctuation is unified with the punctuation marks in the Arabic and Latin blocks, and Thaana also uses the digits from the ASCII or Arabic blocks.

CHAPTER 9 *Scripts of India and Southeast Asia*

[Compositor: I had a lot of trouble getting the examples into this chapter, and a bunch of them are still missing. I don't have a font for Sinhala, and I don't have software that can correctly render the contextual forms in some of the other scripts. For many of the scripts, the font's okay (at least for some things), but Word 2000 has a bug that prevents correct rendering: you get weird spacing behavior: combining marks show up after, rather than on top of, the characters they combine with, there's extra space after Indic syllables or some individual characters, and sometimes (especially in Tamil) there isn't enough space and I had to insert extra spaces to keep things from colliding. Most of these problems are said to be fixed in Office XP (at least when it's running on Windows XP, but I haven't had the money or time to upgrade. There are also a few cases where I had to switch fonts for the same script because the font I was using didn't form the combinations I was looking for.]

[I've placed red notes in the text indicating what's wrong with the corresponding examples. Please work with me to figure out how to do these.]

In this chapter, we'll look at the third major group of writing systems, usually referred to as the "Indic" group. This is the biggest single group of scripts in Unicode, comprising a whopping nineteen scripts. Fortunately, they have a lot in common.

The geographic area covered by the Indic scripts extends from Pakistan in the west all the way east to Cambodia and Laos, and from Tibet in the north to Sri Lanka in the south.⁴⁹ Several different language groups are represented in this area, although the writing systems are all related.

Just as the European alphabetic scripts are all more or less descended from Greek and the Middle Eastern bidirectional scripts are all more or less descended from Aramaic, the Indic scripts are all derived from the ancient Brahmi script, which dates from the middle of the third century. In the succeeding centuries, it spread across most of southern Asia, gradually developing different forms as it spread. The Brahmi script is said to have spawned over 200 different scripts. The fifteen scripts we look at here were more or less frozen into place as printing caught on in India in the eighteenth and nineteenth centuries.

The origins of the Brahmi script are rather unclear. There are schools of thought that say it developed from the early Semitic scripts, probably Phoenician or Aramaic, and schools of thought that say it was an indigenous Indian development. Although there isn't totally conclusive evidence, there seems to be considerable circumstantial evidence that Brahmi did indeed evolve in some way or another from a Semitic source, probably Aramaic, although the Indic scripts and the Middle Eastern scripts are pretty distant cousins now.

The diversity of scripts in South Asia seems to stem at least in part from the desire of the different ethnic groups to differentiate themselves from one another, particularly within India, in the absence of either international boundaries or major geographical features to separate them. A language, and thus the people who spoke it, was considered "major" if it had its own writing system. The religions predominant in this part of the world were (with some important exceptions) much more heavily dependent on oral tradition than on written tradition, the opposite of the religious groups in the other parts of the world, meaning that the imperative to keep a writing system stable so that its scriptures remain legible and have an unambiguous interpretation wasn't really there to countervail this tendency to diversify the scripts.

Perhaps the most interesting example of diversification of scripts is the issue of Hindi and Urdu. Urdu, the main language of Pakistan, is written using the Arabic alphabet (see Chapter 8). Hindi, one of the main languages of India and a very close cousin of Urdu, is written using the Devanagari script, which we'll look at in this chapter. In fact, Hindi and Urdu are so closely related that the difference in scripts is one of the main things that differentiates them.

Instead of literature and religion being the main things that drove the development of writing in this part of the world, administration and commerce seem to have been (although the scholars, priests, and literati were quick to take advantage as well). Even today, the bias in much of South Asia is toward spoken rather than written communication.

This chapter covers nineteen scripts from the Indic family:

- Devanagari, used for classical Sanskrit and modern Hindi, Marathi, and Nepali, plus a lot of less-common languages, is used in most of northern India and in Nepal.
- Bengali, used for writing Bengali, Assamese, and a number of less-prevalent languages, is used in eastern India and Bangladesh.

⁴⁹ My sources for the information in this introductory section are Colin P. Masica, "South Asia: Coexistence of Scripts," in *The World's Writing Systems*, pp. 773-776; Richard G. Salomon, "Brahmi and Kharoshthi" and the accompanying section heading, *op. cit.*, pp. 371-383; the introduction to Chapter 9 of the Unicode standard, pp. 209-210; and a few bits from Nakanishi, pp. 45-47.

- Gurmukhi is used for writing the Punjabi language spoken in the Punjab in northern India.
- Gujarati is used for writing Gujarati, the language of Gujarat in western India.
- Oriya is used to writing Oriya, the language of Orissa in eastern India.
- Tamil is used to write Tamil, spoken in southern India and Singapore.
- Telugu is used to write Telugu, spoken in southern India.
- Kannada is used to write Kannada (or Kanarese), also spoken in southern India.
- Malayalam is used to write Malayalam, spoken in Kerala in southern India.
- Sinhala is used to write Sinhala (or Sinhalese), the main language of Sri Lanka.
- Thai is used to write Thai, the language of Thailand.
- Lao is used to write Lao, the language of Laos.
- Khmer is used to write Cambodian, the language of Cambodia.
- Myanmar is used to write Burmese, the language of Burma.
- Tibetan is used to write Tibetan, the language of Tibet and Bhutan.
- Tagalog isn't used anymore, but was once used to write the Tagalog language spoken in the Philippines, as well as a few other Philippine languages.
- Hanunóo is used to write the language of the Hanunóo people on the island of Mindoro in the Philippines.
- Buhid is used to write the language of the Buhid people on the island of Mindoro in the Philippines.
- Tagbanwa is used to write the language of the Tagbanwa people on the island of Palawan in the Philippines.

There's a fair amount of variation in the Indic scripts, but they generally share a number of important characteristics:

- All are written from left to right.
- All are uncased: The letters of each only come in one flavor.
- Most, but not all, use spaces between words.
- Each script has its own set of numerals.
- All of these writing systems are somewhere between alphabetic and syllabic. A single "character" represents a whole syllable, but it has constituent parts that consistently represent the same phonemes from syllable to syllable.

The system is related to that used for the Middle Eastern languages. You have a basic character representing the combination of a consonant sound and a basic vowel sound. You change the vowel sound by attaching various marks to the basic character. If a syllable has more than one initial consonant, the basic characters for the consonant sounds involved usually combine together into something called a conjunct consonant. (Sometimes the basic consonant shapes don't combine, but instead a mark called a *virama*, or "killer," is attached to the first one to indicate the absence of the basic vowel sound.)

Without further ado, we'll delve into the characteristics of the individual scripts. In the section on Devanagari, we'll explore in detail just how this formation of syllable clusters works, and then we'll highlight only the differences between Devanagari and the other scripts in the subsequent sections.

Devanagari

The Nagari or Devanagari script (pronounced “deh-vuh-NAH-guh-ree,” not “duh-vah-nuh-GAH-ree”) is used across a wide swath of northern and western India and Nepal.⁵⁰ It’s used to write classical Sanskrit and modern Hindi, Marathi, and Nepali, as well as a large number of less-common languages (the Unicode standard lists Awadhi, Bagheli, Bhatneri, Bhili, Bihari, Braj Bhasha, Chhattisgarhi, Garhwali, Gondi, Harauti, Ho, Jaipuri, Kachchhi, Kanauji, Konkani, Kului, Kuamoni, Kurku, Kurukh, Marwari, Mundari, Newari, Palpa, and Santali). It developed from the northern variety of the Brahmi script and reached its current form somewhere around the twelfth century.

Devanagari, like the other writing systems in this chapter, is basically phonetic, with some kind of mark for each sound, but they can combine in complex and interesting ways. The word *purti* [find out what this means], for example, looks like this⁵¹:



The box marked “1” in the figure surrounds the “p,” the “2” marks the “u,” “3” and “4” together mark the “i” (we’ll see later why there are two numbers on this mark), “5” marks the “t,” and “6” marks the “i.”

The best way to think about Devanagari writing is to understand it as consisting of clusters of characters, each cluster representing a vowel sound, optionally preceded by a sequence of consonant sounds (the Unicode standard refers to these clusters as “orthographic syllables”). A cluster can be thought of as roughly analogous to a syllable, but they don’t correspond exactly: a consonant sound that’s spoken at the end of a syllable is actually written at the beginning of the cluster representing the next syllable. An Indic orthographic syllable is one kind of grapheme cluster, according to Unicode 3.2.

Devanagari includes thirty-four basic consonant characters:

क ख ग घ ङ
च छ ज झ ञ

⁵⁰ Most of this material is taken directly from chapter 9 of the Unicode standard, pp. 211- 223, but some is taken from William Bright, “The Devanagari Script,” in *The World’s Writing Systems*, pp. 384-390, and from Nakanishi, pp. 48-49.

⁵¹ The picture is taken from the Unicode standard, p. 14.

ट ठ ड ढ ण
त थ द ध न
प फ ब भ म
य र ल व
श ष स
ह

Every character has a horizontal line, called the *headstroke*, running across the top, and the character basically “hangs” from the headstroke. The headstrokes of the letters in a word all connect together into a single horizontal line. Indians, when they’re writing on lined paper, use the lines of the paper as the headstrokes of the letters (when writing quickly on lined paper, they often leave the headstrokes off the letters altogether). Each consonant also has a vertical stroke on the right-hand side (or sometimes in the center) that marks the core of the syllable. This is called the *stem*.

One interesting consequence of the connected headstrokes is that the “baseline” of Devanagari writing is the headstroke. Text of different point sizes is lined up so the headstrokes all connect:

[insert example]

When Devanagari is combined on the same line of text as Latin, or text in some other writing system where the baseline is under the characters, the text-rendering process (or the user) has to do some interesting adjustments to make the line come out looking good, especially if the point sizes on the line vary.

Unlike all of the other writing systems we’ve looked at so far, however, these symbols don’t merely represent consonant sounds. Instead, they’re said to carry an *inherent vowel sound*; in other words, they represent syllables. In Devanagari, the inherent vowel is *a*, so this symbol...

क

...represents the syllable *ka*. (This is actually pronounced closer to “kuh.”)

You change the vowel by adding an extra sign to the consonant symbol called a *dependent vowel*. Although some of the dependent vowels look more like letters, they basically behave more like diacritical marks. Some dependent vowels attach to the right-hand side of the consonant and just look like “the next letter” in the character stream, but one attaches to the left-hand side, giving the appearance that the characters have reordered from the order in which they were typed. There are several dependent vowel signs that attach to the bottom of the consonant and several more that attach to the top, and some of the vowels that appear to one side of the consonant also overhang it on top.

There are fourteen dependent vowel signs, here shown attached to the letter *ka*:

का कि की कु कू कृ कृ कँ को के कै काँ
को को कौ

The sound of the dependent vowel replaces the consonant's inherent vowel sound. For example, *ka* plus *i* gives you *ki*:

क + ि = कि

There are also fourteen *independent* vowel characters. The dependent vowel signs are basically abbreviated forms of the independent vowels. The independent vowels are used for vowel sounds that occur at the beginnings of words or follow other vowel sounds:

Independent vowels: अ आ इ ई उ ऊ ऋ ॠ ए ऐ ओ औ

Dependent vowels: ा ि िी ु ू ृ ृ ँ ै े ै ँ ो ो ौ

You can also cancel the consonant's inherent vowel sound by adding a mark called a *virama*. Thus, this represents the sound *k* by itself:

क्

Thus, if you just write the characters *sa* and *ta* next to each other, like this...

स त [the headstrokes of these characters should touch]

...you get "sata." But if you put a *virama* on the *sa*...

स्त

...you get *sta*.

Actually, you don't see the *virama* a lot, however. This is because when two consonants appear next to each other without an intervening vowel sound, they usually combine to form something called a *conjunct consonant*. The *virama* does show up on the character representing the final consonant in words that end with a consonant sound (actually, Hindi spelling isn't always phonetic: often, the word is spelled as though it ends with a consonant's inherent vowel, but the inherent vowel sound isn't actually spoken.)

The exact form a conjunct consonant takes depends on which consonants are joining together. Sometimes they stack on top of each other:

क + क = क्क

Often the first consonant (the one without the vowel, often referred to as the “dead” consonant) loses the stem and attaches to the next consonant, as in our “sta” example:

स + त = स् + त = स्त [these characters should connect]

The consonant symbol without its stem is called a *half-form*.

Sometimes, the two characters combine into a whole new glyph that doesn’t obviously resemble either of the characters that make it up:

क + ष = क्ष

The character र , representing the *r* sound, combines with other consonants in even more exotic ways. When the *r* is the first sound in a conjunct consonant, it appears as a hook on top of the second character, above the syllable’s stem:

र + प = र्प

This hook is called a *repha*.

When the *r* is the second sound in a conjunct consonant, it appears as a chevron underneath the second character:

घ + र = घ्र

Sometimes, the chevron representing the *r* sound combines graphically with the character it’s supposed to appear under:

क + र = क्र

Sometimes, र also combines with a dependent vowel attached to it:

र + ु = रु

र + ू = रू

More complex combinations can also occur. If *three* consonant sounds occur in a row, you might have two half-form consonants preceding the main consonant:

स + क + ल = ः + क + ल = स्कल

If the first two characters in the conjunct consonant would normally form a ligature, that ligature may turn into a half-form that combines with the third character:

क + ष + य = क्ष + य = ङ + य = ङ्य

And sometimes, all three characters will form a ligature:

स + त + र = स + त्र = स्त्र

There's also one case of a nonstandard half-form. The half-form of र sometimes looks like this: ः. This is called the “eyelash ra”:

र + न = ः + न = ः न [these glyphs should connect]

When a conjunct consonant combines with a dependent vowel, there are also strict rules for where the dependent vowel goes. If it attaches to the top or bottom of the consonant, it goes over or under the syllable's stem:

स्त्र िं = स्त्रे

If the dependent vowel attaches to the left-hand side of the consonant, it attaches to the extreme left of the conjunct consonant:

स्त्र िं = स्त्रिं

This can get a little weird when, for example, you have a top-joining vowel sign attaching to a conjunct consonant with a *repha*:

क िं = के

Most of the above examples are taken directly out of Chapter 9 of the Unicode standard, which gives a rigorous set of minimal rendering rules for Devanagari that cover all of the things we've looked at above.

Exactly which conjunct consonant forms you get when—that is, whether you use a ligature, a half-form, or just the nominal forms with viramas—depends on various factors, such as language, spelling, and font design. Depending on these things, you might see “kssa” look like this...

क्ष

...or this...

क ष **[again, these glyphs should connect]**

...or this...

क् ष **[yet again, these glyphs should connect]**

...and any of these is correct based on the rules of the script. The choice is purely one of font design.

There are a few marks that are used with the Devanagari characters. The *nukta*, a dot below the basic character, is used to modify the sounds of some of the consonants, widening the available selection of sounds that can be represented to cover all the sounds in modern Hindi. The *candrabinḍu* (ँ) and *anusvara* (ं) are used to indicate nasalization. The *visarga* (ः) represents the *h* sound.

Two punctuation marks are used with traditional texts: the *danda* (।), which functions more or less like a period, and the *double danda* (॥), which marks the end of a stanza in traditional texts. With modern texts, European punctuation is usually used.

As with all of the Indic writing systems, there is a distinctive set of digits used for writing numbers:

० १ २ ३ ४ ५ ६ ७ ८ ९

The Devanagari block

The Unicode Devanagari block runs from U+0900 to U+097F. It's based on the Indian national standard, ISCII (the Indian Script Code for Information Interchange). The ISCII standard includes separate encodings for Devanagari, Bengali, Gurmukhi, Gujarati, Oriya, Tamil, Telugu, Kannada, and Malayalam, switching between them with escape sequences. All of these encodings share the same structure, just as the various scripts do. The same code point value in each code page represents the same basic sound: for example, the value 0xB3 represents the character for ka in each of the various scripts, making transliteration between them simple.

The corresponding blocks in Unicode are all the same size, and the characters have the same positions within these blocks (and thus the same relative ordering) that they do in the ISCII standard, again facilitating transliteration between them.

[The Unicode blocks are based on the 1988 version of ISCII. ISCII has since been revised, bringing it slightly out of sync with Unicode, but converting between Unicode and ISCII is still pretty simple.]

Unicode supplements the ISCII characters with various additional characters. As much as possible, Unicode follows the rule of coding characters that correspond between the related scripts in the same positions in each of their blocks.

Unicode follows the characters-not-glyphs rule religiously with regard to the Indic scripts. Each character is given a single code point value, even though it can have many shapes.

The nominal consonants, complete with their inherent vowels, are given single code point values in Unicode. Thus, the syllable *ka* is represented as a single code point. That is, the sequence...

```
U+0915 DEVANAGARI LETTER KA
```

...shows up like this:

क

A dependent-vowel code is used to change the vowel of a syllable. Thus, *ko* is encoded like this...

```
U+0915 DEVANAGARI LETTER KA  
U+094B DEVANAGARU VOWEL SIGN O
```

...and looks like this:

को

The characters are stored in logical order, the order they're typed and pronounced. Thus, *ki*, which looks like this...

कि

...is still encoded like this...

```
U+0915 DEVANAGARI LETTER KA  
U+093F DEVANAGARI VOWEL SIGN I
```

...even though the vowel sign appears to the left of the consonant.

Conjunct-consonant combinations are represented using U+094D DEVANAGARI SIGN VIRAMA after the nominal consonant characters, even when the virama doesn't actually show up in the rendered text. Thus,

क

...is represented as...

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+0937 DEVANAGARI LETTER SSA
```

...and this sequence...

स्त

...is represented as the following:

```
U+0938 DEVANAGARI LETTER SA
U+094D DEVANAGARI SIGN VIRAMA
U+0924 DEVANAGARI LETTER TA
```

Conjunct consonants involving three consonant sounds require a virama after both of the first two sounds.

स्त्र

...is represented as...

```
U+0938 DEVANAGARI LETTER SA
U+094D DEVANAGARI SIGN VIRAMA
U+0924 DEVANAGARI LETTER TA
U+094D DEVANAGARI SIGN VIRAMA
U+0930 DEVANAGARI LETTER RA
```

The zero-width joiner and non-joiner characters we looked at in the previous chapter can be used to control the display of conjunct consonants. As we saw earlier, the following sequence...

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+0937 DEVANAGARI LETTER SSA
```

...is usually rendered as...

क

The zero-width non-joiner can be used to split the conjunct consonant up, causing it to appear as two nominal consonants with a virama on the first one. In other words...

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+200C ZERO WIDTH NON-JOINER
U+0937 DEVANAGARI LETTER SSA
```

...gets rendered as...

क् ष [these glyphs should connect]

You can surround the non-joiner with joiners to split up the ligature but keep the conjunct-consonant rendering. This causes the syllable to be drawn using a half-form for the *ka* sound:

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+200D ZERO WIDTH JOINER
U+200C ZERO WIDTH NON-JOINER
U=200D ZERO WIDTH JOINER
U+0937 DEVANAGARI LETTER SSA
```

...get rendered like this:

क् ष [these glyphs should connect]

When a consonant with a *nukta* is represented as a combining character sequence (a couple precomposed consonant-nukta combinations are encoded), the *nukta* follows the consonant immediately, in the manner of all other Unicode combining marks. The *candrabindu* and *anusvara*, on the other hand, apply to a whole syllable. They follow the whole syllable in memory.

To tie things together, then, the word “Hindi” would be typed and encoded like this...

```
U+0939 DEVANAGARI LETTER HA
U+093F DEVANAGARI VOWEL SIGN I
U+0928 DEVANAGARI LETTER NA
U+094D DEVANAGARI SIGN VIRAMA
U+0926 DEVANAGARI LETTER DA
U+0940 DEVANAGARI VOWEL SIGN II
```

...and looks like this:

हि न्दी [*sigh* these pieces should also connect]

This breaks down into two syllable clusters. The first is *hi*:

हि

The *ha* (ह) combines with the *i* sign, which attaches on the left. The second cluster is *ndii*:

न्दी

The *na* (न) turns into a half-form, which attaches to the *da* (द) on the left. The *ii* vowel (ी) then attaches to the combined consonant on the right. The virama doesn't appear as a separate glyph; instead, it simply tells the system to combine the *na* and *da* into a conjunct consonant.

The example we opened the chapter with, *purti*, is encoded like this...

```
U+092A DEVANAGARI LETTER PA
U+0942 DEVANAGARI VOWEL SIGN UU
U+0930 DEVANAGARI LETTER RA
U+094D DEVANAGARI SIGN VIRAMA
U+0924 DEVANAGARI LETTER TA
U+093F DEVANAGARI VOWEL SIGN I
```

...and looks like this:

पू ति **[again, the two clusters should connect]**

Again, it breaks down into two clusters. The first is *puu*:

पू

The *uu* sign (ू) attaches to the bottom of the *pa* (प). The second cluster is *rti*:

ति

The *ra* (र) turns into a *repha* and attaches to the top of the *ta* (त). The *i* vowel (ि) then attaches to the left-hand side of the combined consonant. Again, the virama doesn't show up as a glyph, but just tells the text renderer to treat *ra* and *ta* as a conjunct consonant.

[it'd be nice to add a longer example that shows an initial vowel, a three-sound conjunct consonant, and a conjunct consonant that forms a ligature rather than using a half-form]

Note, by the way, that the clusters don't correspond directly to spoken syllables. "Hindi," for example is spoken "hin-di," but written as "hi-ndi".

Bengali

The Bengali, or Bangla, script is a close cousin of Devanagari and bears a close resemblance to it. Like Devanagari, the characters have a horizontal headstroke that connects from character to character in a continuous horizontal line, and like Devanagari, a vertical stem marks the core of each syllable.

Bengali script is used in Bangladesh and northeastern India to write not only Bengali, but also Assamese, Daphla, Garo, Hallam, Khasi, Manipuri, Mizo, Munda, Naga, Rian, and Santali.

Bengali has thirty-two consonants...

ক খ গ ঘ ঙ চ ছ জ ঝ ঞ ট ঠ ড ঢ ণ ত থ দ ধ ন প ফ
ব ভ ম য র ল শ ষ স হ

...and twelve independent vowels:

অ আ ই ঈ উ ঊ ঋ ঌ এ ঐ ও ঔ

One interesting difference between Bengali and Devanagari is that Bengali has two *split vowels*. These are vowel signs that consist of two glyphs, written on either side of the consonant character.

Here are the twelve dependent vowel signs, attached to the letter *ka* (ক):

কা কি কী কু কূ ক্ ক় কে কৈ কো কৌ

[This approximates the desired appearance, but the left-joining and right-joining marks should connect to the consonants, and the bottom-joining vowels should appear under the consonants rather than to the right.]

The Bengali virama looks like this (attached to ক):

ক্

[again, the virama should display under the consonant, not to its right]

The Bengali anusvara looks like this: ঁ The Bengali candrabindu looks like this: ং The Bengali visarga looks like this: ঃ

Bengali forms conjunct consonants in the same way as in Devanagari, with many of the conjunct consonants having similar forms. The special ways in which the consonant representing the *r* sound combines with other consonants also apply in Bengali, although the shapes are a little different.

The Bengali digits look like this:

০ ১ ২ ৩ ৪ ৫ ৬ ৭ ৮ ৯

In addition to the regular digits, there are a series of special characters used to write fractions, a couple currency symbols, and the *ishar* (ॐ), which is used to write the name of God.

The word “Hindi” written using Bengali characters looks like this:

[example—can’t do without proper shaping behavior]

The Bengali block

The Bengali block in Unicode runs from U+0980 to U+09FF. Like the Devanagari block, it follows the ISCII order, which means corresponding characters in the two blocks are found at corresponding positions within their blocks.

The split vowels can be represented in Unicode either as single code points or as pairs of code points, one representing the left-hand glyph and the other representing the right-hand glyph. The representation that uses one code point is preferable, but has a compatibility decomposition to the two-code-point representation. (When two code points are used for a split vowel, they both appear in memory after the consonant.)

Gurmukhi

The Gurmukhi, or Punjabi, script is used to write the Punjabi language of the Punjab region of northern India.⁵² It evolved from the now-obsolete Lahnda script, of which someone once said, “It is convenient, with only one fault; it is seldom legible to anyone except the original writer.”⁵³ A revision was undertaken under the second Guru of the Sikhs, Guru Angad, in the sixteenth century, producing the modern Gurmukhi script. The name Gurmukhi attests to this effort: it means “from the mouth of the Guru.”

Gurmukhi resembles Devanagari; like Devanagari, the characters all hang from a continuous horizontal line, although the vertical stem that characterizes most Devanagari consonants isn’t as frequent in Gurmukhi. Gurmukhi has thirty-three basic consonants...

ਸ ਸ਼ ਹ ਕ ਖ ਗ ਘ ਙ ਚ ਛ ਜ ਝ ਵ ਟ ਠ ਡ ਢ ਣ ਤ ਥ ਦ
ਧ ਨ ਪ ਫ ਬ ਭ ਮ ਯ ਰ ਲ ਲ਼ ਵ

...and ten independent vowels:

ਅ ਆ ਇ ਈ ਉ ਊ ਐ ਐ ਓ ਔ

52 My sources for information on Gurmukhi are Harjeet Singh Gill, “The Gurmukhi Script,” in *The World’s Writing Systems*, pp. 395-398, and Nakanishi, pp. 50-51.

53 The quote is from Nakanishi, p. 50, but doesn’t carry an attribution there.

The ten independent vowels here aren't completely separate forms as they are in Devanagari and Bengali; instead, there are three “vowel-bearers” for the three classes of vowels, and the regular dependent vowel signs are attached to them. There are nine dependent vowels, here shown on **ਕ** (as in the other scripts, the consonants carry an inherent *a*):

ਕਾ ਕਿ ਕੀ ਕੁ ਕੂ ਕੇ ਕੈ ਕੌ ਕੇ

Conjunct consonants aren't as common in Gurmukhi as they are in Devanagari and Bengali. The symbols **ੳ**, **ੲ**, and **ੳ** are attached to the other consonant in the cluster as subscripts when they appear as the second consonant in the cluster:

ਖ + ਰ = ਖੁ

The character **ਯ** takes an abbreviated form when it appears as the second consonant in the cluster:

ਦ + ਯ = ਦੁ

[this font uses a half-form rather than forming a ligature, as the example describes: see TWWS, p. 395]

Generally, in other consonant combinations an explicit virama (◌◌) is used. Gurmukhi spelling isn't always phonetic; in many cases, the virama is simply understood and not written.

For comparison, this is the word “Hindi” written in Gurmukhi characters:

ਹਿ ਹੁਦੀ

[as in the other examples, the syllables should touch]

There are a couple of extra marks that are used with Gurmukhi characters: the *bindi* (◌◌) and *tippi* (◌◌) indicate nasalization; they're roughly analogous to the *anusvara* and *candrabindu* in

Devanagari. There's also a special sign, the *addak* (◌◌), which indicates the doubling of a consonant sound. Like the other scripts, Gurmukhi also uses a *nukta* (a dot below) as a diacritic with some consonants to allow for more consonant sounds.

Punjabi is a tonal language, but there aren't any special characters to represent tone. Instead, some of the consonant letters do double duty as tone indicators.

There is also a unique set of Gurmukhi digits:

੦ ੧ ੨ ੩ ੪ ੫ ੬ ੭ ੮ ੯

The Gurmukhi block

The Gurmukhi block in Unicode runs from U+0A00 to U+0A7F and follows the same ISCII order followed by the Devanagari and Bengali blocks, with analogous characters in analogous positions in all the blocks. The additional consonants used in Punjabi, which are the regular consonants with *nuktas* added, also get their own code points in this block; they have canonical decompositions to the consonant-*nukta* forms.

The Gurmukhi block includes codes for the unadorned vowel-bearers, meaning an independent vowel could be represented using a vowel-bearer and a dependent vowel sign. This isn't recommended, however, as Unicode doesn't include a decomposition from the regular independent-vowel representation to the representation including the vowel-bearer.

Gujarati

The Gujarati script is used to write the Gujarati and Kacchi languages of the Gujarat region of western India.⁵⁴ The earliest known document using Gujarati characters dates from 1592. It's the most closely-related of the Northern Brahmi-derived scripts to Devanagari. Most of the letterforms are quite similar, but the script at first appears strikingly different from Devanagari (and the other scripts we've looked at so far) because the characters lack the horizontal headstroke connecting the letters of a word (a line of Gujarati type of different sizes still has the letters aligned at the top, however).

Gujarati has thirty-four consonants, which carry an inherent a...

ક ખ ગ ઘ ઙ ચ છ જ ઝ ઞ ઠ ડ ઢ ણ ત થ દ ધ ન પ ફ
ભ મ ય ર લ વ શ ષ સ હ

...and thirteen independent vowels:

અ આ ઈ ઈ ઉ ઊ ઋ ઋ ઌ ઌ ઓ ઓ ઔ ઔ

As in the other scripts we've looked at, the independent vowel forms are used only at the beginnings of words or when they follow other vowels. The rest of the time, the dependent vowels, which attach themselves to the consonants, are used. There are thirteen dependent vowel signs (here shown

attached to ક):

ક િ ક િ ક િ ક િ ક િ ક િ ક િ ક િ ક િ ક િ ક િ
ક િ

54 My source for Gujarati is P. J. Mistry, "Gujarati Writing," in *The World's Writing Systems*, pp. 391-398.

ଅ ଥା ଇ ଈ ଉ ଊ ଋ ଌ ଏ ଐ ଓ ଔ

As with the other scripts, Oriya consonants carry an inherent *a* sound, and the other vowels are represented as diacritical marks when they follow a consonant. The ten dependent vowel signs are as follows (on କ):

କା କି କିଂ କିଂ କିଂ କିଂ ଋ କିଂ ଋ କିଂ କିଂ

[same spacing problems: the bottom- and top-joining vowels should be closer to, or combined with, the consonants, and the left- and right-joining vowels should be closer to the consonants]

Oriya also has conjunct consonants, and they take a wide variety of forms. In addition, the dependent vowel signs often form ligatures with the consonants they join to, giving rise to a very wide variety of Oriya letterforms.

Here's "Hindi" written using Oriya characters:

[example—need font that does the shaping right]

The same diacritical marks are used with Oriya as are used with Bengali. Oriya also has a set of distinctive digits:

୦ ୧ ୨ ୩ ୪ ୫ ୬ ୭ ୮ ୯

The Oriya block

The Oriya block in Unicode runs from U+0B00 to U+0B7F. Like the other ISCII-derived blocks, it follows the ISCII order.

Like some other scripts we've looked at, Oriya has split vowels, including one with components that appear to the left, to the right, *and* above the consonant. Unicode includes two characters, U+0B56 ORIYA AI LENGTH MARK and U+0B57 ORIYA AU LENGTH MARK, for the pieces of the split vowels that aren't dependent vowels themselves (*au* is the vowel sign with three components: the *au* length mark comprises both the top and right-hand components). All of the split vowels have canonical decompositions to their pieces and so can be represented using either single code points or pairs of code points.

Tamil

The scripts we've looked at so far—Devanagari, Bengali, Gurmukhi, Gujarati, and Oriya, form the North Brahmi group of scripts. We'll now turn to the South Brahmi scripts, which descend from the ancient Brahmi script along a line that diverged fairly early from the North Brahmi scripts. The scripts in this group are Tamil, Telugu, Kannada, Malayalam, and Sinhala. With the exception of

Sinhala, the main languages these scripts are used to write belong to the Dravidian language group, which is completely unrelated to the Indic language group to which the others belong—their adoption of the same basic writing system as the Indic languages stems from their geographic proximity, and they have some differences from the North Brahmi scripts that reflect the differences in their underlying languages.

The Tamil script (rhymes with “camel,” not “Camille” **[double-check pronunciation]**) is used to write the Tamil language of the Tamil Nadu region of southern India, in addition to several minority languages, including Badaga.⁵⁵ It follows the same basic structure as the other Indic scripts. Tamil has fewer sounds, so the Tamil script has fewer characters than the others. There are twenty-two consonants...

க ங ச ஐ ஞ ட ண த ந ன ப ம ய ர ற ல ள ழ வ ஷ
ஸ ஹ

[bsaically correct, but the spacing between the characters is uneven]

...and twelve independent vowels:

அ ஆ இ ஈ ஐ உ ஊ எ ஏ ஐ ஒ ஓ ஔ

As with the other scripts, the independent vowels are used at the beginnings of words, but they are also sometimes used within a word to represent an especially prolonged vowel sound. Most of the time, however, the eleven dependent vowel signs (here shown with **ப**) are used when the vowel sound follows a consonant:

பா பி பீ பு பூ பெ பே பை பொ போ
பெள

Like many of the other scripts, it’s worth noticing that several of the dependent vowel signs are drawn to the *left* of their consonants, and that three of them have components that appear on *both sides* of the consonant.

As with the other scripts, there are only eleven dependent vowel signs because the consonant characters have an inherent *a* sound. Like the other scripts, the absence of the inherent vowel is represented using a virama, which takes a different form than in the other scripts:

ப்

⁵⁵ Most of the information on Tamil comes straight out of the Unicode standard, but has been supplemented by information from Sanford B. Steever, “Tamil Writing,” in *The World’s Writing Systems*, pp. 426-430.

One important thing that sets Tamil apart from the other writing systems we've looked at so far in this chapter is that it generally doesn't use conjunct consonants. Instead, the normal consonant forms are usually used, and the virama is actually visible. There is one exception to this rule: When க் is followed by ஷ, instead of getting this...

க்ஷ

...you get this:

க்ஷ

This form behaves like conjunct consonants in the other scripts we've looked at, but it's the only conjunct consonant in Tamil.

However, the fact that Tamil doesn't have conjunct consonants doesn't mean that Tamil letters don't form ligatures. It's just that the ligatures are formed between consonants and vowels, rather than between successive consonants. Tamil actually has a considerable number of consonant-vowel ligatures. A few examples:

ண + ா = ணா

ன + ா = நா

ற + ா = ரா

ட + ி = டி

ல + ி = லி

ல + ீ = லீ

க + ற = க்ற

க + ள = க்ள

ஐ + ற = ஐ்ற

ஐ + ள = ஐ்ள

...and so on...

There are also two forms of the vowel றை, which changes shape in front of certain consonants (this is a left-joining vowel):

க + ீ = கை ...

...but ண + ீ = னை

[Not totally sure I have the right glyph here (didn't get it automatically): should this ligature have another loop?]

Like the other scripts in this chapter, there's a unique set of Tamil digits:

௧ ௨ ௩ ௪ ௫ ௬ ௭ ௮ ௯

However, they don't form a normal set of decimal digits used with positional Arabic notation. Instead, they're used with multiplicative notation in conjunction with signs for 10 (௩), 100 (௩௩), and 1,000 (௪௪):

௩ = 3

௩௩ = 13

௩௩௩ = 30

௩௩௩௩ = 33

This system of numeration is actually nearly obsolete; European numerals are generally used now.

The Tamil block

Like the other ISCII-derived blocks, the Unicode Tamil block (U+0B80 to U+0BFF) follows the ISCII order. The three split vowels all have canonical decompositions to their constituent parts. The right-hand side of ீ, which isn't used as a vowel sign by itself, is encoded as U+0BD7 TAMIL AU LENGTH MARK.

Telugu

The Telugu script (pronounced "TELL-a-goo," not "te-LOO-goo" [double-check pronunciation]) is used to write the Telugu language of the Andhra Pradesh region of India, as well as minority languages such as Gondi and Lambadi. It follows the same basic structure as the other Indic scripts, and the letters are characterized by a chevron-shaped headstroke.⁵⁶ There are thirty-five consonants...

⁵⁶ Most of the examples in this section are taken from Nakanishi, pp. 60-61.

క ఖ గ ఘ జ చ ఛ జ ఝ ఞ ఠ డ ఢ ణ త థ ద ధ న ప
ఫ బ భ మ య ర ల ల ళ వ శ ష స హ

...and fourteen independent vowels:

అ ఆ ఇ ఈ ఉ ఊ ఋ ౠ ఎ ఏ ఐ ఒ ఓ ఔ

As always, the independent vowels are used only at the beginnings of words. In the middles of words, there are thirteen dependent vowel signs, here shown on క :

కా కి కీ కు కూ కృ కౄ కౌ క్షా క్షా క్షా
కా క్షా

[again, shaping isn't happening: the right-joining vowels should be closer to the consonant, the top- and bottom-joining vowels need to be shifted to the left, and the consonant should lose the V on top when it combines with a top-joining vowel]

As always, the consonants carry an inherent *a* sound. It can be canceled using a virama...

క[◌]

[again, the virama should appear on top of the consonant, and the consonant should lose the V]

...but multiple consonant sounds with no intervening vowel sounds are usually represented using conjunct-consonant forms instead. Generally, the second consonant is written as a subscript under the first consonant...

[can't do these examples: need a font that can do the necessary shaping]

[ssa] + [ttha] = [ssttha]
[da] + [da] = [dda]

...but sometimes they take other forms:

[ka] + [ka] = [kka]
[sa] + [ta] + [ra] = [stra]

Notice that the consonants generally lose their headstrokes when they appear in combination. This also generally happens when they're combined with a dependent vowel sign.

For comparison, here's "Hindi" written using Telugu characters:

[example]

And like the other scripts, Telugu has its own set of digits:

౦ ౧ ౨ ౩ ౪ ౫ ౬ ౭ ౮ ౯

The Telugu block

Like the other ISCII-derived blocks, the Unicode Telugu block (U+0C00 to U+0C7F) follows the ISCII order. There's one split vowel, although the component parts appear above and below the consonant. The bottom-joining part is encoded separately as U+TELUGU AI LENGTH MARK, and the split vowel has a canonical decomposition that includes this character.

Kannada

The Kannada script (**[pronunciation?]**) is used to write the Kannada, or Kanarese, language of the Karnataka region of southern India, as well as minority languages such as Tulu. It is very closely related to the Telugu script; the two both derive from Old Kannada script, and diverged around 1500, with the two different forms being hardened into place by the advent of printing in the nineteenth century.⁵⁷

The basic letterforms of Kannada are extremely similar to those of Telugu, with the main difference being the shape of the headstroke. Like Telugu, Kannada has thirty-five consonants...

ಕ ಖ ಗ ಘ ಙ ಚ ಛ ಜ ಝ ಞ ಠ ಡ ಢ ಣ ತ ಥ ದ ಧ ನ ಪ
ಫ ಬ ಭ ಮ ಯ ರ ಱ ಲ ಳ ವ ಶ ಷ ಸ ಹ

...and fourteen independent vowels:

ಅ ಆ ಇ ಈ ಉ ಊ ಋ ೠ ಎ ಏ ಐ ಒ ಓ ಔ

The shapes of the thirteen dependent vowels differ more from their counterparts in Telugu. Here they are attached to ಕೆ :

⁵⁷ This information comes from William Bright, "Kannada and Telugu Writing," in *The World's Writing Systems*, pp. 413-419.

కా కే కై కృ కౌ క్క క్కే క్కొ క్కొ
కా

[Again, the font's not doing the right shaping: the top-joining vowels aren't showing up on top of the consonant (and the consonant should lose the headstroke when they do), and the right-joining vowels are too far from the consonant. Because of this, I can't do the other examples in this section without locating a better font.]

Kannada forms conjunct consonants in the same basic way as Telugu, and some of the dependent vowel signs form ligatures with certain consonants, just as they do in Telugu. "Hindi" looks like this in Kannada letters:

[example]

One interesting wrinkle of Kannada that's different from Telugu is the behavior of ರ, which represents the *r* sound (which takes irregular forms in so many of the other Indic scripts). When this character is the *first* character in a conjunct consonant, it changes shape and appears *after* the rest of the syllable:

kaa = ಕ + ీ = [kaa]
kraa = ಕ + ರ + ీ = [kraa]
rkaa = ರ + ಕ + ీ = [rkaa]

[Does it really come after the whole syllable, or between the consonant and the vowel?]

As with the other scripts we've looked at, there's a unique set of Kannada digits:

೦ ೧ ೨ ೩ ೪ ೫ ೬ ೭ ೮ ೯

The Kannada block

The Unicode Kannada block runs from U+0C80 to U+0CFF and follows the ISCII order, just like the other scripts we've looked at so far. There are several split vowels that have canonical decompositions to pairs of vowels. The Kannada block has two special code point values that represent the right-hand components of the split vowels.

The Kannada block includes one misnamed character. U+0CDE KANNADA LETTER FA should really have been called KANNADA LETTER LLLA, but since the Unicode standard never renames a character, we're stuck with the wrong name. Fortunately, this is an obsolete character only used in historical writings.

The *anusvara* (◌̣) is used to represent the *m* sound at the end of a syllable, instead of using a conjunct-consonant form.

The word “Hindi” looks like this in Malayalam:

[example]

A script reform movement in the 1970s and '80s has led to a simplified version of the script with fewer conjunct forms and the modification of the vowel signs that appear below the letters so that the script is easier to print, consisting of a series of independent glyphs that can be written linearly. However, it never completely caught on, giving rise to a bunch of hybrid printing styles that use the simplified forms for some things and the traditional forms for others. As a result, the appearance of the same piece of Malayalam text can vary fairly significantly depending on font design.

Unlike the other scripts we've looked at, spaces aren't always used between words in Malayalam. Some documents do this, but some run sequences of words together without spaces and only use spaces between phrases.

Like the other scripts, there is a set of Malayalam-specific digits:

൭ ൧ ൨ ൩൪ ൫ ൬ ൭ ൮൯

The Malayalam block

The Unicode Malayalam block runs from U+0D00 to U+0D7F. It follows the ISCII order and the same encoding principles as the other scripts we've looked at in this chapter. Choice between traditional and simplified rendering is purely a matter of font design and not reflected in the encoding, although it's possible to break up a conjunct consonant and use the virama forms of the consonants by using the zero-width non-joiner (see the “Devanagari” section earlier in this chapter for more information).

Sinhala

The Sinhala script is used to write the Sinhala (or Sinhalese) language of Sri Lanka, and is also used in Sri Lanka to write Sanskrit and Pali. Although Sinhala is an Indic language, the Sinhala script belongs to the South Brahmi family of scripts used to write the Dravidian languages.⁵⁹ It follows the same general structure as the other Indic scripts, with a set of consonants, and separate sets of dependent and independent vowels. There are forty-one consonants...

[example]

⁵⁹ Information in this section comes from James W. Gair, “Sinhala Writing,” in *The World's Writing Systems*, pp. 408-412, and from Nakanishi, pp. 66-67.

...although only twenty-two of them are used to write regular Sinhala. The others are used to write loanwords, or to write Pali or Sanskrit. Likewise, there are eighteen independent vowels...

[example]

...but only twelve of them are used to write native Sinhala words.

As with the other scripts we've looked at, the independent vowels are used only at the beginnings of words. After consonants, vowel sounds are represented using the seventeen dependent vowel signs (here shown on **[ka]**) attached to the consonants:

[example]

As with the other scripts, the consonants carry an inherent *a* sound, which can be canceled using the virama (called *al-lakuna* in Sinhala):

[ka + virama]

The *anusvara* (**[glyph]**) is used in Sinhala not only to indicate nasalization, but in modern spelling to indicate an actual *n* sound at the end of a syllable.

Like many of the other Indic scripts, Sinhala has a fairly complex system of conjunct consonants, and many of the dependent vowel signs form ligatures with certain consonants. These include special forms of the *r* consonant when it appears at the beginning or end of a conjunct-consonant cluster similar to the ones used in Devanagari and some of the other North Brahmi scripts.

In traditional texts, spaces were not used between words, and the only punctuation used was the *kunddaliya* (**[glyph]**), which was used where we'd use a period. In modern writing, spaces are used between words and European punctuation is used. There is also a traditional set of Sinhala digits; these, too, have passed into obsolescence: today, European digits are used. Some modern printing practices eschew many of the conjunct-consonant clusters in favor of forms explicitly using the *al-lakuna*.

The Sinhala block

The Unicode Sinhala block runs from U+0D80 to U+0DFF. It does *not* follow the ISCII order. This is partially because the ISCII standard doesn't include a code page for Sinhala and partially because Sinhala includes a lot of sounds (and, thus, letters) that aren't in any of the Indian scripts. The basic setup of the block is the same: *anusvara* and *visarga* first, followed by independent vowels, consonants, dependent vowels, and punctuation. Unlike in the ISCII-derived blocks, the *al-lakuna* (virama) precedes the dependent vowels, rather than following them.

The same basic encoding principles used for the ISCII-derived scripts are used by Unicode for encoding Sinhala: The character codes encode characters, not glyphs, so there are no special code-point values for the conjunct ligatures or contextual forms of the vowels. The characters are stored in logical order, as they are pronounced, meaning vowel signs always follow their consonants, even when they're drawn to their left. The code points for the consonant characters include the inherent *a* sound; the *al-lakuna* is used to cancel the inherent vowel and signal the formation of conjunct-consonant clusters (when conjunct consonant clusters are formed, the *al-lakuna* will not show up as a

separate glyph, even though it's typed and stored). The joiner and non-joiner characters can be used to control the formation of ligatures and consonant conjuncts.

Consonant forms carrying a *nukta* are used for writing Tamil using the Sinhala script; these aren't separately encoded. Interestingly, the Sinhala block doesn't include a code-point value for the *nukta*. The Unicode standard points this out, but doesn't say what to do about it. Since the other Indic blocks each include their own *nuktas*, it's unclear which one to use with Sinhala, if indeed any of them is appropriate. The actual answer will depend on individual implementations and fonts. The Unicode Sinhala block includes a code point for the *kunddaliya*, even though it's not used in modern writing, but doesn't include code point values for the Sinhala digits.

Thai

So much for the Indian scripts. The next four scripts we'll look at in this chapter are used in Southeast Asia for various Asian languages. These scripts all descend from the South Brahmi forms, but show greater diversification, reflecting the greater diversity of languages represented.

The Thai script is used to write the Thai language of Thailand.⁶⁰ The Thai and Lao scripts both have their roots in a script developed in 1283 by King Ramkhamhaeng of Sukhothai. They had diverged into separate scripts, with different spelling rules, by the sixteenth century.

Thai has a whopping forty-six consonants:

ก ข ฃ ค ฅ ฉ ง จ ฉ ช ซ ฌ ญ ฎ ฏ ฐ ฑ ฒ ณ ด ต ถ ท ธ น บ ป ผ ฝ พ ฟ ภ ม ย
ร ฤ ล ฦ ว ศ ษ ส ห พ อ ฮ

Interestingly, Thai only has twenty-one consonant sounds. The extra letters represent sounds that aren't distinguished anymore in modern pronunciation, but still carry etymological information and are used to distinguish tone.

Thai doesn't have independent vowel signs like the other Indic scripts. Instead, the letter อ, which represents the glottal stop, is used with the dependent vowel signs to write an initial vowel sound.

Thai consonants are generally considered to carry an inherent *o* sound, with the other vowel sounds represented with the twenty-one dependent vowel signs, here shown on อ:

อะ ั อ่า อี้ อี้ อี้ อี้ อู่ อู๋ เออะ เอ เอ็ แอะ แอ โอะ โอ เออะ อ่า โอ เออ

Thai, like many Asian languages, is a tonal language. Different consonants carry different inherent tones, but there's a set of four tone marks as well:

⁶⁰ The information in this section and the next comes from Anthony Diller, "Thai and Lao Writing," in *The World's Writing Systems*, pp. 457-466, with some additional info coming from Nakanishi, pp. 76-79.

อ ออ ออ ออ

When a consonant has both a dependent vowel sign and a tone mark on it, the tone mark moves up and gets smaller to make room for the vowel sign:

อ

Thai doesn't have conjunct consonants. There's a virama (๑) that's used to cancel the inherent vowel when the Thai script is used to write Pali, but normal Thai doesn't use a virama—successive consonants are just juxtaposed and you have to know when to pronounce the inherent vowel and when not to.

This isn't a big deal, since Thai spelling in general isn't phonetic—there are only three sounds that can occur at the end of a syllable, but most of the letters can be used in this position, for example. Sometimes there are extra letters at the end of a word that aren't pronounced; sometimes the *thanthakhat* (๑) is written over a silent letter to indicate that it's silent.

There are two special symbols that are used like letters in the middle of words: the *maiyamok* (๑) is used to represent repetition of the syllable that precedes it, and the *paiyannoi* (๑) is used to indicate elision of letters from long words or phrases. The sequence ๑๑๑ means “etc.”

There are also a bunch of special Thai punctuation marks: the *fongman* (๑) is used the way we use a bullet: to mark items in lists, or at the beginnings of verses or paragraphs. The *angkhankhu* (๑) is used at the ends of long segments of text, the sequence ๑๑ marks the end of an even longer segment of text, such as a verse in poetry. The *khomut* (๑) may follow this sequence at the end of a paragraph or document.

Spaces don't generally separate words in Thai text (although space is sometimes added between words when necessary to justify text between the margins); instead, spaces are used between phrases or sentences, where you'd use a comma or period in English.

Finally, there's a set of Thai digits:

๐ ๑ ๒ ๓ ๔ ๕ ๖ ๗ ๘ ๙

The Thai block

The Unicode Thai block runs from U+0E00 to U+0E7F. It's based on TIS 620, the Thai national standard.

In Lao, as in Thai, there are multiple consonants with the same sound but different inherent tone.

Some consonants that don't have counterparts with the opposite tone are written with ຫ, which is silent but changes the tone. Sometimes, this letter forms a ligature with the letter carrying the actual consonant sound:

ຫ + ນ = ຫນ

ຫ + ມ = ຫມ

The symbol ັ is used to represent the *l* or *r* sound (the letter ັ) when it occurs after another consonant:

ພ + ັ = ພ ັ [the "lo" mark should appear underneath the consonant, not beside it]

These forms are the only remaining remnants of an earlier system of conjunct consonants; other than these forms, Lao is written like Thai: the characters for successive consonant sounds are just juxtaposed, with the reader understanding when the inherent *o* is not pronounced. Like Thai, Lao doesn't use a virama. Like Thai, Lao doesn't generally use spaces between words, using them instead where we'd normally use a comma.

Lao, like Thai, has an ellipsis character (⋯) and a repetition character (⋮). There is also a distinctive set of Lao digits:

໐ ໑ ໒ ໓ ໔ ໕ ໖ ໗ ໘ ໙

Although both the languages and writing systems are similar, the Laotian government instituted a set of spelling reforms in the 1960s, with the result being a divergence in spelling between Lao and Thai. Thai spelling preserves etymological details at the expense of phonetics. Lao spelling is phonetic at the expense of etymology.

With the exception of the conjunct-consonant forms, all of the letters in the Lao script have direct counterparts in the Thai script, making the two scripts freely convertible (at least when used to write Lao, since Thai has more letters). There is a sizable Lao-speaking minority in Thailand; they generally use the Thai script to write their language.

The Lao block

Since the Thai and Lao scripts are so closely related and convertible, the Unicode Lao block (U+0E80 to U+0EFF) is organized with the characters in the same relative positions within their block as the analogous characters in the Thai block. Lao follows the same encoding principles as Thai, with characters stored in visual order rather than logical order. Thus, left-joining vowel signs precede their consonants in memory, even though they're spoken after the consonants. Split vowels are represented using two separate character codes: one, representing the left-joining part of the vowel, precedes the consonant in memory, and the other, representing the right-joining part of the

[as always, Word isn't honoring the overhang on the vowel signs: the right-joining vowels should be closer to the consonants and the top- and bottom-joining vowels should actually be above and below the consonants]

...but these represent different sounds depending on whether they're used with a first-series or second-series consonant. Sometimes it's necessary to put a second-series vowel sound on a first-series consonant, or vice-versa. Diacritical marks are used to do this: The *muusikatoan* (្ក្រ) converts a second-series consonant to a first-series consonant, and the *triisap* (្ក្រ) converts a first-series consonant to a second-series consonant.

Like most of the other Indic scripts, Khmer does have conjunct consonants. Normally, the second consonant in a pair is written as a subscript under the first one...

[ka] + [ta] = [kta]

...but sometimes the subscript has an ascender that sticks up to the right of the main consonant...

[tho] + [yo] = [thyo] [Is this a good example?]

...and sometimes the subscript has a completely different form from its normal form:

[ko] + [vo] = [kvo] [Is this a good example?]

The letter representing the *r* sound (្ក្រ), when used as a subscript, picks up an ascender that sticks up to the *left* of the main consonant:

[ko] + [ro] = [kro]

The *r* character thus appears to *precede* the main consonant even though its sound *follows* it (this is analogous to the left-joining dependent vowels).

The vowel-shifter characters *muusikatoan* (្ក្រ) and *triisap* (្ក្រ) can change shape in certain circumstances. If the consonant has an ascender that sticks up into the space these characters normally occupy, or if there's a vowel sign in that position, the character takes a different shape and moves *under* the consonant:

[sa] + [ii] + [triisap] = [sii] + [triisap] = [sii (subscript triisap)]

The *nikahit* (្ក្រ) represents nasalization, the *reahmuk* (្ក្រ) represents a final glottal stop, and the *bantoc* (្ក្រ) shortens the preceding vowel.

Like Thai and Lao, Khmer is written without spaces between words; instead, spaces are used more or less where we'd use a comma. The *khan* (្ក្រ) is used at the ends of sentences or groups of sentences. The *baryoosan* (្ក្រ) marks the end of a section. There is a distinctive set of Khmer digits:

០ ១ ២ ៣ ៤ ៥ ៦ ៧ ៨ ៩

The Khmer block

The Unicode Khmer block runs from U+1780 to U+17FF. It includes not only the regular Khmer letters, but a series of rare letters and symbols used only in transcribing Sanskrit and Pali. It also includes a wide selection of punctuation marks.

The encoding principle is the same as the Indic scripts (and different from Thai and Lao): the vowel signs always follow the consonants, even when the consonant glyph appears before or around the consonant. Conjunct consonants are formed by using U+17D2 KHMER SIGN COENG, the Khmer analogue to the Indic virama. This code point has no visual presentation (the glyph shown in the Unicode standard is arbitrary); its function is only to signal the formation of a conjunct consonant. Even though the following (subscript) form of *r* appears to the left of the preceding consonant, it's still represented in memory following the main consonant (and the *coeng*).

If there's a series-shifter in a syllable, it follows all the consonants in the cluster, and precedes the vowel. Any other marks on the cluster come last.

Myanmar

The Myanmar script, which is used to write Burmese, the language of Burma, dates back to the early twelfth century. It is descended from the Mon script, which ultimately derives from the South Brahmi scripts.⁶² As with the other scripts in this chapter, it follows the same basic structure as the other Indic scripts. There are thirty-four consonants:

က ခ ဂ ဃ င စ ဆ ဇ ဈ ည ဋ ဌ ဍ ဎ ဏ တ ထ ဒ ဓ န
ပ ဖ ဖ ဘ မ ယ ရ လ ဝ သ ဟ ဌ

There are seven dependent vowel signs (here shown on အ), each of which represents a combination of a vowel and a tone:

အာ အိ အီ အု အူ ဧ အဲ

[This is closer than I've gotten with most of the other combining-vowel examples, but still not there: the one left-joining vowel should be closer to the consonant, and the top- and bottom-joining vowels should be centered on their consonants, not skewed to the right]

62 My sources for Myanmar are Julian K. Wheatley, "Burmese Writing," in *The World's Writing Systems*, pp. 450-456, and Nakanishi, pp. 72-73.

The letter ခ is a neutral consonant that is used to represent independent vowel sounds, but there are also a few independent vowel signs, used generally in Indian loanwords.

The dependent vowel sign ဘ takes a different shape when its normal shape would cause the consonant-vowel combination to be confused with another consonant:

[pa] + **[a]** = **[paa]**

Myanmar has conjunct consonants. Most of the time, the second consonant is written as a subscript under the first one...

[ma] + **[ma]** = **[mma]**

...but the consonants **[ya]**, **[ra]**, **[wa]**, and **[ha]**, representing the *y*, *r*, *w*, and *h* sounds, take special diacritic forms when they're used in conjuncts:

[kha] + **[ya]** = **[khya]**
[ma] + **[ra]** = **[mra]**
[ma] + **[wa]** = **[mwa]**
[ra] + **[ha]** = **[rha]**

These characters can also participate in three-consonant clusters:

[la] + **[ha]** + **[ya]** = **[lhya]**

The *anusvara* (◌̃) and *visarga* (◌̃) from the Indic languages, along with a subscript dot (◌̣) are used as tone marks. A virama, or “killer” (◌̣), cancels a consonant’s inherent *a* sound when a word ends with a consonant sound and occasionally in the middles of words.

Burmese is written with spaces between phrases rather than words. For punctuation, the danda (၊) and double danda (။), corresponding roughly to the semicolon and period, are borrowed from the Indic languages. There is also a unique set of Myanmar digits:

၀ ၁ ၂ ၃ ၄ ၅ ၆ ၇ ၈ ၉

The Myanmar block

The Unicode Myanmar block (U+1000 to U+109F) follows the same basic principles as the other Indic blocks: the dependent vowel signs always follow their consonants in memory, even when the vowel is drawn to the left of the consonant; each letter gets only one code point value, even if its shape changes depending on context; and conjunct consonants are formed by placing the virama between the consonants. (The zero-width non-joiner can be used to break up a conjunct in the middle of a word and force the virama to be visible.)

[ra] + [ka] = [rka]

The consonants ར , ལ , འ , ཡ , and ཅ adopt an abbreviated shape when they appear below certain other consonants:

[ka] + [la] = [kla]
[ka] + [ya] = [kya]
[ka] + [ra] = [kra]
[ka] + [wa] = [kwa]

Certain consonants just go before a consonant stack, and certain ones just go after it. A consonant after the main consonant of the syllable might represent a word-final consonant rather than part of the stack, in which case the vowel signs go on the consonant before it. If the main consonant of the syllable keeps its inherent *a* (i.e., it doesn't have a vowel sign), another consonant after it is treated as a syllable-ending consonant. If it's not supposed to—that is, if the syllable ends with the inherent *a* sound—the syllable will end with ཅ . For example,

ར ལ ར is *dag...*

...but ར ལ འ ར is *dga*.

[spacing is broken again—the horizontal headstrokes on all the characters should be touching]

A syllable can consist of up to four consonants in a row, with the main consonant optionally being a consonant stack. (In most words like this, most of the letters are actually silent; this reflects the fact that Tibetan spelling hasn't changed in centuries.)

Syllables (words, basically) are separated from one another with dots called *tshegs*, which are aligned along the baseline like periods in English text (Tibetan is written hanging from the baseline, like Devanagari, rather than sitting on it, like English, though, so the *tsheg* appears on the upper-right shoulder of the last character in the syllable). When text is justified between margins, extra *tshegs* are added at the end of the line to pad it out to the margin.

The rules for how consonant stacks form are very regular for Tibetan, but a much wider variety of conjuncts can be formed in Sanskrit and other languages written using Tibetan characters. There is also a complex system of abbreviation used in Tibetan writing that can break most of the normal character-joining rules.

Tibetan is normally written without spaces, although the *tsheg* serves the same basic purpose. There is a wide variety of punctuation marks used with Tibetan, the most important of which are the *shad* (།) and the *nyis shad* (༎), which derive from the Devanagari danda and double danda. They're very roughly equivalent to the comma and period: the *shad* marks the end of an expression, and the *nyis shad* marks a change in topic. There is also a unique set of Tibetan digits:

༠ ༡ ༢ ༣ ༤ ༥ ༧ ༨ ༩

The Tibetan block

The Unicode Tibetan block is quite large, running from U+0F00 to U+0FFF. In addition to the basic Tibetan letters, it includes extra letters used for transliterating various other languages and a very wide variety of punctuation marks, cantillation marks, and astrological signs.

The Tibetan block follows a radically different encoding practice from the other Indic blocks. Instead of using a virama character to signal the formation of conjunct consonants, each consonant is encoded twice: The first series of code point values is to be used for consonants when they appear by themselves or at the top of a conjunct stack, and the second series is to be used for consonants that appear in a conjunct stack anywhere but at the top (only certain combinations are legal in Tibetan, but all the consonants are encoded this way to allow for non-Tibetan conjuncts, such as those you get in Sanskrit). Thus, **[kra]** (*kra*), which would be encoded like this in Devanagari (and most of the other Indic scripts)...

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+0930 DEVANAGARI LETTER RA
```

...is instead encoded like this in Tibetan:

```
U+0F40 TIBETAN LETTER KA
U+0FB2 TIBETAN SUBJOINED LETTER RA
```

Since the letter འ can change shape when it appears at the top of a conjunct stack, something you don't always want in some non-Tibetan languages, Unicode includes a separate code point value, U+0F6A TIBETAN LETTER FIXED-FORM RA, to represent a འ that keeps its full size and shape all the time. Likewise, since ར , ལ , and ཤ change shape when they appear at the bottom of a conjunct stack, Unicode provides separate code point values—U+0FBA TIBETAN SUBJOINED LETTER FIXED-FORM WA, U+0FBB TIBETAN SUBJOINED LETTER FIXED-FORM YA, and U+0FBC TIBETAN SUBJOINED LETTER FIXED-FORM RA—that also keep their normal shape no matter what. These alternate code point values should never be used in encoding native Tibetan.

The Philippine Scripts

Unicode 3.2 adds four additional scripts to write various languages of the Philippines. These four scripts are sandwiched into a comparatively small area of the Unicode encoding space running from U+1700 to U+177F.⁶⁴

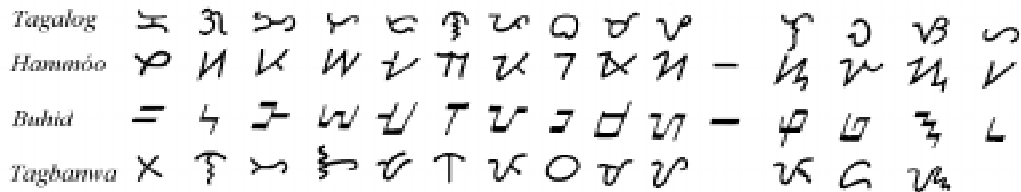
⁶⁴ Most of the small amount of information on the Philippine scripts here comes from section 9.16 of PDUTR #28, supplemented by some material from Joel C. Kuipers and Ray McDermott, "Insular Southeast Asian Scripts," *The World's Writing Systems*, pp. 474-484.

The four Philippine scripts are related to each other, probably coming from a common ancestor. They're clearly related to the Southeast Asian scripts, and probably came to the Philippines from Southeast Asia by way of Indonesia, possibly as late as the 14th century.

The Tagalog script fell into disuse in the 1700s under European colonization, superseded by the Latin alphabet. The other three scripts, Hanunóo, Buhid, and Tagbanwa, are all still used, but to varying degrees. Hanunóo is still relatively popular, used most often (interestingly enough) to write love poetry. Buhid and Tagbanwa are less often used.

The four Philippine scripts all have a similar structure. Each has three independent vowels...

...and from twelve to fourteen consonants:



[Will need to fix this drawing so the rows are straight]

As with most of the other Indic scripts, the consonants all carry an inherent *a* sound, and the vowel sound can be changed through the addition of a vowel mark. There are only two dependent vowel marks in the Philippine scripts (there are more vowel sounds; the vowel marks each represent more than one sound). In each script, the two marks are the same (depending on the scvript, a dot, dash, or chevron); whether the vowel mark is drawn above or below the consonant determines which sound it represents.

In all four languages, syllables have a consonant-vowel-consonant format. Most of the time, syllable-ending consonants are simply not written, although there are virama-like marks in Tagalog and Hanunóo. The Hanunóo virama (called the *pamudpod*) is used somewhat more often. In both Tagalog and Hanunóo, the virama is always visible; the consonants do not form conjuncts. In Hanunóo and Buhid, however, certain consonant-vowel combinations form ligatures.

The Philippine scripts are generally written from left to right, as the other Indic scripts are, but you also see them written *from bottom to top*, with the lines running from left to right. This seems to be a byproduct of the instruments usually used to write these scripts.

Punctuation marks similar to the danda and double danda in Devanagari are used with all four scripts.

In Unicode, the four scripts are encoded in four successive blocks in the range from U+1700 to U+177F. The Tagalog block runs from U+1700 to U+171F, the Hanunóo block from U+1720 to U+173F, the Buhid block from U+1740 to U+175F, and the Tagbanwa block from U+1760 to U+177F. Characters that correspond between the scripts are encoded at analogous positions in their respective blocks. The punctuation marks are encoded one for all scripts, in the Hanunóo block.

CHAPTER 10 *Scripts of East Asia*

Now we come to the last of the four major groups of modern scripts. This last group comprises the East Asian scripts. In particular, it covers the various writing systems used to write Chinese, Japanese, Korean, and (prior to the 1920s) Vietnamese. These scripts are often referred to collectively as the “CJK” scripts (for “Chinese, Japanese, and Korean”), or sometimes the “CJKV” scripts (for “Chinese, Japanese, Korean, and Vietnamese”).

This group of scripts is fundamentally different from the others. If you’ll remember, the European alphabetic scripts all either descended from the Greek alphabet or arose under its influence. The bi-di scripts of the Middle East all either descended from the ancient Aramaic alphabet or arose under its influence. The Indic scripts used in southern Asia all descended from the ancient Brahmi script. The Brahmi script is generally thought to have its roots the Aramaic script, and the Greek and Aramaic alphabets are both direct descendants of the ancient Phoenician alphabet. So the Phoenician alphabet can properly be thought of as the ancestor of all the writing systems we’ve looked at so far.

Not so the CJK scripts. These scripts all either descended from or arose under the influence of the Chinese characters. The Chinese characters have no known antecedent. They’re not known to have evolved from anybody else’s writing system. Just as the ancient Phoenician alphabet (well, actually its earlier ancestors) arose spontaneously in Mesopotamia with no outside influence, the Chinese characters arose spontaneously in China with no outside influence.

So the writing systems of East Asia are related to each other, but they’re not related to any of the other scripts we’ve looked at. They share some unique characteristics:

- The Chinese characters are used in almost every system we'll look at. Unlike the other scripts, where the characters stand for sounds, Chinese characters stand for whole words or ideas⁶⁵, so there are orders of magnitude more Chinese characters than there are characters in any of the other writing systems—they number in the tens of thousands. Needless to say, the individual characters can get quite complex. Moreover, there are multiple variant forms of many Chinese characters, and new Chinese characters are still being coined.
- All of the CJK languages, except for Yi, make use of some kind of system of phonetic characters to supplement the Chinese characters. How important this system is varies with language—the Chinese system is rarely used, the Korean phonetic system is used almost exclusively, and the Japanese phonetic systems are somewhere in the middle. The modern Yi script is completely phonetic, and doesn't actually use the Han characters, but arose under their influence.
- All of these writing systems make use of characters that are all the same size and shape (that is, the characters are all designed to fit into a square, and successive characters in a passage of text all fit into the same-size squares). The characters do not interact typographically—they don't change shape depending on context, they don't connect, and they don't form ligatures. Each character stands alone. A page of CJK text often looks like the characters were laid out on graph paper.
- All of these writing systems are traditionally written vertically. That is, the characters in a line of text start at the top of the page and work their way down, with lines of text starting on the right-hand side of the page and working their way to the left. Under the influence of Western writing and Western printing and computing equipment, however, all of these scripts are now also written like the Latin alphabet: on horizontal lines running from left to right and progressing from the top to the bottom of the page.
- Spaces are not used between words (in fact, the whole question of what constitutes a “word” in Chinese writing can be somewhat complicated), except sometimes in Korean.
- Many of these writing systems make occasional use of a system of interlinear annotation (i.e., notes between lines of text) to clarify the pronunciation or meaning of unfamiliar characters.

The Han characters

Since the Chinese characters form the cornerstone, in one way or another, of all the scripts we'll be talking about, we'll start by looking at them.

The Chinese characters are called *hanzi* in Mandarin Chinese, *kanji* in Japanese, *hanja* in Korean, and *chu Han* in Vietnamese. All of these expressions mean “Han characters,” the term we'll also use for them from here on out. Han, from the Han Dynasty, is used to refer to traditional Chinese culture. In other words, in Asia, “Han” more or less means “Chinese.”

⁶⁵ Not entirely true, but a decent approximation for now (and, in fact, the approximation that leads to their being called “ideographs”). Bear with me, and we'll get into the meat of the matter in the next section.

The Han characters have a very ancient history.⁶⁶ The earliest examples of Chinese writing date back to about 1200 BC. The Han characters in use today are directly descended from these early examples. The shapes of the characters have changed over time, and many more characters have been added, but the general principles underlying the system have remained constant through more than three millennia.

The Han characters are variously referred to as “ideographs,” “logographs,” and “pictographs,” all of which are sort of right and sort of wrong. “Ideographs,” the term the Unicode standard uses, implies that the characters simply stand for things and concepts and not for sounds. In practice, this is more or less true, especially for non-Chinese speakers, but for Classical Chinese (and, to a lesser degree, for modern Chinese languages such as Mandarin), the Han characters also have a strong phonetic function. “Logographs” implies that the characters stand for words. Again, this is close enough, although the simple “syllable = word = character” idea that many have of Chinese doesn’t always hold in practice. Polysyllabic words in Chinese are usually written with multiple characters, so it may be more accurate to think of the characters as “syllabographs.” “Pictographs” is just plain too simplistic. It implies that all the Han characters are little pictures of the things they represent. Some may have started out that way, and many of the simplest characters still work like that, but for most characters, the truth is more complicated.

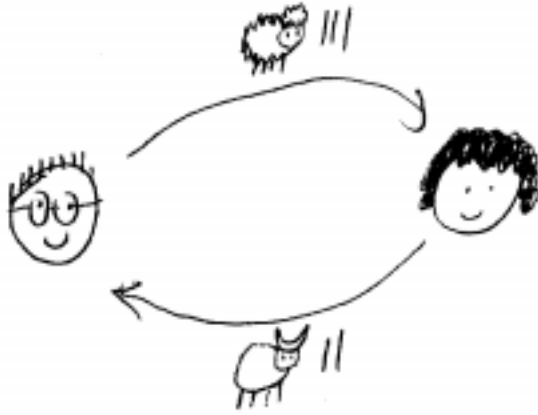
Western writing and Chinese writing probably started out in more or less the same way. Writing wasn’t invented to tell stories; it was invented to keep records.⁶⁷ Say you have ten sheep and three head of cattle and you want to make a record of how many pieces of livestock you have. So on some handy object you have, maybe a leaf or a piece of bark or a piece of clay, you draw a little picture of a sheep and follow it with ten scratch marks and a little picture of a bull and follow it with three scratch marks:



Now say you want to record a business transaction. Say you give your neighbor Joe three of your sheep in exchange for two of his cattle. You might draw something like this:

⁶⁶ My sources for information on the Han characters are William G. Boltz, “Early Chinese Writing,” in *The World’s Writing Systems*, pp. 191-199; Victor H. Mair, “Modern Chinese Writing,” *op. cit.*, pp. 200-208; Ken Lunde, *CJKV Information Processing* (Cambridge: O’Reilly, 1999), pp. 27-65; and Nakanishi, pp. 82-87.

⁶⁷ See, for example, Georges Ifrah (trans. David Bellos, E. F. Harding, Sophie Wood, and Ian Monk), *The Universal History of Numbers* (New York: Wiley, 2000), p. 77ff.



Here we have the basics of pictographic writing. The simplest Han characters are pictographs just like this:

人	man
木	tree
山	mountain
川	river
口	mouth

Of course, some of them you have to use your imagination with a little more than others:

女	woman
日	sun
月	moon

As often happens, the pictographs tend to become simpler and more stylized over time, looking less and less like the things they represent. The straight lines in the “sun” and “moon” characters, for example, reflect the fact that the tools these characters were originally written with didn’t do circles well.

In many cases, it’s fairly simple to come up with a picture to describe an abstract concept, giving rise to characters that can properly be called “ideographs”:

上	up
下	down

中 **middle**

一 **one**

二 **two**

三 **three**

You can also do variations on a single character to mean different things. For example, if you add an extra stroke to the “tree” character, you get “root”:

本 **root**

You can also use characters to mean one of a couple related concepts. For example, the “sun” and “moon” characters also represent “day” and “month.” The character for “root” also means “origin.”

Then you can combine characters. For example, take two trees and put them next to each other...

林

...and you get the character for “woods.” Add another tree...

森

...and you get the character for “forest.” Put “sun” and “moon” together...

明

...and you get the character for “bright.”

So far we haven’t talked about how the Han characters can be used to represent sounds; all the discussion so far as revolved around the semantic values of the characters—that is, how they are used to represent things or concepts. To understand how phonetics get involved, consider this example:



These three pictures represent the word “carpenter.” To get “carpenter,” you just take the names of the three objects and string them together: you have a picture of a car, a picture of a pen, and a picture of a tear running down someone’s cheek. This, as you probably know, is called a *rebus*, and examples abound. A favorite old chestnut is this one...

STAND

I

...which means “I understand” (that is, “I” under “stand”). Rebuses involving only the names of the letters and numbers abound on personalized license plates these days. Another oldie but goodie is

OU812

...which, of course, means “Oh, you ate one too?” and was the title of an old Van Halen album (yeah, I’m really dating myself here).

So if we follow the “rebus principle,” we see that a Han character can be used not only for the word the character originally meant to symbolize, and not just other words with similar or related meanings, but also for words that *sound like* the original word. For example, consider the character

王

This character means “king.” No, there’s nothing about it that really leaps out at you and says “king.” This is another example of the stylization these shapes undergo over time. The original version of this character⁶⁸ looks like a human figure standing on a line. This doesn’t really say “king,” either (maybe the line is a carpet or dais or something), but it’s closer.

Anyway, I’m digressing. The word for “king” in Chinese is “wang.” Given this, you could use this character for other words that are pronounced the same way. In this way, the same character can be used to mean “to go toward,” which is also pronounced “wang” (but with a different intonation). There’s no semantic or etymological connection between the two words; they just happen to be pronounced similarly. But because we have a character for “wang” meaning “king,” we can also use it for “wang” meaning “to go toward,” which doesn’t lend itself well to inventing an ideograph.

Another example is this character:

貝

This character originally meant “cowrie shell” (these were once used as currency), and was pronounced “bei.” This character also came to be used to represent the word “bai,” which meant “defeat.”

68 See the Boltz article in *The World’s Writing Systems*; table 14.1 on p. 192 and the discussion of the character on pp. 192-193, from which this example, and all of the following examples, were lifted.

So between using a character to represent a group of semantically related words (such as “root” and “origin”) and using a character to represent a group of words with similar pronunciations, you can greatly increase the expressive power of these characters.

This is the point where the development of Chinese writing goes off in a different direction from the direction the development of Western writing took. The Phoenicians took existing pictographs and began using them to represent their initial sounds. (This is kind of like the “Alpha, Bravo, Charlie” system used by the U.S. military for spelling out words.) For example, the word *beth* meant “house” and was represented by a little picture of a house. Over time, the character came to represent the “b” sound at the beginning of “beth,” instead of the whole word “beth” or the concept of “house,” and its shape also gradually changed to look less and less like a house (or anything else). This process eventually led to the Phoenician alphabet (and the picture of a house eventually evolved into the Latin letter B and practically every other alphabetic letter representing the “b” sound).

Chinese writing occasionally made feints in the same direction, but the Han characters never lost their semantic and etymological baggage. The characters, right down to modern times, have always represented both sounds and ideas.

Now if you’re continually making the characters do double, triple, or quadruple duty, you’re introducing ambiguity into the system. Up to a point, you can use context to determine which of several different words is intended for a particular character in a particular spot, but it gets unwieldy after a while.

The system that evolved with the Han characters involved using an extra character to resolve the ambiguity. If we go back to our “cowrie shell” example, we could set apart the “defeat” meaning by adding another character that clarified the meaning. The character chosen was this one...

女

...which means “strike” and is pronounced “pu.” Combine the two of them together and you get this:

敗

[fix to use Traditional Chinese form—the left-hand part should be the same as the example above]

The figure on the left represents the sound “bei,” and the figure on the right represents the concept of “strike,” and together you get “bai,” or “defeat.”

The same thing can also work in the other direction. Consider the character...

口

...which normally means “mouth” and is pronounced “kou.” It also came to be used to represent the word “ming,” which means “to call out.” This time, you have a semantic similarity with a difference in pronunciation. To resolve the ambiguity, they added this character...

夕

...which means “brighten” and is also pronounced “ming.” This gives us this:

名

Again, you have two components, one representing the sound “ming” and a second representing the concept “mouth.” Together they represent the word “ming” meaning “to call out.”

This process is actually recursive. You can now take the compound character and do the same thing with it. So we take our new compound character...

名

...representing the word “ming” and meaning “to call out” and use it to represent yet *another* word that’s pronounced “ming,” this time meaning “inscription (as on a bronze vessel)”.

Now to set it apart from the meaning “to call out,” you add another component:

金

This means “metal.” Add it to the other two...

銘

...and you get a unique character meaning “inscription.”

This process can theoretically continue ad infinitum, although in practice the characters generally become unreadable if they have more than five or six components. Characters with more than five or six components are quite rare.

So this is basically how the Han characters work. As you might imagine, this gives rise to a few interesting problems. First, what’s to keep different people from coming up with totally different characters for the same words? Two different people might follow the same basic rules, but use

different component characters, leading to two different, but both perfectly valid, characters for the same word. Who decides which one wins?

This problem is akin to the problem to standardizing spelling of words across a population that uses an alphabetic script, and the standard solution is the same: publish a dictionary. (Of course, this works best if you're a powerful emperor and can force people to follow the dictionary's prescriptions.) The first known dictionary of Han characters was compiled by Xu Shen in about AD 100, and it had the effect (over time, anyway) of standardizing on one way of writing each word. Xu Shen's dictionary, the *Shuo wen jie zi*, contained 9,353 characters. It began by dividing them into unit characters (those that couldn't be broken down into component parts) and compound characters (those that could). With the compound characters, Xu Shen put together a list of 540 component elements, at least one of which occurred in each compound character. These component elements are called "radicals," and one of them was classified as the main radical in every compound character. In this way, a hierarchy is imposed on the component parts of a compound character, and this both helps standardize the characters for each word and to ensure that the system that had developed by that time, using components that represent both sound and meaning, remains the standard way of writing Chinese words.

By the early eighteenth century, the number of characters had grown to almost 50,000, and the Kangxi emperor commissioned a new dictionary, the *Kangxi Zidian*. It was completed in 1716 and contained 47,021 characters. The Kangxi dictionary, like the many others before it, preserved the basic classification of characters according to radical first seen in the *Shuo wen jie zi*, but despite the huge growth in the number of characters, was able to classify them all using just 214 radicals. Unlike the *Shuo wen jie zi*, the Kangxi dictionary organizes characters within each radical group according to the number of additional pen strokes they contain beyond those used to write the radical.

The Kangxi dictionary remains the definitive authority for classical Chinese literature, and its system of 214 radicals continues to be used today to classify Chinese characters. In fact, this system of radicals is used in the radical-stroke index in the Unicode standard itself. (This system of organizing the characters by radicals and strokes works pretty well; even a non-Chinese speaker such as yours truly was able to find most characters in the Unicode radical-stroke index with only a modicum of effort.)

One of the beauties of the Han characters is that since they represent words, and by extension ideas, the same text can (with proper education, of course) be read and written by speakers of many different languages. (Of course, the farther the spoken language gets from classical Chinese, the less help the phonetic cues in the characters are.) In a phonetic writing system, someone who doesn't speak a particular language but uses the same script can still look at the writing and have at least a vague idea of how the words are pronounced, but probably little or no idea what they mean. With the Han characters, someone can look at something other than his native language written using the Han characters and get a vague idea of what it says, but have little or no idea how it's pronounced. It's interesting to compare the situations of China and India. Both are vast countries consisting of hundreds of millions of people of many different ethnic groups speaking many different languages. But in India, even relatively closely-related languages are frequently written with their own scripts. In China (with a number of important exceptions such as Inner Mongolia and Tibet) all the different spoken languages use the same written language.

(Of course, standard written Chinese consists of characters designed either based on Classical Chinese pronunciation or modern Mandarin pronunciation, and the characters are arranged according to Madarin grammar, so speakers of other Chinese languages, such as Cantonese, are effectively reading and writing a different language from the one they speak. There actually *are* Han characters

specific to Cantonese and other Chinese languages, but up until recently, there's been no comprehensive attempt to put them into any encoding standards. This began to change with Unicode 3.1, which includes many Cantonese-specific characters, but there's still much to do.)

Just as the various Indian scripts spread outside of India, the Han characters have come to be used outside of China. The Han characters were the first writing system used in both Japan and Korea, although both have since supplemented the Han characters with their own scripts. Until recently, the Han characters were also used in Vietnam. In all of these places, people also coined new characters using the same techniques the original ones were constructed with (but presumably according to Japanese or Korean pronunciation). In fact, the Vietnamese developed a whole system of characters, the *chu Nom*, for writing Vietnamese. These characters look and behave like Chinese characters, but are indigenous and were used instead of the traditional Chinese characters. The Japanese and Koreans didn't go this far; they used their own characters along with the traditional Chinese characters.

This points to one of the other great things about the Han characters—the system is very malleable, and it's easy to coin new characters for new words. Because of this, the number of Han characters has steadily grown over the years. This chart shows the numbers of characters in various Chinese dictionaries down through the years⁶⁹:

Year (AD)	Number of Chinese Characters
100	9,353
227-239	11,520
480	18,150
543	22,726
751	26,194
1066	31,319
1615	33,179
1716	47,021
1919	44,908
1969	49,888
1986	56,000
1994	85,000

Because of this, it's proper to think of the Han characters as an open-ended writing system. It's impossible to say just how many Han characters there are, because some occur so rarely, and new ones are constantly being invented.

Of course, this doesn't mean everybody knows every character—this would clearly be impossible. 2,400 characters is generally considered the lower limit for basic literacy, with most educated people having a working repertoire of four or five thousand and some especially literate people up to eight thousand. Studies have repeatedly shown that 1,000 characters cover 90% of written Chinese, 2,400 cover 99%, 3,800 cover 99.9%, 5,200 cover 99.99%, and 6,600 character cover 99.999% of written Chinese. Thus, the vast majority of Chinese characters are extremely rare, many so rare that neither their pronunciation nor their meaning is known—only their shape. Some may have only been used

69 This table is taken directly from Lunde, p. 58.

once or twice in all of history.⁷⁰ Unfortunately, many of these characters do occur, however rarely, in modern writing and it must be possible to represent (and thus input and print) them.

(It's worth pointing out that this is fairly analogous to words in English. The average English speaker has a working vocabulary roughly equivalent to the number of characters the average literate Chinese reader has in his working vocabulary. And new characters get invented at roughly the same rate as new English words get invented. No one regularly uses *all* the words in Webster's Third, and even though two people may have vocabularies of about the same size, they won't share all the same word.)

Variant forms of Han characters

Of course, despite efforts at standardization, variant forms of the characters do creep in over time. This tends to happen most often with more complicated characters: people come up with quicker and less-complicated ways of writing them, and some of these catch on and come into widespread use. Which forms catch on will tend to vary according either to a particular writer's preference, or, more often, according to regional convention.

As an extreme example, there are at least six variant forms of the character for "sword" in Japanese:

劍 劍 劒 劒 劒 劒

More important than this kind of thing, however, are national variants. Over time, some characters have developed different forms in Japan and Korean than they have in China, for example. In the example above, the first character is a Japanese-specific variant of the second character.

Of particular importance is the difference between Traditional Chinese and Simplified Chinese. In the 1950s, the Mao government in the People's Republic of China undertook a systematic effort to simplify written Chinese. Among the outcomes of this effort were the Pinyin system of writing Chinese words with Latin letters, and new versions of several thousand Han characters, many of which were so drastically simplified as to be unrecognizable given only their original forms.

Today, the new forms are used exclusively in the People's Republic of China, and the old forms are still used in Taiwan. The term "Simplified Chinese" has come to refer to Chinese as written in mainland China, and "Traditional Chinese" has come to refer to Chinese as written in Taiwan. These days the difference between Simplified and Traditional Chinese goes beyond merely using characters with different shapes; the choice of which ideographs to use in the first place also tends to be different between the locales. The traditional forms of the characters also tend to be used in Hong Kong

⁷⁰ These statistics and the following sentences are taken almost verbatim from the Mair article in *The World's Writing Systems*, pp. 200-201.

and among Chinese outside of China, while the simplified forms of the characters tend to be used in Singapore.

Japan and Korea tend to have their own forms, which sometimes match the Traditional Chinese form, sometimes match the Simplified Chinese form, and often are completely different. This means the same basic character can have as many as four different shapes, depending on the place where it's being used (the line between saying this and saying that the four different places have different, but historically related, characters for expressing the same concept is actually quite thin).

Here are some examples of characters with variant forms in Traditional Chinese, Simplified Chinese, and Japanese⁷¹:

Gloss	T.C.	S.C.	J
broad	廣	广	広
country	國	国	国
infant	兒	儿	児
east	東	东	東
open	開	开	開
see	見	见	見
picture	畫	画	画
old	舊	旧	旧
not/without	無	无	無
deficient/vacancy	缺	欠	缺
hall	廳	厅	庁
surround/enclose	圍	围	囲
abundant	豐	丰	豊

The mapping between Simplified and Traditional Chinese isn't always straightforward, either, as in many cases the language reformers mandated the *same* simplified character for several different words that originally had *different* traditional characters. For example, all of these characters in Traditional Chinese...

干 幹 乾 榦

⁷¹ These examples are taken from "Comparative Table of Sinitic Characters," *The World's Writing Systems*, pp. 252-258

...have the same representation in Simplified Chinese:

干

The four Traditional Chinese characters have widely varying meanings, but their pronunciations are all variations on “gan.” Here, rather than simplifying the more-complicated characters, they simply fell back on rebus writing.⁷²

Korean-specific variants are harder to find—generally the Korean forms of the Han characters follow the Traditional Chinese forms, but Korean font-design principles (things like stroke length and termination) tend to follow Japanese font-design principles rather than Chinese ones.

Han characters in Unicode

Unicode includes a truly huge number of Han characters: 70,195 different characters are included in Unicode 3.1. The Han character set in Unicode represents a truly huge effort by a lot of East Asian-language experts over a long period of time.

⁷² This example is taken from Jack Halpern, “The Pitfalls and Complexities of Chinese to Chinese Conversion,” an undated white paper published by the CJK Dictionary Publishing Society.

Unlike pretty much all of the other modern scripts encoded in Unicode, there's no set number of Han characters.⁷³ They number well into the tens of thousands, but there's no definitive listing of *all* of them. Not only are new characters being created all the time, but there are numerous characters that appear maybe only once or twice in all East Asian writing. Some of these might just be idiosyncratic ways of writing more common characters, but some are new coinages created for some ad-hoc purpose in some piece of writing (often called “nonce forms”). There are characters in some Chinese dictionaries whose pronunciations and meanings have both been lost.

So the idea of official character sets goes back well before the introduction of information technology into the Han-character-using world. Not only do you have sets of characters collected in dictionaries, but various governmental entities publish official lists of characters that are to be taught in school to all children, or that restrict the set of characters that can be used in official government documents, etc. This is more or less akin to states in the U.S. publishing vocabulary lists containing words that every child at certain grade levels should know.

These lists form the basis of the various encoding standards used in the various Han-character-using countries, but all of the character encoding standards go beyond these lists of basic characters to include rarer characters. Each of the main locales using the Han characters—the People's Republic of China, Taiwan, Japan, North and South Korea, Vietnam, Hong Kong, and Singapore—have developed their own sets of character-encoding standards, generally for their own use.

Unicode uses many of these as the basis for its collection of Han characters. The CJK Unified Ideographs area, the main body of Han characters in Unicode, is based on no fewer than eighteen different source standards, containing approximately 121,000 characters total. Yet the CJK Unified Ideographs area contains only 20,902 characters.

How is this possible? Well, consider the various ISO 8859 standards. Together, the first ten parts of ISO 8859 comprise 1,920 printing characters. However, this number is misleading because each part of ISO 8859 contains the ASCII characters (in fact, the various extended Latin encodings included in ISO 8859 also contain considerable overlap). Thus, 960 characters of that 1,920-character total represent ten copies of the 96 ASCII characters. Unicode only encodes these once.

In the same way, there's considerable overlap between the eighteen source standards that the CJK Unified Ideographs block are based on. Unfortunately, determining just which characters are duplicates and which aren't isn't all that straightforward.

Unicode's designers followed several rules in deciding which characters were duplicates and could therefore be weeded out (or “unified”):

- If one of the encoding standards that went into the original set of Han characters (the ones designated as “primary source standards” by the committee) encodes two characters separately, then Unicode does too, even if they would otherwise have been unified. This is called the “Source Separation Rule,” and is analogous to the round-trip rule used in most of the rest of the Unicode standard. For characters added to Unicode more recently, the Source Separation Rule is not followed—round-trip compatibility is only guaranteed for the original primary source standard (this is also true from some more recent source standards in other parts of Unicode). The various different versions of the “sword” character we looked at earlier are encoded separately in the JIS

⁷³ Most of the information in this section comes either from the Unicode standard itself, pp. 258-267, or from Lunde, pp. 66-137.

X 0208 standard, which is a primary source standard, and so are also encoded separately in Unicode.

- If two characters look similar but have distinct meanings and etymologies, they aren't unified. This is called the “Noncognate Rule.”

For example, this character, which means “earth”...

土

...looks a lot like this one...

士

...but this second character is completely unrelated to the first character. It means “warrior” or “scholar.” These two characters each get their own code point values in Unicode.

For characters that *are* related historically and have similar shapes, closer examination is required to decide whether or not to unify them. The basic principle is that if the characters have the same “abstract shape”—that is, if the differences between them can be attributed merely to differences in typeface design—they can be unified. Otherwise, they're considered to be distinct.

To figure out whether the differences between two characters are merely typeface-design differences, the characters are decomposed into their constituent parts. If they have different numbers of components, or if the components are arranged differently, the characters are not unified. For example, these two characters have different numbers of components⁷⁴...

崖 ≠ 厓

...and these two characters have the same components, but arranged differently:

峰 ≠ 峯

If corresponding components in the two characters have a different structure, or are based on different radicals, the characters are not unified. For example, these two characters have one corresponding component that has a different structure in the two characters...

擴 ≠ 擴

⁷⁴ The following set of examples are all taken from the Unicode standard, p. 265.

...and these two characters are similar, but based on different radicals:

祕 ≠ 秘

This leaves us with relatively minor differences, such as stroke overshoot...

雪 ≈ 雪
鉅 ≈ 鉅

...stroke initiation and termination...

朱 ≈ 朱
父 ≈ 父

...and stroke angle:

丕 ≈ 丕
八 ≈ 八

[these would look a lot better if we could get access to the original EPS files from TUS]

These differences (and thus the above pairs of characters) *are* considered differences in typeface design, and these pairs of characters *are* unified in Unicode.

One result is that, for the most part, the Simplified and Traditional Chinese versions of most characters are encoded separately in Unicode, as are the Japanese variants when they're different. Still, there are cases where this didn't happen because the shapes were close enough to be unified, but different typeface designs are still preferred in different locales. For example, the Traditional Chinese character for "bone"...

骨

...looks like this in Simplified Chinese:

骨

The differences between the two were small enough for them to be unified. In these cases, the choice of typeface is considered outside the scope of Unicode, just as the choice of typeface is outside the

scope with the other scripts in Unicode (such as when historical and modern versions of letters are unified even when they have different shapes, for example Old Cyrillic and modern Cyrillic letters).

The result of the rigorous application of these principles is the Unified Han Repertoire, the collection of Han characters in Unicode. The Unified Han Repertoire is often nicknamed “Unihan.” This is how the 110,000 characters in the original eighteen source standards were boiled down to the 20,902 characters in the original Unihan collection (this original collection of characters, the ones in the CJK Unified Ideographs area, was originally called the Unified Repertoire and Ordering, or URO). The same principles have also been applied to the collections of characters that were added after the original version of Unicode.

There is one other interesting problem that had to be solved: which order to put the characters in. Unicode’s designers didn’t want to specifically follow the ordering of any of the source standards, since this would appear to favor one country over the others. They opted instead for a culturally-neutral ordering based on the ordering of the Kangxi dictionary. This is regarded as authoritative in all the countries that use the Han characters, and contains most of the characters in Unicode. For characters that aren’t in the Kangxi dictionary, there’s a list of three other dictionaries (one each for Japanese, Chinese, and Korean) that were used instead. The character in question was located in one of these dictionaries and placed in the Unicode order after the nearest character in that dictionary that was also in the Kangxi dictionary.

The CJK Unified Ideographs area

The CJK Unified Ideographs area runs from U+4E00 to U+9FAF. This is the original set of Han characters in Unicode. It contains 20,902 characters drawn from eighteen different source standards.

The source standards are divided into groups that are given letters: the “G source” consists of seven national standards from mainland China, the “T source” consists of three sections from the Taiwanese national standard, the “J source” consists of two Japanese national standards, and the “K source” consists of two Korean national standards. Eight additional sources which aren’t official national standards and are collected from all the preceding locales constitute the “U source,” but characters that exist only in the U source aren’t included in this block.

The CJK Unified Ideographs Extension A area

The CJK Unified Ideographs Extension A area runs from U+3400 to U+4DBF. It adds 6,582 additional characters drawn from thirteen additional sources. These include three extra G-source standards (plus extra characters from some of the original sources), six additional T-source standards (basically, the rest of the Taiwanese national standard, which is huge), another J-source standard, two more K-source standards, two Vietnamese standards (the “V source”), and a bunch of extra supplemental (“U source”) materials (again, characters *only* in the U source don’t actually go in this area). This collection follows the same unification and ordering rules as the original Unihan collection (except that they dropped the Source Separation Rule), but is encoded separately because this collection came along after the other collection had already been encoded. Extension A was added in Unicode 3.0.

The CJK Unified Ideographs Extension B area

Unicode 3.1 opened up the non-BMP space for colonization, and the Ideographic Rapporteur Group, the group responsible for maintaining the Unihan collections, took to it with a vengeance. Plane 2,

the Supplementary Ideographic Plane (or “SIP”), now contains 42,711 Han characters. These make up the CJK Unified Ideographs Extension B area (sometimes called “Vertical Extension B”), which runs from U+20000 to U+2A6D6.

The characters are drawn from twenty new source standards: eight new G-source standards, most of which are dictionaries instead of character-interchange standards (the KangXi dictionary was one of them), a collection of characters used in Hong Kong (the new “H source”), five new T-source standards (characters added to the Taiwanese national standard in its most recent incarnation), two new J-source standards, a new K-source standard, and three new V-source standards. Extension B doesn’t add any U sources (supplemental works consulted by the IRG but not formally submitted by any of the member countries).

Most of the characters in Extension B are quite rare, although this is also where new characters from non-Mandarin languages go, so there are a lot of Cantonese-specific characters, many of which you see in common things like Hong Kong place names, in here.

The CJK Compatibility Ideographs block

The CJK Compatibility Ideographs block runs from U+F900 to U+FAFF. The South Korean national encoded standard, KS X 1001 (originally KS C 5601), includes duplicate encodings for a bunch of Han characters that have multiple pronunciations in Korean (the Unicode standard likens this to having a separate code point value for each different pronunciation of the letter *a*). These characters are pure duplicates: they look identical and have the same meaning; they just have variant pronunciations. However, KS X 1001 is a primary source standard, so the Source Separation Rule dictates that Unicode also had to give these characters multiple code points. The designers of Unicode compromised, giving them separate code-point values, but putting them into a separate zone, the CJK Compatibility Ideographs area, rather than interspersing them with the other Han characters in the main Unihan area. They did the same thing with duplicate characters that were inadvertently included in a few other source standards. All of these characters have *canonical*—not compatibility—decompositions to the original characters in the CJK Unified Ideographs area.

There are actually twelve characters in this block that *aren’t* duplicates and can be considered part of the main Unihan collection. They’re here rather than in the main Unihan block because they were drawn from the supplemental sources (the “U source”) rather than from the national standards. These characters don’t decompose and aren’t compatibility characters, despite the name of the block they’re in.

The CJK Compatibility Ideographs Supplement block

The SIP also contains the CJK Compatibility Ideographs Supplement block, which runs from U+2F800 to U+2FA1D. These characters are also duplicates, but were added to maintain round-trip compatibility with the most recent version of CNS 11643, the Taiwanese national standard. These characters also have canonical decompositions to characters in the main Unihan area.

The Kangxi Radicals block

The Kangxi Radicals block, which runs from U+2F00 to U+2FDF, contains the 214 radicals used in the Kangxi dictionary to classify the Han characters. Most of these radicals can also be used as ideographs, and are already encoded in the main Unihan block. The characters here aren’t Han characters; they’re just radicals and they’re not supposed to be used as Han characters. Unicode

classifies them as symbols rather than letters to emphasize this. All of these characters have compatibility decompositions to the ideographs that consist of just the radical.

The CJK Radicals Supplement block

The set of radicals used to classify characters varies from dictionary to dictionary, and some radicals have variant shapes depending on the character that includes them (or a different shape in Simplified Chinese than in Traditional Chinese). The CJK Radicals Supplement block, which runs from U+2E80 to U+2EFF, includes a bunch of these supplemental radicals and supplemental shapes.

Ideographic description sequences

It hardly seems like, with more than 70,000 characters, any more would be necessary. But the set of Han characters is an open-ended set, with more being created every day (a common source, for example, is proud parents wanting to give their children unique names). As if that weren't enough, not all of the variant forms of the encoded characters are given separate encodings. When there are big differences, they are, but as we saw earlier, smaller differences that can be attributed to font design are unified. The problem is that there are many cases where there are definite regional preferences for particular font designs, meaning it can sometimes be difficult to get precisely the glyph you want, even when the character is encoded. Unicode provides a couple of ways of dealing with these problems.⁷⁵

Let's start by looking at a rare glyph we might need these techniques for. Consider this character:



[should probably get John to send me the EPS of this character... might also be possible to put it together from the pieces using Illustrator]

This is an archaic form of this character...



...which means “dry”. This passage from the writings of Mengzi, the Second Sage of Confucianism, uses it⁷⁶:

⁷⁵ Most of the material in this section is drawn from John H. Jenkins, “New Ideographs in Unicode 3.0 and Beyond,” *Proceedings of the Fifteenth International Unicode Conference*, session C15. In particular, all of the examples in this section are ripped off from that source.

⁷⁶ Jenkins supplies this English translation for the excerpt, attributing it to D.C. Lau: “When the sacrificial animals are sleek, the offerings are clean, and the sacrifices are observed at due times, and yet floods and droughts come, then the altars should be replaced.”

犧牲既成，粢盛既潔，祭祀以時；然而旱乾水溢，則變置社稷。

A modern edition would just use the modern glyph, but a scholar seeking an accurate transcription of the original text might need an accurate way to represent the ancient glyph. There are a number of ways of dealing with this in Unicode.⁷⁷

The oldest and simplest solution is the geta mark (U+3013), which looks like this: ■. This mark is used by Japanese printers to reserve space on the line for a character they can't print. Using it, the passage above would look like this:

犧牲既成，粢盛既潔，祭祀以時；然而旱■水溢，則變置社稷。

This gives you a good visual tip-off that there's a character in this passage that can't be rendered, but it doesn't tell you anything about that character. Unicode 3.0 introduces a refinement of this idea that gives you a little more information. It's called the ideographic variation indicator (U+303F) and looks like this:

𠄎

The idea is to use this in conjunction with some character that *is* encoded in Unicode. You precede that character with the ideographic variation indicator to say that the character you *really* want is somehow related to the character following the ideographic variation mark. It might be a glyphic variant of that character, look kind of like it, mean the same thing, or whatever. If we use it in our example, it winds up looking like this:

犧牲既成，粢盛既潔，祭祀以時；然而旱𠄎乾水溢，則變置社稷。

This may not be a whole lot better than using the geta mark, but it at least provides some information as to what's really desired. You could, of course, just use the substitute glyph, and in this example at least, the meaning of the passage would still get across. But the ideographic variation indicator gives a visual signal that it isn't the glyph you really wanted to use.

Unicode 3.0 also introduced a much richer way of getting across which character you mean when it isn't available. This is called the *ideographic description sequence*, and it makes use of the characters in the Ideographic Description Characters block (U+2FF0 to U+2FFF). The basic idea

⁷⁷ Actually, this character *is* encoded in Unicode, in the new CJK Unified Ideographs Extension B area, so this example may not be the best. Nevertheless, the character was added in Unicode 3.1, so if you were dealing with a Unicode 3.0 or earlier implementation, or didn't have a font that actually contained this character, you might still have to use the techniques described here.

behind ideographic description sequences is simple: Most Han characters are simply combinations of two or more other Han characters. You could get at least a decent approximation of the character you want by using the characters that make it up along with some kind of character indicating how they are combined together. The ideographic description characters fulfill that role of indicating how other ideographs are combined together.

For example, let's say you want to combine this character...

井

...which means “well,” with this one...

蛙

...which means “frog,” to produce a “frog in a well” character:⁷⁸

井蛙

[can probably put together a decent-looking version of this glyph in Illustrator]

You could use the character U+2FF1 IDEOGRAPHIC DESCRIPTION CHARACTER ABOVE TO BELOW (𠄎) to indicate you want to stack the other two characters on top of each other, yielding something that looks like this:

𠄎井蛙

Perhaps not the most beautiful sequence of characters ever designed, but it does a good job of saying what you mean.

More complicated sequences are also possible. You could, for example, break down the “frog” character even further, since “frog” is made up of individual characters meaning “insect” and “pointed jade” (this is obviously the phonetic). Either of these sequences could also be used to describe our “frog in a well” character:

𠄎井 虫圭

𠄎井 虫 土土

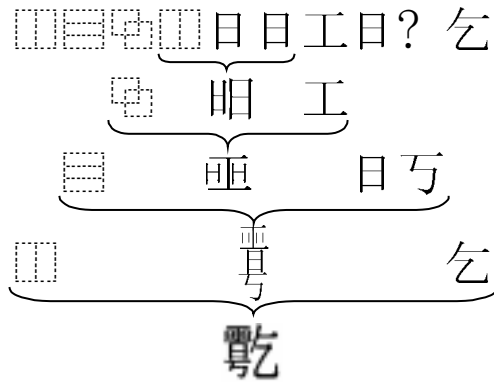
⁷⁸ “Like a frog at the bottom of a well” is a common expression in Chinese for a narrow-minded person, which is where this example comes from.

What this shows is that the ideographic description characters are recursive. The basic rule is that you start an ideographic description sequence with an ideographic description character, and then follow that character with the two (or three) characters that are to be composed, in order from left to right (or top to bottom). These characters may be regular Han characters, radicals, or ideographic description sequences themselves.

This can go to fairly wild extremes. Let's go back to the old character from the Mengzi example. It can be described like this:

☐☐☐☐☐ 日日工日 乞

Crystal clear, huh? This points up one of the problems with ideographic description sequences. They're designed to be machine-readable, and some human readability is sacrificed in the process. The sequence can be parsed like this:



[There appears to have been an error importing this drawing—the question mark should be the same character as the one directly below it: the one that looks kind of like the number 5.]

Of course, this takes a bit of work on the part of the reader. It also gets kind of messy when you place a sequence like this into actual text:

犧牲既成，粢盛既潔，祭祀以時；然而旱☐☐☐☐☐
日日工日 乞水溢，則變置社稷。

Of course, what you *really* want is for a sufficiently sophisticated text-rendering engine to give you this when it sees that sequence:

犧牲既成，粢盛既潔，祭祀以時；然而旱乾水溢，
則變置社稷。

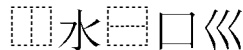
Unicode gives implementations the freedom to do this if they want. But it's important to note that this isn't *required* of systems that know about the ideographic description characters. It's perfectly legal for an implementation to just draw them as regular visible glyphs, as in the preceding example. (It is *not* legal, by the way, to make them invisible and just show the Han characters they operate on.) The idea is that these characters can be used to *describe* a missing character, *not* that they can be used to *represent* that character. They're more akin to the phrase "an e with an acute accent over it" than they are to the sequence U+006E U+0301.

Part of the reason for this is that they're really still only approximate descriptions of what is intended. For example, this character...



[Composer: This character is U+23CE7. I don't have any fonts with the Plane 2 Chinese characters in them]

...can be described like this:



But it might also theoretically be represented like this:



This is because the "water" character (水) usually adopts an abbreviated shape (𠂆) when it's used as part of another character. But there's always the possibility that it shouldn't do that in your particular case. There's no way to say which you mean using the facilities that Unicode gives you.

It's important to remember that ideographic description sequences are explicitly *not* combining character sequences. The ideographic description characters are classified as symbols, not as combining marks or formatting characters. Multiple methods of describing the same character using ideographic description sequences are possible and are explicitly not equivalent. If an ideographic description sequence is used to describe a character that's actually encoded in Unicode, it's explicitly not equivalent to that character. An ambitious Unicode implementation can parse an ideographic description sequence and attempt to render the character it describes, but it doesn't have to. Ideographic description sequences are basically intended as a way of representing some character until it's actually officially encoded in Unicode.

To make life easier on implementations that want to parse ideographic description sequences and try to display the characters they describe, the Unicode standard prescribes a grammar for their use and specifies that an individual ideographic description sequence may be no longer than sixteen characters, may not contain more than six Han characters or radicals in a row with no intervening ideographic description characters (this is to minimize how far back a process examining the text from a random spot in the middle has to scan to figure out whether it's in an ideographic description

Bopomofo are rarely used alone; usually they're used to annotate regular Chinese text with pronunciations. (For more information on how this works, see the “Ruby” section later in this chapter.) You generally see them in children's books or other educational materials.

The Bopomofo block

The Unicode Bopomofo block, running from U+3100 to U+312F, contains the basic set of Bopomofo characters used to write Mandarin, plus three characters used to write non-Mandarin dialects. The tone marks are unified with Western accent marks in the Spacing Modifier Letters block.

The Bopomofo Extended block

The Bopomofo Extended block, running from U+31A0 to U+31BF (there wasn't enough room to put these characters in the same block as the other Bopomofo characters) contains a bunch of less-universally-recognized Bopomofo characters used to write various non-Mandarin Chinese languages. A few additional tone marks are also unified with characters in the Spacing Modifier Letters block.

Japanese

Japanese may have the most complicated writing system of any language spoken today. It's written in a combination of three complete scripts, sometimes supplemented with Latin letters and various other marks.⁷⁹

The history of written Japanese goes back to around the third century, when the Han characters were first used to write Japanese (in Japanese, the Han characters are called *kanji*). But Japanese isn't actually linguistically related to Chinese (although it contains a lot of words borrowed from Chinese), and so the Han characters never worked as well for Japanese. Most of the phonetic cues in the Han characters don't help in Japanese (one estimate holds that only 25% of the commonly-used Kanji characters contain useful clues as to their Japanese pronunciations). Furthermore, while Chinese is generally monosyllabic and noninflecting (or at least has a lot fewer polysyllabic words than most other languages, including Japanese and Chinese), Japanese is polysyllabic and has a fairly complex inflectional system. Writing Japanese exclusively with Kanji characters generally means that the grammatical endings on words aren't written and must be supplied by the reader, or that they have to be represented using Kanji characters with the right pronunciations, which can be confusing (especially since almost all Kanji characters have at least two pronunciations in Japanese).

Eventually a standardized set of Kanji characters came to be used for writing the grammatical endings and for other types of phonetic representation. These came to be known as the *man'yōgana* (“Chinese characters used phonetically to write Japanese”), and their shapes gradually became more simplified and differentiated from their shapes when used as regular Kanji characters. This process culminated in the two modern Japanese syllabic scripts: Hiragana and Katakana (collectively called “Kana”), which were finally officially standardized by the Japanese government in 1900. The

⁷⁹ My sources for the section on Japanese are Janet Shibamoto Smith, “Japanese Writing,” in *The World's Writing Systems*, pp. 209-217; Lunde, pp. 42-47; and Nakanishi, pp. 94-95.

Hiragana were derived from more cursive, calligraphic forms of the original characters, while the Katakana were very simplified variants of the printed forms of the characters.⁸⁰

Both forms of Kana are what is known as “pure syllabaries,” that is, each Kana character represents an entire syllable of spoken Japanese. Unlike so-called “alphasyllabaries,” such as Hangul, which we’ll look at next, Kana characters can’t be broken down into component parts that represent the individual sounds in the syllable.

The more cursive form, Hiragana, is generally used to write native Japanese words and grammatical endings. There are forty-eight basic Hiragana characters:

	–	K	S	T	N	H	M	Y	R	W
A	あ	か	さ	た	な	は	ま	や	ら	わ
I	い	き	し	ち	に	ひ	み	り	る	
U	う	く	す	つ	ぬ	ふ	む	ゆる		
E	え	け	せ	て	ね	へ	め	れ	ゑ	
O	お	こ	そ	と	の	ほ	も	よ	ろ	を
–					ん					

(The characters ゐ and ゑ are rarely used nowadays.)

Each character (with the exception of あ, い, う, え, お, and ん) represents the combination of an initial consonant sound and a vowel sound. The vowels by themselves are represented by あ, い, う, え, and お, and are used not just at the beginnings of words but as the second characters in diphthongs. There is only one syllable-final consonant, and it’s represented by ん. Thus, the system isn’t a 100% perfect syllabary—there are a number of syllables that are actually represented using more than one Hiragana character.

Not all of the characters have the pronunciation shown in the chart above. For example, つ is actually pronounced “tsu,” ち is pronounced “chi,” ふ is pronounced “fu,” and し is pronounced “shi.” ん is actually pronounced “m” before “b” and “p” sounds.

⁸⁰ Lunde provides a very interesting table (pp. 46-47) showing which Kana characters derived from which Kanji characters.

In addition to the basic sounds shown in the table above, diacritical marks are used to widen the palette of available sounds. Two strokes on the upper right-hand shoulder of the character, called *nigori* or *dakuten*, are used with certain characters to turn their initial consonant sound from an unvoiced sound into a voiced sound. For example, さ (“sa”) becomes “za” when the dakuten is added (ざ). The characters in the h-series (i.e., は, ひ, ふ, へ, and ほ) are special: a small circle on the upper-right shoulder, called *maru* or *handakuten*, gives them an initial *p* sound. The dakuten gives them an initial *b* sound. That is, は is pronounced “ha,” ぱ is pronounced “pa”, and ば is pronounced “ba.” Finally, the dakuten can be added to う to get “vu” (うゝ). These combinations of basic Hiragana characters and diacritical marks give us another twenty-five syllables:

	G	Z	D	B	P	V
A	が	ざ	だ	ば	ぱ	
I	ぎ	じ	ぢ	び	ぴ	
U	ぐ	ず	づ	ぶ	ぷ	うゝ
E	げ	ぜ	で	べ	ぺ	
O	ご	ぞ	ど	ぼ	ぽ	

Smaller versions of some of the Hiragana characters are also used more or less like diacritics. For example, Japanese has palatalized syllables (syllables with a *y* sound between the consonant and the vowel). These are represented using a basic character representing the initial consonant and the *i* sound, followed by a small version of や, ゆ, or よ representing the *y* sound and the final vowel.

For example, *nya* is represented with に (‘‘ni’’) followed by a smaller version of や (‘‘ya’’): にゃ.

The one other case of two initial consonants, *kwa*, is represented similarly: く (‘‘ku’’) plus a small version of わ (‘‘wa’’): くわ. This gives us another thirty-six possible combinations:

	K	S	T	N	H	M	R	G	Z	D	B	P
YA	きゃ	しゃ	ちゃ	にゃ	ひゃ	みゃ	りゃ	ぎゃ	じゃ	ぢゃ	びゃ	ぴゃ
YU	きゅ	しゅ	ちゅ	にゅ	ひゅ	みゅ	りゅ	ぎゅ	じゅ	ぢゅ	びゅ	ぴゅ
YO	きよ	しよ	ちよ	によ	ひよ	みよ	りよ	ぎよ	じよ	ぢよ	びよ	ぴよ

The small version of つ is used to indicate consonant gemination (i.e., the doubling or prolonging of a consonant sound): Thus, ぶつだ is “butsuda,” but ぶつだ is “Buddha” (pronounced “budda”).

The small versions of the bare vowels (あ, い, う, え, and お) are used with certain characters for borrowed sounds. Generally, the resulting syllable is a combination of the initial consonant sound from the main character and the vowel sound denoted by the small character. For example, ふあ is “fa,” とう is “tu,” and ちあ is “cha.”

Long vowel sounds are usually represented by using two of the regular vowel characters in a row. For example, “ā” is written as ああ. The one exception is ō, which is written おう.

When Hiragana is written vertically, the character く is occasionally used to indicate that the preceding two characters are to be repeated. Adding the *dakuten* changes the initial consonant on the repeated sequence to a voiced consonant. This example...

て
ん
ト
\

[need to correct example so that the two halves of the repeat mark touch]

...is pronounced “tenden.” Notice, by the way, that the repeat mark takes up two display cells, visually indicating that it’s repeating two characters.

Katakana, the more angular form of Kana, is used for foreign words and onomatopoeia. Sometimes, native Japanese words are written in Katakana for emphasis or to indicate euphemism or irony (in other words, using Katakana for a word that’s normally written in Kanji or Hiragana is similar to putting quotation marks around an English word that’s not actually part of a quotation). Katakana works almost exactly the same way as Hiragana. There are forty-eight basic characters:

	–	K	S	T	N	H	M	Y	R	W
A	ア	カ	サ	タ	ナ	ハ	マ	ヤ	ラ	ワ
I	イ	キ	シ	チ	ニ	ヒ	ミ		リ	ヰ
U	ウ	ク	ス	ツ	ヌ	フ	ム	ユ	ル	

E	エ	ケ	セ	テ	ネ	ヘ	メ	レ	エ	
O	オ	コ	ソ	ト	ノ	ホ	モ	ヨ	ロ	ヲ
-					ン					

(Again, the characters 𛄁 and 𛄂 are rarely used nowadays.)

Dakuten and *handakuten* work exactly the same way with Katakana (with a few additional syllables possible):

	G	Z	D	B	P	V
A	ガ	ザ	ダ	バ	パ	ヴァ
I	ギ	ジ	ヂ	ビ	ピ	ヴィ
U	グ	ズ	ヅ	ブ	プ	ヴ
E	ゲ	ゼ	デ	ベ	ペ	ヴェ
O	ゴ	ゾ	ド	ボ	ポ	ヴォ

Smaller versions of certain Katakana characters also work the same way as their Hiragana counterparts:

	K	S	T	N	H	M	R	G	Z	D	B	P
YA	キ	シ	チャ	ニ	ヤ	ミ	リ	ギ	ジャ	ヂ	ビ	ピ
YU	ク	シュ	チュ	ニ	ユ	ミ	リ	グ	ジュ	ヂ	ビ	ピ
YO	キ	ョ	チャ	ニ	ヨ	ミ	リ	ギ	ジョ	ヂ	ビ	ピ

ブツダ = "Buddha"

ファ = "fa"

トゥ = "tu"

チャ = "cha"

Long vowels are generally represented in Katakana using a dash (“choon”) instead of a doubled vowel. That is, “ā” is written as *アー* rather than *アア*.

With Katakana, the *ヽ* mark is occasionally used to indicate repetition. Unlike the Hiragana repetition mark, the Katakana repetition mark only repeats the last character: for example, *コヽ* is “koko.”

Like Chinese, Japanese is traditionally written in vertical columns that run from the right-hand side of the page to the left-hand side, but it’s becoming more common to see Japanese written in horizontal lines with the same reading order as English. Spaces are not used between words, although word boundaries are usually fairly detectable due to the mixture of scripts and use of punctuation in most Japanese text. Japanese text is not word-wrapped: there are certain punctuation marks that can’t appear at the beginning of a line and others that can’t occur at the end of a line, but basically line breaks can occur anywhere, including the middles of words. Latin letters are used for certain foreign words and expressions, and numbers are written using both Kanji characters and Western numerals, depending on context.

Anyway, if you put all these writing systems together, you get something that looks like this:⁸¹

EUC等のエンコーディング方法は日本語と英語が混交しているテキストをサポートします。

It breaks down as follows:

Original text	Meaning	Pronunciation	Script
EUC	EUC	EUC	Romaji
等	such as...	nado	Kanji
の	possessive marker	no	Hiragana
エンコーディング	encoding	enkōdingu	Katakana
方法	method	hōhō	Kanji
は	topic marker	wa	Hiragana
日本語	Japanese	nihongo	Kanji
と	and	to	Hiragana
英語	English	eigo	Kanji
が	subject marker	ga	Hiragana
混交	to mix	konkō	Kanji
している	doing	shite-iru	Hiragana
テキスト	text	tekisuto	Katakana

⁸¹ This example is ripped off from Lunde, p. 3.

を	object marker	o	Hiragana
サポート	support	sapōto	Katakana
します	do	shimasu	Hiragana

Put it all together and you get “Encoding methods such as EUC can support texts that mix Japanese and English.”

The Hiragana block

The Unicode Hiragana block runs from U+3040 to U+309F. The basic syllables, the syllables with *dakuten* and *handakuten*, and the small versions of the syllables all get individual code points. In addition, the *dakuten* is encoded as U+3099 COMBINING KATAKANA-HIRAGANA VOICED SOUND MARK and the *handakuten* as U+309A COMBINING KATAKANA-HIRAGANA SEMI-VOICED SOUND MARK. These are regular Unicode combining characters, and the syllables that use them have canonical decompositions to their unadorned forms and the combining marks. However, the standard practice in native Japanese encoding standards is to use the precomposed versions of the characters. The analogous characters in Japanese encodings don’t have combining semantics, so Unicode also has non-combining versions of these characters.

The Hiragana repetition marks are also encoded in this block as U+309D HIRAGANA ITERATION MARK and U+309E HIRAGANA VOICED ITERATION MARK. A special version of the Hiragana repetition mark for use in vertical text is encoded in the CJK Symbols and Punctuation block as U+3031 VERTICAL KANA REPEAT MARK. (A clone with the *dakuten* is encoded at U+3031, and since this character normally takes up two display cells, the upper and lower halves are encoded at U+3033, U+3034, and U+3035.)

The Katakana block

The Unicode Katakana block runs from U+30A0 to U+30FF and encodes the Katakana characters in the same relative positions as the Hiragana characters are encoded in the Hiragana block, along with a few extra characters that don’t have counterparts in Hiragana. The combining voiced and semi-voiced sound marks from the Hiragana block are also used with Katakana, and the appropriate characters used them in their canonical decompositions. The choon mark is encoded as U+30FC KATAKANA-HIRAGANA PROLONGED SOUND MARK. The Katakana repetition marks are also included here, as is a centered dot that is used to separate words in a sequence of foreign words written in Katakana.

The Katakana Phonetic Extensions block

Unicode 3.2 adds a new Katakana Phonetic Extensions block running from U+31F0 to U+31FF. It contains a bunch of additional “small” Katakana characters used for phonetic transcription of Ainu and some other languages.

The Kanbun block

The Kanbun block (U+3190 to U+319F) contains a small set of characters that are used in Japanese editions of classical Chinese texts to indicate the Japanese reading order. These are small versions of certain ideographs and are usually written as annotations to the side of a line of vertically-arranged

Chinese text. All of these characters have compatibility mappings to regular Han characters, since they're just smaller versions of them used as annotations.

Korean

Like Japanese, the Han characters (called *hanja* in Korean) were originally used to write Korean, with the oldest example dating to about 414.⁸² Korean, which is distantly related to Japanese⁸³ but not to Chinese or the other Asian languages, isn't really any better suited to writing with Hanja than Japanese is, and various phonetic systems similar to Japanese Kana were devised to supplement the Han characters. Several of these systems are obsolete, but one of them—Hangul—has come to be the main writing system used for modern Korean. Interestingly, even though Hangul was invented in 1446, the Chinese characters persisted as the preferred method of writing Korean right up to the beginning of this century; Hangul was accorded a kind of second-class status as the writing of the uneducated. Since Korea was liberated from Japan in 1945, however, Hangul has risen steadily in importance and the Han characters have increasingly fallen into disuse. In North Korea today, Hangul is used pretty much exclusively, although a small set of Hanja is still taught in school. In South Korea, Hanja are still used, but they tend to make up a small minority of the characters in most modern texts.

Hangul (also spelled Hankul or Hangeul, and also called by a variety of other names, including [before 1910] *changum*, *enmun*, or *kwukmun*, or [in North Korea] *chosengul* or *wuli kulcha*) was invented in 1446 by the sage King Sejong. Although it bears a superficial resemblance to the Han characters (probably not accidentally) and borrows some concepts from other writing systems, it's pretty much a pure invention rather than something that kind of grew organically from other origins, as almost every other writing system we've looked at did. The design is based on sound linguistic principles, and it's quite well-suited to the writing of Korean.

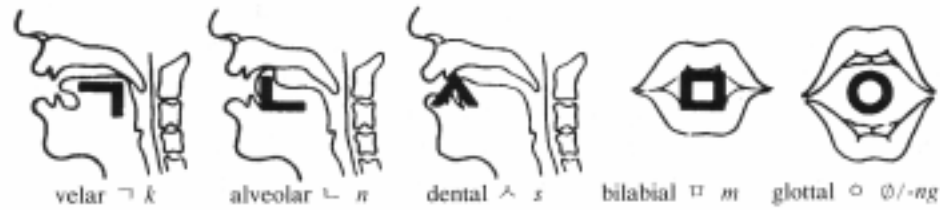
Hangul is basically an alphabetic script, consisting of letters called *jamo*, which represent both consonant and vowel sounds (unlike the Middle Eastern and Southeast Asian scripts, the vowels are neither left out nor written as marks that attach to the consonants; they're independent characters in their own right like the vowels in the Latin, Greek, and Cyrillic alphabets). However, instead of being written sequentially like letters in Western alphabets are, Hangul jamo are arranged into blocks representing whole syllables. The blocks fit into square display cells, like the Han characters, and look kind of like them. There have been proposals to write the jamo sequentially on a line of text, but they've never caught on. Even though it requires more different pieces of movable type, the arrangement into syllable blocks makes more efficient use of space and is (arguably, at least) more efficient to read as well.

There are five basic consonant characters, whose shapes originally were designed to resemble the positions of the speech organs while making the sounds⁸⁴:

⁸² The material on Hangul, especially the examples and tables, comes from Ross King, "Korean Writing," in *The World's Writing Systems*, pp. 218-227.

⁸³ At least according to some linguists; there's a lot of controversy on this point.

⁸⁴ The illustration comes from King, *op. cit.*, p. 220.



(The shapes have changed a little since the inventions of the characters, but you can still see the basic principles at work in the modern shapes.) The other consonants were then organized into families according to the five basic ones, with extra strokes added⁸⁵:

ㄱ <i>k</i>	ㄴ <i>n</i>	ㄷ <i>s</i>	ㅁ <i>m</i>	ㅇ <i>ng</i>
	ㅌ <i>t</i>	ㅊ <i>c</i>	ㅍ <i>p</i>	
ㅋ <i>kh</i>	ㅍ <i>th</i>	ㅌ <i>ch</i>	ㅍ <i>ph</i>	ㅎ <i>h</i>
ㄲ <i>kk</i>	ㅍ <i>tt</i>	ㅌ <i>ss</i>	ㅍ <i>pp</i>	
		ㅌ <i>cc</i>		
	ㄹ <i>l</i>			

The construction of the vowels isn't quite as concrete: Three basic shapes—the horizontal line, representing earth, the dot, representing the heavens, and the vertical line, representing man—were used in various combinations to represent the vowels⁸⁶:

ㅣ *i* ㅡ *u* ㅏ *wu*
 ㅑ *ey* ㅓ *e* ㅗ *o*
 ㅕ *ay* ㅛ *a*

(The dot has since morphed into a short stroke.) These also combine in various ways to form diphthongs:

ㅑ *uy* ㅗ *wuy*
 ㅓ *oy*
 ㅏ *ywu*
 ㅑ *yey* ㅓ *ye* ㅗ *yo*
 ㅕ *yay* ㅛ *ya*
 ㅑ *wey* ㅗ *we*

⁸⁵ This table is reproduced directly from King, *op. cit.*, p. 222.

⁸⁶ *Ibid.*, p. 221

ㅘ way ㅚ wa

A syllable consists of an initial consonant, a vowel, and, optionally, a final consonant. If the vowel is based on a vertical line, the initial consonant goes to its left. If the vowel is based on a horizontal line, the initial consonant goes above it. If the vowel is a combination of vertical and horizontal strokes, the initial consonant goes in the upper left-hand corner. The final consonant, if there is one, goes underneath the initial consonant and the vowel. Syllables beginning with a vowel sound use ㅇ as their initial consonant (it's silent at the beginning of a syllable). There are certain combinations of final consonants that can go at the end of a syllable; the first sound goes to the left of the second. Thus, the reading order within a syllable block is left to right, top to bottom:

ㅍ (p) + ㅏ (a) = ㅑ (pa)

ㅍ (p) + ㅏ (a) + ㅁ (m) = ㅑㅁ (pam)

ㅅ (s) + ㅑ (o) = ㅓ (so)

ㅅ (s) + ㅑ (o) + ㄴ (n) = ㅕ (son)

ㅇ (silent) + ㅏ (a) = ㅑ (a)

ㅇ (silent) + ㅚ (wa) = ㅜ (wa)

ㅇ (silent) + ㅏ (a) + ㅇ (ng) = ㅇㅑ (ang)

ㅍ (p) + ㅏ (a) + ㄹ (l) + ㅍ (p) + ㅌ (t) + ㅏ (a) = ㅑㅑㄹㅑ (palpta)

ㅇ (silent) + ㅣ (i) + ㄹ (l) + ㄱ (k) + ㅌ (t) + ㅏ (a) = ㅇㅣㄹㅑ (ilkta)

Like Chinese and Japanese, Korean is traditionally written in vertical columns running from right to left across the page, but horizontal writing is becoming more common, especially in North Korea. In South Korea, as noted before, Hangul is mixed with Hanja, while Hangul is used exclusively in North Korea. Spaces are frequently, but not always, used between words. When spaces are not used between words, lines can be broken anywhere (except for certain punctuation marks), but when spaces are used, line boundaries must come between words, as in Western writing. There are no native Korean numerals—numbers are written using Han characters or Western numerals.

The Hangul Jamo block

There are two basic ways of representing Hangul in Unicode: unique code point values can be given to each possible syllable block, or unique code point values can be given to the individual jamo, which would then be used in combination to build up syllable blocks. The designers of Unicode opted to have it both ways.

The Hangul Jamo block, which runs from U+1100 to U+11FF, gives a code point value to each jamo. The basic encoding approach is based on the Johab encoding from the Korean national standards. The Johab encoding classifies the jamo into three categories: initial consonants (or *choseong*), vowels (*jungseong*), and final consonants (*jongseong*). There are 19 possible initial consonants, 21 possible vowels and vowel combinations, and 27 possible final consonants and

combinations of final consonants. Each of these gets its own code point value (in Johab, each of these categories gets a unique five-bit combination, which are then concatenated, with a 1 bit on the front, to form a 16-bit character code, but in Unicode, each jamo combination is given a full 16-bit code point value). This approach means that a jamo that can appear either at the beginning or the end of a syllable has a different code point value depending on its position within the syllable. It also means that combinations of vowel or consonant jamo get their own code point values. In this way, every syllable can be represented using a sequence of three jamo (for those that only consist of two jamo, invisible “filler” characters are also provided).

Unicode doesn’t define decompositions for any of the jamo, even those that are made up of combinations of smaller jamo. This is consistent with standard Korean encodings and with the way Koreans think of the jamo. The characters in this block have combining semantics: a sequence of initial-consonant, vowel, final-consonant represents a single syllable block, not three isolated jamo. (Any arbitrary sequence of initial consonants, followed by any arbitrary sequence of vowels, followed by any arbitrary sequence of final consonants, are defined by the standard to be a single combining character sequence representing one syllable, although any sequence consisting of something other than exactly one of each category is considered to be malformed.)

In addition to the 67 modern Hangul jamo, the Hangul Jamo block includes a bunch of ancient jamo that aren’t used in modern Korean, plus two invisible “filler” characters that can be used to make each syllable consist of three code points even when it doesn’t contain three jamo.

The Hangul Compatibility Jamo block

The Hangul Compatibility Jamo block, which runs from U+3130 to U+318F, contains duplicates of the characters from the Hangul Jamo block, but these characters are intended to be used to represent individual jamo, rather than being used in combination to represent whole syllables (that is, they don’t have combining semantics like the characters in the Hangul Jamo block). This block doesn’t make a distinction between initial and final consonants, since the characters aren’t intended to be used in combination.

All of the characters in this block have compatibility decompositions to characters in the regular Hangul Jamo block. This has an interesting side effect: The compatibility jamo don’t conjoin, but the regular jamo do. This means that if you convert a passage of compatibility jamo to Normalized Form KD or KC, they combine together into syllables. The conjoining jamo can be prevented from clustering together by inserting non-Korean characters (the zero-width space works nicely for this).

The Hangul Syllables area

The Korean Johab encoding allows for an encoding of every possible combination of the 67 modern jamo, even though only about 2,500 actually occur in written Korean. This is a grand total of 11,172 possible syllables. Even though the characters in the Hangul Jamo block work the same way as the Korean Johab encoding, a single syllable takes up six bytes in Unicode, but only two bytes in Johab. So Unicode, as of Unicode 2.0, also provides a unique 16-bit code point value for every one of the 11,172 syllables representable using Johab. These occupy the Hangul Syllables area, which runs from U+AC00 to U+D7FF. The ordering is the same as Johab, but the relative positions aren’t—Johab contains holes, while the Unicode Hangul Syllables area keeps all the characters contiguous, squeezing out the holes.

Unicode defines a syllable represented using a code point from this area to be equivalent to the same syllable represented using three code points from the Hangul Jamo block—in effect, even though the Unicode Character Database doesn’t specifically spell them all out, the characters in this block all have canonical decompositions to sequences of characters from the Hangul Jamo block. The Hangul Jamo block and the Hangul Syllables area are arranged in such a way that converting from one representation to the other can be done algorithmically (rather than requiring a translation table of some kind).

The Hangul Syllables area only contains combinations of the 67 modern jamo. Syllables using the archaic jamo can only be represented using sequences of characters from the Hangul Jamo block.

Halfwidth and fullwidth characters

Now that we’ve taken a look at the characters involved and how they work, we can go back and look a little more closely at some of the interesting issues involved in dealing with these characters.

You’ll often hear the terms “halfwidth” and “fullwidth” (or their Japanese equivalents, *hankaku* and *zenkaku*) used in conjunction with various characters. These terms have their origin in the variable-width character encoding standards used for most East Asian languages: Japanese variable-length encodings such as SHIFT-JIS or EUC-JP included the characters from both the JIS X 0201 and JIS X 0208 standards, but these standards overlap, resulting in a bunch of characters with two different encodings: a one-byte version and a two-byte version. This collection included the ASCII characters and a basic set of Katakana.

The double encoding of these collections of characters led to a distinction being made by many implementations in how they should be displayed. In a monospaced East Asian typeface, all of the regular East Asian characters (Kanji, Kana, Bopomofo, Hangul, etc.) fit into a square design space: they’re as wide as they are high. The foreign characters (in particular, the Latin letters, digits, and punctuation) would be sized to take up half of one of these display cells: they were half as wide as they were high. This came to be identified with the encoding length: one-byte characters occupied half a display cell, and two-byte characters occupied an entire display cell. For those characters that had both one-byte and two-byte encodings, this was the difference: the two-byte version would fit into a single display cell, and the one-byte version would fit two to a display cell. For the Latin letters, this meant the two-byte versions were shown using glyphs that were stretched horizontally and surrounded with extra white space. For the Katakana characters, this meant the one-byte versions would be squooshed horizontally to fit in half of a display cell (the original encodings of halfwidth Katakana only included the basic 48 characters: distinctions between regular and small kana were lost, and the *dakuten* and *handakuten* were treated as separate spacing marks that appeared in their own half-width display cells after the characters they modified).

Even in these days of proportional fonts and sophisticated typesetting, the distinction between halfwidth and fullwidth occasionally remains a useful one. East Asian word-processing software, for example, often still systematically treats fullwidth and halfwidth characters differently for the purposes of line breaking and line layout.

In order to retain backwards compatibility with the various East Asian character encoding standards that give certain characters redundant one- and two-byte encodings, Unicode does the same for these characters, segregating the redundant versions off in a separate block within the Compatibility Zone (see below). These are the only characters in the Unicode standard that are officially designated as either “halfwidth” or “fullwidth.”

Unicode Standard Annex #11, “East Asian Width,” extends these definitions to the entire Unicode character repertoire, allowing applications that treat “halfwidth” and “fullwidth” characters differently to know how to treat every character in Unicode. It does this by classifying the Unicode characters into six categories. An application resolves these six categories down to two broad categories, “narrow” and “wide,” depending on context. The six categories are:

- Fullwidth (F). These are the characters in the compatibility zone with “FULLWIDTH” in their names. These always get treated as “wide.”
- Halfwidth (H). These are the characters in the compatibility zone with “HALFWIDTH” in their names. They always get treated as “narrow.”
- Wide (W). These are the other East Asian characters. This category includes those characters that have counterparts in the H category (the regular Katakana), plus all the other East Asian characters: Han, Hangul, Hiragana, Katakana, Bopomofo, all of the symbols that are designed for use in East Asian typography (or taken exclusively from East Asian encoding standards), and so on. These characters, of course, always get treated as “wide.”
- Narrow (Na). These are the characters that have counterparts in the F category. Basically, this includes the ASCII characters. They’re obviously always treated as “narrow.”
- Ambiguous (A). These are the characters that have two-byte encodings in the legacy East Asian encodings and thus usually get treated in East Asian systems as “fullwidth,” but that normally are treated by other systems the same way Latin characters are. This category includes the Cyrillic and Greek alphabets, plus some math symbols. These characters get treated as either “wide” or “narrow” depending on context. In an East Asian context, they’re “wide”; otherwise, they’re “narrow.” (UAX #11 doesn’t really define what an “East Asian context” means, leaving that up to implementation; the basic idea is that an “East Asian context” would be a document or other piece of text consisting predominantly of East Asian characters, or a piece of text set using an East Asian font, even for the Western characters.)
- Neutral. This contains all characters that don’t occur in any East Asian legacy character encoding standards and, by extension, don’t occur in East Asian typography. Officially, they’re neither “wide” nor “narrow,” but for all practical purposes can be treated as “narrow.” This category includes the Arabic and Devanagari alphabets.
- Combining marks don’t have an East Asian width property at all—they’re defined to take on the East Asian width property of the character they’re applied to, just as they take on most of the other properties of the characters they’re applied to.

The East Asian width properties of all Unicode characters are given in the `EastAsianWidth.txt` file in the Unicode Character Database (see Chapter 5).

It’s important to keep in mind that the classification of characters into “halfwidth” and “fullwidth” (or, as UAX #11 does it, into “wide” and “narrow”) categories doesn’t mean that these characters are always typeset this way. Purely monospaced typefaces aren’t used any more in East Asian typesetting than they are in Western typesetting. It’s true that the Han characters (and the characters that are used with them, such as Bopomofo, Kana, and Hangul) are generally all set in the same-size display cells (which aren’t always square, by the way—newspapers, for example, often use fonts with rectangular display cells in order to fit more text on a page), punctuation and symbols often aren’t. Latin letters and digits are usually set using proportional typefaces as well, and there are various justification techniques that are used in different situations that fix it so that even when most of the characters are the same size, you don’t get that “graph paper” look. Furthermore, there are proportional fonts that actually do use different-sized display cells for certain Han and Han-related characters. For an extensive treatment of all these issues, see Chapter 7 of Ken Lunde’s *CJKV Information Processing*.

The Halfwidth and Fullwidth Forms block

Unicode’s Halfwidth and Fullwidth Forms block, which runs from U+FF00 to U+FFEF, is the dumping ground for all the characters that have redundant one- and two-byte encodings in legacy East Asian encoding standards. All of the characters in this block have compatibility mappings to characters elsewhere in the standard. It includes fullwidth versions of all the ASCII printing characters, halfwidth versions of the basic Katakana characters, halfwidth versions of the modern Hangul jamo, halfwidth versions of certain symbols and punctuation marks, and fullwidth versions of several currency symbols.

Vertical text layout

As mentioned repeatedly throughout this chapter, East Asian text is traditionally written vertically, with characters proceeding in a line from the top of the page to the bottom, and with lines of text proceeding from the right-hand side of the page to the left. Books in Chinese or Japanese, like books in right-to-left languages such as Hebrew, appear to English speakers to be bound on the wrong side: they begin at what we normally think of as the back.⁸⁷

Thanks to Western influence and the inability of early computer systems to handle vertical text, most East Asian languages are also now written horizontally, with characters running from the left to the right and lines running from the top to the bottom, as in English text.⁸⁸ The main East Asian characters don’t really behave any differently when they’re arranged horizontally as opposed to vertically, but many other characters do. Consider this example⁸⁹:

良治は、「僕には二十三にんの子供がいる」と言って、笑った。

(This is basically the same example from our discussion of the bi-di algorithm translated into Japanese: “Ryoji said, ‘I have 23 children,’ and smiled.”)

This is what it looks like laid out vertically:

⁸⁷ I’m relaying rather heavily for this section and the next section on Lunde, pp. 336-386.

⁸⁸ Although right-to-left writing also happens: Nakanishi (p. 114) gives an interesting example of a single newspaper page from Taiwan containing Chinese text written vertically, left to right, and right to left.

⁸⁹ Many thanks to Koji Kodama for his help translating the vertical-text examples.

良治は、「僕には二十三にんの子供がいる」と言っ
て、笑った。

Notice what happened to the Japanese quotation marks (「」): In the vertical text, they were rotated 90 degrees clockwise. This happens with a lot of punctuation marks: parentheses, ellipses, dashes, and so forth all get rotated. Some undergo additional transformations. A more subtle transformation can be observed if you look closely at the period and comma (。、). Their shapes stay the same, but they occupy a different position relative to the other characters: In horizontal text, the period and comma appear in the lower left-hand corner of its display cell, like a Western period. But in vertical text, they appear instead in the upper right-hand corner. Various other characters, such as the small kana in Japanese, keep their shapes but get drawn in different parts of their display cells depending on the directionality of the surrounding text.

Except for a few characters in the compatibility zone (see Chapter 12), Unicode adheres to the characters-not-glyphs rule for those characters whose glyphs change depending on the directionality of the surrounding text. An implementation that supports vertical text must be smart enough to make the appropriate glyph selections when drawing vertical text.

In some cases, exactly how the character is transformed depends on the locale: The colon (:), for example, looks the same in horizontal Japanese and Chinese text, but in Japanese vertical text it turns sideways. In Chinese text, on the other hand, it remains upright but gets smaller and moves to the upper right-hand corner of its display cell:

Chinese:	Japanese:
良 治 は :	良 治 は ..

Unicode doesn't give these language-specific transformations separate code point values, just as it doesn't for language-specific variants of certain Han characters. The assumption is that you'll use a font designed for the language you're writing in, and that the font will be selected either manually by

the user and specified using a higher-level protocol, or selected automatically by the system based on analysis of the text (or, in a pinch, use of the language-tagging characters).

Things get even more fun when you start mixing non-Asian text with Asian text. Let's say, for example, that we used Western digits for the "23" in our pervious example:

良治は、「僕には23にんの子供がいる」と言って、笑った。

When you write this vertically, there are three basic ways you can treat the "23." You can keep the digits upright and stack them vertically:

良	良
治	治
は	は
、	、
「	「
僕	僕
に	に
は	は
2	2
3	3
に	に
ん	ん
の	の
子	子
が	が
い	い
る	る
」	」
と	と
言	言
っ	っ
て	て
、	、
笑	笑
っ	っ
た	た
。	。

You can rotate the number ninety degrees clockwise:

良治は、「僕には23にんの子
供がいる」と言つて、笑つた。

Or you can keep the “23” upright and put the digits side-by-side in a single vertical display cell:

良治は、「僕には23にんの子
供がいる」と言つて、笑つた。

The third approach (often referred to in Japanese as *tate-chu-yoko*, or “horizontal in vertical”) is generally reserved for abbreviations and other very short snippets of text, where you’ll even see it done with native characters (you’ll see month names, which are written in Japanese with a Western numeral and the “moon” character, written in a single display cell (10月), and will see some abbreviations consisting of as many as four or five characters stuck into a single display cell (キリスト)). The second approach, rotating the horizontal text, works best for long snippets of foreign text.

Since Unicode allows for mixing of languages that were much harder to mix before, such as Arabic and Japanese, you also run into the possibility of mixing right-to-left text with vertical text. There are

no hard-and-fast rules for how to do this. Generally, especially with Arabic, you'd do it by rotating the text. It makes sense to keep the overall top-to-bottom writing direction, so you'd probably rotate right-to-left text ninety degrees counterclockwise:

と 良
言 治
つ は
て、 、
笑 「
つ 笑
た。 っ
た。」

Of course, you could also mix left-to-right, right-to-left, *and* vertical text in the same document. In this case, you'd probably rotate all of the horizontal text the same direction (probably ninety degrees clockwise) and use the Unicode bi-di algorithm to lay out the individual pieces of the horizontal text relative to each other (which means the right-to-left text would actually read from the bottom of the page up).

良 良
治 治
は は
、 、
「The first two books of the Bible are
「
笑 笑
つ っ
て、 た。
、」
and
笑 笑
つ っ
て、 た。
、」

Unicode doesn't set forth any specific rules for mixing vertical and horizontal text. The correct behavior depends on the situation. (It's unclear that there's a word-processing system out there that can even handle all of the above examples.)

Ruby

Even native speakers of the languages that use the Han characters don't know *all* of them. This means that in East Asian typography (especially Japanese) you'll see annotation techniques used to

help readers deal with unfamiliar characters. The most common of these techniques is something called *interlinear annotation*. This is usually referred to, at least in Japanese, as “ruby” or *furigana*.⁹⁰

The basic idea is that you adorn a character with one or more smaller characters that clarify its meaning or, more commonly, its pronunciation. In Japanese, you might see an unusual Kanji character adorned with a few small Hiragana characters that give the Kanji character’s pronunciation. The rubies go above the character being annotated in horizontal text, and (generally) to the right in vertical text. The effect is something like this, using English:

SLO-buh-dahn mi-LO-she-vitch
Slobodan Milosevic was defeated by
VOY-slav ko-SHTOO-nit-sa
Vojislav Kostunica in the Yugoslav
general election.

Ruby occurs relatively rarely in most text, but frequently in children’s books or books written for non-native speakers learning the language. It’s most common in Japanese, but also occurs in Chinese and Korean.

The Interlinear Annotation characters

Ruby is one of those things generally best left to a higher-level protocol, but Unicode includes some very rudimentary support for ruby for those situations where it absolutely has to be exchanged using plain text. (They can also be useful as an internal implementation detail, giving implementations a way to store the annotations in with the actual text without getting them mixed up.) There are three characters in the Specials block that are used to represent ruby:

```
U+FFF9 INTERLINEAR ANNOTATION ANCHOR
U+FFFA INTERLINEAR ANNOTATION SEPARATOR
U+FFFB INTERLINEAR ANNOTATION TERMINATOR
```

This is how they work: U+FFF9 marks the beginning of a piece of text that is to be annotated. U+FFFA marks the end of the piece of text to be annotated and the beginning of the annotation. U+FFFB marks the end of the annotation. The above example would be represented as follows (with the anchor, separator, and terminator characters represented by $\boxed{\text{IA}}$, $\boxed{\text{S}}$, and $\boxed{\text{T}}$ respectively):

$\boxed{\text{IA}}$ Slobodan $\boxed{\text{S}}$ SLO-buh-dahn $\boxed{\text{T}}$ $\boxed{\text{IA}}$ Milosevic $\boxed{\text{S}}$ mi-LO-she-vitch $\boxed{\text{T}}$ was defeated by
 $\boxed{\text{IA}}$ Vojislav $\boxed{\text{S}}$ VOY-slav $\boxed{\text{T}}$ $\boxed{\text{IA}}$ Kostunica $\boxed{\text{S}}$ ko-SHTOO-nit-sa $\boxed{\text{T}}$ in the Yugoslav general
election.

Unicode doesn’t give you any way to control how exactly the ruby will be laid out, and it also stipulates that paragraph separators can’t occur inside annotations. A paragraph separator inside an annotation terminates the annotation.

⁹⁰ Ken Lunde, the Japanese text-processing expert whose book provided a lot of the information in this chapter, has a daughter named Ruby, which somehow seems quite appropriate.

You can attach more than one annotation to a piece of annotated text by separating the annotations with extra instances of U+FFFA INTERLINEAR ANNOTATION SEPARATOR.

Generally you only want annotations to be significant when you're laying out or otherwise rendering text; for other operations, such as sorting and line breaking, you want them to be transparent. So most applications can safely ignore all three annotation characters, plus all characters that occur between U+FFFA and U+FFFB (or between U+FFFA and the end of a paragraph). Typically, unless you know you're interchanging text with something you know understands these characters, they, along with any annotations, should be filtered out.

Yi

Finally, before we leave East Asia, let's take a quick look at one more writing system. The Yi (or Lolo or Nuo-su) people are one of the largest minority groups in China. Most live in various parts of southwestern China, but they're scattered all over southeastern Asia. The Yi language is related to Tibetan and Burmese and is written with its own script, called, not surprisingly, the Yi script, but also called Cuan or Wei.⁹¹

Classical Yi is an ideographic script, like the Chinese characters, and although it probably arose under the influence of the Han characters, it's not directly related to them, and the characters have a notably different look from Han characters. The earliest examples of Yi writing date back about 500 years, although linguists suspect it may have an actual history going back as much as 5,000 years. It's estimated that there are somewhere between 8,000 and 10,000 different characters in the surviving examples of classical Yi writing. Because the Yi have been scattered geographically, the same characters look quite different from group to group; there was never really any standardization.

In the 1970s, in an effort to increase literacy in Yi, a new writing system for Yi was developed. It takes some of the old Yi ideographs and removes their semantic value, creating a purely phonetic script. Like the Japanese Kana scripts, modern Yi writing is a pure syllabary: each character represents a whole syllable, and the characters can't be broken down into smaller components representing the original sounds (as, for example, Korean Hangul syllables can).

There are a lot more possible Yi syllables than there are Japanese syllables, so the Yi syllabary is considerably larger than the Japanese syllabaries. Each Yi syllable consists of an optional initial consonant and a vowel. There are forty-four initial consonants and ten vowels, for 440 possible syllables. Each syllable can have four tones, so each syllable has four different symbols, one for each tone. This means that there are 1,760 possible syllables in Yi. Of course, not all of these syllables actually exist in Yi: the standard Yi syllabary actually has 1,165 characters. (Compare Korean, where there are over 11,000 possible Hangul syllables, but only 2,500 or so actually occur in Korean.)

The characters representing the various Yi syllables don't have anything in common (they're derived instead from earlier ideographic characters with those sounds), with one exception: each syllable actually has *three* unique characters. The fourth tone (the middle-rising tone) is actually written using the character for the middle-high tone and putting an arch (an inverted breve) over it.

⁹¹ My sources for this section are the Unicode standard itself and Dingxu Shi, "The Yi Script," in *The World's Writing Systems*, pp. 239-243.

Thus there are actually 819 unique Yi characters; the rest are duplicates with the arch over them. For example, here are the four characters for the syllable *pi* (these represent that syllable with, respectively, the high, middle-high, low-falling, and middle-rising tones):

𐩈 𐩉 𐩊 𐩋

Yi is written horizontally, from left to right.

The Yi Syllables block

Currently, Unicode encodes only the modern Yi syllabary, not the classical Yi ideographic script. The Unicode Yi Syllables block runs from U+A000 to U+A48F and is based on the Chinese national standard for Yi, GB 13134-91, and contains separate code point values for each of the 1,165 syllables. The characters representing syllables with the middle-rising tone get their own code point values, and even though they're the same as the characters representing the syllables with the middle-high tone with an added arch, they *don't* decompose and *can't* be represented with combining character sequences. Unlike the Han characters and the Hangul syllables, each Yi syllable has an assigned name rather than an algorithmic name. The names are an approximation of the syllable's sound, with an extra character added to the end to indicate the tone: That is, the four characters shown in the previous example have the following code point values and names:

```
U+A038 YI SYLLABLE PIT
U+A03A YI SYLLABLE PI
U+A03B YI SYLLABLE PIP
U+A039 YI SYLLABLE PIX
```

The Yi Radicals block

The Yi Radicals block (U+A490 to U+A4CF) contains a group of radicals which function similarly to the radicals in the Han characters. The Yi syllables are divided into categories according to certain visual elements (radicals) they have in common, and these radicals are used to facilitate lookup in dictionaries. This block contains fifty radicals.

CHAPTER 11 *Scripts from Other Parts of the World*

After four chapters and more than 150 pages, we've taken a good in-depth look at the four major script groups still in use around the world today. To recap:

- In Chapter 7, we looked at the five European alphabetic scripts, which are descended from the ancient Phoenician alphabet by way of the ancient Greek alphabet. These scripts are all alphabetic in nature, consists of upper-/lower-case pairs, are written from left to right, and use spaces between words. There is little or no typographic interaction between characters in the printed forms of these scripts, but most of these scripts make moderate to heavy use of combining diacritical marks. These scripts are used throughout Europe and other parts of the world colonized by the Europeans.
- In Chapter 8, we looked at the four Middle Eastern scripts, which are descended from the ancient Phoenician alphabet by way of the ancient Aramaic alphabet. These scripts are consonantal in nature, generally using letters for consonant sounds and writing most vowel sounds, if at all, with combining marks above or below the consonants. They're written from right to left. Two of these scripts, Arabic and Syriac, are cursive, with the characters in a word generally being connected and with the shapes of the letters varying greatly depending on the surrounding letters. These scripts are used in the Middle East (i.e., southwestern Asia and northeastern Africa) and throughout the Muslim and Jewish worlds.
- In Chapter 9, we looked at the fifteen Indic scripts, which are descended from the ancient Brahmi script, which in turn is probably descended from the ancient Phoenician alphabet by way of the Aramaic alphabet. These scripts are alphasyllabic in nature, with individual characters for consonant and vowel sounds reordering, reshaping, and combining in interesting ways to form clusters representing syllables. They're written from left to right and most, but not all, of them use spaces between words. These scripts are used throughout south and southeastern Asia.
- Finally, in Chapter 10, we looked at the Han characters, which have an independent history going back more than 3,000 years and didn't descend from anything. They're ideographic or

logographic in nature, numbering in the tens of thousands, have numerous regional and historic variants of some characters, and are an open-ended set to which new characters continue to be added. They're supplemented in various ways by the other three scripts we looked at, two of which are syllabic (Japanese and Yi) and one of which is alphabetic (Korean). These scripts are traditionally written vertically, but can also be written from left to right, spaces are *not* used between words (except in Korean), and the characters generally don't interact typographically. These scripts are used in East Asia.

All in all, we've looked at twenty-eight different writing systems, but you may have noticed something interesting: We've really only looked at Europe and Asia. What about the other continents?

The Europeans colonized Africa, Australia, and the Americas, and so their languages (along with Arabic in northern Africa) have come to be spoken on those continents, and their writing systems (especially the Latin alphabet) have come to be the chief writing systems used in those places. Among the indigenous languages used in most of those places, most of those that are still spoken either don't have a written language or historically didn't and now have a recently-developed written language based on one of the European scripts (usually the Latin alphabet). But there are exceptions.

In this chapter, we'll take a look at the four major exceptions. These four scripts don't fit well into any of the other categories we've looked at, although some of them are related to the others. Three of the four are indigenous to areas outside Europe and Asia.

Mongolian

The exception is the Mongolian alphabet, which is (not surprisingly) used in Mongolia. Mongolian is written using this alphabet in Inner Mongolia (the Mongolian district in China); in Outer Mongolia (the Mongolian People's Republic), Mongolian was written using the Cyrillic alphabet for much of the twentieth century, but the Mongolian alphabet was restored by law in 1992. (Unlike other biscriptal languages, such as Serbo-Croatian, there is no direct one-to-one mapping between the two versions of written Mongolian; spelling of words in the two scripts is independent.)

The Mongolian alphabet is at best only very distantly related to the other Asian and Southeast Asian scripts; it's actually a cousin of the Arabic alphabet. It has been used in Mongolia since the early thirteenth century, when it evolved from the old Uighur script, which had evolved from the ancient Aramaic alphabet. (Modern Uighur is written using the Arabic alphabet.)

Like Arabic, Mongolian letters are cursively connected and change their shape depending on the surrounding letters. But there are a couple of important differences:

- The Mongolian script is fully alphabetic, using full-fledged letters for both vowels and consonants, instead of using combining diacritical marks for vowel sounds.
- The Mongolian script (like the ancient Uighur script before it) is written *vertically*. Letters proceed from the top of the page to the bottom, and lines of text proceed from left to right (unlike the East Asian scripts, where lines proceed from right to left).

The effect is similar to what you'd get if you took a page of Arabic writing and rotated it ninety degrees counterclockwise. In fact, this is probably how the difference happened: Right-handers often write Arabic vertically in lines that go from left to right across the page, even though it's read horizontally—this helps keep from smudging the page. It's not a big leap from writing with the page held like this to reading it that way as well.

The Mongolian alphabet has thirty-five letters:

[example]

(In this example, the same letterforms are used as are used in the Unicode code charts—this is generally the initial or independent form of the letter. Some letters have the same initial forms, and so for some letters some other form of the letter is used so that they all appear different in the example.)

The Mongolian alphabet is also used for a few other languages, including Todo, Sibe, and Manchu. The Mongolians also use their alphabet for religious and classical texts in Sanskrit and Tibetan; extra letters are added for each of these languages.

As with Arabic, the letters of the Mongolian alphabet change shape depending on their context, with most letters having initial, medial, and final forms. Some letters have similar appearances, with the initial form of one letter serving as the medial form of another, or the final form of one letter being the initial form of another, etc. Context and the principle of vowel harmony, which specifies which combinations of letters in a word are legal, allow similar-looking characters to be differentiated. Certain pairs of letters also form ligatures.

The shaping rules for Mongolian are actually more complicated than those for Arabic. Some letters take on special shapes in certain contexts, and these can't always be determined algorithmically; sometimes the correct shape for a letter is something other than the default choice, and the exact shape to use is a matter of spelling (i.e., different in that position depending on the actual word). Some letters have up to ten different contextual glyph shapes.

Spaces are used between words in Mongolian, but white space also appears within some words. Many words consist of a stem word with one or more suffixes attached to it; these suffixes are separated from the basic word with a small space. Sometimes the presence of the space causes the letters on either side of it to take on special shapes. Also, sometimes when the letters **[a]** (*a*) or **[e]** (*e*) appear at the end of the word, they're separated from the rest of the word with a break in the continuity of the cursive stroke; this usually affects the shape of the preceding letter.

Like many writing systems, Mongolian has its own period (**[glyph]**), comma (**[glyph]**), colon (**[glyph]**), and so forth. A special symbol (**[glyph]**) is used to mark the end of a chapter or a text. In Todo, hyphenation is indicated with a small hyphen at the *beginning* of the line containing the second half of the divided word, rather than with a hyphen at the end of the line containing the first half, as in English. There is also a set of Mongolian digits:

[example]

There aren't any really hard-and-fast rules governing what happens when Mongolian text is combined with text written in other scripts. If words from normally horizontal scripts are included in the middle of a Mongolian text, they're generally rotated ninety degrees clockwise, as they would be in Chinese or Japanese (right-to-left scripts would be rotated ninety degrees counterclockwise, to preserve the top-to-bottom ordering, unless the included excerpt consisted of a mixture of left-to-right and right-to-left text). If a few snippets of Mongolian are included in a generally left-to-right text, the Mongolian can be rotated ninety degrees counterclockwise, preserving the left-to-right reading order of the text (this generally only works with short excerpts, though, as the relative order

of successive lines of Mongolian text would be backwards when you do this—for long blocked quotes using the Mongolian script, it'd be better to preserve the vertical directionality).

[This section would probably be better with some examples and a bit more specificity on the shaping behavior. Unfortunately, my normal sources are failing me: the article on Mongolian in TWWS really sucks and doesn't give any of the needed detail. Need to find a better source on Mongolian to supplement the information in TUS itself.]

The Mongolian block

The Unicode Mongolian block runs from U+1800 to U+18AF. It includes not only the basic Mongolian alphabet, but extensions necessary for Todo, Sibe, Manchu, Sanskrit, and Tibetan; the digits; and various Mongolian- and Todo-specific punctuation marks.

As with all of the other scripts that have contextual glyph shaping as one of their basic features, Unicode only encodes the letters themselves, not the individual glyphs. There's no “presentation forms” block containing code point values for the individual glyph shapes for the Mongolian letters.

There are a few special formatting characters in the Mongolian block: the three characters U+180B, U+180C, and U+180D, MONGOLIAN FREE VARIATION SELECTOR ONE, TWO, and THREE, are used to affect the contextual shaping of the letters. They're intended for the situations where a letter has a special shape only in certain positions or words and you can't get the correct shape using the zero-width joiner and non-joiner (which work the same way on Mongolian as they do on Arabic). The free variation selectors have no visible representation themselves; they simply serve as a signal to the text rendering process to use an alternate glyph for the preceding letter. The three free-variation selectors allow for up to three alternate glyphs for each letter. The characters have no effect when they follow something other than a Mongolian letter.

U+180E MONGOLIAN VOWEL SEPARATOR is similar in function to the free variation selectors, but is intended specifically to be used before word-final **[a]** and **[e]**. Unlike the free variation selectors, it does have a visual presentation: it appears as a thin space. But it also causes the preceding character to take on a special shape.

Finally, U+202F NARROW NO-BREAK SPACE has a special function in Mongolian: it's used for the other situations where suffixes are set off from the base word with a small space. This character counts as part of the word; it's used for word-internal space. It can also cause special choices of contextual glyph selection.

The narrow no-break space works differently in Mongolian than does the regular non-breaking space (U+00A0). The regular no-break space is a normal-size space rather than a narrow one, so it's not suitable for use as a word-internal space in Mongolian; it also doesn't cause any special Mongolian glyph selection to happen. It can be used to glue two words together on the same line, but shouldn't be used for the word-internal spaces.

Finally, special mention should be made of U+1806 MONGOLIAN TODO SOFT HYPHEN. This works like the regular soft hyphen (see Chapter 12)—that is, it marks a legal hyphenation position in a Todo word and is invisible unless the word is actually hyphenated there, but it appears at the beginning of the second line rather than at the end of the first line.

At the time version 3.0 of the Unicode standard was published, the definition of the Mongolian block was incomplete. In particular, it lacked full information on the Mongolian shaping rules, especially

the definition of exactly which glyphs the free variation selectors caused when used in conjunction with each of the letters. The book anticipated publication of the full Mongolian shaping rules as a Unicode Standard Annex and mentions a special committee made up of experts from China, Mongolia, and other interested parties similar to the Ideographic Rapporteur Group that is working on this document. As I write this in late February 2001, that annex has not yet appeared on the Unicode Web site, even in draft form. The updated information was also not included in UAX #27, Unicode 3.1, so apparently the definition of Mongolian is still incomplete in Unicode 3.1.

[Need to do research to find out what the current status of this work is, probably updating this section closer to publication time after finishing the rest of the chapters. If it's possible to get my hands on a copy of "Users' Convention for System Implementation of the International Standard on Mongolian Encoding," alluded to on p. 289 of TUS (I'm assuming this will either be the eventual UAX, or the UAX will be a summary of this document), I should do so and include the preliminary information here.]

Ethiopic

The Ethiopic script was originally developed to write the Ge'ez language spoken in ancient Ethiopia. Ge'ez today is limited to liturgical use, but its writing system is also used to writing several different modern languages spoken in Ethiopia and Eritrea, including Amharic, Oromo, Tigre, and Tigrinya.⁹²

Ge'ez and Amharic are southern Semitic languages, cousins of Arabic and Hebrew. The writing system is also of Semitic origin, more distantly related to the Arabic and Hebrew alphabets. It doesn't actually descend from the Phoenician alphabet, like so many other modern scripts, but it shares a common ancestor with the Phoenician alphabet.⁹³ The Ethiopic script was adapted from an earlier script, the Sabeian script, which was used for Sabeian, another southern Semitic language. A distinct writing system for Ge'ez developed sometime in the early fourth century. Several letters have since been added to represent sounds in Amharic that aren't also in Ge'ez (the added characters can frequently be identified by a barbell-shaped stroke at the top of the character—compare **ከ** and **ኸ**, for example).

Originally, Ethiopic was a consonantal script, like Hebrew and Arabic, with letters representing only consonant sounds. The basic Amharic consonantal alphabet has thirty-three letters:

ሀ ለ ሐ መ ሠ ረ ሰ ሸ ቀ ባ ተ ቸ ጎ ነ ኘ ኡ ከ ኸ ወ ዐ ዘ ዠ የ ደ ጀ ገ ጠ ጪ ጰ ጸ ፀ ፈ ፐ

This was fairly quickly deemed unsuitable for Ge'ez, and a system of vowel marks was developed to supplement the basic consonantal script. Unlike Hebrew and Arabic, where the vowel sounds are represented (if at all) with various dots placed near and around the basic letter, the Ethiopic consonants sprouted various appendages to represent the vowels, and these became mandatory. There are seven vowel sounds, so there are seven combinations of base consonant and vowel mark. Here's a typical example, on the basic letter *qaf*:

⁹² My sources for the section on Ethiopic are Getatchew Haile, "Ethiopic Writing," in *The World's Writing Systems*, pp. 569-576; as well as the Unicode standard itself, pp. 284-286; and a couple facts from Nakanishi, pp. 102-103.

⁹³ See chart on p. 89 of *The World's Writing Systems*.

ቀ ቁ ቂ ቃ ቄ ቅ ቆ

A few consonants also have labialized forms, where a *w* sound is interpolated between the main consonant sound and the main vowel sound. The consonant sprouts extra appendages to represent these combinations:

ቈ ቑ ቒ ቓ ቔ

The formation of some of these different consonant-vowel combinations can be irregular. Sometimes the character sprouts a leg to provide a place for the vowel mark to attach...

ወ ዑ ዒ ዓ ዔ ዖ

...and the marks for some vowels aren't always consistent. Compare, for example, these three versions of *lawe*...

ሌ ል ሎ

...with the corresponding versions of *hawt*:

ሐ ከ ኮ

Because of all this variation in the basic forms, today the Ethiopic script is generally thought of as a syllabary, with each character representing a combination of consonant and vowel sounds (sometimes with the interpolated *w*). The syllabary is normally shown as a 43 by 8 matrix: 43 consonants (there are ten extra consonants for writing languages other than Amharic) and eight vowels. This gives you 344 possible characters, although the actual number is smaller (not every syllable actually exists in the languages represented).

Ethiopic is written from left to right, and the syllable characters don't interact typographically. Traditionally, a colon-like mark (፥) was used to separate words, but the ordinary Western space is becoming common. Other punctuation consists of the word-space character with other marks added to it: ፡ is a period, ፣ a comma, ፤ a semicolon, and ። a question mark, for example.

Ethiopic has its own system of writing numbers. It uses an additive-multiplicative system similar to that used in Chinese, rather than Western positional notation (for more info, see Chapter 12). These are the basic number characters, representing the numbers from 1 to 10 and 100:

፩ ፪ ፫ ፬ ፭ ፮ ፯ ፰ ፱ ፲

The Ethiopic block

The Ethiopic block (U+1200 to U+137F) separately encodes all of the syllables, arranging them in the traditional consonant-vowel matrix, with gaps left for combinations of sounds that don't have corresponding characters. The Ethiopic syllables don't decompose; despite their historical relationship, they're all treated as separate, independent characters in Unicode. This block also contains the Ethiopic wordspace characters and the other punctuation marks, as well as a complete set of Ethiopic numerals.

Cherokee

Unlike all of the other scripts we've looked at so far, the Cherokee script doesn't have a long history. It was invented in 1821 by Sequoyah, a monolingual Cherokee, and was quickly picked up by the Cherokee nation in general, having spread as early as 1824.⁹⁴ Sequoyah apparently got the idea from seeing written English, and then went on to invent his script based on the general idea. The modern forms of many of the characters resemble Latin, Greek, or Cyrillic letter and Arabic digits, although there's no relationship between the sounds of the characters in Cherokee and the sounds of similar-looking characters in other scripts. (Sequoyah's original letterforms are said to have been unique, but the forms were later modified by missionaries to look more like European letters in an effort to make the language easier to print.)

The Cherokee script is a syllabary, with eighty-five characters representing various combinations of thirteen basic consonants and six vowels:

D	R	T	Ꭰ	Ꭱ	i
Ꭲ	Ꭳ	Ꭴ	Ꭵ	Ꭶ	E
Ꭷ	Ꭸ	Ꭹ	Ꭺ	Ꭻ	Ꭼ
Ꭽ	Ꭾ	Ꭿ	Ꮀ	Ꮁ	Ꮂ
Ꮃ	Ꮄ	Ꮅ	Ꮆ	Ꮇ	
Ꮈ	Ꮉ	Ꮊ	Ꮋ	Ꮌ	Ꮍ
Ꮎ	Ꮏ	Ꮐ	Ꮑ	Ꮒ	Ꮓ
Ꮔ	Ꮕ	Ꮖ	Ꮗ	Ꮘ	Ꮙ
Ꮚ	Ꮛ	Ꮜ	Ꮝ	Ꮞ	Ꮟ
Ꮠ	Ꮡ	Ꮢ	Ꮣ	Ꮤ	Ꮥ
Ꮦ	Ꮧ	Ꮨ	Ꮩ	Ꮪ	Ꮫ
Ꮬ	Ꮭ	Ꮮ	Ꮯ	Ꮰ	Ꮱ
Ꮳ	Ꮴ	Ꮵ	Ꮶ	Ꮷ	Ꮸ
Ꮹ	Ꮺ	Ꮻ	Ꮼ	Ꮽ	Ꮾ
Ꮿ	Ᏸ	Ᏹ	Ᏺ	Ᏻ	Ᏼ

94 My source is Janine Scancarelli, "Cherokee Writing," in *The World's Writing Systems*, pp. 587-592.

These 85 characters actually represent a far wider selection of actual sounds, as can be seen by the fact that there are multiple characters in certain cells of the chart, representing variants of the basic consonant sounds in those syllables. (Other cells have similar variant sounds, but they're all represented with the same character.) With a few exceptions (often involving ou , which represents the *s* sound without a vowel), successive consonant sounds with no intervening vowel sounds are written using an arbitrary syllable with the right consonant sound—the vowel sound is understood to be silent.

Cherokee is written from left to right, with Western punctuation and spaces between words (in handwritten Cherokee, words are sometimes separated by raised dots). Cherokee writing is basically caseless, but some printed materials in Cherokee use larger forms of the syllables in a way analogous to capital letters in English. The Cherokee syllables don't interact typographically and aren't used with combining diacritical marks. Sequoyah also invented a Cherokee numeration system, but it never caught on—today numbers are either written out in words or written using regular Western digits.

The Cherokee block

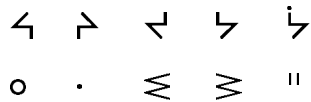
The Unicode Cherokee block runs from U+13A0 to U+13FF. It just contains the basic Cherokee syllables. Normal Western spaces, numerals, and punctuation are used with these characters. Unicode treats Cherokee as uncased; it doesn't include separate code point values for the larger versions of the letters used in some printed Cherokee materials; these must be represented using out-of-band styling information.

Canadian Aboriginal Syllables

The Canadian aboriginal syllabary is used to write various Native American languages in Canada.⁹⁵ It was originally invented in 1841 by the Wesleyan missionary James Evans to write Cree and Ojibwe, and the letterforms were originally based on a form of British shorthand. The characters generally represent syllables (combinations of initial consonant and vowel), but additional characters are used to represent syllable-initial vowels and syllable-final consonants. The basic Cree syllabary consisted of sixty-three characters, as follows:

▽	△	▷	◁	◁̇	
▽̇	△̇	▷̇	◁̇	◁̇̇	
∨	∧	>	<	<̇	ı
U	∩	∪	∩̇	∩̇̇	ı̇
q	p	d	b	ḃ	ı̇̇
ŋ	ɾ	j	ɬ	ɬ̇	-
ɫ	ɫ̇	ɫ̇̇	ɫ̇̇̇	ɫ̇̇̇̇	ı̇̇̇̇
ɔ	σ	ɔ̇	ɔ̇̇	ɔ̇̇̇	ı̇̇̇̇̇
ɥ	ɥ̇	ɥ̇̇	ɥ̇̇̇	ɥ̇̇̇̇	ı̇̇̇̇̇̇

⁹⁵ My source is John D. Nichols, "The Cree Syllabary," in *The World's Writing Systems*, pp. 599-611.



The vowels alone were represented with triangles pointing in various directions; the syllables were then represented with various other character shapes, with the orientation specifying the vowel sound. Small superscripted forms of the various syllables were generally used to represent syllable-final consonant sounds.

This script spread fairly rapidly to the different Native American communities in Canada, and it's now used for a variety of languages in the Algonquian, Inuktitut, and Athapascan languages families. Among others, the script is used to write Eastern and Western Algonquian, Cree, Ojibwe, Naskapi, Inuktitut, Inuktitun, Athapascan, Chipewyan, and Carrier.

As the script was adopted by each new group, it was modified as necessary to suit that language, with various new characters being added and dots and other marks added to many of the original characters. The end result is that virtually every language using this writing system has its own version with its own characters and its own mapping of the characters to sounds. Generally speaking, the spelling of the various languages isn't standardized, and various dots and marks used with the characters may be omitted by some writers.

The Canadian syllables are written from left to right, spaces are used between words, and (with one exception) English punctuation marks are used. The exception is the period, which looks like this: ×

The Unified Canadian Aboriginal Syllabics block

The Unified Canadian Aboriginal Syllabics block (U+1400 to U+167F) contains all of the characters necessary to write all the languages that are written with scripts based on the original Cree syllabary. The unification is based on glyph shape: if the same syllable has a different glyph in different languages, they're still encoded separately, and if different languages use the same shape for different sounds, they're still unified. The basic idea is to allow the mixture of different Canadian aboriginal languages in the same document without forcing a change of fonts.

Many of the characters in this block consist of a basic letterform with various added dots or other marks. These characters don't decompose—they're treated as fully independent first-class syllables rather than combinations of a basic character and one or more marks.

The Canadian aboriginal period and the chi sign (×), used to write the name of Christ in some texts, are also included in this block.

Historical scripts

With that, we've rounded out the set of 32 modern scripts encoded in Unicode. With the characters in these blocks, all of the modern written languages in common business or official use can be represented. There are certainly other modern scripts that are in use and could be encoded (and

most eventually will be), but the current Unicode set is sufficient to write at least one language written by practically every literate person on Earth.

However, there are many more writing systems that have fallen into disuse but still occur in scholarly writing about them. In many cases, there are also communities of hobbyists that use them. Unicode 3.1 includes encodings for five such scripts, and we'll take a brief look at each of them next.

Runic

The Runic script was used in various parts of Europe from the first century to the nineteenth century.⁹⁶ Its origins are uncertain, but the earliest inscriptions using runes were found in Denmark and Germany. The traditional Runic alphabet (or “futhark,” a name derived from the sounds of the first six letters) consisted of twenty-four letters:

ƿ ᚠ ᚢ ƿ ᚦ ᚨ < X ƿ H † | > J C Y 4 ↑ B M M
Γ ◊ M X

The use of runes spread north and west into Scandinavia and the British Isles, and as it spread to new peoples, various changes occurred both in the alphabet itself (addition and removal of letters, changes in their sounds) and in the shapes of the individual letters. In Scandinavia, the alphabet was shortened to sixteen letters, but in Britain it was extended to twenty-eight letters.

Runic was a full alphabetic script, with letters for both consonant and vowel sounds. Like many early writing systems, the writing direction varied: some inscriptions were written from left to right, some from right to left, and some in alternating directions (“boustrophedon”). Spaces weren't used between words, but sometimes words were separated by various dots. Certain pairs of letters form ligatures in certain inscriptions.

The Runic block in Unicode (U+16A0 to U+16FF) contains 75 letters. Letters with common historical origins are generally unified (there are a few exceptions, notably the so-called “long branch” and “short twig” runes), even when their shapes differ significantly. For this reason, a user must use a font designed especially for the variety of Runic he wants to represent. The Unicode Runic block also includes a few Runic punctuation marks whose exact meaning is unknown, and a few special runes that were only used for writing numbers. One interesting quirk of the Runic block is that the name of each character is actually the concatenation of *several* names for the character, in the cases where the same character was used in different places and times and had different names. This helps figure out which character you want when it's been unified with several other characters that had different names.

⁹⁶ The material from the Unicode standard is supplemented by material from Ralph W. V. Elliott, “The Runic Script,” *The World's Writing Systems*, pp. 333-339.

Ogham

Ogham (pronounced “Ohm”) was used in the fifth and sixth centuries for inscriptions on stone in early Irish (the surviving examples are generally territory markers and tombstones).⁹⁷ Its origins are uncertain, but there’s some consensus that Ogham was inspired by the Latin script. It’s alphabetic in character, with characters for both consonants and vowels. There are twenty basic marks, organized into four groups of five:

├ ┆ ┆ ┆ ┆ ┆├ ┆ ┆ ┆ ┆ ┆┆ ┆ ┆ ┆ ┆ ┆ ┆ ┆┆ ┆ ┆ ┆ ┆ ┆ ┆┆ ┆ ┆ ┆ ┆ ┆ ┆ ┆ ┆┆ ┆ ┆ ┆ ┆ ┆ ┆ ┆ ┆┆

[spacing problems again—the center line in all these marks should connect]

These characters originally were series of notches cut along the edge of the stone: the first group would be cuts in the stone along one face, the second group along the other face, the third group crossing both faces, and the fourth group would be notches just in the edge between the faces. Inscriptions would begin on the lower left-hand corner of the stone and continue up the left side, along the top, and down the right side. There was no punctuation or word division.

In later periods, Ogham was also written on paper, where a central line would often be used to symbolize the edge of the stone (hence the forms shown above). Six additional letters (“forfeda”) were added during this time. The normal left-to-right writing direction is usually used when writing Ogham on paper, and spaces are more common, at least in scholarly work.

The Unicode Ogham block runs from U+1680 to U+169F and includes the twenty-six Ogham letters, plus a space character consisting of just the center line with no marks crossing it (in fonts that don’t use the center line, this looks the same as a regular space) and “feather marks,” chevron-like marks used to mark the beginning and end of a text.

Old Italic

The Latin alphabet evolved from the ancient Etruscan alphabet, which evolved from an ancient form of the Greek alphabet. There are a bunch of related alphabets used around Italy in that same period that derived from the same version of the Greek alphabet. They may be thought of as siblings of the Etruscan alphabet. The Unicode Old Italic block, which runs from U+10300 to U+1032F in the SMP (i.e., Plane 1), is a unification of these various Etruscan-related alphabets used in Italy beginning in roughly the eighth century BC. In particular, the Etruscan, Umbrian, North and South Picene, Faliscan, Sabellian, Oscan, Volscian, Elimian, Sicani, and Siculan alphabets are unified into the Old Italic block (various other alphabets from that area and time period, such as Gallic, Venetic, Rhaetic, and Messapic, belong to different families and aren’t unified here. They remain unencoded in Unicode 3.1). The unified letters are related historically and phonetically, but don’t necessarily look the same in each of the unified alphabets; use of a font customized for the particular alphabet you want to represent is therefore required.

These alphabets were generally written from right to left, with *boustrophedon* sometimes occurring, but since modern scholarly usage generally turns them around and writes them left to right (the better to include glosses and transcriptions in modern Latin letters), Unicode encodes

⁹⁷ The material from the Unicode standard is supplemented with material from Damian McManus, “Ogham,” *The World’s Writing Systems*, pp. 340-345.

them as strict left-to-right characters. The directional-override characters can be used to display these characters in their original directionality; rendering engines should use mirror-image glyphs when these characters are written from right to left.

Punctuation was not used in any of these alphabets, and spaces were not used between words, although in later usage dots were sometimes used between words. Unicode doesn't include special codes for these dots, unifying them instead with similar symbols elsewhere in the standard. The numbering system used was not well attested, but some fairly well-attested numeral characters are included in this block; they work like Roman numerals.

Gothic

The language of the Goths, from a now-extinct branch of the Germanic languages, has come down to us only in very few examples, but is important to scholars because of the historical importance of the Goths and its status as a representative of a now-extinct branch of a major language group.⁹⁸ The few examples of written Gothic are written in a distinctive script that tradition holds was invented by the Gothic bishop Wulfila in the late fourth century, primarily to allow for a translation of the Bible into Gothic. Indeed, almost all examples of written Gothic are Bible fragments, and these fragments are important for Biblical scholarship. Many sources have been proposed for Wulfila's script—it probably derived at least partially, if not completely, from the Greek alphabet.

The Gothic alphabet has twenty-seven letters:

ᚠ ᚢ ᚦ ᚨ ᚱ ᚷ ᚹ ᚻ ᚾ ᚿ ᚰ ᚱ ᚴ ᚷ ᚸ ᚹ ᚰ ᚱ ᚴ ᚷ ᚸ ᚹ
ᚺ ᚻ ᚼ ᚽ ᚾ

Two of these (ᚱ and ᚾ), as far as we know, were only used for writing numerals. The letter ᚷ is written with a diaeresis (ᚷ̈) when it occurs at the beginning of a word or at certain grammatically-significant positions within a word. Contractions are written with a bar over the word (\overline{xus} , *xus*, was the abbreviation for “Christ,” for example).

As with many ancient alphabets, the letters were also used to write numbers (as we saw, there were two letters that, as far as we know, were only used for this purpose)—numerals were set off with a dot on either side ($\cdot\eta\epsilon\cdot$) or with lines above and below ($\overline{\eta\epsilon}$).

Gothic was written from left to right; spaces were not used between words, but spaces, centered dots, and colons were all used between phrases and for emphasis.

The Unicode Gothic block runs from U+10330 to U+1034F in the SMP and contains only the twenty-seven Gothic letters. ᚷ̈ is represented using the code point for ᚷ (U+10339), followed by U+0308 COMBINING DIAERESIS. Similarly, U+0305 COMBINING OVERLINE is used for

⁹⁸ The information from the Unicode standard is supplemented by Ernst Ebbinghaus, “The Gothic Alphabet,” *The World's Writing Systems*, pp. 290-293.

the bar that indicates contractions. Numerals are set off using either U+00B7 MIDDLE DOT before and after, or with a combination of U+0304 COMBINING MACRON and U+0331 COMBINING MACRON BELOW attached to each character. The normal space, middle-dot, and colon characters are used for phrase division and emphasis.

Deseret

Unlike the other historical scripts in the current version of Unicode, all of which are of considerable antiquity, the Deseret alphabet is a modern invention. It was developed in the 1850s at the University of Deseret, now the University of Utah, and promoted by the Church of Jesus Christ of Latter-Day Saints (the Mormons).⁹⁹ The word Deseret comes from the Book of Mormon, which defined it to mean “honeybee.” (The beehive was used as a symbol of cooperative industry by the early Mormons, who used “Deseret” as the name of the region where they were settling.)

The interesting thing about the Deseret alphabet is that it was invented for the writing of *English*, not some exotic language only spoken historically or by some small group of people. The idea was to provide a completely phonetic way of writing English in an effort to make it easier for Mormon immigrants from Europe to learn English. It was heavily promoted by Brigham Young over a period of some fifteen years, but never really caught on. On the other hand, it’s never completely died out, either—there is still interest among some Mormons, and some new materials continue to be published from time to time using this alphabet.

The Deseret alphabet has thirty-eight letters, sixteen vowels and diphthongs...

Ɔ (long i) Ǝ (long e) Ɔ (long a) Ɔ (long ah) Ɔ (long o) Ɔ (long oo)
† (short i) † (short e) † (short a) † (short ah) † (short o) † (short oo)
ɔ (ay) ɔ (ow) ɔ (wu) ɔ (yee)

...and twenty-two consonants:

Ɔ (h) Ɔ (pee) Ɔ (bee) Ɔ (tee) Ɔ (dee) Ɔ (chee) Ɔ (jee) Ɔ (kay)
Ɔ (gay) Ɔ (ef) Ɔ (vee) Ɔ (eth) Ɔ (thee) Ɔ (es) Ɔ (zee) Ɔ (esh)
Ɔ (zhee) Ɔ (er) Ɔ (el) Ɔ (em) Ɔ (en) Ɔ (eng)

The letters are written from left to right, spaces are used between words, punctuation and numerals are the same as in English, and no diacritical marks or other complex typography are required. The letters come in capital/small pairs, but they differ only in size, not in shape. One interesting quirk of Deseret writing was that when the name of a letter was pronounced the same as an actual word,

⁹⁹ The information from the Unicode standard is supplemented by John H. Jenkins, “Adding Historical Scripts to Unicode,” *Multilingual Computing & Technology*, September 2000.

CHAPTER 12 *Numbers, Punctuation, Symbols, and Specials*

It's been a long trip, but we've finally covered all of the writing systems in Unicode 3.2. But we're not done yet with our guided tour of the character repertoire. At this point, we've looked at basically all of the characters in Unicode that are used to write words: the various letters, syllables, ideographs, diacritical marks, and so on that are used to write the various languages of the world. In this chapter, we'll look at everything else:

- The various characters used to write numeric values, and the various characters used with them.
- The punctuation marks.
- The special Unicode characters and non-character code points, including control characters, invisible formatting characters, non-characters, and other characters with special properties.
- Various symbols.
- Various miscellaneous characters, such as the dingbats, the box-drawing characters, the geometric shapes, and so on.

We've checked out some of these characters as we've gone through the writing systems, mainly the ones that are used only with particular scripts. In this chapter, we'll look at those characters briefly again in the context of a more comprehensive treatment of their character categories.

Numbers

Written numbers are probably older than written words, and it's just as import to have characters for the writing of numeric values as it is to have characters for the writing of words. There are lots of digit and numeral characters in Unicode, as well as many characters that do double duty as digits or

numerals and as something else (usually letters or ideographs). In this section, we'll look at all of them.

Western positional notation

The numeration system we use in the United States and Europe is called positional, or place-value notation. When used to write decimal numbers, positional notation makes use of nine basic “digit” characters that represent not only the values from one to nine, but also those values times the various powers of ten. The exact value of a digit (i.e., which power of ten you multiply by) is determined by the digit’s position within the numeral. The rightmost digit in an integer keeps its nominal value, and each succeeding digit to the left has its nominal value multiplied by a higher power of 10. A special placeholder digit, called a “zero,” is used to mark digit positions that aren’t used, so that the meanings of the surrounding digits are unambiguous.

The positional system of numeration and the shapes of the characters we use as digits have their origins in India somewhere around the fifth century AD.¹⁰⁰ In the West, the characters are commonly called “Arabic numerals” because the Europeans got this system of numeration from the Arabs. (As if this weren’t confusing enough, the Arabs call their numerals “Hindi numerals” because they got them from the Indians.)

Although pretty much everybody using positional notation uses digits derived from the original Indian ones, the shapes have diverged quite a bit in the intervening centuries. Unicode includes sixteen different sets of decimal digits. Each set of digits is used to make up numerals in the same way our familiar Western digits are. In some countries, both the Western digits and the native digits are used. Unicode encodes each set of digits in the block for the script that uses those digits (the Western digits are in the ASCII block), and in every case puts them in a contiguous block arranged in ascending order from 0 to 9. Number-formatting code (i.e., code that converts a number from its internal representation into characters for display to the user) can thus be adapted to use a locale’s native digits rather than the European digits simply by telling it which code point to use for 0—the others follow automatically.

Here’s a table of the various sets of decimal digits in Unicode **[still need to fill in Mongolian]**:

European	U+0030	0	1	2	3	4	5	6	7	8	9
Western Arabic	U+0660	٠	١	٢	٣	٤	٥	٦	٧	٨	٩
Eastern Arabic	U+06F0	۰	۱	۲	۳	۴	۵	۶	۷	۸	۹
Devanagari	U+0966	०	१	२	३	४	५	६	७	८	९
Bengali	U+09E6	০	১	২	৩	৪	৫	৬	৭	৮	৯
Gurmukhi	U+0A66	੦	੧	੨	੩	੪	੫	੬	੭	੮	੯
Gujarati	U+0AE6	૦	૧	૨	૩	૪	૫	૬	૭	૮	૯
Oriya	U+0B66	୦	୧	୨	୩	୪	୫	୬	୭	୮	୯

¹⁰⁰ Ifrah, p. 420.

Numbers

Telugu	U+0C66	౦	౧	౨	౩	౪	౫	౬	౭	౮	౯
Kannada	U+0CE6	೦	೧	೨	೩	೪	೫	೬	೭	೮	೯
Malayalam	U+0D66	൦	൧	൨	൩	൪	൫	൬	൭	൮	൯
Thai	U+0E50	๐	๑	๒	๓	๔	๕	๖	๗	๘	๙
Lao	U+0ED0	໐	໑	໒	໓	໔	໕	໖	໗	໘	໙
Tibetan	U+0F20	༠	༡	༢	༣	༤	༥	༦	༧	༨	༩
Myanmar	U+1040	၀	၁	၂	၃	၄	၅	၆	၇	၈	၉
Khmer	U+17E0	០	១	២	៣	៤	៥	៦	៧	៨	៩
Mongolian	U+1810										

There are other characters in Unicode that have the “decimal digit” property, but they’re either presentation forms, or the series isn’t complete because a different system of numeration is used (typically, the zero is missing and there are additional characters for 10, 100, and so on).

Alphabetic numerals

Many cultures, notably the Greeks and Hebrews, although there are others, have used the letters of their alphabets as digits as well as letters. Typically the way this works is that the first nine letters of the alphabet represent the values 1 to 9, the next nine letters represent the multiples of 10 from 10 to 90, the third group the multiples of 100 from 100 to 900, and so on. In some languages, such as Greek and Gothic, some of the letters eventually fell into disuse as letters but continued to be used as digits. In Hebrew, the alphabet only had twenty-four letters, so 700, 800, and 900 had to be written using combinations of characters. In most languages that used letters to write numbers, there was some special mark that was used to distinguish numerals from words.

None of the cultures that used alphabetic notation for writing numbers still does, except sometimes for things like list numbering or dates (similar to when we might use Roman numerals). Unicode, in keeping with its principle of unification, just designates the regular letters as doing double duty as digits when alphabetic notation is used for historical reasons. In all cases, letters that came to be used only as digits, and the special distinguishing marks used to mark sequences of letters as numerals instead of words are also encoded.

Roman numerals

Just as we’re all familiar with Arabic numerals, we’re all familiar with Roman numerals, where the letters I, V, X, L, C, D, and M are used in various combinations to represent numeric values. Roman numerals used a basically additive notation, with the values of the characters just being added up to get the value of the whole numeral (XVIII = 10 + 5 + 1 + 1 + 1 = 15) [the subtractive notation that causes IV to be interpreted as 4 was a later development, and doesn’t really change the basic nature of the notation]. The I, V, and X characters probably evolved out of early methods of representing

numeric values by making notches or scratch marks in things—similar methods exist in lots of ancient cultures.

The basic way of representing Roman numerals in Unicode is the same as for representing alphabetic notation: to use the normal Latin-letter characters. There’s no need for special I, V, X, etc. characters for writing Roman numerals. There are, however, exceptions.

The M for 1,000 was a fairly late development; before it the symbol **ↀ** was used in Roman numerals to represent 1,000. You could represent higher powers of 10 by adding concentric circles: for example, **ↁ** represents 10,000. You cut these symbols in half to represent half the value:

ↂ represents 5,000, and this is also where D for 500 comes from. Unicode includes code points for these three characters in the Number Forms block at U+2180, U+2182, and U+2181 respectively.

In more recent centuries, you’ll see the symbols for the large numbers represented using whole characters. In situations where printers didn’t have type for the special characters, they’d fake them using C, I and an upside-down C: **ↀ** would be written as Cↀ and **ↁ** as CCↀↀ. The upside-down C is in the Number Forms block at U+2183. **[MS Word bug is causing bad spacing between upside-down Cs]**

A more modern practice represents large numbers using bars over the traditional letters: a bar over a

letter multiples its value by 1,000: **ↆ** **[fix alignment]**, for example, represents 10,000. You can use U+0305 COMBINING OVERLINE to do this.

Finally, even though Unicode doesn’t need special code-point values for the other Roman-numeral characters, it has them. This is a legacy issue: Several Japanese and Chinese character-encoding standards have code points for the Roman numerals (where they’re useful for such things as making a whole Roman numeral show up in a single display cell in vertical text), and Unicode adopts them as well for full round-trip compatibility with these standards. For this reason, the Unicode Number Forms block includes two sets of Roman numerals, representing the values from 1 to 12, plus 50, 100, 500, and 1,000. One series uses capital letters and the other uses small letters. Thanks to this, you can represent “VIII” either with four code points, for the letters V, I, I, and I, or with a single code point, U+2167, representing the concept “Roman numeral 8.” These characters occupy the range U+2160 through U+217F in the Number Forms block, and all of them have compatibility decompositions to sequences of regular letters.

Han characters as numerals

Certain Han characters represent basic numeric values and are combined using an additive-multiplicative system to write the other numeric values. These are the basic characters representing the numbers from 1 to 9:

一	二	三	四	五	六	七	八	九
1	2	3	4	5	6	7	8	9

They're used in conjunction with these characters that represent various powers of 10:

十	百	千	万	億	兆
10	100	1,000	10,000	100,000,000	1,000,000,000,000

The additive-multiplicative principle is pretty simple: if one of the basic-number characters appears before a power-of-10 character, their values are multiplied together. If it appears after it, their values are added together. So you get this:

二 = 2

十 = 10

十二 = 12

二十 = 20

二十二 = 22

二百二十二 = 222

For values above 10,000 you get a new character every power of 10,000 instead of every power of 10. This is analogous to how in English, for values above 1,000, you introduce a new word for every power of 1,000 instead of every power of 10. Just as 123,456 is “one hundred twenty-three thousand four hundred fifty-six,” so 1,234,567 in Han characters is

百二十三万四千五百六十七

The last part is simple, representing 4,567 with pairs of characters for 4,000, 500, and 60, and a single character for 7:

四千五百六十七

The first part follows the same principle for 123 (a single character for 100, a pair of characters for 20, and a single character for 3) and follows this sequence with the character for 10,000, which multiplies the whole sequence by 10,000. Thus, this sequence represents 1,230,000:

百二十三万

Other Han characters are also used in special situations to represent numbers, particularly in accounting. For example, let's say somewhere in the amount on a check you have...

二

...representing 2. All you have to do is add another stroke to it to get 3:

三

In fact, you can add two additional strokes to get 5:

五

For this reason, alternate characters that can't be altered just by adding a couple strokes are used on checks or in other places where it's necessary to deter counterfeiting. Exactly which characters are substituted varies from region to region, but here's a fairly complete selection:¹⁰¹

一	二	三	四	五	六	七	八	九	十	百	千	万
1	2	3	4	5	6	7	8	9	10	100	1,000	10,000
壹	貳	參	肆	伍	陸	柒	捌	玖	拾	佰	仟	萬
壹	貳	參			陸					陌		
貳	貳	叁										
	貳	叁										

Another Han character is sometimes used to represent the value 0:

零

101 My source for the information in this chart is the Unicode standard, p. 97.

Sometimes, you'll also see the Han characters used with Western positional notation. In these situations, a special zero character is used. 1,203 would normally be written like this...

千二百三

...but you may occasionally see it written like this, using positional notation:

一 二 〇 三

All of the characters shown above are regular Han characters in the CJK Unified Ideographs block, where some have other meanings and some are used exclusively for writing numbers. The one exception is 〇, which is in the CJK Symbols and Punctuation block.

Also in the CJK Symbols and Punctuation block are the so-called “Hangzhou-style” numerals. They come to Unicode via a number of different East Asian encoding standards, but their provenance is rather cloudy. Some of them are similar to the “regular” Han characters for numbers, but based on vertical lines instead of horizontal lines:

丨			ㄨ	ㄣ	⊥	≡	≡	ㄨ	+	++	+++
1	2	3	4	5	6	7	8	9	10	20	30

It appears these are another set of “commercial” numerals, used by shopkeepers in parts of China (and maybe also Japan) to mark prices.¹⁰²

Finally, the Japanese and Koreans, and less frequently the Chinese, also use the Western digit characters to write numbers.

Other numeration systems

Tamil uses the same additive-multiplicative system used with the Han characters. The characters for 1 through 9 are...

௧ ௨ ௩ ௪ ௫ ௬ ௭ ௮ ௯

...and the following characters represent 10, 100, and 1,000, respectively:

¹⁰² My information here comes from Thomas Emerson, “On the ‘Hangzhou-Style Numerals in ISO/IEC 10646-1,” a white paper published by Basis Technology. Emerson cites the source encodings for these characters, his best information as to what they’re used for, and goes on to say that the term “Hangzhou” isn’t attested in any source he could find—he proposes the term “Suzhou,” which was attested in a number of sources, instead. (Thanks, by the way to Tom for forwarding his paper to me.)

ω π Ϡ

This system appears to be dying out in favor of Western numerals.

Ethiopic uses a system that's kind of a hybrid of the additive system underlying Greek and Hebrew alphabetic notation and the additive-multiplicative system used in Japanese and Chinese. The following nine characters represent the values from 1 to 9:

፩ ፪ ፫ ፬ ፭ ፮ ፯ ፰ ፱

The following nine characters represent the multiples of 10 from 10 to 90:

፲ ፳ ፴ ፵ ፶ ፷ ፸ ፹ ፺

The other values from 1 to 99 are represented using pairs of these characters. 100 is represented with this character:

፺፻

Multiples of 100 are represented by putting pairs of the other characters before the ፺፻. So 1,234 is written like this:

፲፪፻፳፻፳፻፺፻

The same pattern is followed for values over 10,000: each power of 100 is represented using an appropriate number of ፺፻ characters. Thus, 123,456 is

፲፪፻፳፻፳፻፺፻፲፪፻፳፻፳፻፺፻፲፪፻፳፻፳፻፺፻

A few languages also have special ways of writing fractions. Bengali has several special characters for writing the numerator of a fraction (✓, ২, ৩, ৪, and ৫, representing 1, 2, 3, 4, and one less than the denominator) and one special character for writing the denominator (৬, representing a denominator of 16).

Tibetan has special characters representing the halves. Each is a regular Tibetan digit with a hook that subtracts 1/2 from its value: For example,

༣

...is the Tibetan digit 3, and

༣

...represents 2½ (or “half three”).

Numeric presentation forms

Unicode also includes a bunch of numeric presentation forms. These are all compatibility characters, included for round-trip compatibility with some national or vendor standard, and most have compatibility decompositions. Most represent a number with some sort of styling added, but some are simply single-code-point representations of a sequence of characters. (These latter combinations pretty much all come from various East Asian encoding standards, where they’re used for things like list numbering where you want all the digits and punctuation to show up in a single display cell in vertical text.)

Superscripts and subscripts. Unicode includes superscript and subscript versions of all the digits, as well as some math operators and the letter n (as in “2ⁿ”). The superscripted versions of 1, 2, and 3 are in the Latin-1 block at U+00B9, U+00B2, and U+00B3 respectively, and the other characters are in the Superscripts and Subscripts block, which ranges from U+2070 to U+209F.

Circled. Unicode includes versions of the numbers 1 through 20 with circles around them in the range U+2460 to U+2473 in the Enclosed Alphanumerics block. Several other versions of the numbers from 1 to 10 (sans-serif versions with circles around them, and black circles with white serifed and sans-serif numbers on them) are in the Dingbats block from U+2776 to U+2793. The Han characters for 1 to 10 with circles around them are encoded in the Enclosed CJK Letters and Months block from U+3280 to U+3289.

Parenthesized. Unicode also includes single code points representing the numbers from 1 to 20 surrounded with parentheses in the Enclosed Alphanumerics block from U+2474 to U+2487.

With a period. Unicode also includes single code points representing the numbers from 1 to 20 followed by a period. These are also in the Enclosed Alphanumerics block, from U+2488 to U+249B.

Fractions. Finally, Unicode has several single code points representing fractions (the so-called “vulgar fractions”). ½, ¼, and ¾ are in the Latin-1 block at U+00BC, U+00BD, and U+00BE. Several more, for the thirds, fifths, sixths, and eighths, are in the Number Forms block from U+2153 to U+215E. The Number Forms block also includes a single code point representing just the numerator 1 and the slash (⅓, U+215F). This can be used with the subscript digits to make other fractions.

National and nominal digit shapes

The ISO 8859-6 standard didn’t have enough room to encode both the Western digits and the native Arabic digits, so it had the code points in the ASCII range that normally represent the Western digits do double duty and represent both the Western digits and the native Arabic digits. When transcoding

from this encoding to Unicode, you're not guaranteed, for example, that U+0032 represents "2"; it might actually represent "٢", depending on the source material. The original version of ISO 10646 included two characters, U+206E NATIONAL DIGIT SHAPES and U+206F NOMINAL DIGIT SHAPES, invisible formatting characters that would control whether the digits in the ASCII block get rendered with their normal glyph shapes or with the local (usually Arabic) glyph shapes. These characters came into Unicode with the merger with 10646.

These characters were deprecated in Unicode 3.0. The preferred approach to this problem is for the process that converts from 8859-6 to Unicode to be smart enough to know whether to convert 0x30 through 0x39 to U+0030 through U+0039 or to U+0660 through U+0669. The digit characters in the ASCII block should now always be treated as representing the European digit shapes.

Punctuation

The other thing you need for a writing system to be complete (generally speaking; there are exceptions) is punctuation. As with all the other character categories, Unicode includes a wide variety of punctuation marks, some specific only to one writing system, and some common to several. As we've seen, there are script-specific punctuation marks scattered throughout the Unicode encoding range. There are also two encoding blocks reserved exclusively for punctuation marks. We'll take a brief look at the whole universe of Unicode punctuation and then take a closer look at a few of the more interesting cases.

Script-specific punctuation

A lot of punctuation marks have been scattered throughout the Unicode encoding range. Some of these, mainly those in the ASCII and Latin-1 blocks, are where they are for historical or compatibility reasons; most of the others are specific to one or two writing systems and are kept in the same block as the other characters in that script. Most punctuation marks are used to delimit ranges of text in various ways. Here's a table showing most of these marks. The rows show which block they're in (which generally corresponds to which scripts they're used with). The columns, on the other hand, attempt to categorize the marks according to the amount of text they delimit. The first column contains marks that are used within words or in place of words (abbreviation or elision marks, "etc." or "and" signs, etc.). The second column contains marks that are used between words. The third column contains marks that are used between groups of words (e.g., phrases or list items). The fourth column contains marks that are used between sentences (or roughly-equivalent groupings of words). The fifth column contains marks that are used between groups of sentences (e.g., paragraphs or topics). And the last column contains marks that generally appear at the beginning or end of an entire text.

[The table is missing glyphs for some characters I couldn't find a font for. The code point values for these glyphs are shown in red and need to be filled in with the actual glyphs before publication.]

Special Characters

	<word	word	>word	sentence	>sentence	text
ASCII	' &	/ -	, ; :	. ! ?		
Latin-1	- .			ı İ		
Greek			•	;		
Armenian	՛ ՛	՛	՛	՛ օ :		
Hebrew	' "					
Arabic			‘ ’	؟ -		
Syriac	ⲥ	ⲥ	ⲥ ⲥ ⲥ ⲥ	ⲥ ⲥ ⲥ	ⲥ ⲥ	
Devanagari	◌					
Sinhala				[0DF4]		
Thai			๐		๑	๒
Tibetan		། །	། །	། ། །	། ། །	། ། ། །
Myanmar			၊	။		

Georgian					∴	
Ethiopic		፥	፡፡ ፡። ፡፣ ፡፤	፡። ፡፣	፡፤	
Canadian Aboriginal				×		
Ogham						ᚿ ᚻ
Khmer	ឱ ្ក ្ង		[17F6] ្ក	្ក	្ក	្ក
Mongolian	[1801] [1806] [1807]		[1802] [1804] [1808]	[1803] [1809]		[1800] [1805]
CJK		～ ～	、	。		

There are a lot of marks that don't fit well into this taxonomy (parentheses and quotation marks, for instance); we'll look at most of these below.

The General Punctuation block

The General Punctuation block, running from U+2000 to U+206F, is just what its name suggests: It contains a bunch of punctuation marks that are intended to be used with many different scripts. For historical reasons, most of the most common punctuation marks, such as the period, comma, question mark, exclamation points, and so on are actually in the ASCII block, so the General Punctuation block tends to contain more-specialized or less-ambiguous versions of characters that appear in the ASCII block and a miscellany of less-common marks.

In particular, it contains a large collection of spaces, dashes, quotation marks, dot leaders, and bullets, all of which we'll look at below, as well as various other marks and symbols. It also contains most of Unicode's special invisible formatting characters, which we'll also look at below.

The CJK Symbols and Punctuation block

The CJK Symbols and Punctuation block, which runs from U+3000 to U+303F, is analogous to the General Punctuation block, but contains various marks drawn from various East Asian encodings and specific to East Asian text. These include the period and comma used with CJK characters, a space that's correctly sized for use with CJK characters, various parentheses, dashes, and quotation marks, Kana repetition marks, tone marks, various miscellaneous signs and symbols, and the Ideographic Variation Indicator, which we looked at in Chapter 10.

Spaces

Now we'll look at some of the more interesting and important categories of characters in Unicode. Unicode has a number of groups of characters with similar appearances and semantics, and it's worth it to take a closer look at these and understand just what their differences are.

The basic space character in Unicode, the one representing the spaces between these words, is U+0020 SPACE, Unicode's equivalent to the ASCII space character.

Then there are a bunch of other space characters that basically differ from the regular space in width only. In rough order from widest to narrowest, the fixed-width Unicode spaces are:

- U+2003 EM SPACE, also known as the “em quad”, is a space one em wide. An em is a unit of horizontal space on a line of type that is exactly as wide as the font in use is high. So a space one em wide (the em quad) is square, at least in most fonts. The term “em” comes from the fact that in old Roman inscriptions the capital letter M was square, so the em quad is the width of the capital letter M (in most modern fonts, the capital M actually isn't an em wide, but the name has stuck). The term “quad” comes from the fact that the space is square. U+2001 EM QUAD was a mistaken duplicate encoding of the same thing, and now has a singleton canonical decomposition to this character.
- U+3000 IDEOGRAPHIC SPACE is a space that's the same width as the Han characters. Since in most fonts, the Han characters occupy a square space, this can be thought of as basically the same thing as the em quad, but intended for use in Japanese, Chinese, and Korean text. However, unlike the em quad, which has a fixed width, the ideographic space can vary in width in justified text.
- U+2002 EN SPACE, also known as the “en quad,” is a space one half of an em wide. In other words, it defines a space one half as wide as it is high. Two en quads make an em quad. The term “en” probably comes from the idea that the en is the width of the capital letter N, but I doubt even in old Roman inscriptions that the letter N was ever really half the width of the letter M. U+2000 EN QUAD is also a mistaken double encoding, and has a singleton canonical decomposition to this character.
- U+2007 FIGURE SPACE is a space that's the same width as the digits in a typeface that makes all the digits the same width. It's useful for lining up figures in columns, as the normal space is usually narrower than the digits.
- U+2004 THREE-PER-EM SPACE is a space that's one-third of an em in width, or one-third as wide as it is high. Three three-per-em spaces (or “three-em spaces” for short) make an em quad.
- U+2005 FOUR-PER-EM SPACE is a space that's one-fourth of an em in width, or one-fourth as wide as it is high. Four four-per-em spaces (or “four-em spaces” for short) make an em quad.
- U+2006 SIX-PER-EM SPACE is (all together now) a space that's one-sixth of an em in width, or six times higher than it is wide. Six six-per-em spaces (or “six-em spaces”) make an em quad.
- U+2009 THIN SPACE is more analogous to the regular U+0020 SPACE in that it doesn't have a set width. It's narrower than the regular space, but like it can vary in width when it's used in justified text.

- U+205F MEDIUM MATHEMATICAL SPACE (new in Unicode 3.2) is a thin space for mathematical typesetting, defined to be four-eighteenths of an em wide.
- U+2008 PUNCTUATION SPACE is defined to be the same width as the period. Like the figure space, it's used for lining up columns of numbers.
- U+200A HAIR SPACE is the narrowest space of all. In most fonts, it's only a point or so wide.
- U+200B ZERO WIDTH SPACE is a space that takes up no horizontal space at all. The name is kind of an oxymoron (a "space" that doesn't take up any space?). The zero-width space, despite its name, is more properly thought of as an invisible formatting character. It's used to mark word boundaries in text without actually causing a visible break. For instance, Thai is written without spaces between words (spaces in Thai are used in a manner roughly analogous to commas in English), but you're supposed to respect the word boundaries when laying Thai text out on a line (i.e., unlike in Japanese text, it's not okay to break a line of Thai text in the middle of a word). One way to deal with this problem is to put zero-width spaces in the text to mark the word boundaries. This will cause word wrapping to do the right thing, but you won't see any spaces between the words.

So where does the normal space we all know and love fit into this equation? Unlike most of the space characters listed above, U+0020 SPACE can vary in width as necessary to produce properly justified text. The base width before any adjustments are applied varies from font to font but is usually the same size as the three-per-em or four-per-em space.

In addition to the various spaces listed above, Unicode also defines a class of so-called non-breaking spaces. The non-breaking characters in Unicode (they're not all spaces) behave normally except that code that wraps lines of text isn't allowed to put the characters on either side of the non-breaking character on different lines. In other words, the non-breaking characters are used to keep other characters together on one line of text. Unicode has four non-breaking spaces:

- U+00A0 NO-BREAK SPACE (often abbreviated "NBSP") is the basic non-breaking space. It's the same width as the regular space (U+0020), and like it can vary in width in justified text, but unlike the regular space, it doesn't mark a word boundary for the purposes of word wrapping. In fact, it specifically prevents word wrapping. For example, in some languages (French, for example), the space is used as the thousands separator in large numbers instead of a period or comma, but you don't want the number split across two lines. Another example is in some languages, some punctuation is separated by a space from the words they follow. In French again, there's usually a space between the last word of a question and the question mark, unlike in English, where the question mark appears right next to the last word of the question. But again, you don't want the question mark to appear at the beginning of the next line. Sophisticated word-wrapping algorithms will do this correctly, but you can use the NBSP to get the correct effect with naive word-wrapping algorithms. Because of the NBSP's use as a thousands separator in French and some other languages, the Unicode bi-di algorithm treats it as a number separator.
- U+202F NARROW NO-BREAK SPACE is roughly analogous to the thin space in the way that the regular NBSP is analogous to the regular space. It's narrower than a regular space, may or may not be variable-width depending on the language and font, and doesn't mark a word boundary. In Mongolian, this character is used to separate grammatical particles from their stem words while keeping them semantically part of the same word. (In Mongolian, this character also causes special shaping behavior, but this isn't true with any other script.)
- U+FEFF ZERO WIDTH NO-BREAK SPACE is again not technically a space, despite its name. In Unicode 2.x, 3.0, and 3.1, it was an invisible formatting character used to control word wrapping. Starting in Unicode 3.2, however, this use has been deprecated. Now, use U+2060 WORD JOINER instead.

- U+2060 WORD JOINER is the new zero-width non-breaking space. It “glues” the two characters on either side of it together, forcing them to appear on the same line. Any character can be turned into a non-breaking character by putting U+2060 before and after it.

You could also consider U+180E MONGOLIAN VOWEL SEPARATOR to be a Unicode non-breaking space, but its use is restricted to Mongolian. In other contexts, it’s either ignored or treated the same as U+202F NARROW NO-BREAK SPACE, but in Mongolian it causes special shaping behavior when it precedes certain letters.

Dashes and hyphens

Unicode also has a lot of dashes and hyphens. The original ASCII character (0x2D) did triple duty as a hyphen, dash, and minus sign. Since ASCII was originally used mostly on computer terminals and teletype machines where everything was monospaced and these three things looked the same anyway, it wasn’t a big deal originally. But in real typesetting, these characters look and act differently and it helps tremendously if the different appearances and semantics have different code point values.

The ASCII hyphen was retained at its same position in Unicode as U+002D HYPHEN-MINUS and continues to have ambiguous semantics. Several unambiguous characters have been added to replace it (although U+002D is still the most commonly used):

- U+2010 HYPHEN is the hyphen character (-). It’s the character used in hyphenated English words and phrases, such as “day-to-day.” In most typefaces, you get the same glyph for U+002D HYPHEN-MINUS and U+2010 HYPHEN, but this is the preferred representation for that glyph.
- U+2013 EN DASH is a dash an en (i.e., half of an em) wide (–). It’s used for things like separating ranges of values (“1921–1991”).
- U+2012 FIGURE DASH has the same ambiguous semantics as U+002D HYPHEN-MINUS and is, like U+002D, in Unicode for compatibility with older standards. It’s defined to have the same width as the digits in fonts where the digits are all the same width (in many fonts, this means it’s the same as the en dash).
- U+2212 MINUS SIGN is the minus sign, used to indicate subtraction or negative numbers (“–23” or “75 – 27 = 48”). This character usually has the same glyph as the en dash, but it behaves differently. In particular, most word-wrapping algorithms will wrap words on a hyphen or dash, but not on a minus sign. The rule of thumb should be to use U+2013 when it’s okay to break the line after it (such as when it’s used to separate the numbers in a range) and U+2212 when it’s not (such as in “–23”).
- U+2014 EM DASH is a dash that’s an em wide (—). In other words, the length of the dash is the same as the height of the font. The em dash is often used in English to indicate a break in continuity or to set off a parenthetical expression (“Bill would—if *I* had anything to say about it—be in a lot of trouble right now”).
- U+2015 HORIZONTAL BAR (aka the “quotation dash”) is an even wider dash (⎯) used in some older typographical styles to set off quotations.

There are also a couple of dash characters with special behavior:

- U+00AD SOFT HYPHEN functions in a manner analogous to the zero-width space—it’s used to indicate to a word-wrapping algorithm places where it’s okay to hyphenate a word. A word-wrapping routine can treat it the same as any other dash or hyphen for the purposes of figuring out where to break a line, but it’s invisible unless it actually appears at the end of a line. At the end of a line (i.e., when the word containing it actually needs to be hyphenated), it has the same glyph as the regular hyphen (U+2010).

- U+2011 NON-BREAKING HYPHEN functions in a manner analogous to the non-breaking space and opposite to the regular and soft hyphens: you use it when you want a hyphen in a word but you still want the entire hyphenated word or expression to be on the same line. A word-wrapping routine isn't allowed to break lines on either side of the non-breaking hyphen. Again, the glyph is the same as U+2010.

There are also a few language-specific hyphen characters:

- U+058A ARMENIAN HYPHEN (֊) is just what its name suggests: it works the same way in Armenian as the regular hyphen works in English, but has an Armenian-specific appearance.
- U+1806 MONGOLIAN TODO SOFT HYPHEN (᠎ᠠ) is used when Todo is written using the Mongolian alphabet (it's a vertical line because the Mongolian script is written vertically). It works in the same way as the regular soft hyphen, except that when a line is broken at the hyphen position, the hyphen itself appears at the *beginning* of the line *after* the break, as opposed to Western hyphens, which always appear at the *end* of the line *before* the break.
- U+301C WAVE DASH and U+3030 WAVY DASH are dashes used in Japanese. When Japanese text is laid out vertically, these two characters rotate.

There are also various other characters that are basically horizontal lines and are occasionally confused with the dash characters. Notable among them is U+30FC KATAKANA-HIRAGANA PROLONGED SOUND MARK. You'll occasionally run across Japanese text that uses one of the dash characters for the Katakana length mark instead of U+30FC.

You also periodically run into a “swung” dash (“~”). This character isn't specifically encoded in Unicode, but since U+007E TILDE generally isn't used as a tilde anymore (in the old days of teletype machines, this character would be used in conjunction with the backspace character to form “ñ”, for example, something you don't see anymore), most fonts these days use a glyph that's vertically centered, rather than one that's raised above the font's x-height, for U+007E TILDE. Thus, this character can be thought of as the Unicode swung-dash character. (Theoretically, you could use the CJK wave dash, but the metrics are usually wrong for use with Western characters and it's usually not in Western fonts.)

Quotation marks, apostrophes, and similar-looking characters

Quotation marks. Even though most languages that use the Latin alphabet use the same marks for such things as the period, comma, question mark, and exclamation point, they all seem to use different marks for quotation marks. The situation with quotation marks in Unicode is rather complicated and ambiguous.

For starters, there's U+0022 QUOTATION MARK. This is the character that's familiar to all of us from typewriter keyboards (and, later on, computer keyboards). It's used for both opening and closing quotation marks, and so its glyph shape is a compromise. It usually looks like this:

''

Like the “hyphen-minus” character, this character has ambiguous semantics both because all of the semantics corresponded to a single key on a typewriter keyboard and because of the limited encoding space ASCII had. We're all kind of used to seeing this mark from older computer-generated or typewritten material, but with the advent of laser printers and computer-generated typography, we'd

generally prefer to see the quotation marks the way we write them or the way they've always appeared in printed material:

“English”

Unicode provides unambiguous code-point values for these characters as U+201C LEFT DOUBLE QUOTATION MARK and U+201D RIGHT DOUBLE QUOTATION MARK. Most modern word processors now map from the U+0022 character you get from the keyboard to either U+201C or U+201D depending on context, a feature which can work quite well (or course, it doesn't always work well, leading to stupidities such the upside-down apostrophe in “the ‘60s” that you see all too often).

French and a number of other languages use completely different-looking quotation marks. Instead of the comma-shaped quotation marks used with English, French and a number of other languages use chevron-shaped quotation marks known in French as *guillemets*:

« Français »

The guillemets were included in the Latin-1 standard and are thus in the Latin-1 block as U+00AB LEFT-POINTING DOUBLE ANGLE QUOTATION MARK and U+00BB RIGHT-POINTING DOUBLE ANGLE QUOTATION MARK.

Now here's what it gets interesting. These same marks don't mean the same thing in all languages. For example, in German, U+201C is actually the *closing* quotation mark rather than the opening quotation mark. A quotation is usually punctuated like this in German:

„Deutsch“

The German opening double-quote mark is included as U+201E DOUBLE LOW-9 QUOTATION MARK. (The “9” refers to the fact that the marks look like little 9s.) The single-quote version of this character has the same glyph form as the comma—systems should be careful not to use the comma for a single quote.

You'll also sometimes see the guillemets used with German text, but in German they point *toward* the quoted material, rather than away from it as in French:

»Deutsch«

Swedish and some other Scandinavian languages sidestep the problem of paired quotation marks altogether, using the same character for both the opening and closing quotation mark. Depending on typographic style, they can use either the guillemets or (more commonly) the curly quote marks, and you get the same effect either way:

”svenska” or »svenska»

There are numerous other linguistic and regional variants in quotation-mark usage, but you get the basic idea. The practical upshot is that while quotation-mark characters come in pairs, exactly *which* pair is used depends on language, and whether a given character (with a few exceptions) is an opening or closing quote also depends on language.

Unicode deals with this by encoding each glyph shape only once and declaring that different languages should use different code-point values for their quotation marks, rather than just having abstract “opening quotation mark” and “closing quotation mark” code-point values that have different glyph shapes depending on language. This kind of goes against the “characters, not glyphs” rule, but it eases confusion in most situations. Encoding the abstract semantics and depending on some outside process to pick the right glyphs in this case would have led to different fonts being necessary for each language that uses (for example) the Latin alphabet, just because their quotation-mark characters look different, or would have required language tagging along with the regular abstract characters in order to get the right glyph shapes. Unicode’s designers have gone for this kind of approach when an *entire alphabet* is semantically the same but different font styles are preferred in different languages, but doing it for *just the quotation-mark characters* is kind of ridiculous.

This leads to a problem: Quotation marks are inherently paired punctuation, like parentheses, but which character is which member of the pair depends on language. For this reason, the quotation marks generally can’t be put into the regular Ps and Pe (starting and ending punctuation) general categories.¹⁰³ Beginning in Unicode 2.1, two new categories, Pi and Pf (for “initial-quote punctuation” and “final-quote punctuation” respectively) were created expressly for these characters, and the quotation marks were grouped rather arbitrarily into one bucket or the other. The idea is that a process that cares about paired punctuation would know the language the text was in and be able to map Pi and Pf to Ps and Pe in an appropriate manner for that language.

The initial-quote category consists of the following marks:

```
«      U+00AB LEFT-POINTING DOUBLE ANGLE QUOTATION MARK
‘      U+2018 LEFT SINGLE QUOTATION MARK
ˆ      U+201B SINGLE HIGH-REVERSED-9 QUOTATION MARK
“      U+201C LEFT DOUBLE QUOTATION MARK
”      U+201D SINGLE RIGHT QUOTATION MARK
”      U+201E DOUBLE HIGH-REVERSED-9 QUOTATION MARK
‹      U+2039 SINGLE LEFT-POINTING ANGLE QUOTATION MARK
```

The final-quote category consists of the following marks:

```
»      U+00BB RIGHT-POINTING DOUBLE ANGLE QUOTATION MARK
’      U+2019 RIGHT SINGLE QUOTATION MARK
”      U+201D RIGHT DOUBLE QUOTATION MARK
›      U+203A SINGLE RIGHT-POINTING ANGLE QUOTATION MARK
```

These two marks always appear at the beginning of a quote and are thus in the normal “starting punctuation” category:

¹⁰³ There are a few exceptions to this rule: The low-9 marks used in German and some other languages always appear at the beginning of a quotation, for example.

Special Characters

, U+201A SINGLE LOW-9 QUOTATION MARK
,, U+201E DOUBLE LOW-9 QUOTATION MARK

East Asian quotation marks. Unicode includes three different pairs of characters used in East Asian languages:

┌ U+300C LEFT CORNER BRACKET	┐ U+300D RIGHT CORNER BRACKET
┌ U+300E LEFT WHITE CORNER BRACKET	┐ U+300F RIGHT WHITE CORNER BRACKET
“ U+301D REVERSED DOUBLE PRIME QUOTATION MARK	” U+301F LOW DOUBLE PRIME QUOTATION MARK

Within these pairs, the glyph shape *does* vary depending on context. The brackets turn sideways when used with vertical text¹⁰⁴:

ㄣ
字
內
ㄣ

[need to fix spacing in this example: characters are too far apart]

The “double-prime” quotation marks may either point towards or away from the quoted text depending on font design:

“字內” or “字內”

[Need to fix second example: quote marks need to be closer to text]

They’ll also change direction and move to different quadrants of their display cells when used with vertical text. (There’s an additional character, U+301E DOUBLE PRIME QUOTATION MARK, which is used only when Western text is quoted inside East Asian text; it’s not used with East Asian quotations.)

Single quotes and the apostrophe. The single-quote character in ASCII was even more overloaded than the double-quote character. The same code point, 0x27, was used for both the opening and

¹⁰⁴ This example and the one that follows are taken from the Unicode standard, p. 153.

closing single quotes, the apostrophe, the acute accent, the prime mark, and occasionally other things. It's normally rendered with a direction-neutral glyph like the one used for the double quote...

,

...although sometimes you'll see it rendered with a left-slanted or curly glyph on the assumption that 0x60, the grave-accent character, would be used as the opening single quote. As with the double-quote, these characters and their original numeric values have been adopted into Unicode with their original muddled semantics and extra characters with unambiguous semantics have been added. We've already talked about how quotation marks work, and the rules apply equally to single and double quotes. Most languages will alternate between double and single quote marks for quotes within quotes; different languages have different preferences for whether single or double quote marks should be used for the outermost pair or quotation marks; in British usage, for example, single quotes are usually used for the outermost pair instead of the double quotes you usually see in American usage.

There's the additional problem of the apostrophe, which has exactly the same glyph as the (English) closing single quote. In many languages, and in Latin-alphabet transcriptions of many languages that don't use the Latin alphabet, an apostrophe is used as a letter, usually to represent the glottal stop, as in "Hawai'i"; in particular the Arabic letter alef and its counterparts in languages such as Hebrew are usually written with an apostrophe when transliterated into Latin letters.

Unicode draws a distinction between the apostrophe being used as a letter, where U+02BC MODIFIER LETTER APOSTROPHE is the preferred representation, and the apostrophe being used as a punctuation mark, where U+2019 RIGHT SINGLE QUOTATION MARK is the preferred representation (in other words, Unicode doesn't distinguish between this usage of the apostrophe and the single quotation mark). Both of these characters have the same glyph shape, but may be treated differently by some text-analysis processes (such as word wrapping).

As with the double quotes, you'll usually see a single key that generates U+0027, and the system automatically maps it to U+2018 or U+2019 as appropriate. This is another argument for not having a separate character code for the apostrophe—keyboards would probably still generate the right-single-quote character code. For this reason, U+2019 will probably show up in a lot of places where the truly appropriate character code would be U+02BC. This is just one example where code that operates on Unicode text probably can't assume everything is represented the way it's supposed to be.

Modifier letters. In addition to U+02BC, mentioned above, the Spacing Modifier Letters block includes a number of other apostrophe-like characters. These include U+02BB MODIFIER LETTER TURNED COMMA, which has the same glyph shape as U+2018 LEFT SINGLE QUOTATION MARK, and U+02BD MODIFIER LETTER REVERSED COMMA, which has the same glyph shape as U+201B SINGLE HIGH REVERSED-9 QUOTATION MARK. The same rule applies here: the modifier-letter characters are to be used when the glyph is used as a letter or diacritical mark, and the quotation-mark characters are to be used when the glyph is being used as a quotation mark (or other punctuation mark).

Other apostrophe-like characters. Further adding to the confusion, Unicode includes a whole host of other characters that look kind of like apostrophes or quotation marks but have distinct semantics (usually, they also have distinct, but often subtly so, glyph shapes, but this isn't always the case).

For example, there are the half-ring characters in the Spacing Modifier Letters block, which are very similar to the apostrophe and reversed apostrophe. There are the various acute-accent and grave-accent characters in the ASCII and Spacing Modifier Letters blocks. There are the prime, double prime, reversed prime, etc. characters in both the Spacing Modifier Letters block, where they're used as diacritical marks, and the General Punctuation block, where they're usually used as math or logical symbols or as abbreviations for "feet" and "inches" or "minutes" and "seconds." There are language-specific marks, such as the Armenian half-ring and apostrophe, which look like the characters in the Spacing Modifier Letters block or the Hebrew geresh and gershayim marks, which look like the prime and double-prime marks.

These marks are generally extensively cross-referenced in the Unicode standard, and it's usually fairly apparent which of the various similar-looking marks to use if you look in the Unicode standard, but real life doesn't always match theory. You'll very often see some code point value used in a spot where some other code point value is more appropriate. There are a number of reasons for this. Sometimes the wrong character is used because only one of several similar-looking characters is available on the keyboard. Sometimes it's because fonts only support one code point value. Sometimes it's because of usage in a legacy encoding where more appropriate character codes aren't available and appropriate adjustments aren't made when converting to Unicode (sometimes this is a flaw in the conversion algorithm, and sometimes it's because Unicode draws a semantic distinction the original source encoding didn't draw).

The sad consequence of all this is that code that operates on Unicode text can't always assume that the correct code point value has been used to represent some mark; often, a code point value representing a similar-looking character has been used instead. I mention this here because there are so many characters that look like apostrophes and quotation marks and often get represented with U+0022 and U+0027 even in Unicode, but this also happens with many other groups of characters (see the section below on line-breaking conventions, for example). One can hope that this situation will improve as Unicode use becomes more widespread, but there'll probably never be anything close to 100% purity in use of Unicode code point values.

Paired punctuation

In addition to the quotation marks, there are a bunch of other punctuation marks in Unicode that come in pairs. Most of them are various types of parentheses and brackets, and most of them are used in a variety of languages and thus are encoded in the ASCII and General Punctuation blocks. Here's a complete list of the non-quotation paired punctuation marks:

(U+0028) U+0029	PARENTHESIS
[U+005B] U+005D	SQUARE BRACKET
{ U+007B	} U+007D	CURLY BRACKET
། U+0F3A	༎ U+0F3B	TIBETAN GUG RTAGS
། U+0F3C	༎ U+0F3D	TIBETAN ANG KHANG
᳚ U+169B	᳛ U+169C	OGHAM FEATHER MARK
[U+2045] U+2046	SQUARE BRACKET WITH QUILL
⟨ U+2329	⟩ U+232A	ANGLE BRACKET

⟨ U+3008 ⟩	U+3009	ANGLE BRACKET
《 U+300A 》	U+300B	DOUBLE ANGLE BRACKET
【 U+3010 】	U+3011	LENTICULAR BRACKET
〔 U+3014 〕	U+3015	TORTOISE SHELL BRACKJET
⏏ U+3016 ⏏	U+3017	WHITE LENTICULAR BRACKET
⏏ U+3018 ⏏	U+3019	WHITE TORTOISE SHELL BRACKET
⏏ U+301A ⏏	U+301B	WHITE SQUARE BRACKET
❧ U+FD3E ❧	U+FD3F	ORNATE PARENTHESIS

The two Tibetan pairs of characters function like parentheses. The Ogham feather marks generally bracket a whole text. The various brackets in the U+3000 block originally come from various CJK encodings and are intended specifically for use with CJK text. The “ornate parentheses” are intended for use specifically with Arabic—they’re in the Arabic Presentation Forms A block even though they’re neither presentation forms nor compatibility characters.

Dot leaders

The em and en dashes and the curly quotation marks and apostrophes are three leading examples of instances where the way we type something and the way it’s printed in quality typography are different from one another. Another notable exception is the ellipsis—three periods in a row generally don’t come out spaced right, so there’s something called a “typographer’s ellipsis,” a single piece of type with three periods that are spaced correctly. In Unicode this is U+2026 HORIZONTAL ELLIPSIS. (The “horizontal” distinguishes it from the vertical and diagonal ellipses, which are used sometimes to abbreviate matrices and are in the Mathematical Operators block.)

You’ll also see other numbers of periods used in a row for various other purposes, such as lines of dots in tables of contents, and Unicode also includes U+2025 TWO DOT LEADER, which has two periods, and U+2024 ONE DOT LEADER, which only has one. The one-dot leader looks exactly the same as the regular period, but when a bunch of them are used in a row, they space optimally.

These three characters all have compatibility mappings to series of regular periods of appropriate length.

Bullets and dots

Unicode also has a bunch of centered dots for various purposes. Many of these are designed as list bullets. The canonical bullet is U+2022 BULLET (•), but the General Punctuation block also includes these specially-shaped bullets...

- U+2023 TRIANGULAR BULLET (◦)
- U+2043 HYPHEN BULLET (◡)

- U+204C BLACK LEFTWARDS BULLET (◀)
- U+204D BLACK RIGHTWARDS BULLET (▶)

...and these language-specific characters which are usually used in their respective languages as list bullets:

- U+0E4F THAI CHARACTER FONGMAN (◌◌)
- U+0F09 TIBETAN MARK BSKUR YIG MGO (◌)
- U+17D9 KHMER SIGN PHNAEK MUAN (◌)

In addition, many characters from the Dingbats and Geometric Shapes blocks are often used as bullets.

There are also a number of other centered-dot characters with distinct semantics but similar appearances:

- U+000B7 MIDDLE DOT has multiple semantics, but one of the main ones is as the dot that’s used in Catalan between a pair of *ls* to indicate that they should be pronounced distinctly rather than treated as a single letter and pronounced like “ly”.
- U+2219 BULLET OPERATOR is a mathematical operator that looks like a bullet. **[What’s it used for?]**
- U+22C5 DOT OPERATOR is a mathematical operator used to indicate multiplication.
- U+2027 HYPHENATION POINT is a centered dot used in dictionaries to indicate places within a word where it may be hyphenated.

Special characters

Not every Unicode code point value represents something most people would actually call a “character”—some Unicode code point values have no glyph shape at all, but exist solely to influence the behavior of some process operating on the surrounding text. In addition to the “control characters” we’re all familiar with from ASCII, Unicode includes “characters” whose purpose is to influence line breaking, glyph selection, bidirectional text layout, and various other text processes. We’ll take a look at these characters next.

Line and paragraph separators

Arguably the most important special characters in Unicode are those that are used to mark a line or paragraph division. Unfortunately, the situation in Unicode with regard to line and paragraph divisions is rather complicated.

There are two ASCII characters that deal with line breaks: LF, or Line Feed (0x0A), and CR (0x0D), or Carriage Return. On a teletype machine, an LF would cause it to move the platen roller forward one line without moving the carriage, and a CR would cause it to return the carriage to the beginning of the line without advancing the platen roller. Starting a new line of text, as you would on a manual typewriter by hitting the carriage-return lever and pushing the carriage all the way back, or on an electric typewriter by hitting the carriage-return key, required two signals, a CR and an LF. You could use CRs by themselves to perform overprinting effects, such as underlining or adding accents

(as we discussed earlier, BS (Backspace) was also used for these purposes), but LF without CR wasn't terribly useful.

The carriage-return/line-feed combination (usually abbreviated as "CRLF") is still used as the end-of-line signal in most DOS and Windows applications, but other platforms have diverged. Since LF doesn't really do anything interesting without CR, you can interpret it as carrying an implicit CR and signaling a new line all by itself. Most UNIX and UNIX-derived systems (including the Java programming language) do this. The Macintosh platform, on the other hand, went the opposite direction: Noting that the "carriage return" button on a typewriter both returns the carriage and advances the platen, they give the CR character the same semantic—it effectively carries an implicit LF and signals the beginning of a line all by itself.

To complicate matters even further, the EBCDIC encoding actually has a specific "new line" character (NL, either 0x15 or 0x25 depending on system), and the ISO 6429 standard, which defines the control characters used with many other encoding standards (including the ISO 8859 family), adopted it into the C1 space as 0x85. There are both EBCDIC and 6429-based systems that use NL as their new-line signal, although this is rarer than using CR, LF, or CRLF.

As if all this wasn't bad enough, different applications treat the platform-specific new-line sequences in different ways, regardless of which specific character codes are involved. Originally, as the name suggests, the new-line sequences signaled the beginning of a new line of text, analogous to hitting the carriage-return key on a typewriter. Many applications, especially simple text editors, still use the new-line sequences in this way. The standard Internet mail protocols also use new-line sequences between lines of an email message, even though most email programs don't actually require you to type a new-line sequence at the end of each line (they wrap the lines automatically and supply the extra new-line sequences at transmission time).

But with the advent of word-processing software and other programs that wrap lines automatically, having a special character to signal the beginning of a new line became superfluous. Instead, these programs generally treated the new-line sequence as a new-paragraph signal.

Then, complicating things even *more*, you have word-processing programs that allow explicit line divisions within a paragraph. But since they were already using the regular new-line sequence to signal a new *paragraph*, they had to invent a *new* character to use to signal an explicit line break within a paragraph. Microsoft Word uses the ASCII VT (Vertical Tab, 0x0B) for this purpose, but this is even less standardized than the regular new-line sequences.

The problem is that because Unicode includes all the characters from the ASCII and ISO 8859-1 standard in their original relative positions, it includes all the characters that are used to signal line and paragraph divisions in these encodings, and inherits their ambiguous meanings.

Unicode includes two new characters in the General Punctuation block, U+2028 LINE SEPARATOR and U+2029 PARAGRAPH SEPARATOR (often abbreviated as LS and PS respectively), which have unambiguous semantics, as a way of trying to solve this problem. Unfortunately, most systems that use Unicode still persist in using the old characters to signal line and paragraph boundaries. The Java programming language, for example, follows the UNIX convention and uses the LF character (in Unicode, this is U+000A LINE FEED) as its new-line sequence. This is the character that gets output at the end of the line when you call `System.out.println()`, for example, and it's the character represented by `\n`. Similarly, Windows NT and the other Unicode-based versions of Windows still follow the old DOS convention of using CRLF at the ends of lines or paragraphs.

This means that software that operates on Unicode text has to be prepared to deal with *any* new-line convention. It should always interpret LS as a line separator and PS as a paragraph separator, but should also treat CR, LF, CRLF, and NL in Unicode text the same as either LS or PS depending on the application (generally, you'd treat them all the same as PS except in program editors and a few other specialized applications).

For code that doesn't know whether the legacy new-line sequences should be treated as LS or PS (library code, for example), it should make whatever the safest assumption is. The boundary-finding code in Java, for example, treats CR, LF, and CRLF the same as PS when looking for character or word boundaries (where they're really only whitespace) and as LS when looking for sentence boundaries (where a new-line sequence might be the end of a paragraph or it might just be whitespace, but where you know that there'll generally be punctuation at the end of a paragraph anyway that tips you off you're at the end of a sentence).

Character encoding conversion is also complicated by the new-line mess. When going to Unicode from some other encoding, if you know whether the platform new-line sequence is being used in the source text to represent line breaks or paragraph breaks, you can map the platform new-line sequence to LS or PS. If you don't, you can either just map to PS and hope you didn't get it wrong, or leave the characters alone and depend on the process that's actually going to do something with the text to interpret the new-line sequences correctly.

Converting *to* another encoding from Unicode is similarly complicated. You can simply leave CR, LF, and NL alone, which at least doesn't mess them up any more than they're already messed up, but you're still stuck coming up with something to do with LS and PS. Assuming you know what platform you're running on, the simplest course of action, unless you know the ultimate destination of the text, is to map both PS and LS to the platform new-line sequence, whatever it happens to be. (It may be a good idea to map any other new-line sequences in the text to the platform sequence as well.) If you don't know what platform you're running on, you're pretty much stuck. You could play the odds and convert LS and PS to LF or CRLF, but you may well get it wrong.

This issue is sufficiently complicated that the Unicode Consortium issued a Unicode Standard Annex (UAX #13) specifically to discuss these issues and nail down some recommendations. I've basically parroted back most of UAX #13 here.

Segment and page separators

Related to the issue of line and paragraph separators are the issues of page and segment separators. Fortunately, there isn't the same level of ambiguity here that there is with line and paragraph separators.

Consider the page separator, for example. The ASCII FF (Form Feed, 0x0C) character originally signaled the teletype or printer to eject the page and begin a new one. This makes it the natural choice for a character to signal the beginning of a new page, and in fact it's pretty much universally used for this purpose. Since Unicode already has all the ASCII control characters, this character is in Unicode (as U+000C FORM FEED), and it's considered the correct character to use for this purpose even in Unicode. (A page break obviously is also a line break, so while FF shouldn't be messed with in conversion, processing code that deals with explicit line breaks should generally treat FF as having all the semantics of LS, plus the additional semantic of signaling an explicit page break.)

There are also breaks that occur within a single line, usually individual pieces of text that are to be lined up in columns. In ASCII text, the TAB or HT (Horizontal Tab, 0x09) character is used for this purpose, analogous to the Tab key on a typewriter. Again, there isn't any ambiguity here, so the analogous Unicode character (U+0009 HORIZONTAL TABULATION) means the same thing in Unicode.

The ASCII standard includes a number of other control characters for use in structured text: for example the FS, GS, RS, and US (File Separator, Group Separator, Record Separator, and Unit Separator, respectively) characters were once used to delimit different-sized units of structured text. These characters are rarely used nowadays, although their analogues do exist in Unicode. Today, in very simple structured-text documents, TAB is used as a field separator and the platform new-line sequence is used as a record separator. But more sophisticated structured-text interchange is usually considered the province of a higher-level protocol, rather than the job of the character encoding. These days, the higher-level protocol of choice for this kind of thing in XML, which uses markup tags as field and record delimiters.

Control characters

Unicode adopts all of the control characters from the ISO 6429 standard. These occur in two blocks, the C0 block, which contains the original ASCII control characters, and the C1 block, an analogous encoding area in the 0x80-0xFF range containing an extra 32 control codes. Since Unicode is designed to be freely convertible to ISO 8859-1 (Latin-1, which includes all the ASCII characters and is based on ISO 6429), Unicode also has all the same control characters. The Unicode standard doesn't formally assign any semantics to these characters, but they're generally agreed to have the same interpretations as they do in their original source standards. In particular, U+0009 HORIZONTAL TABULATION and U+000C FORM FEED are the tab and page-break characters of choice in Unicode—Unicode's designers didn't see a reason to encode “official” Unicode characters with these semantics.

In addition to U+0009 and U+000C, mentioned above, you'll also frequently see U+000A LINE FEED, U+000D CARRIAGE RETURN, and U+000B VERTICAL TABULATION in Unicode text, even though U+2028 LINE SEPARATOR and U+2029 PARAGRAPH SEPARATOR are the officially-sanctioned characters to representing line and paragraph breaks. See the section on line and paragraph separators above for more information on this.

Characters that control word wrapping

It's a little misleading to call this section “Characters that control word wrapping” because almost every character in Unicode has an effect on word wrapping. For example, the line and paragraph separators always cause a line break, and the various spaces, hyphens, and dashes always produce a place in the text where a line break can happen. Other characters affect word wrapping in more complicated, language-defined ways. We'll come back and examine this subject in greater depth in Chapter 16.

But there are a couple of special characters whose purpose is solely to affect word wrapping. We looked at them earlier in the section on spaces, but it's worth look at them briefly again.

U+200B ZERO WIDTH SPACE (of “ZWSP” for short) isn't really a space at all, but an invisible formatting character that indicates to a word-wrapping process that it's okay to break a line between the two characters on either side of it. For example, let's say you have the following sentence:

I work at
SuperMegaGiantCorp.

The margins are just a little too tight for “SuperMegaGiantCorp.” to fit on the same line as the rest of the sentence, leaving a big white gap at the beginning of the first line. Maybe it’d be okay, however, to split “SuperMegaGiantCorp.” where the capital letters are. You could do this by inserting ZWSP before each of the capital letters. Then you’d get this...

I work at SuperMegaGiant
Corp.

... with a reasonable amount of text on the first line. Now if the text changed or the margins were widened, you’d get this:

I work at SuperMegaGiantCorp.

When it all fit on one line, you’d still see “SuperMegaGiantCorp.” without any spaces.

Of course, you wouldn’t see this kind of thing in English very much, but in other languages this happens frequently. In Thai, for example, spaces aren’t used between words, but word wrapping is still supposed to honor word boundaries. Most Thai word processors employ sophisticated language-sensitive word-wrapping processes to figure out where the word boundaries are, but they sometimes guess wrong. You can use the ZWSP character as a hint to the word wrapping process as to where the word boundaries are.

In English, such as in our SuperMegaGiantCorp. example, what you’d probably actually want is for “SuperMegaGiantCorp.” to be hyphenated. Related to the ZWSP is U+00AD SOFT HYPHEN (or “SHY” for short). Like the ZWSP, it normally functions as an invisible formatting character that indicates to the word-wrapping routine a place where it’s okay to put a line break. The difference is that it tells the word-wrapping routine that if it actually puts a line break at that position, it should precede it with a hyphen. In other words, it’s an invisible formatting character that indicates places where it’s okay to hyphenate a word. Like ZWSP in Thai, you might be working with a word processor that hyphenates words automatically. But if the hyphenation process hyphenates some word wrong (for example, a proper name that’s not in its dictionary, such as “SuperMegaGiantCorp.”), you can use the SHY as a hint to tell it where the right places are. In our example, if the margins are wide enough, you see the sentence the way it looks in the previous example. But if the margins are smaller, you get this:

I work at SuperMegaGiant-
Corp.

There's also U+1806 MONGOLIAN TODO SOFT HYPHEN, which works the same way with Mongolian text. (See the previous section on dashes and hyphens.)

You may also want the reverse behavior. Unicode generally says that spaces and dashes indicate legal places to break a line, but there are times when you don't want that. In French, for example, spaces are used to separate powers of 1,000 in large numbers. So let's say you have the following sentence:

You owe me 100 000 000 F.

(We'll pretend it's in French.) If the margins were tighter, you might get this instead:

You owe me 100 000
000 F.

But of course you don't want the number to be split across two lines like this. In fact, you probably want the currency symbol to appear on the same line as the number as well. U+2060 WORD JOINER (or "WJ" for short), like its cousin ZWSP, isn't really a space, but an invisible formatting character that tells word-wrapping routines that it's *not* okay to put a line break between the two characters on either side. It's often called the "glue" character because it "glues" two characters together and forces them to appear on the same line. If you were to put a WJ after each real space in the number above, the same example would wrap like this:

You owe me
100 000 000 F.

The number and the currency symbol that follows it are all kept together on a single line.

In versions of Unicode prior to Unicode 3.2, the "glue" function was represented with U+FEFF ZERO WIDTH NO-BREAK SPACE ("ZWNBS" for short), which also did double duty as the Unicode byte-order mark. This is one of those "seemed like a good idea at the time" things: The semantics would seem to be unambiguous: At the beginning of a file, U+FEFF is a BOM, and anywhere else, it's a ZWNBS. But you can then have a particular U+FEFF shift meanings when

you concatenate two files together or extract characters from the middle of a file. Unicode 3.2 deprecated the ZWNBSPP meaning of U+FEFF, making it just a byte order mark again.

Most of the time, word-wrapping routines treat spaces and hyphens as marking legal break positions, and this is the behavior you're trying to suppress by using WJ. Because of this, Unicode also includes special space and hyphen characters that also have non-breaking semantics. U+00A0 NO-BREAK SPACE (or "NBSP" for short) and U+2011 NON-BREAKING HYPHEN look and act exactly the same as their regular counterparts U+0020 SPACE and U+2010 HYPHEN, but they also suppress line breaks both before and after themselves. That is, using NBSP produces the same effect as putting a WJ on either side of a regular space—it glues together the space and the characters before and after it.

In addition to the regular no-break space, U+202F NARROW NO-BREAK SPACE and U+2077 FIGURE SPACE have non-breaking semantics, as does U+0F0C TIBETAN MARK DELIMITER TSHEG BSTAR, a non-breaking version of the Tibetan tsheg mark, which is used between words and normally indicates a legal line-break position.

Characters that control glyph selection

Unicode also includes invisible formatting characters that control glyph selection and character shaping.

The joiner and non-joiner. The two most important of these characters are U+200C ZERO WIDTH NON-JOINER ("ZWNJ" for short) and U+200D ZERO WIDTH JOINER ("ZWJ" for short).

The purpose of the non-joiner is to break a connection between two characters that would otherwise be connected. To take a trivial example:

```
U+0066 LATIN SMALL LETTER F  
U+0069 LATIN SMALL LETTER I
```

...would normally be rendered like this...

fi

...but if you insert ZWNJ...

```
U+0066 LATIN SMALL LETTER F  
U+200C ZERO WIDTH NON-JOINER  
U+0069 LATIN SMALL LETTER I
```

...you get this instead:

fi

ZWNJ breaks up both cursive connections and ligatures. For example, in Arabic this example...

U+0644 ARABIC LETTER LAM
U+0647 ARABIC LETTER HEH
U+062C ARABIC LETTER JEEM

...normally looks like this...

لهج

...but if you add in the non-joiner...

U+0644 ARABIC LETTER LAM
U+200C ZERO WIDTH NON-JOINER
U+0647 ARABIC LETTER HEH
U+200C ZERO WIDTH NON-JOINER
U+062C ARABIC LETTER JEEM

...the letters break apart, giving you this:

ل ه ج

In the Indic scripts, the ZWNJ can be used to break up conjunct consonants, forcing the explicit virama form. For example, this sequence...

U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+0937 DEVANAGARI LETTER SSA

...is normally rendered as a conjunct consonant:

क्ष

But if you stick the ZWNJ into the sequence...

U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+200C ZERO WIDTH NON-JOINER
U+0937 DEVANAGARI LETTER SSA

...the conjunct is broken, and the virama becomes visible:

क् ष [these glyphs should connect]

The zero-width *joiner* has the opposite effect: it causes cursive connections where they wouldn't otherwise exist. Its behavior, however, is a little more complicated than the ZWNJ. The ZWJ doesn't actually work on pairs of characters; it independently affects the characters on either side of it, both of which may not necessarily change shape based on context. In other words, if you put it between two characters A and B, it'll cause the rendering engine (assuming properly-designed fonts and so forth) to use the right-joining version of A, if there is one, and the left-joining version of B, if there is

one. If A has a right-joining form and B has a left-joining form, they appear cursively connected, but if only one of them has the appropriate form, it still takes it on even though there's no cursive connection (of course, if *neither* character connects cursively, ZWJ has no effect).

For example, this character...

```
U+0647 ARABIC LETTER HEH
```

...normally looks like this...

ه

...but if you precede it with the ZWJ...

```
U+200D ZERO WIDTH JOINER
U+0647 ARABIC LETTER HEH
```

...you get the final form of the character...

ﻩ

...and if you follow it with ZWJ...

```
U+0647 ARABIC LETTER HEH
U+200D ZERO WIDTH JOINER
```

...you get the initial form...

ﻪ

...and if you flank it with ZWJs...

```
U+200D ZERO WIDTH JOINER
U+0647 ARABIC LETTER HEH
U+200D ZERO WIDTH JOINER
```

...you get the medial form:

ﻪ

This is true no matter what the character on the other side of the joiner is; the characters on either side of the joiner effectively join to the joiner, not to each other. For example, there's no left-joining form of alef, so the sequence...

```
U+0627 ARABIC LETTER ALEF
U+0647 ARABIC LETTER HEH
```

...would normally be drawn with both characters' independent forms:

ا

But if you put ZWJ between them...

```
U+0627 ARABIC LETTER ALEF
U+200D ZERO WIDTH JOINER
U+0647 ARABIC LETTER HEH
```

...it'll force the final (i.e., right-joining) form of heh to be used, even though it can't actually join to the alef (which has no left-joining form):

ا

There's an exception to this normal rule about the ZWJ affecting the characters on either side independently: If the characters on either side can form a ligature, ZWJ causes the ligature to happen. If it's a mandatory ligature, the ZWJ is superfluous, but if it's an optional ligature, you can use the ZWJ, rather than out-of-band styling information, to cause it to happen.

For example, many Latin-alphabet fonts (especially italic fonts) include an optional ligature for the letters *c* and *t*. So while this sequence...

```
U+0063 LATIN SMALL LETTER C
U+0074 LATIN SMALL LETTER T
```

...would normally be rendered like this...

ct

...if you put a ZWJ between the letters...

```
U+0063 LATIN SMALL LETTER C
U+200D ZERO WIDTH JOINER
U+0074 LATIN SMALL LETTER T
```

...you might get this instead:

[ct ligature]

Sometimes there's an intermediate form between fully joined and fully broken apart that some pairs of characters can take. For example, two adjacent characters can either form a ligature, join cursively without forming a ligature, or just be drawn independently. To get the middle form (joining cursively without forming a ligature), you use both ZWJ and ZWNJ. You use ZWNJ to break up the ligature and then flank it with ZWJs to cause the cursively-joining versions of the characters to be used. Consider the Arabic letters lam and alef. Represented in their normal way...

```
U+0644 ARABIC LETTER LAM
U+0627 ARABIC LETTER ALEF
```

...they form a ligature...

ﻻ

...but if you put the ZWNJ between them...

```
U+0644 ARABIC LETTER LAM
U+200C ZERO WIDTH NON-JOINER
U+0627 ARABIC LETTER ALEF
```

...the ligature breaks apart and you get the independent forms of both letters:

ﻻ

If you surround the ZWNJ with ZWJs, however...

```
U+0644 ARABIC LETTER LAM
U+200D ZERO WIDTH JOINER
U+200C ZERO WIDTH NON-JOINER
U+200D ZERO WIDTH JOINER
U+0627 ARABIC LETTER ALEF
```

...the characters join together without forming the ligature:

ﻻ

Both joiners are required, because each independently affects one of the letters (each joiner effectively joins its letter to the non-joiner, which keeps them from forming the ligature). If you only have one joiner...

```
U+0644 ARABIC LETTER LAM
U+200D ZERO WIDTH JOINER
U+200C ZERO WIDTH NON-JOINER
U+0627 ARABIC LETTER ALEF
```

...you still get the independent form on one of the letters:

ﻻ

A similar effect can be seen in Devanagari and some other Indic scripts. Two consonants that normally form a conjunct...

```
U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+0937 DEVANAGARI LETTER SSA
```

क

...instead get rendered independently with an explicit virama if the ZWNJ is used:

U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+200C ZERO WIDTH NON-JOINER
U+0937 DEVANAGARI LETTER SSA

क् ष [these glyphs should connect]

But if you surround the ZWNJ with ZWJs, you actually get the *half-form* of ka instead of the explicit virama:

U+0915 DEVANAGARI LETTER KA
U+094D DEVANAGARI SIGN VIRAMA
U+200D ZERO WIDTH JOINER
U+200C ZERO WIDTH NON-JOINER
U+200D ZERO WIDTH JOINER
U+0937 DEVANAGARI LETTER SSA

क् ष [again, these glyphs should connect]

As with the other formatting characters we look at in this section, the joiner and non-joiner have no appearance of their own, and they don't affect any process on the text except for glyph selection. As far as other processes on the text (searching, sorting, line breaking, etc.) are concerned, the ZWJ and ZWNJ are transparent.

Variation selectors. Unicode 3.2 adds sixteen new characters that affect glyph selection, the sixteen characters in the new Variation Selectors block, which runs from U+FE00 to U+FE0F. For certain characters, you can follow the character with a variation selector. The variation selector acts as a “hint” to the rendering processing that you want the preceding character to be drawn with a particular glyph. The variation selectors never have appearances of their own, and they're transparent to all other processes that operate on Unicode text.

The variation selectors exist for things like math symbols, where 90% of the time, two similar glyphs (for example, the \geq sign with either a horizontal “equals” line or an “equals” line parallel to the lower half of the $>$) are just variants of the same character, but where you sometimes encounter writers that use the two versions to mean different things. In situations where the glyph shape matters, you might follow \geq with U+FE00 VARIATION SELECTOR-1 to indicate you want the version with the slanted line, or with U+FE01 VARIATION SELECTOR-2 to indicate you want to horizontal line. (If you didn't use either variation selector, it means you don't care which glyph is used, and the font or the rendering process is free to use whichever one it wants.)¹⁰⁵

The Unicode Character Database, starting in Unicode 3.2, includes a new file, StandardizedVariants.html, which lists specifically which combinations of a normal character and a variation selector are legal, and which glyphs they represent. You're not free to use unassigned combinations for your own uses: like unassigned code point values, all unassigned variation-selector combinations are reserved for future standardization. If a variation selector follows a character that isn't listed in StandardizedVariants.html as taking a variation selector, the variation selector simply

¹⁰⁵ This isn't actually the best example: The version of \geq with the slanted line is actually one of the new math symbols in Unicode 3.2.

has no effect. There are no “private-use” variation selectors, although you could designate some of the regular private-use characters to work that way if you wanted.

Variation selectors are generally treated by non-rendering processes as combining characters with a combining class of 0, which gives the right behavior most of the time with no special-casing code. This does mean, however, that variation selectors can only ever be used with non-combining characters (you’ll never see a variation selector used to select acute accents with different angles, for example), non-composite characters.

In addition to the sixteen generic variation selectors, there are three more variation selectors in the Mongolian block. These work exactly the same way (including having the variants listed in `StandardizedVariants.html`), but their use is restricted to Mongolian. (This a historical quirk: the Mongolian variation selectors were added in Unicode 3.0, before the generic variation selectors.)

The grapheme joiner

Unicode 3.2 adds another interesting character, U+034F COMBINING GRAPHEME JOINER. Like the variation selectors, it’s treated as a combining character and has no visible presentation. But (with one important exception) it doesn’t affect the way text is rendered.

The grapheme joiner is more or less analogous to the word joiner, except that instead of suppressing word breaks, it suppresses a grapheme cluster break, “gluing” the grapheme clusters on either side together into a single grapheme cluster. This causes the entire sequence to be treated as a unit for such things as cursor movement, searching, and sorting. In cases where a word-wrapping process has to wrap on character boundaries rather than word boundaries, it also suppresses a line break.

Enclosing marks are defined in Unicode 3.2 to enclose everything up to the most recent grapheme cluster boundary, so the grapheme joiner can be used to cause an enclosing mark to surround more text than it might otherwise. For example, the sequence

```
U+0031 DIGIT ONE
U+0032 DIGIT TWO
U+20DD COMBINING ENCLOSING CIRCLE
```

...gets drawn like this:

1②

The 1 and the 2 are separate grapheme clusters, so the circle only surrounds the 2. But if you put a grapheme joiner between them...

```
U+0031 DIGIT ONE
U+034F COMBINING GRAPHEME JOINER
U+0032 DIGIT TWO
U+20DD COMBINING ENCLOSING CIRCLE
```

...you get this:

⑫

The grapheme joiner “glues” the 1 and the 2 together into a single grapheme cluster, causing the circle to surround both of them.

For more information on grapheme clusters and the grapheme joiner, see Chapter 4.

Bidirectional formatting characters

Unicode includes seven characters whose purpose is to control bidirectional reordering:

```
U+200E LEFT-TO-RIGHT MARK
U+200F RIGHT-TO-LEFT MARK
U+202A LEFT-TO-RIGHT EMBEDDING
U+202B RIGHT-TO-LEFT EMBEDDING
U+202D LEFT-TO-RIGHT OVERRIDE
U+202E RIGHT-TO-LEFT OVERRIDE
U+202C POP DIRECTIONAL FORMATTING
```

Chapter 8 includes an in-depth treatment of all of these characters, so there’s no point in repeating the whole thing here, but here’s a brief recap:

- The left-to-right mark (LRM) and right-to-left mark (RLM) have no visual appearance, but look to the Unicode bi-di algorithm like strong left-to-right and strong right-to-left characters, respectively. These characters affect the interpretation of characters with weak or neutral directionality that appear in ambiguous positions in the text.
- The left-to-right and right-to-left override characters (LRO and RLO) cause all of the characters between themselves and the next PDF (or another override or embedding character) to be treated by the bi-di algorithm as characters with strong directionality in the direction specified, no matter their natural directionality.
- The left-to-right and right-to-left embedding characters (LRE and RLE) like LRO and RLO, are also used with the PDF to delimit a range of text, but they cause the enclosed run of text to be *embedded* in the enclosing run of text. For example, if you have a Hebrew quotation in an English sentence, but the Hebrew quotation includes an English word, you can use LRE to cause the English word to be treated as part of the Hebrew quotation rather than as part of the enclosing English sentence (which would cause the Hebrew parts of the quotation to appear out of order).
- The pop-directional-formatting character (PDF) terminates the effect of a preceding LRO, RLO, LRE, or RLE.

Again, for a much more in-depth discussion of these characters, complete with examples, see Chapter 8.

As with the other formatting characters in this section, the bi-di formatting characters have no visual appearance of their own and are transparent to all processes operating on the text except for bidirectional reordering.

Deprecated characters

The General Punctuation block includes six more invisible formatting characters that were deprecated in Unicode 3.0. These characters shouldn’t be used anymore, but since characters can’t be

removed from the Unicode standard, they're still there. The Unicode Consortium strongly discourages users from using these characters, and strongly discourages developers implementing the Unicode standard from supporting them. But for completeness, here's a quick explanation of what they were originally intended to do:

- U+206A INHIBIT SYMMETRIC SWAPPING basically turns off mirroring. Normally, characters such as the parentheses have a different glyph depending on whether they appear in left-to-right text or right-to-left text. The starting parenthesis appears as (in left-to-right text but as) in right-to-left text, maintaining its semantic identity as the starting parenthesis regardless of text direction.

U+206A turns off this behavior, causing the mirrored characters to take on their left-to-right glyphs even when they appear in right-to-left text. U+206B ACTIVATE SYMMETRIC SWAPPING turns mirroring back on, returning things to the default state.

This was included for backward compatibility with older standards that didn't have the concept of mirrored characters. When converting from such a standard to Unicode, the conversion process should be smart enough to understand the bi-di state and exchange the code point values for the mirroring characters instead of relying on being able to turn mirroring on and off.

- U+206D ACTIVATE ARABIC FORM SHAPING causes the same glyph-selection process that's used for the Arabic letters in the U+0600 block to be used with the characters in the Arabic Presentation Forms blocks. In other words, all of the different presentation forms of the letter *heh* in the presentation-forms blocks get treated the same as the nominal version of *heh* in the Arabic block—their glyph shapes are determined by context, rather than always being the glyph shapes specified in the Unicode standard. U+206C INHIBIT ARABIC FORM SHAPING turns this behavior off, returning to the default state where the Arabic presentation forms always have their specified glyph shapes.

These codes came into Unicode in the 10646 merger. They shouldn't be used; if you want contextual glyph selection, use the nominal forms of the Arabic letters instead of the presentation forms.

- U+206E NATIONAL DIGIT SHAPES causes the Unicode code point values from U+0030 to U+0039 to take on the glyph shapes for the digits in the script of the surrounding text. For example, in Arabic text, U+206E causes U+0031, U+0032, U+0033 to appear as ١, ٢, ٣ instead of as 1, 2, 3. U+206F NOMINAL DIGIT SHAPES returns things to their default state, where the digits in the ASCII block are always drawn as European digits.

This was included for compatibility with ISO 8859-6, which had the ASCII digit codes do double duty, representing both the European digits and the Arabic digits depending on context. Code converting to Unicode from 8859-6 should be smart enough to map the ASCII digit characters to the appropriate code points in the Arabic block when they're supposed to show up as the Arabic digits.

Interlinear annotation

In Chinese and Japanese text, you'll periodically see characters adorned with little annotations between the lines of text. Usually these annotations are there to clarify the pronunciation of unfamiliar characters. In Japanese, this practice is called *furigana* or *ruby*.

Text containing rubies is normally considered to be structured text, in the same way that text containing sidebars or footnotes would be structured text. Usually Unicode would be used to represent the characters in the annotations themselves, but some higher-level protocol, such as HTML or XML, would be used to associate the annotation with the annotated text and specify where in the document the text of the annotation would go.

But occasionally it's necessary to include the annotations in otherwise plain, unstructured text, and so Unicode includes characters for interlinear annotation.

- U+FFF9 INTERLINEAR ANNOTATION ANCHOR marks the beginning of a piece of text that has an annotation.
- U+FFFA INTERLINEAR ANNOTATION SEPARATOR marks the end of a piece of text that has an annotation and the beginning of the annotation itself. A piece of text might have multiple annotations, in which case the annotations are separated from one another with additional U+FFFA characters.
- U+FFFB INTERLINEAR ANNOTATION TERMINATOR marks the end of the annotation and the resumption of the main text.

A full treatment of interlinear annotations, complete with examples, is included in Chapter 10.

The object-replacement character


Fully-styled text may often include non-text elements, such as pictures or other embedded elements. How such elements are included in a document is the province of a higher-level protocol, but Unicode includes one thing to make it simpler.

In a markup language such as HTML, embedded elements are simply included via appropriate pieces of markup text (such as the “” tag) in appropriate places in the text. But in many word processors and similar programs, the non-text information is kept in a separate data structure from the text itself and associated with it by means of a *run-array mechanism*, a data structure that associates ranges of text with records of styling information (for example, it includes records that say, in effect, “characters 5 through 10 are in boldface”). Embedded elements are usually stored in the style records in the run array.

But unlike normal styling information, embedded elements aren't really associated with a character; they associated with a position *between* two characters, something a normal run-array mechanism doesn't handle very well. Unicode provides U+FFFC OBJECT REPLACEMENT CHARACTER to help out. In most styled-text situations, this character is invisible and transparent to operations operating on the text. But it provides a spot in the text to associate embedded elements with. In this way, an embedded element like a picture can be stored in a style record and the run array has a character in the text to associate it with.

In purely plain Unicode text, U+FFFC can be used to mark the spot where a picture or other embedded element had once been or is to be placed.

The general substitution character

U+FFFD REPLACEMENT CHARACTER is used to represent characters that aren't representable in Unicode. The idea is that if you're converting text from some legacy encoding to Unicode and the text includes characters that there's no Unicode code point value for, you convert those characters to U+FFFD. U+FFFD provides a way to indicate in the text that data has been lost in the translation. There's no standardized visual representation for this character, and in fact it doesn't have to have a visual representation at all, but  is the glyph shown in the Unicode standard and is frequently used.

The ASCII SUB character (U+001A SUBSTITUTE in Unicode) is often used in ASCII-based encodings for the same purpose. Even though this character also exists in Unicode, it shouldn't be used for this purpose in Unicode text.

Since the designers of Unicode used most of the other encoding standards out there as sources for characters, there won't be many situations where it's necessary to convert something to U+FFFD. Up until Unicode 3.0, you'd mainly see this character used for CJK ideographs that weren't yet in Unicode. Unicode 3.0 added more-specialized character codes for representing otherwise-unrepresentable CJK ideographs. These include the geta mark (U+3013), which functions like U+FFFD and merely marks the position of an unknown ideograph (but which has a standard visual representation), the ideographic variation indicator (U+303E), which is used in conjunction with a character that *is* representable to say "a character kind of like this one," and the ideographic description characters, which provide a whole grammar for describing an otherwise-unencodable character. For a full treatment of these characters, see Chapter 10.

Tagging characters

One longtime point of controversy in Unicode has been the doubling up of certain blocks of characters to represent multiple languages where the characters have historical relationships but distinct glyph shapes. Such pairs include Greek and Coptic; modern and old Cyrillic; Traditional Chinese, Simplified Chinese, Japanese, and Korean; and (possibly) Arabic and Urdu. To get the right appearance for the characters in these ranges, you need some additional out-of-band information either specifying the font to be used or the language the text is in.

Various processes operating on the text might also work differently depending on language. Spell-checking and hyphenation are obvious examples, but there are many others.

The position of the Unicode Consortium on things like specifying the language of a range of text is that it's the proper province of a higher-level protocol, such as HTML or MIME, and such information doesn't belong in plain text. But there are rare situations where language tagging might be useful even in plain text.

Unicode 3.1 introduces a method of putting language tags into Unicode text. The tagging mechanism is extensible to support other types of tagging, but only language tagging is supported right now.

Many of the members of the Unicode Technical Committee approved this while holding their noses, and the Unicode tagging characters come with all sorts of warnings and caveats about their use. The bottom line is that if you can use some other mechanism to add language tags to your text, do it. If you can get by without language tags at all, do it. Only if you absolutely have to have language tags in Unicode plain text should you use the Unicode tag characters, and even then you should be careful.

The tagging characters are in Plane 14, the Supplementary Special-Purpose Plane. There are 95 main tagging characters, corresponding to the 95 printing characters of ASCII and occupying the same relative positions in Plane 14 that the ASCII characters occupy in Plane 0, from U+E0020 to U+E007E. For example, if U+0041 is LATIN CAPITAL LETTER A, U+E0041 is TAG LATIN CAPITAL LETTER A.

The basic idea is that a tag consists of a string of ASCII text in a syntax defined by some external protocol and spelled using the tag characters rather than the ASCII characters. A tag is preceded by an identifier character that indicates what kind of tag it is. Right now, there's one tag identifier character: U+E0001 LANGUAGE TAG. There might eventually be others.

The nice thing about this arrangement is that processes that don't care about the tags can ignore them with a very simple range check, and that you don't need a special character code to terminate a tag: the first non-tag character in the text stream terminates the tag.

For language tags, the syntax is an RFC 3066 language identifier: a two-letter language code as defined by ISO 639, optionally followed by a hyphen and a two-letter region or country code as defined by ISO 3166. (You can also have user-defined language tags beginning with "x-".)

For example, you might precede a section of Japanese text with the following:

```
U+E0001 LANGUAGE TAG
U+E006A TAG LATIN SMALL LETTER J
U+E0061 TAG LATIN SMALL LETTER A
U+E002D TAG HYPHEN-MINUS
U+E006A TAG LATIN SMALL LETTER J
U+E0070 TAG LATIN SMALL LETTER P
```

This, of course, is "ja-jp" spelled out with tag characters and preceded with the language-tag identifier. All of the text following this tag is tagged as being in Japanese (of course, there's nothing to guarantee the tag is right: you could easily have Spanish text tagged as French, for example).

A language tag stays in scope until you reach the end of the text ("the end of the text" being an application-defined concept), see another language tag, or reach U+E007F CANCEL TAG. U+E007F in isolation cancels all tags in scope at that point. U+E0001 U+E007F cancels just the language tag in effect (right now, there's no difference because language tags are the only kinds of tags there are, but this may change in the future). Language tags don't nest.

If other tag types are defined in the future, tags of different types will be totally independent. Except for U+E007F canceling all tagged text in scope (a feature allowed basically to permit concatenation of files without tags bleeding from one to another), the scope of one kind of tag doesn't affect the scope of any other kind of tag that may also be in effect.

It's completely permissible for a Unicode implementation to ignore any tags in the text it's working on, but implementations should try not to mess up the tags themselves by, for example, inserting regular characters into the middle of a tag. Tag-aware applications should, however, be prepared to deal with just such a thing in some reasonable way. (Unicode 3.1 prescribes that malformed language tags should be treated the same as CANCEL TAG by conforming applications—text after a malformed tag is simply considered untagged.)

Since language tags are stateful, they pose a problem for interpretation. If you care about the language tag in effect at some point in the text, you might have to scan backward from that point an indefinite distance to find the appropriate language tag. A good policy is to read the file sequentially when it's loaded into your program, strip out all the language tags, and put the information into an out-of-band data structure such as a style run array.

The Unicode standard specifically doesn't say that conforming applications that care about the tags have to do anything special with tagged text other than interpreting the text. In particular, conforming implementations don't have to do anything different with text that's tagged as being in some language than they'd do with the same text if it were untagged. They're *permitted* to treat tagged text differently, of course, but this isn't *required*.

Non-characters

Unicode also includes a number of code point values that are not considered to be legal characters. These include:

- **U+FFFF.** Having this not be a character does two things: It allows you to use an `unsigned short` in C++ (or corresponding type in other programming languages) for Unicode characters and still be able to loop over the entire range with a loop that'll terminate (if U+FFFF was a legal character, a naive loop would loop forever unless a larger data type were used for the loop variable or some extra shenanigans are included in the code to keep track of how many times you've seen a particular value, such as 0). It also provides a non-character value that can be used for end-of-file or other sentinel values by processes operating on Unicode text without their having to resort to using a larger-than-otherwise-necessary data type for the characters (in spite of this, the relevant Java I/O routines return `int`—go figure).
- **U+FFFE.** The byte-swapped version of this is U+FEFF, the Unicode byte-order mark. Having U+FFFE not be a legal character make it possible to identify UTF-16 or UTF-32 text with the wrong byte ordering.
- **U+xxFFFE and U+xxFFFF.** The counterparts to U+FFFE and U+FFFF in all of the supplementary planes are also set aside as noncharacters. This is a holdover from the early design ISO 10646, when it was expected that the individual planes would be used independently of one another (i.e, the UCS-2 representation could be used for planes other than the BMP when, for example, the text consisted entirely of Han characters).
- **U+FDD0 to U+FDEF.** These additional code points in the BMP were set aside starting in Unicode 3.1. The idea here was to give implementations a few more non-character values that they could use for internal purposes (without worrying they'd later be assigned to actual characters). For example, a regular-expression parser might use code point values in this range to represent special characters like the `*` and `|`.

The noncharacter code points are never supposed to occur in Unicode text for interchange, and it's not okay to treat them as characters—they're not just extensions of the Private Use Area. They're expressly reserved for internal application use in situations where you need values you know don't represent characters (such as sentinel values).

Symbols used with numbers

Of course, Unicode also has a vast collection of signs and symbols. We'll examine these in two separate sections. This section deals with those symbols that are used with numbers or in mathematical expressions.

Numeric punctuation

Many of the cultures that use Arabic positional notation make use of punctuation marks in numerals. For example, Americans use the period to separate the integral part of the number from the fractional part of the number, and in large numbers, we use a comma every three digits to help make the

numeral easier to read. Many other cultures also use punctuation marks for these purposes, but they use different punctuation marks. In much of Europe, for example, a comma, rather than a period, is used as the decimal-point character, and various characters from the period (Germany) to the space (France) to the apostrophe (Switzerland) are used as thousands separators.

The one thing that almost all of these cultures have in common is that they use regular punctuation marks for these purposes. Unicode follows this and doesn't separately encode special numeric-punctuation characters. There's no "decimal point" character in Unicode, for example; the regular period or comma code point values are used as decimal points.

The one exception is Arabic, which has special numeric punctuation. When using native Arabic digits, you also use U+066B ARABIC DECIMAL SEPARATOR and U+066C ARABIC THOUSANDS SEPARATOR. 1,234.56, thus, looks like this when written with Arabic digits:

١٬٢٣٤٫٥٦

Currency symbols

Unicode includes a wide selection of currency symbols. Most of these are in the Currency Symbols block, which runs from U+20A0 to U+20CF. Notably, the Euro sign (€) is in this block, at U+20AC (U+20A0 EURO-CURRENCY SIGN is a historical character that's different from the modern Euro sign; be careful not to confuse them!).

But there are also a number of currency symbols, including many of the most common, that aren't in the Currency Symbols block for one reason or another. They include the following:

Character	Code point value	Units
\$	U+0024	Dollars (many countries), plus some other units
¢	U+00A2	Cents (U.S.)
£	U+00A3	British pounds
₹	U+00A4	<i>(see note below)</i>
¥	U+00A5	Japanese yen
f	U+0192	Dutch guilder
৳	U+09F2	Bengali rupee decimal separator
৳	U+09F3	Bengali rupees
฿	U+0E3F	Thai baht
៛	U+17DB	Khmer riels

U+00A4 CURRENCY SIGN is a holdover from the international version of the ISO 646 standard. Rather than enshrine the dollar sign (or any other country's currency symbol) at 0x24 (where ASCII put the dollar sign), ISO 646 declared this to be a national-use character and put this character, which isn't actually used anywhere as a currency symbol, in that position as the default. The later International Reference Version of ISO 646 restored the dollar sign as the default character for 0x24, and ISO 8859-1 set aside another code point value, 0xA4, for the generic currency symbol. It comes into Unicode by way of ISO 8859-1. It still isn't used as an official currency symbol anywhere, but

is occasionally used as a “generic” currency symbol (the Java currency formatting APIs, for example, use it to mark the position of a currency symbol in a formatted number).

Unit markers

A lot of characters are used with numerals to specify a measurement unit. A lot of the time, the unit is simply an abbreviation in letters (such as “mm”), but there are a lot of special unit symbols as well. Among them are the following:

In the ASCII block: # [pounds, number] (U+0023), % [percent] (U+0025)

In the Latin-1 block: ° [degrees] (U+00B0)

In the Arabic block: ٪ [Arabic percent sign; used with native Arabic digits] (U+066A)

In the General Punctuation block: ‰ [per thousand] (U+2030), ‰‰ [per ten thousand] (U+2031), ′ [minutes, feet] (U+2032), ″ [seconds, inches] (U+2033)

The **Letterlike Symbols block** (U+2100 to U+214F) also includes many characters that are symbols for various units and mathematical constants. Most of the characters in this block are regular letters with various stylistic variants applied to them. In a few cases, they’re clones of regular letters, but they’re categorized as symbols (in the case of the few clones of Hebrew letters encoded in this block, the symbols are also strong left-to-right characters).

The **CJK Compatibility block** (U+3300 to U+33FF) consists almost entirely of various unit symbols used in East Asian typography. These all consist of several Latin letters or digits, Katakana or Kanji characters, or combinations thereof arranged to fit in a standard East Asian character cell. For example, both of the of the following characters are used in vertical Japanese typography as a “kilometer” symbol (the one on the left is the word “kilometer” written in Katakana):

キロメ
ートル km

Math symbols

There’s also a whole mess of mathematical symbols in Unicode. Again, the greatest concentration of them is in the Mathematical Operators block, which runs from U+2200 to U+22FF. Many math symbols are also unified into the Arrows and Geometric Shapes blocks. There’s also a fair amount of unification that has taken place within the Mathematical Operators block itself: when a particular glyph shape has multiple meanings, they’re all unified as a single Unicode character, and when a symbol with a single semantic has several variant (but similar) glyph shapes (such as \geq with a horizontal line versus the same thing with the line parallel to the lower half of the $>$), they’re unified into a single code point value.

Further increasing the number of possible math symbols, there’s the Combining Diacritical Marks for Symbols block, running from U+20D0 to U+20FF, which, as its name suggests, contains a bunch of combining marks that are used specifically with various symbols to change their meaning (some of

the characters in the regular Combining Diacritical Marks block, such as U+0338 COMBINING LONG SOLIDUS OVERLAY, are also frequently used with symbols).

There are also a few math symbols in the Letterlike Symbols block.

In addition to the math symbols in these blocks, there are a number of math symbols in other Unicode blocks. Among them:

In the ASCII block: + (U+002B), < (U+003C), = (U+003D), > (U+003E), | [absolute value] (U+007C) [the hyphen/minus, slash, and tilde characters are often used as math operators, but these have counterparts in the Mathematical Operators block that are specifically intended for mathematical use]

In the Latin-1 block: ¬ [not sign] (U+00AC), ± (U+00B1), ÷ (U+00F7) [the centered dot character at U+00D7 is also often used as a math operator, but also has a more appropriate counterpart in the Mathematical Operators block]

In the General Punctuation block: ‖ (U+2016) [norm of a matrix], plus the various prime marks from U+2032 to U+2037.

In the Miscellaneous Technical block: Γ ∟ ⊥ ⊓ (U+2308 to U+230B)

On top of all this, Unicode 3.2 adds a whopping 591 new math symbols, one product of something called the STIX project (Scientific and Technical Information Exchange), an effort by a consortium of experts from the American Mathematical Society, the American Physical Society, and a number of mathematical textbook publishers.¹⁰⁶ This resulted in five new blocks being created in the Unicode encoding space:

The **Miscellaneous Mathematical Symbols–A** (U+27C0 to U+27EF) and **Miscellaneous Mathematical Symbols–B** (U+2980 to U+29FF) blocks contain a wide variety of operators, symbols, and delimiters. Of particular note in these blocks are a number of brackets that are disunified from similar-looking counterparts in the CJK Miscellaneous Area. The CJK brackets are supposed to be used as punctuation in East Asian text and have the “wide” East-Asian-width property, which may make them space funny in mathematical text, so new brackets specifically for mathematical use have been added.

The **Supplemental Arrows–A** (U+27F0 to U+27FF) and **Supplemental Arrows–B** (U+2900 to U+297F) blocks contain a variety of new arrows and arrowlike symbols used as math operators. Of note here are a few arrows that are longer versions of arrows in the regular Arrows block. The longer arrows actually have different semantics from their shorter counterparts and can appear in the same expressions.

The **Supplemental Mathematical Operators block** (U+2A00 to U+2AFF) contains more mathematical operators.

¹⁰⁶ As always with comments about Unicode 3.2, which is still in beta as I write this, the details may have shifted a bit by the time you read this.

Together, these five blocks, along with the math symbols already in Unicode 3.1, cover pretty much everything that's in reasonably common use as a math symbol in a wide variety of publications. Of course, math symbols are like Chinese characters: new ones are always being coined. Of course, if one researcher invents a symbol just for his own use, that's generally not grounds to encode it (it's a "nonce form" and can be handled with Private Use Area code points), but if it catches on in the mathematical community, it becomes a candidate for future encoding.

Unicode 3.2 also adds a bunch of "**Symbol pieces**" to the Miscellaneous Technical block (adding to a couple that were already there). These are things like the upper and lower halves of brackets, braces, parentheses, integral signs, summation and product signs, and so forth, as well as extenders that can be used to lengthen them. You see these used to put together extra-large (or at least extra-tall) versions of the characters they make up, such as brackets to enclose arrays and vectors that are several lines long. Sometimes they also get pressed into service as independent symbols in their own right. As a general rule, the symbol pieces shouldn't appear in interchange text, but they're useful as internal implementation details for drawing extra-large math symbols.

Finally, the **General Punctuation block** includes several special characters for mathematical use:

U+2044 FRACTION SLASH can be used with regular digit characters to form typographically-correct fractions: This sequence, for example...

```
U+0035 DIGIT FIVE
U+2044 FRACTION SLASH
U+0038 DIGIT EIGHT
```

...should be rendered like this:

$\frac{5}{8}$

The basic behavior is that any unbroken string of digits before the fraction slash gets treated as the numerator and any unbroken string of digits after it gets treated as the denominator. If you want to have an integer before the fraction, you can use the zero-width space:

```
U+0031 DIGIT ONE
U+200B ZERO WIDTH SPACE
U+0031 DIGIT ONE
U+2044 FRACTION SLASH
U+0034 DIGIT FOUR
```

...should be rendered as...

$1\frac{1}{4}$

Unicode 3.2 adds a couple new "invisible operators" to the General Punctuation block. These include U+2061 INVISIBLE TIMES, U+2062 FUNCTION APPLICATION, and U+2063 INVISIBLE SEPARATOR. None of these characters has a visible presentation; they exist to make life easier for code that's parsing a mathematical expression represented in Unicode. Consider this expression:

A_{ij}

This could either be the j th element in the i th row of A , a two-dimensional matrix, or the i -times- j th element of A , a one-dimensional vector. You could put U+2063 INVISIBLE SEPARATOR between i and j to let the parser know you mean the first meaning, and U+2061 INVISIBLE TIMES between them to indicate you mean the second meaning. Similarly, if you have...

$f(x + y)$

...this could either be the application of function f to the sum of x and y , or the variable f times the sum of x and y . Again, you could use U+2062 FUNCTION APPLICATION to indicate the first and U+2061 INVISIBLE TIMES to indicate the second.

In all of these cases, the invisible operator doesn't specifically show up in the rendered expression, although it might cause a sufficiently-sophisticated math renderer to space things a little differently.

All of these issues and more are covered in Proposed Draft Unicode Technical Report #25, "Unicode Support for Mathematics," which will probably be a real Technical Report by the time you read this. It goes into great detail on the various math symbols in Unicode, proper design of fonts for mathematical use, and tips for the design of Unicode-based mathematical-typesetting software.

Among the more interesting parts of PDUTR #25 are a list of proposed new character properties to help mathematical typesetters parse Unicode text (these might get adopted into future versions of the standard) and an interesting and detailed proposal for using plain Unicode text for expressing structured mathematical expressions, a method that is potentially both more concise and more readable than either MathML or TeX.

Mathematical alphanumeric symbols

Finally, Unicode 3.1 added the **Mathematical Alphanumeric Symbols block** (U+1D400 to U+1D7FF in the SMP), which, like the Letterlike Symbols block, contains clones of various letters and digits with various stylistic variations added. Unlike the Letterlike Symbols block, this block systematically encodes every combination of a prescribed set of characters and a prescribed set of stylistic variations (except those that are already encoded in the Letterlike Symbols block, for which corresponding holes have been left in the Mathematical Alphanumeric Symbols block). Also unlike the characters in the Letterlike Symbols block, the characters in the Mathematical Alphanumeric Symbols block don't generally have well-established semantics.

The idea behind this block is that all of these characters can appear as symbols in mathematical equations and that the different styles of the same letter are actually *different* symbols meaning different things and can appear in the same equation. Using normal styled-text markup for this kind of thing can be both clumsy and heavy-weight, and the letters in the equation might get treated as real letters by some text-processing operations. The characters in this block are unambiguously to be treated as *symbols*, and because they all have different code point values, a search for, say, a bold **h** won't accidentally find an italic *h*. You also don't want a case-mapping routine to accidentally convert the symbols from one case to another (which would effectively turn them into completely different symbols). If you transmit the equation as plain text, you wouldn't want an equation

containing both h and **h**, obviously meaning different things, to have both turn into a regular h , which might mean yet a third thing in addition to blurring the original distinction between **h** and h .

A classic example of this would be the Hamiltonian formula¹⁰⁷:

$$H = \int d\tau (\epsilon E^2 + \mu H^2)$$

If you convert this to plain text and lose the distinction between the H s, you get a simple integral equation over the variable H :

$$H = \int d\tau (\epsilon E^2 + \mu H^2)$$

It's important to note that it's *not* okay to use the characters in the Mathematical Alphanumeric Symbols block to fake styled text in a plain-text environment. That's not what these characters are designed for, and that's part of the reason why they're classified as symbols rather than letters.

Other symbols and miscellaneous characters

Finally (puff puff puff), Unicode encodes a whole host of characters that don't readily fit into any of the other categories. These include both various non-mathematical symbols and various other types of characters.

Musical notation

Unicode 3.1 adds the new **Byzantine Musical Symbols block** (U+1D000 to U+1D0FF) and **Musical Symbols block** (U+1D100 to U+1D1FF), which include various musical symbols. Full-blown Western musical notation is often considered to be the most comprehensive and complex written language even devised, and Unicode doesn't attempt to provide the means to fully specify a musical score. In particular, the specification of pitch is left to a higher-level protocol, as are most of the more complex layout issues that you get into when laying out music.

What Unicode does provide is a complete set of symbols that can be used by a higher-level protocol for representing musical notation. In plain text, these symbols can be used just as simple symbols and are useful for writing *about* music.

Some characters in the Musical Symbols block require special treatment even in plain text. In order to allow flexibility in things like the selection of noteheads, notes are actually built out of pieces that get treated as combining character sequences. Precomposed characters are also provided for the most common cases. For example, a simple eighth note...

¹⁰⁷ This example is lifted straight from the text of UAX #27.



...can be represented as a single character...

```
U+1D160 MUSICAL SYMBOL EIGHTH NOTE
```

...but this actually has a canonical decomposition down to...

```
U+1D158 MUSICAL SYMBOL NOTEHEAD BLACK  
U+1D165 MUSICAL SYMBOL COMBINING STEM  
U+1D16E MUSICAL SYMBOL COMBINING FLAG-1
```

There's also a set of combining marks for the common articulation marks and for augmentation dots. So we can make our eighth note into a dotted eighth note and add an accent to it...



...by adding the appropriate characters to the combining character sequence:

```
U+1D160 MUSICAL SYMBOL EIGHTH NOTE  
U+1D16D MUSICAL SYMBOL COMBINING AUGMENTATION DOT  
U+1D17B MUSICAL SYMBOL COMBINING ACCENT
```

The Musical Symbols block also includes several invisible formatting characters that, in combination with the musical symbols, signal the beginning and endings of beamed groups of notes, slurs, ties, and phrase marks. For example, the following sequence...

```
U+1D173 MUSICAL SYMBOL BEGIN BEAM  
U+1D160 MUSICAL SYMBOL EIGHTH NOTE  
U+1D161 MUSICAL SYMBOL SIXTEENTH NOTE  
U+1D175 MUSICAL SYMBOL BEGIN TIE  
U+1D161 MUSICAL SYMBOL SIXTEENTH NOTE  
U+1D174 MUSICAL SYMBOL END BEAM  
U+1D173 MUSICAL SYMBOL BEGIN BEAM  
U+1D161 MUSICAL SYMBOL SIXTEENTH NOTE  
U+1D176 MUSICAL SYMBOL END TIE  
U+1D161 MUSICAL SYMBOL SIXTEENTH NOTE  
U+1D160 MUSICAL SYMBOL EIGHTH NOTE  
U+1D174 MUSICAL SYMBOL END BEAM
```

...would be rendered like this by a sufficiently-sophisticated rendering engine:



The Byzantine Musical Symbols block encodes the symbols used by the Eastern Orthodox Church (especially the Greek Orthodox Church) for writing hymns and other liturgical music. In plain text, all of these symbols can just be treated as regular spacing symbols, but music-notation software can use them as primitives for building up real Byzantine musical notation.

Braille

Unicode 3.0 added the **Braille Patterns block** (U+2800 to U+28FF), which encoded the 256 possible eight-dot Braille patterns. Generally speaking, Braille can be thought of as a font variant of a regular script in Unicode, and text in Braille can be represented in Unicode using the regular code point values for the characters being represented in Braille. But there are many different standards for representing writing in Braille and they're all different. The purpose of this block is to allow an abstract representation of the dot patterns themselves, without worrying about what actual characters they represent. In certain situations, the ability to represent the abstract dot patterns saves systems from having to do extra round-trip conversions back to standards that assign actual meaning to the dot patterns. Most of the time, this is an internal representation; users would usually deal with text represented using the normal Unicode characters and an internal process would convert the normal Unicode characters to the characters in this block for output on a Braille printing device (or a Braille input device would generate the abstract characters in this block, which would then be converted by a separate process into the appropriate regular Unicode characters for the language the user speaks).

Other symbols

There are various other miscellaneous-symbols blocks in Unicode, including:

- The **Letterlike Symbols block** (U+2100 to U+214F) includes a collection of symbols that either are clones of regular letters with some special stylistic variant applied or are special combinations of multiple letters. Some of these are mathematical symbols (the Riemann integral \int), the Weierstrass elliptic function \wp), some represent constants (the Planck constant h), the first transfinite cardinal \aleph), some represent units of measure (degrees Fahrenheit $^{\circ}\text{F}$), ohms Ω), Angstrom units \AA), and some are just abbreviations (care-of [%], trademark $^{\text{TM}}$, versicle $\text{\textcircled{V}}$).
- The **Arrows block** (U+2190 to U+21FF) includes various arrows. Most of these are mathematical or technical symbols of one kind or another, but all arrowlike symbols have been collected into this one block. Multiple semantics with the same glyph shape have been unified.
- The **Miscellaneous Technical block** (U+2300 to U+23FF) includes a variety of different symbols that didn't fit into any of the other blocks. These include various keycap symbols used in computer manuals, and a complete set of functional symbols from the programming language APL, as well as various other symbols from various technical disciplines.
- The **Control Pictures block** (U+2400 to U+243F), which includes glyphs that provide visible representations for the various ASCII control characters (plus two different visible representations for the space).

- The **Optical Character Recognition block** (U+2440 to U+245F) includes various symbols from the OCR-A character set that don't correspond to regular ASCII characters, as well as the MICR characters (magnetic-ink character recognition) that are used to delimit the various numbers on the bottom of a check.
- The **Miscellaneous Symbols block** (U+2600 to U+26FF) includes a bunch of nontechnical symbols that didn't fit anywhere else. It includes such things as weather symbols, astrological symbols, chess pieces, playing-card suits, die faces, I-Ching trigrams, recycling symbols, and so forth.

Presentation forms

Unicode also includes a bunch of blocks consisting of nothing but special presentation forms of characters whose normal encoding is elsewhere in the standard. These characters generally exist solely for round-trip compatibility with some legacy encoding standard. Among the blocks are:

- The **Superscripts and Subscripts block** (U+2070 to U+209F), which consists of various digits and mathematical symbols in superscripted or subscripted form (a couple of superscripts and subscripts are also encoded in the Latin-1 block)
- The **Number Forms block** (U+2150 to U+218F), which consists of various precomposed fractions and Roman numerals.
- The **Enclosed Alphanumerics block** (U+2460 to U+24FF), which consists of various letters and numbers with circles around them, with parentheses around them, or with periods after them.
- The **Hangul Compatibility Jamo block** (U+3130 to U+318F) contains non-conjoining clones of the characters in the regular Hangul Jamo block.
- The **Enclosed CJK Letters and Months block** (U+3200 to U+327F) contains various CJK characters with circles or parentheses around them, plus precomposed symbols for the twelve months in Japanese.
- The **CJK Compatibility block** (U+3300 to U+33FF) contains a wide variety of abbreviations arranged so as to fit in single CJK character cells.
- The **CJK Compatibility ideographs block** (U+F900 to U+FA5F) contains clones of characters in the CJK Unified Ideographs block. Most of these characters are from the Korean KS X 1001 standard and are identical to their counterparts in the CJK Unified Ideographs block. KS X 1001 had duplicate encodings for a bunch of Han characters that had more than one Korean pronunciation—each pronunciation got its own code point. A number of other characters that got double encodings in other East Asian standards for other reasons are also in this block. The characters in this block are all included here solely for round-trip compatibility with their original source standards—they would have been unified with their counterparts in the Unihan block otherwise.
- The **Alphabetic Presentation Forms block** (U+FB00 to U+FB4F) contains some common Latin and Armenian ligatures, as well as precomposed Hebrew-letter-with-point combinations used in Yiddish.
- The **Arabic Presentation Forms A block** (U+FB50 to U+FDFF) contains a whole mess of precomposed Arabic ligatures which are virtually never used.
- The **Combining Half Marks block** (U+FE20 to U+FE2F) contains clones of the double diacritics (diacritics that appear over two characters) from the Combining Diacritical Marks block. The clones in the Combining Half Marks block have been split in half: to represent a double diacritic over two letters using the half marks, you follow the first letter with the first half of the desired mark and follow the second letter with the second half of the desired mark. The use of these characters is discouraged.

- The **CJK Compatibility Forms block** (U+FE30 to U+FE4F) contains clones of various CJK punctuation marks rotated sideways. Normally in vertical text you use the regular punctuation marks and the rendering engine automatically turns them sideways; these characters always represent the glyphs used in vertical text for those punctuation marks.
- The **Small Form Variants block** (U+FE50 to U+FE6F) contains small versions of the ASCII symbols. These are included for round-trip compatibility with the Taiwanese CNS 11643 standard.
- The **Arabic Presentation Forms B block** (U+FE70 to U+FEFF) contains characters representing the contextual glyphs used for the Arabic letters (e.g., the initial form of *heh*). Unlike the regular Arabic letters, these characters represent specific glyph shapes and aren't subject to contextual glyph selection. Some rendering engines use these characters for rendering the regular Arabic letters, and they can also be used to represent specific glyphs in text that discusses the Arabic letters. The byte-order mark/zero-width no-break space is also encoded in this block.
- The **Halfwidth and Fullwidth Forms block** (U+FF00 to U+FFEF) contains fullwidth variants of the ASCII characters and halfwidth variables of various CJK characters.

Miscellaneous characters

Finally, Unicode also includes these blocks of miscellaneous characters:

- The **Box Drawing block** (U+2500 to U+257F) contains a bunch of characters that can be used to “draw” lines on a character-oriented display. This block is derived from the box-drawing characters in the old IBM PC code pages that were used to draw Windows-like user interfaces in old DOS programs, and from a similar collection of box-drawing characters in the Japanese JIS standards.
- The **Block Elements block** (U+2580 to U+259F) contains a bunch of glyphs that represent different fractions of a monospaced display cell being colored in. These also come from the old DOS and terminal code pages.
- The **Geometric Shapes block** (U+25A0 to U+25FF) contains a bunch of geometric shapes. Some of these are used as symbols or bullets of various kinds, and they come from a variety of sources.
- The **Dingbats block** (U+2700 to U+27BF) contains the characters in the Zapf Dingbats typeface. Since this typeface has been included on so many laser printers, these characters (various pictures, bullets, arrows, enclosed numbers, etc.) have become something of a de facto standard, and the Unicode designers decided they deserved unique code points in Unicode.

And that (finally) concludes our guided tour of the Unicode character repertoire! Starting in the next chapter, we'll examine how Unicode text is actually handled in software.

SECTION III

Unicode in Action

*Implementing and Using the
Unicode Standard*

CHAPTER 13 *Techniques and Data Structures for Handling Unicode Text*

We've come a long way. We've looked at the general characteristics of Unicode, the Unicode encoding forms, combining character sequences and normalization forms, and the various properties in the Unicode Character Database. We've taken a long tour through the entire Unicode character repertoire, looking at the various unique requirements of the various writing systems, including such things as bidirectional character reordering, Indic consonant clusters, diacritical stacking, and the various interesting challenges associated with representing the Han characters. Now we've finally gotten to the point where we can talk about actually doing things with Unicode text.

In this section, we'll delve into the actual mechanics of performing various processes on Unicode text. We'll look at how you actually do such things as search and sort Unicode text, perform various types of transformations on Unicode text (including such things as converting to other character sets, converting between encoding forms, performing Unicode normalization, mapping between cases, performing various equivalence transformations, and more generalized transformations), and draw Unicode text on the screen. We'll also look at how Unicode is used in databases and on the Internet, and at the Unicode support provided in various operating systems and programming languages.

Most of the time, you don't need to know the nitty-gritty details of how to do various operations on Unicode text; the most common operations have already been coded for you and exist in various operating systems and third-party libraries. You just need to know the support is there and how to use it. Chapter 17 summarizes what's in the various Unicode support libraries and how to use them. If you're not actually in the business of implementing pieces of the Unicode standard, you may just want to read that chapter.

However, you might find you need to do things with Unicode text that aren't supported by the libraries you have at your disposal. This chapter will give you the ammunition to deal with your particular situation, as well as a grounding in the techniques that are used in the commercial libraries to perform the more common operations on Unicode text.

Useful data structures

In this chapter, we'll look at a series of useful data structures and techniques for doing various things with Unicode text. Most operations on Unicode text boil down to one or more of the following:

- Testing a character for membership in some collection of characters.
- Mapping from a single character to another single character, or to some other type of value.
- Mapping from a variable-length sequence of characters to a single character (or, more commonly, some other single value).
- Mapping from a single character (or other type of value) to a variable-length sequence of characters (or other values).
- Mapping from a variable-length sequence of characters or values to another variable-length sequence of characters or values.

More simply put, you may need to do one-to-one, one-to-many, many-to-one, or many-to-many mappings, with the universe of possible Unicode characters being the set on either side (or both sides) of the relation, and with various mixes of different-length sequences and different repetition characteristics in the sets. We'll look at a variety of different data structures designed to solve different variations on these problems.¹⁰⁸

Testing for membership in a class

A very common operation, and one that can't always be accounted for in out-of-the-box software, is testing a character for membership in a particular category of characters. For those categories that are defined by the Unicode standard, such as lowercase letters or combining diacritical marks, the standard Unicode libraries almost always provide APIs you can use. In addition, many APIs provide additional functions to do character-class-membership queries that are specific to their environment—for example, the Java class libraries provide methods for testing whether a character is legal in a Java identifier.

If the test you're looking for is available in an API library you're using (such as testing whether a character is a letter, digit, or whitespace character), you're usually best off using that library's functions. But you may run into situations where the category or categories of characters you're interested in is specific to your application and doesn't match the predefined categories set up by the designers of the APIs you're using.

Say, for example...

¹⁰⁸ Except where otherwise noted, the techniques described in this chapter are either from Mark Davis, "Bits of Unicode," *Proceedings of the Eighteenth International Unicode Conference*, session B11, April 27, 2001, or from my own research.

- ...you're interested in whitespace characters, but want to include not just the Unicode space-separator characters, but the line separator, paragraph separator, and the ASCII horizontal tab, vertical tab, line feed, form feed, and carriage return characters
- ...you're interested in letters, but only in certain scripts (say, Latin, Greek, Cyrillic, Arabic, and Hebrew, but not Armenian, Georgian, Syriac, Thaana, or any of the Indic characters)
- ...you're interested in Indic vowel signs only
- ...you're filtering out any character whose use in XML files is discouraged, and you don't have access to an API that'll do this for you
- ...you're interested in valid filename or identifier characters and your application or system defines this set in some way other than the Unicode standard (a personal favorite example is the Dylan programming language, where identifiers may not only consist of letters and digits, but also a wide variety of special characters, including !, &, *, <, >, |, ^, \$, %, @, _, -, +, ~, ?, /, and =.)
- ...you're defining a grammar (e.g., a programming-language grammar) and the categories of characters that are significant to the grammar don't line up exactly with the Unicode categories (a particularly interesting example of this is scanning natural-language text for linguistic boundaries, something we'll talk about in depth in Chapter 16)
- ...you're interested only in the subset of Unicode characters that have equivalents in some other standard (only the Han characters in the Tongyong Hanzi, only the Hangul syllables in the KS X 1001 standard, only the Unicode characters that aren't considered "compatibility characters" according to some appropriate definition, only the Unicode characters that are displayable by some font you're using, etc.)

In all of these cases, you need a way to check an arbitrary character to see whether it falls into the set you're interested in (or, if you're interested in multiple sets, which of them it falls into).

There are a number of ways to do this. If the category is relatively close to one of the Unicode categories (or the union of several), you might use a standard API to check for membership in that category and then just special-case the exceptions in code. For instance, you could do the whitespace check above with something like this in Java:

```
public static boolean isWhitespace(char c) {
    // true for anything in the Zs category
    if (Character.getType(c) == Character.SPACE_SEPARATOR)
        return true;

    // true also for any of the following:
    switch (c) {
        case '\u2028': // line separator
        case '\u2029': // paragraph separator
        case '\n':     // line feed
        case '\t':     // tab
        case '\f':     // form feed
        case '\r':     // carriage return
        case '\u000b': // vertical tab
            return true;

        // everything else is a non-whitespace character
        default:
            return false;
    }
}
```

(This, of course, is kind of an unnecessary example, as `Character.isWhitespace()` does essentially the same thing, but it's useful for illustrative purposes.)

For some sets of characters, this kind of thing is sufficient, but it can be kind of clunky and it doesn't perform well for complicated categories.

There are other simple alternatives: If the characters all form one or two contiguous ranges, you can, of course, do a simple range check...

```
public static boolean isASCIILetter(char c) {
    if (c >= 'A' && c <= 'Z')
        return true;
    else if (c >= 'a' && c <= 'z')
        return true;
    else
        return false;
}
```

...but this also doesn't scale well.

If the characters you're interested in all fall within, say, the ASCII block, you could use a range check followed by a simple array or bit-map lookup...

```
public static boolean isASCIIConsonant(char c) {
    if (c > 0x7f)
        return false;
    else
        return consonantFlags[c];
}

private static boolean[] consonantFlags = {
    false, false, false, false, false, false, false, false,
    false, false, false, false, false, false, false, false,
    false, false, false, false, false, false, false, false,
    false, false, false, false, false, false, false, false,
    false, false, false, false, false, false, false, false,
    false, false, false, false, false, false, false, false,
    false, false, true, true, true, false, true, true,
    true, false, true, true, true, true, true, false,
    true, true, true, true, true, false, true, true,
    true, true, true, false, false, false, false,
    false, false, true, true, true, false, true, true,
    true, false, true, true, true, true, true, false,
    true, true, true, true, true, false, true, true,
    true, true, true, false, false, false, false, false
};
```

...but arrays and bitmaps gets unwieldy very fast if the characters you're interested in cover more than a fairly small subset of the Unicode characters.

The inversion list

A nice, elegant solution for representing any arbitrary set of Unicode characters in the *inversion list*. It's a simple array, but instead of having an entry for each Unicode character, or simply listing all the Unicode characters that belong to the category, the entries define *ranges* of Unicode characters.

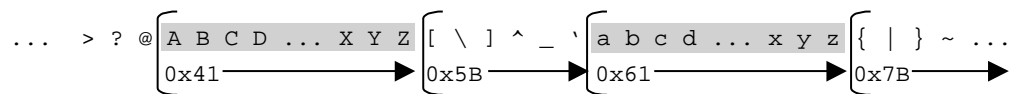
It works like this: The first entry in the array is the first (that is, the lowest) code point value that's in the category. The second entry in the array is the first code point value that's higher than the first entry and *not* in the set. The third entry is the first code point value that's higher than the second entry and *is* in the set, and so on. Each even-numbered entry represents the lowest code point value in a contiguous range of code point values that are in the set, and each odd-numbered entry represents the lowest code point value in a contiguous range of code point values that are *not* in the set.

The set of basic Latin letters would thus be represented like this:

U+0041 ('A')
U+005B ('[')
U+0061 ('a')
U+007B ('{')

The first entry represents the first character in the set (the capital letter A). The set then consists of all the characters with code point values between that one and the next entry (the opening square bracket). The second entry (the opening square bracket) is the beginning of a range of characters that isn't in the set. That range is terminated by the third entry (the small letter a), which marks the beginning of the next range, and *that* range is terminated by the fourth entry (the opening curly brace), which marks the beginning of a range of code point values that isn't included in the set. This range extends to the top of the Unicode range, so the last character before this range (the small letter z) is the character with the highest code point value that's included in the set.

Conceptually, you get something like this:



The brackets represent the beginnings of ranges (the entries in the table), and the shaded ranges are the ones that are included in the set.

Testing a character for membership in the set is a simple matter of binary-searching the list for the character (you code the binary search so that if the character you're looking for isn't specifically mentioned, you end on the highest code point value less than the character). If you end up on an even-numbered entry, the character's in the set. If you end up on an odd-numbered entry, the character's not in the set.

The code to check membership would look something like this in Java:

```
public static boolean charInSet(char c, char[] set) {
```

```
int low = 0;
int high = set.length;
while (low < high) {
    int mid = (low + high) / 2;
    if (c >= set[mid])
        low = mid + 1;
    else if (c < set[mid])
        high = mid;
}

int pos = high - 1;

return (pos & 1) == 0;
}
```

Armed with this new routine, we can re-cast our `isASCIIConsonant()` example using a much shorter table:

```
public static boolean isASCIIConsonant(char c) {
    return charInSet(c, asciiConsonants);
}

private static char[] asciiConsonants = {
    'B', 'E', 'F', 'I', 'J', 'O', 'P', 'U', 'V', '[',
    'b', 'e', 'f', 'i', 'j', 'o', 'p', 'u', 'v', '{'
};
```

An empty array represents the null set. An array with a single entry containing U+0000 represents the set containing all the Unicode characters. A set containing a single character would (generally) be a two-element array where element 0 is the code point for the character you're interested in and element 1 is that code point plus one.

Performing set operations on inversion lists

One neat thing about this representation is that set operations on inversion lists are pretty simple to perform. A set inversion is especially simple—if the first entry in the array is U+0000, you remove it; otherwise, you add it:

```
public static char[] invertSet(char[] set) {
    if (set.length == 0 || set[0] != '\u0000') {
        char[] result = new char[set.length + 1];
        result[0] = '\u0000';
        System.arraycopy(set, 0, result, 1, set.length);
        return result;
    }
    else {
        char[] result = new char[set.length - 1];
        System.arraycopy(set, 1, result, 0, result.length);
        return result;
    }
}
```

Testing for membership in a class

Union and intersection are a little more interesting, but still pretty simple. To find the union of two sets, you basically proceed the same way you would to merge two sorted lists of numbers (as you would in a merge sort, for example). But instead of writing every entry you visit to the output, you only write some of them.

The way you decide which entries to write is by keeping a running count. You start the count at 0, increment it every time you visit an even-numbered entry (i.e., the beginning of a range), and decrement it every time you visit an odd-numbered entry (i.e., the end of a range). If a particular value causes you either to increment the running count away from 0 or decrement it to 0, you write it to the output.

Say, for example, you have two sets, one containing the letters from A to G and M to S and the other containing the letters from E to J and U to X. The running counts would be calculated as follows:

```
A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
|----->|----->|----->
|----->|----->
1 1 1 1 2 2 2 1 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 0 0
```

The union of the two sets should contain all the letters where the running count is greater than zero. So when you see A, the running count gets incremented to 1 and you write the A to the result set. When you see the E, the running count gets incremented to 2 and you don't write anything. When you see the H, the running count gets decremented to 1 and you still don't write anything. When you see the K, however, the count gets decremented to 0. This causes you to write the K to the result set. And so on.

There's an important wrinkle to watch out for here: what to do if the entries from the two lists are equal. Say, for example, you have two sets, one containing the letters from A to E and another containing the letters from F to K:

```
A B C D E F G H I J K L M
|----->|----->
1 1 1 1 1 1 1 1 1 1 1 0 0
```

The letter F occurs in both sets: as the endpoint of the first set (the first character after the range it contains) and as the starting point of the other set. When you're merging the two arrays, you'll thus come across a situation where you're looking at two Fs. If both of them are even-numbered entries or both of them are odd-numbered entries, you can process them in any order: only one of them will decrement the running count to 0 or increment it away from 0, and it doesn't matter which one.

But if, as in our example, one of them is an odd-numbered entry and one is an even-numbered entry, you have a problem. If you process the odd-numbered entry first, you might decrement the count to zero and then immediately increment the count away from zero again. This would, in our example, cause F to be written to the result set twice. This gives us a result that consists of two ranges: the range from A to E and the range from F to K. This'll still give us the right answer, but it wastes space. What you really want is a single range, running from A to K. Processing the even-numbered entry first will correctly coalesce the two ranges into a single range and prevent the Fs from being written to the result set.

Anyway, here's what the union function looks like in Java:

```
public static char[] union(char[] a, char[] b) {
    char[] tempResult = new char[a.length + b.length];
    int pA = 0;
    int pB = 0;
    int pResult = 0;
    int count = 0;
    char c;
    int p;

    // go through the two sets as though you're merging them
    while (pA < a.length && pB < b.length) {

        // if the lower entry is in the first array, or if the
        // entries are equal and this one's the start of a range,
        // consider the entry from the first array next
        if (a[pA] < b[pB] || (a[pA] == b[pB] && (pA & 1) == 0)) {
            p = pA;
            c = a[pA++];
        }

        // otherwise, consider the entry from the second array next
        else {
            p = pB;
            c = b[pB++];
        }

        // if the entry we're considering is the start of a range
        // (i.e., an even-numbered entry), increment the running
        // count. If the count was zero before incrementing,
        // also write the entry to the result set
        if ((p & 1) == 0) {
            if (count == 0)
                tempResult[pResult++] = c;
            ++count;
        }

        // if the entry we're considering is the end of a range
        // (i.e., an odd-numbered entry), decrement the running
        // count. If this makes the count zero, also write the
        // entry to the result set
        else {
            --count;
            if (count == 0)
                tempResult[pResult++] = c;
        }
    }

    // figure out how big the result should really be
    int length = pResult;

    // if we stopped in the middle of a range, decrement the count
    // before figuring out whether there are extra entries to write
    if ((pA != a.length && (pA & 1) == 1)
        || (pB != b.length && (pB & 1) == 1))
```

```
        --count;

// if, after the adjustment, the count is 0, then all
// entries from the set we haven't exhausted also go into
// the result
if (count == 0)
    length += (a.length - pA) + (b.length - pB);

// copy the results into the actual result array (this may
// include the excess from the array we hadn't finished
// looking at)
char[] result = new char[length];
System.arraycopy(tempResult, 0, result, 0, pResult);
if (count == 0) {
    // only one of these two calls will do anything
    System.arraycopy(a, pA, result, pResult, a.length - pA);
    System.arraycopy(b, pB, result, pResult, b.length - pB);
}
return result;
}
```

The example above contains an optimization: rather than continue through the loop until we've exhausted both input sets, we drop out of the loop as soon as we've exhausted one of them. This avoids a couple extra checks in that first `if` statement to account for the possibility of one of the input lists being exhausted and speeds up that loop. If, after dropping out of the loop, our running count is greater than 0, the result set ends with an even-numbered entry and everything higher than that is included in the set; we can safely ignore the entries we haven't looked at yet. If, on the other hand, the running count *is* 0, we can just copy all of the entries we haven't looked at yet directly into the result.

The code to perform an intersection is exactly like the code for performing the union, except for these differences:

- Instead of writing a value to the output when the count is incremented away from 0 or decremented to 0, you write a value to the output when the count is incremented to 2 or decremented away from 2.
- If the entries from the two input sets are equal, process the *odd-numbered* entry (i.e., decrement the running count) first.
- In the optimization at the end, instead of copying the extra entries from the non-exhausted list into the result when the count is 0, you do so when the count is 2.

Other operations can be built up from the primitives we've looked at above. A set difference is simply an inversion followed by an intersection. To add a single character, first turn it into an inversion list (to make *c* into an inversion list, create a two-element array consisting of *c* followed by *c + 1*) and then take the union of the new set and the set you're adding to. To remove a single character, turn that character into an inversion list and do a set difference. Optimized implementations of these operations are certainly possible, but are left as an exercise for the reader.

Inversion lists have a number of desirable characteristics. Chief among them is their compactness. The entire set of ideographic characters in Unicode (as defined in `PropList.txt`) can be represented as

```
{ 0x3006, 0x3008, 0x3021, 0x302A, 0x3038, 0x303B, 0x3400, 0x4DB6, 0x4E00,
```

```
0x9FA6, 0xF900, 0xFA2E, 0x20000, 0x2A6D7, 0x2F800, 0x2FA1D }
```

That's a pretty compact way to represent a list of 71,503 characters.

Similarly, the entire set of Latin letters (ignoring IPA characters and fullwidth variants) compresses to

```
{ 0x0041, 0x005B, 0x0061, 0x007B, 0x00C0, 0x00D7, 0x00D8, 0x00F7, 0x00F8,
0x0220, 0x0222, 0x0234, 0x1E00, 0x1E9C, 0x1EA0, 0x1EFA }
```

Again, 662 characters compress to 16 entries.

Inversion lists are also relatively fast to search and easy to modify.

Still, they have a couple of drawbacks. One, they're fairly quick to search unless you're doing a lot of searches quickly, as you might in a very tight loop while iterating over a long string. In this case, you need a data structure that's truly optimized for fast lookup, and we'll look at just such a beast below.

Second, because they don't operate on groups of characters, inversion lists can only operate on the full Unicode range if they store their code point values as UTF-32 units. This means that a Unicode 3.1-compatible inversion list in Java has to be an array of `int`, not an array of `char`, as in the example code we just looked at, and it also means that you may have to perform conversions from UTF-16 to UTF-32 before you can do a lookup. (In Java, for example, strings are sequences of UTF-16 code units; if you encounter a surrogate pair, you have to convert it to a single UTF-32 code unit before you can look it up in an inversion list. This kind of thing can make inversion lists less ideal and force the use of an alternative data structure that can handle variable-length sequences of code units and thus work on UTF-16 in its native form.)

Finally, the benefit of the compression goes away if you have long stretches of code point values that alternate between being in the set and not being in the set. For example, an inversion list representing the set of capital letters in the Latin Extended A block would look like this:

```
{ 0x0100, 0x0101, 0x0102, 0x0103, 0x0104, 0x0105, ... }
```

If the inversion list is only being used in a transient manner (say, to derive character categories for a parser), this may be an acceptable limitation. For longer-lived lists of characters, there are alternative representations that provide better compression.

Mapping single characters to other values

Testing a character for membership in a class is really just a special case of mapping from a single character to a single value—in essence, you're mapping from a character (or, more accurately, a code unit) to a Boolean value.

But most of the time, you want to map to something with more values than a Boolean. For example, instead of checking to see whether a character is a member of some class, you may want to perform

one operation to see which of several mutually-exclusive classes a given character belongs to (as opposed to checking each category one by one for the character, as you might have to do if the categories overlapped).

There are plenty of other uses for this kind of thing: mapping characters from upper case to lower case or to a case-independent representation¹⁰⁹, looking up a character's combining class, bi-di category, numeric value, or (sometimes) converting from Unicode to some other encoding standard.

Inversion maps

One way of doing this is to use an inversion *map*. An inversion map is an extension of an inversion list. With a regular inversion list, the value being mapped to is implicit: even-numbered entries represent ranges of characters that map to “true,” and odd-numbered entries represent ranges of characters that map to “false.”

You can extend this by storing a separate parallel array that contains the values you want to map to. Each entry in the main list contains the first code unit in a range of code units that map to the same output value. The range extends until the beginning of the next range (i.e., the code unit value in the next entry), and the value that all the characters map to is stored in the corresponding entry in the parallel array.

For example, the character categories for the characters in the ASCII range could be represented in two parallel arrays that look like this:

U+0000	Cc (control character)
U+0020	Zs (space separator)
U+0021 (!)	Po (other punctuation)
U+0024 (\$)	Sc (currency symbol)
U+0025 (%)	Po (other punctuation)
U+0028 ((Ps (starting punctuation)
U+0029 ())	Pe (ending punctuation)
U+002A (*)	Po (other punctuation)
U+002B (+)	Sm (math symbol)
U+002C (,)	Po (other punctuation)
U+002D (-)	Pd (dash punctuation)
U+002E (.)	Po (other punctuation)
U+0030 (0)	Nd (decimal-digit number)
U+003A (:)	Po (other punctuation)
U+003C (<)	Sm (math symbol)
U+003F (?)	Po (other punctuation)
U+0041 (A)	Lu (uppercase letter)
U+005B ([)	Ps (starting punctuation)
U+005C (\)	Po (other punctuation)
U+005D (])	Pe (ending punctuation)
U+005E (^)	Sk (modifier symbol)
U+005F (_)	Pc (connector punctuation)
U+0060 (`)	Sk (modifier symbol)
U+0061 (a)	Ll (lowercase letter)
U+007B ({)	Ps (starting punctuation)
U+007C ()	Sm (math symbol)

109 Actually, this isn't always a one-to-one operation, but we'll ignore that for the time being.

```
U+007D (})    Pe (ending punctuation)
U+007E (~)    Sm (math symbol)
U+007F        Cc (control character)
```

It's rather ugly because the punctuation marks and symbols in the ASCII range are scattered across a bunch of Unicode categories, but you still manage to compress a single 128-element array into a pair of 29-element arrays.

The compact array

The inversion map can work fairly well for a lot of tasks, but if you have to do a lot of lookups in the table, this isn't the best way to do it. If the compression leads to a list with only a few entries in it, lookup is pretty fast, and the speed of lookup degrades more and more gradually as the list gets longer. Still, there's no escaping the fact that it's faster to go straight to the entry you want than to binary-search the array for it (even if that, in turn, is a lot faster than linear-searching it), and that constant-time lookup performance is (generally speaking) better than $O(\log n)$ performance. If you're doing a lot of lookups in a very tight loop, as you will when doing almost anything to a large body of text (e.g., Unicode normalization, conversion between Unicode and some other encoding, glyph selection for rendering, bidirectional reordering, searching and sorting, etc.), this difference can become quite large (just how large, of course, depends on what you're doing—both the exact number of lookups the algorithm requires and the exact nature of the table of data you're doing the lookup in make a big difference).

Of course, we come back to the problem that a regular array—the reasonable way to do most of these operations on single-byte non-Unicode text—would be prohibitively huge when you're operating on Unicode text. With Unicode 3.0 and earlier, you'd need an array with 65,536 entries (assuming you don't have to worry about the private-use code points in Planes 15 and 16). With Unicode 3.1, you'd need an array with a whopping 1,114,112 entries. Yeah, hardware is getting cheaper and cheaper and more and more complex, and so memory and storage capacities are getting bigger, but this is ridiculous, especially if you need more than one lookup table.

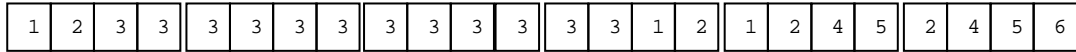
Fortunately, there's a simple and elegant technique for compressing a 65,536- (or 1,114,112-) element array that still gives you fast constant-time lookup performance. Lookup takes a little more than twice the time of a straight array lookup, but this is still blazingly fast, and it's still a constant time for every item in the array. This basic technique (and various others based on it that we'll look at) forms the cornerstone of most of the common algorithms for doing various things on Unicode text, and it's well worth understanding. This technique goes by a number of different names, but for our purposes here, we'll refer to it by the rather unimaginative name of “compact array.”

To understand how the compact array works, let's consider an example array:

1	2	3	3	3	3	3	3	3	3	3	3	3	3	1	2	1	2	4	5	2	4	5	6
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

This array has 24 elements. Obviously, there's no need to compress a 24-element array, but it makes it easy to see what's going on. First, you (conceptually) break the array into a number of equally-sized blocks. For the purposes of this example, we'll use four-element blocks, although again real life this would do you no good at all. Now you have something that looks like this:

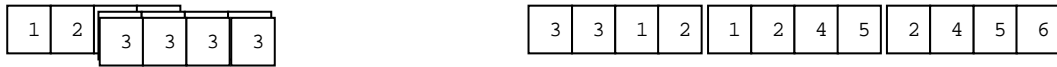
Mapping single characters to other values



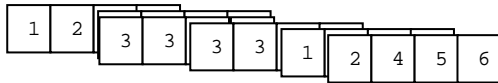
Here's the next part: If a block ends with the same set of elements that the next block begins with (or if the blocks are simply the same), those elements can share the same storage. For example, the first block ends with "3 3" and the second block begins with "3 3," so you can stack them on top of each other, like this:



The second block and the third block are identical, so you can stack them completely on top of each other...

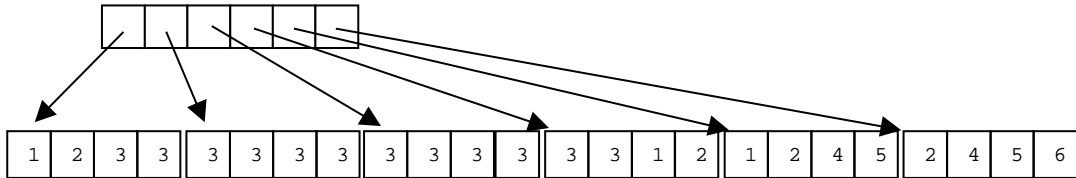


If you continue in this way, stacking adjacent blocks that have elements in common, you end up with something that looks like this:

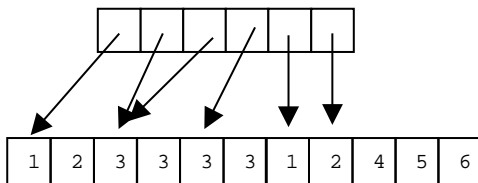


So we've compressed our original 24-element array down to 11 elements.

Obviously, the only way you can access the elements in the array is through the use of an external index that tells you where in the compressed array the beginning of each block is, so the uncompressed data structure really looks like this...

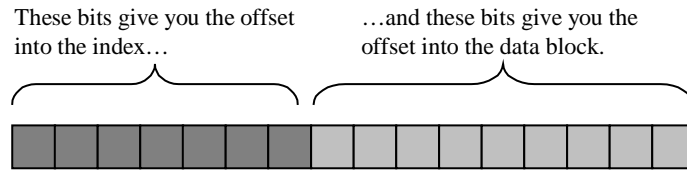


...and the compressed data structure really looks like this:



So in our example array, you save 13 elements in the data array, but you need a new six-element index array, for a net savings of seven array elements. With real Unicode-based tables, the compression is generally much greater.

Looking up a character in a compact array is relatively simple. The most-significant bits of the code point give you the offset into the index array, and the least-significant bits of the code point give you the offset into the data block that the index-array entry points you to.



The code to perform a lookup into a compact array is quite simple:

```
result = data[index[c >> INDEX_SHIFT] + (c & BLOCK_MASK)];
```

In this code snippet, `c` is the character we're looking up a value for, and `index` and `data` are the index and data arrays. `INDEX_SHIFT` is the number of bits in the character that constitute the offset into the data block (and, hence, the number of bits you shift the code point value to the right to get the offset into the index block). `BLOCK_MASK` is a mask used to isolate the offset into the data block, and is always equal to $(1 \ll \text{INDEX_SHIFT}) - 1$. The whole lookup consists of two array lookups, a shift, an AND operation, and an addition. In fact, if you're willing to have the index array have four-byte elements (it's usually only necessary for the index elements to be two bytes), you can store actual pointers in the index array (at least in C or C++) and eliminate the addition. This lookup compiles extremely efficiently on almost all programming languages and machine architectures.

(Note that the `>>` in the above example is a *logical* (i.e., zero-filling) right shift. In C, this means that `c` has to be of an unsigned integral type, or of a type that's at least a bit wider than necessary to hold the code point value, so the sign bit is always zero. In Java, this shouldn't be a problem, since `char` is an unsigned type, and `int` is bigger than necessary to hold a UTF-32 value (which is actually only 21 bits wide—even a full UCS-4 value would only be 31 bits wide). If for some weird reason, you had the code point value stored in a `short`, however, you'd have to remember to use the `>>>` operator to do the shift.)

There are several important things to note about using compact arrays. First, it works best either on sparsely-populated arrays (i.e., arrays where the vast majority of elements are 0), arrays with long stretches of identical values or long stretches of identical patterns of repeating values, or other arrays with high levels of redundancy. Fortunately, virtually all Unicode-related lookup tables either fit these criteria or can be represented in a way that does (for example, mappings between code-point values, such as conversions from Unicode to other encodings, or case-mapping tables, don't tend to lend themselves to compression unless the tables are sparsely populated, and then the populated blocks can't be coalesced. Turning these into tables of numeric offsets that get added to the input value to produce the output value can often turn long stretches of different values into long stretches of the same value, allowing compression).

Second, there's nothing about the design of this data structure that requires that the data blocks appear in the same order in the compressed data array as they did in the original uncompressed array, or that blocks that get coalesced together be adjacent in the original array. Most real-life tables get most of their compression from coalescing identical blocks in the original uncompressed array. The unpopulated space in a sparse array, for example, turns into a single unpopulated block in the compressed representation, even though not all of the unpopulated blocks in the original array were next to each other. (Often, part of the unpopulated block can also be doubled up with the unpopulated space in one or two of the populated blocks.)

Third, it's impossible to make generalizations as to what the optimum value for `INDEX_SHIFT` should be. As `INDEX_SHIFT` gets larger, the individual blocks in the data array get bigger, but the index gets smaller; as it gets smaller, the individual blocks get smaller (possibly permitting more of them to be doubled up in the compressed representation), but the index gets longer, offsetting those gains. If `INDEX_SHIFT` goes either to zero or to the total number of bits in the character code, you're back to doing a simple array lookup, since either the index size or the block size goes to zero. Generally, for a 16-bit code point value, values of `INDEX_SHIFT` between 6 and 10 work the best.

Fourth, there are many methods of compressing the data in a compact array, none of which are particularly fast. They don't work well for tables of data that change a lot at run time. They work best when the tables can be built once at run time and then used repeatedly, or, better yet, can be calculated as part of the development of the software the uses them and simply used at run time (i.e., tables that don't depend on data that's only available at run time, such as tables of data that's defined in the Unicode standard).

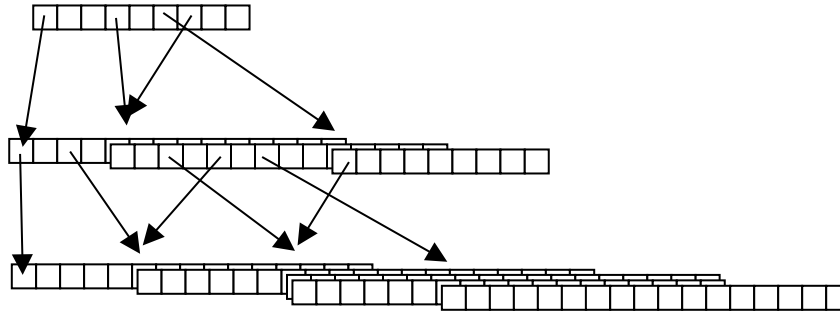
Generally the fastest method of compressing a compact array is to just pick an `INDEX_SHIFT` and worry about coalescing only identical blocks. More compression can be obtained by coalescing half blocks or quarter blocks. Or by examining adjacent blocks for identical heads and tails. The only way to get the absolutely optimal amount of compression for an arbitrary collection of data is to try every possible ordering of the blocks and examine the adjacent pairs for identical heads and tails (this would actually also involve trying all of the possible `INDEX_SHIFT` values as well). Most of the time, this isn't worth the trouble and just coalescing identical blocks or half-blocks (and maybe looking for identical heads and tails in blocks adjacent in the original ordering) gets good results. Fortunately, the algorithm used to compress the data doesn't affect the algorithm used to look up values in the compressed array, nor does it affect the performance.¹¹⁰

Two-level compact arrays

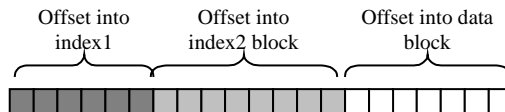
A simple compact array as described above works well with UTF-16 values, but this leaves out all of the new characters in Unicode 3.1 (and all the private-use characters in Planes 15 and 16). If your application can ignore these characters, you can use a simple compact array, but if you need to support the characters in the supplementary planes, a simple compact array doesn't work so well. When you're trying to partition a 21-bit value into two groups, you tend to end up either with very long data blocks, making compression more difficult, or very long indexes, making compression less relevant.

One approach to dealing with this problem is to use a *two-level* compact array. Instead of having a single index array, you have two. In other words, you apply the same compression techniques to the index that you apply to the data array. The first index points you to a block of values in the second index, which in turn points you to a block of values in the data array.

¹¹⁰ Actually, in some processor architectures and programming languages, lookup can be made a little faster if you use an `INDEX_SHIFT` of 8, but the improvement is minuscule.



The code point value gets broken into three groups instead of two: The most-significant bits give you the offset into the first index, the next set of bits gives you the offset into the appropriate block of the second index, and the least-significant bits give you the offset of the result in the appropriate block of the data array.



Performing the lookup into a two-level compact array is little more complicated and a little slower than doing it into a regular compact array:

```
result = data[index2[index1[c >> INDEX1_SHIFT] + (c & INDEX2_MASK) >>
INDEX2_SHIFT] + (c & DATA_MASK)];
```

Instead of two array lookups, a shift, a mask, and an addition, this takes three array lookups, two shifts, two masks, and two additions. Nine operations instead of five, but still quite fast and still constant time.

Mapping single characters to multiple values

Compact arrays are still basically just that—arrays to which compression and indexing has been applied. Like regular arrays, this means that they’re best suited to store a bunch of values that are all the same size. They also work best when the values being mapped to are relatively small; a big array of 32-byte values, for example, gets kind of unwieldy.

But what if you’re looking up a value that is large, or could be of variable length? There are a couple of techniques that can be used.

If the total universe of values that can be mapped to is relatively small, you can save space without imposing too huge a performance penalty by offloading the target values into a separate array (or other data structure). The target values in the compact array would just be indexes into the auxiliary array containing the values.

For large complex data structures, you could even go so far as to have the entries in the compact array’s data array be pointers to the target data structures. The big negative here is that means the data array would have to be of 32-bit values, which could make it pretty big if you couldn’t achieve a

lot of compression. A frequent alternative is to use an auxiliary array to hold the pointers—the data elements in the compact array then need only be big enough to hold the indexes into the auxiliary array: 16 or even eight bits.

If you're mapping to string values, especially if most of the resultant strings are short, it's often useful to store all the target strings in one big array. This cuts down on memory-manager overhead (on most systems, every independent block of memory has some additional amount of memory, usually eight or more bytes, allocated to it as a block header, where housekeeping info for the memory manager is stored). The compact array can then just contain offsets into the big string block, or if you're willing to have the compact array map to 32-bit values, actual pointers to the individual strings in the memory block. This technique works best in C and C++, where pointers can point into the interiors of memory blocks, where character strings are null-terminated, and where a string is just a pointer to an arbitrary run of character data in the first place.

Exception tables

A common situation might be one where you're generally mapping to a fixed-length value, but there are occasional target values that are longer. The most common version of this is where you're mapping to strings and the vast majority of the resulting strings are single characters (you can run into this, for example, doing language-sensitive case mapping, Unicode normalization, language-sensitive comparison, or mapping Unicode to a variable-length target encoding). In these cases, there's an optimization that's worth considering.

The basic idea is that you make the data array elements the size of the most common result values. Then you designate some range of the result values as sentinel values. The non-sentinel values are used as-is as the results; the sentinel values are parsed in some way to derive the result values. Usually this means performing some mathematical transformation on them to get an offset into an auxiliary array containing the large or variable-length results.

Say, for example, you're mapping single Unicode characters to either characters or strings, with characters predominating. You'd set up the compact array to store characters, but set aside a range of character values as sentinel values. An obvious range to use for this purpose is the range from U+D800 to U+F800, the surrogate space and the private-use area. If you look up a character in the compact array and it's outside this range, then it's your result and you're done. If it's inside this range, you subtract 0xD800 from it to get the offset into an auxiliary data structure holding the strings, such as the one we just looked at. If any of the single characters in the private-use area is a result value, you just treat it as a one-character string: it appears in the auxiliary data structure with all the longer strings, and the entry in the compact array is the properly-munged reference to its position in the auxiliary structure.

Even if you're only mapping single characters to single characters, this technique is useful. A generalized mapping to Unicode characters would have to use a data array with 32-bit elements, since some of the result characters might be supplementary-plane characters. If, as is usually the case, BMP characters dominate your result set, you can use a data array with 16-bit elements and store the supplementary characters in a separate array. You can then store the offsets to the auxiliary array in the compact array using the munging technique described above.

Another approach sets aside only one value in the compact array's result set as a sentinel value. If you get something other than this value from the compact-array lookup, it's the result and you're

done. If you get the sentinel value, you repeat the whole lookup in a slower data structure that holds only the exceptional values.

Consider, for example, the Unicode bi-di categories. There are 19 bi-di categories, but there are six categories that consist of only one character. If we collapse these down to a single “category,” we get 14 bi-di categories, and we can represent the bi-di category as a four-bit value. This could let us squish the bi-di categories for two characters into a single byte, potentially halving the size of a table of bi-di categories. Bi-di reordering code would normally have a section of code that looks up the category for each character and takes some action depending on the result. If this code looked up the category and got “single character category,” it could just look at the input character itself to decide what to do.

In other cases, when you get the “exceptional character” value, you’d refer to an “exception table,” a slower data structure, such as a hash table or a sorted array of character-result pairs, that contained the results for the exceptional characters. You pay a performance penalty, but you’d only have to pay it for a small and relatively infrequent subset of the characters.

We’ll see numerous examples of these techniques in the chapters that follow.

Mapping multiple characters to other values

Then there’s the situation where you want to map a key that may consist of *more than one* character to some other value. If, for example, Unicode decomposition is a situation where you’re mapping single characters to variable-length strings, Unicode composition is a situation where you’re mapping variable-length strings to single characters. Even in situations where you’re consistently mapping single characters into single characters (or into some other fixed-length value), a UTF-16-based system (or a system that wants to save table space) might choose to treat supplementary-plane characters as pairs of surrogate code units, effectively treating single characters as “strings.”

As with mapping single characters (or other values) into variable-length strings, mapping variable-length strings into single characters (or other values) is something we’ll run into over and over again in the coming chapters. It’s important, for example, for canonical composition, language-sensitive comparison, character encoding conversion, and transliteration (or input-method handling).

Exception tables and key closure

There are a couple of broad methods of approaching this problem. If the vast majority of your keys are single characters and you have relatively few longer keys, you can use the exception-table approach we looked at in the previous section. You look up a character in a compact array and either get back a real result or a sentinel value. In this case, the sentinel value tells you that the character you looked up is the first character in one or more multiple-character keys. You’ll have to look at more characters to find out the answer, and you’ll be doing the lookup in a separate exception table designed for multiple-character keys.

Again, there are many different approaches you can take for the design of the exception table, the simplest again being a hash table or a sorted list of key-value pairs.

Let’s say for a moment that you use a sorted list of key-value pairs. If you can set aside more than one sentinel value in your result set, you can use the sentinel value to refer you directly to the zone in

the exception table that contains the keys that start with the character you already looked up. Each successive character then moves you forward in the exception table until you get to a sequence that isn't in the table at all (i.e., not only does it not exist as a key, it also doesn't exist as the beginning of a key). At this point, the last entry you were sitting on gives you your answer, and you start over with the last character you were examining, looking it up in the main table again.

This would probably be clearer with an example. Let's say you're implementing code to allow someone to input Russian text with a regular Latin keyboard. As you type, characters will be converted from Latin letters to their Cyrillic equivalents. There are a number of Cyrillic letters that are normally transliterated into Latin using more than one letter. For example, consider this group of letters¹¹¹:

```
c    =>  s
ш    =>  sh
щ    =>  shch
ц    =>  c
ч    =>  ch
```

If I type an "s," we don't know whether I mean for c to appear on the screen, or if I'm going to type "h" next to get ш. For that matter, if I type "sh," I might mean for ш to appear on the screen, or I might follow it with "ch" to get щ.

So if "sh," "ch," and "shch" were the only three multiple-keystroke sequences you had to worry about, you'd wind up with an exception table that looked like this:

c => ц
ch => ч
s => c
sh => ш
shch => щ

The main compact array would contain the mappings for all of the other Latin letters, but would contain pointers into this table for "s" and "c." This is why you need entries for "s" and "c" themselves: their mappings can't appear in the main table; the main table contains the signal that "s" and "c" are exceptional characters.

So I type "s," and the main mapping table points you to line 3 of the exception table:

111 This example is ripped off directly from Mark Davis, "Bits of Unicode," *op. cit.*

c => 𐀀
ch => 𐀁
s => 𐀂
sh => 𐀃
shch => 𐀄

You buffer the “s” in preparation for getting more characters. I type “h,” and you add it to the buffer. The buffer now contains “sh,” and you scan forward in the exception table looking for the first key that begins with “sh.” It’s the next entry, so now you’re pointing there:

c => 𐀀
ch => 𐀁
s => 𐀂
sh => 𐀃
shch => 𐀄

Now I type “a.” You add “a” to the buffer (which now contains “sha”) and scan forward for the first entry that starts with “sha.” The next entry is “shch,” which comes after “sha,” so “sha” isn’t in the table. You output the result for the last key you were on (“sh”), which gives you 𐀃, and then you go back and look up “a” in the main lookup table. “a” isn’t an exceptional character, so you get an immediate result — a — and output that. You’ve now correctly mapped “sha” to “𐀃a”.

But this algorithm breaks down if I type “shca.” After processing “shc,” you have “shc” in your buffer and you’re pointing at the entry for “shch,” the first entry that begins with “shc”:

c => 𐀀
ch => 𐀁
s => 𐀂
sh => 𐀃
shch => 𐀄

If the next letter I type is “a,” you’re already past where it would go in the table, so it’s not in the table. Your code would do one of two things. If it’s designed to always scan forward, it’ll actually behave as though I typed “shcha,” outputting the result for the entry we’re sitting on (“𐀄”) and then looking up “a” again to get “𐀄a”, or it’d know we were already past the entry we were looking for

and use the *last* entry that matched, giving us “ш.” If it does the second, it now has to be smart enough to know that there are *two* input characters that haven’t been accounted for, and to go back and look up *both* of them in the main table.

So the code has to be a lot more complicated to account for the cases where you actually have to go back and re-process more than one character. Either that, or you have to change the table.

Changing the table is the cleaner solution. The principle for changing the table is called *key closure*. The principle of key closure dictates that for every key that’s represented in the exception table, every *prefix* of that key must also appear in the table. For our example, this means that if “shch” is in the table, “s,” “sh,” and “shc” must all also occur in the table. So we solve our bug by including “shc” as a key in the table:

c => ш
ch => ч
s => с
sh => ш
shc => шшч
shch => шч

If I type “shca,” you’ll end up on the entry for “shc” when you realize “shca” isn’t in the table, output шшч, and go back and look up just “a”, adding a to the output. You output the correct answer for that series of characters, which is шшча.

Tries as exception tables

Modeling the exception table as a simple array of pairs of strings is nice and simple, but it’s not terribly efficient. As you walk through the table, you have to keep track of all the characters you’ve seen before and do a bunch of string comparisons to get to the right entry. If the keys and the table are both short, this probably won’t matter much, but for more complicated tables, all these lookups get to be a pain. You can actually change things so that each key is only one character long. The way you do this is by recursively applying the exception-table rule. The entry in the main table has a sentinel character telling you this character may be the first character of a multi-character key and referring you to an exception table. You look up the *next* character from the input in the exception table. This will either produce a real result, or it’ll produce a sentinel value telling you to go to *another* exception table and look up the *next* character from the input in *there*. The process continues until you land on a real result. For our Russian example, this gives us a data structure that looks like this:

Main table:

a => a
b => б
c => ●
...
s => ●
...

Exception tables:

h => ч
∅ => и

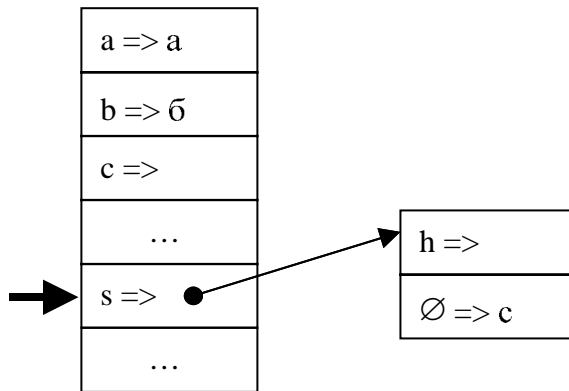
h => ●
∅ => c

c => ●
∅ => ш

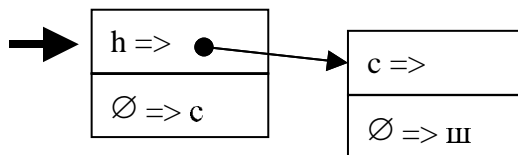
h => шч
∅ => шшч

So each entry either contains a result or a reference to an exception table. If you find a reference to a new exception table, you get the next character and use it to find the entry in the next table. If you don't find the character, you end up on the ∅ entry. This tells you the previous character was the last one in the key and that the character you're on is the beginning of a new key. The result is the result for the key, and you go back and start the lookup over in the main table with the character you're on.

So let's say the input is "shca." You'd read the "s" and look up "s" in the main table.

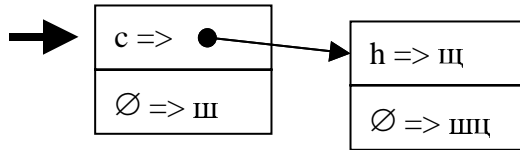


This would refer you to the "s" table. You'd read the "h" and look it up in the "s" table.

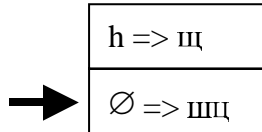


This would take you to the "sh" table. You'd read the "c" from the input and look it up in the "sh" table, taking you to the "shc" table.

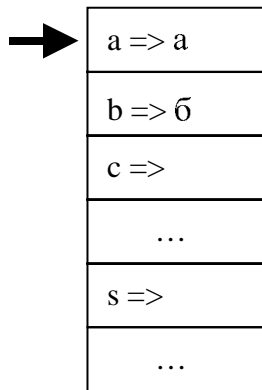
Single versus multiple tables



You read “a” from the input. The “shc” table doesn’t have an entry for “a,” so you use the “∅” entry.

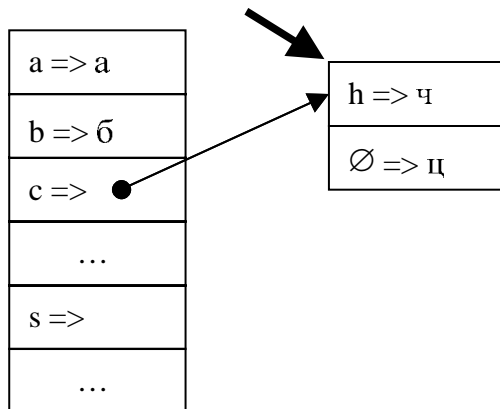


This tells you that “IIII” is the output to use for “shc,” and that you have to go back to the main table and look up “a” again there.



You look up “a” in the main table and get “a” and you’re done.

You only have to look up a character twice if you hit a ∅ entry. Say, for example, the input was “cha.” You’d read “c” from the input, look it up in the main table, and find the reference to the “c” table. You read “h” from the input, look it up in the “c” table, and find “v.”



This is your output, and you've used all the characters you've seen so far. You read "a" from the input and look it up in the main table, where you find "a." This is your output. You've now read all the characters from the input and produced "чa" as your output, which is correct, and you didn't have to look up anything twice.

You can eliminate the \emptyset entries, and the double lookups they cause, by simply having every exception table have an entry for each possible input character (this also lets you use a simple array of some kind, rather than an associative data structure—you basically use the same data structure you're using for the main table for the exception tables—but this will usually get pretty big).

You also don't need them if you don't need key closure. That is, if you never have a key that's a substring of another key, you don't have to do the repeated lookup. For example, if your lookup is UTF-16 based and you see a high surrogate, it's meaningless on its own. You have to look at the next code unit to know what to do. As long as you know you won't have malformed UTF-16 (i.e., unmatched surrogates), you'll never have to do a repeated lookup and you don't need the \emptyset entry.

If you're already experienced in text processing, or you remember back to your data structures class, you've already recognized that the data structure we're talking about here is a trie. A trie is a tree structure where each node in the tree points directly or indirectly to the results for all keys that share a common prefix. In a character-based trie, each level of the tree corresponds to a character in the key, and the number of nodes you traverse to look up the result for a key is equal to the number of characters in the key (actually, as we saw in the previous example, if you can have keys that are substrings of other keys, the number of nodes you visit is actually the number of characters in the key plus one). The word "trie" comes from "information retrieval," but the word "trie" is usually pronounced "try" to distinguish it from "tree."

There are many different ways of representing a trie in memory. In our example, each "exception table" is a node in the trie, so the nodes are represented by key/value arrays that you have to walk linearly. These nodes could be split into linked lists; this makes the overall structure a binary tree. Binary trees work especially well when the data in the trie has to be modified a lot—the array-based representation shown in the examples works more or less the same way as a binary-tree representation but is more compact and useful when the data in the trie won't change. An alternative that gives you faster lookup time and is still good when the trie is being modified is to split each individual node into a binary tree instead of a simple linked list. This turns the overall structure into a *ternary* tree.

You could also expand each node out into a regular flat array with an entry for each possible character. This speeds up lookup by quite a bit (at least when each node has a lot of entries in it), but (generally) at a fairly heavy price in compactness. For a Unicode-based data structure, where each character could possibly be any of more than a million, this obviously isn't practical unless you come up with some kind of compaction scheme for the nodes...

Tries as the main lookup table

...which gets us back to the compact array. You might have noticed that the compact array is essentially a trie itself. You do a lookup based on part of the key, and that takes you to a spot where you do another lookup based on another part of the key, and so on until you've either used up the whole key or gotten to the point where the rest of the key doesn't matter (i.e., where the result will be the same no matter what the rest of the key is). The two things that make the compact array different from a more general trie are:

- Every branch has the same number of levels, so you always know that the result of the first (or, for a two-level compact array, the first and second) lookup is a pointer to the next node of the trie, and that the last lookup returns a real result value.
- Because of this, the chunks of the key that are considered at each node lookup don't have to be the same size. Asymmetrical partitions of the character for lookup in a compact array are actually more the norm than the exception.

This suggests a number of variations on the basic trie idea that might be useful in different situations:

Pruning branches. In a normal compact array, you always do a lookup in the index array followed by a looking into the data array. If the value of every entry in a particular block is the same, you can save space by storing that value directly into the index array and not allocating space for it in the data array. There are a number of things to keep in mind about this technique:

- This technique works best for tables that have a lot of blocks that only consist of one value, but only when those blocks don't all contain the *same* value. If every uniform block contains the same value as the other uniform blocks (i.e., you have a sparsely-populated array where each block is either all zeros or something interesting), this technique only saves you (at most) one block's worth of entries in the data array.
- This technique also works best only when the long stretches of the same value tend to begin and end on block boundaries. If, say, a range of the same value begins halfway into one block, spans three whole blocks, and ends three quarters of the way into another block, pruning the branches only saves you about a quarter of a block's worth of space. (You still need the ends of the long stretch to appear in the data array. You can coalesce them so that you only need as many as are in the longest segment [in our example, three quarters of a block], but since you can do this anyway if the homogeneous blocks are included, you only save the difference between the longest stretch without homogeneous blocks [three quarters of a block] and the longest stretch with them [a whole block].) The savings could be even less if there are several stretches of the same value that don't begin and end on or near block boundaries. You can manage this by fiddling with the block size, but that can introduce other problems.
- If the result values and the pointers to blocks of the data array aren't naturally disjoint, you may have to munge them in some way to make them disjoint, or make the array elements bigger. This could trump any savings in the number of elements.
- The lookup code has to check after each lookup to see whether it's done. This can impose a serious performance penalty on you unless most of your lookups are along the pruned branches.
- Basically, pruning of branches isn't worth the trouble for a simple one-level compact array, unless space is at a premium and you know the data you're trying to compress gets a lot smaller if you prune. It has more utility the more levels the trie has.

Using the same array for the data and index values. This doesn't buy you anything in a single-level compact array. In fact, if the elements in the data array are smaller than the ones in the index, it hurts. As with pruning of branches, it might help with multi-level lookups. In a sparsely-populated two-level compact array, for example, you could double up the zeros in the indexes with the zeros in the data array and save some space. This might be worth it, but generally the win here comes with tables where you traverse a varying number of nodes depending on the key.

Basing lookup on UTF-8 or UTF-16. If you base your lookup on UTF-16 code unit values, you can fix it so that you can do two array lookups for most characters and only have to do a third lookup when you encounter a supplementary-plane character (of course, unless you do branch pruning, you'll actually need *four* lookups when you hit a supplementary-plane character, but since they're pretty rare in most text, this might be a good optimization).

If you do the lookup based on the UTF-8 representation, you can get yourself down to a single lookup when you're dealing with the ASCII characters and get some very nice compaction with little additional cost with the other characters (assuming, of course, that you can either guarantee well-formed UTF-8 or are prepared to get weird results with malformed UTF-8). Because there are a lot of combinations of bytes that don't happen in UTF-8 text, most of the nodes in the tree won't actually need 256 entries in the array, so you get a fair amount of natural compaction (and, potentially, tree pruning). The downside here is that you either have to pay attention to which byte of the character you're on or the result values and pointer values have to be disjoint and you have to pay the cost of checking whether you're done. If the text you're interested in isn't already in UTF-8, you'd also have to waste time converting it (which probably isn't worth the trouble). If the text is already in UTF-8, this technique can be pretty helpful with a lot of sets of data, especially if the majority of the text being processed is in the ASCII range or you can do tree pruning on the branches dealing with Han and Hangul (i.e., get the answer after looking at only the first byte, letting you skip over the other two bytes).

Combining the main lookup table with the multi-character exception table. If you're dealing with variable-length keys, it might make sense to use a unified trie structure for the entire lookup, rather than using a compact array for the first character in the key and some other data structure for the subsequent characters. This can simplify your code and potentially speed it up. But if all of the nodes after the root node are very sparse (i.e., you don't have many multi-character keys, or you don't have many that start with the same character), this can waste a lot of space. If you can safely make assumptions about the input text, this might save space and make this feasible. And if most of the nodes are *not* sparsely populated, this can be a big win.

Single versus multiple tables

One more point worth addressing is whether, if you need to be able to perform several different mappings on Unicode characters, you want to use a completely separate table for each mapping or use one big table for all of them.

It's obviously tough to generalize here, but there are certainly times when a big unified table will save you space and time over a bunch of smaller, more specific tables. The most common case is dealing with the Unicode character properties: Should there be a separate table for each property, or one big one for all of them? The answer is that you can cram most of the Unicode character properties into a single table. The majority of the Unicode properties for a single character can be jammed into a single 32-bit word (the main exception is the canonical and compatibility decompositions, which need their own tables).

The Unicode general category can fit into five bits. As we saw earlier, the bi-di category can fit into four bits, with provision for a few exceptional values. Whether the character is mirrored or not can fit into a single bit. The rest of the bits can be overloaded depending on the character's general category: For non-spacing and enclosing marks, the remaining bits can hold the character's combining class (the combining class is 0 for characters in the other categories). For characters in the "number" categories, the remaining bits can be used for the character's numeric value (some of the numeric values are fractions, and some are especially large, so a bit or two of the numeric-value field have to be set aside as flag bits to indicate which format is being used to encode the actual value into the other bits). And for characters in the cased-letter categories, the other bits can hold the value you add or subtract from the character's code-point value to get its equivalent in the opposite case (or in titlecase). Again, a couple bits may need to be used as flag bits to choose between a number of alternative representations of the case mapping.

Single versus multiple tables

There are, of course, many other ways of storing the data tables that are used to do various things to Unicode characters, but the ones we've looked at in this chapter are some of the most important and useful. In the coming chapters, we'll see how these various techniques are used to perform actual operations on Unicode characters.

CHAPTER 14 *Conversions and Transformations*

Now that we've taken a good look at some of the most common and useful data structures and techniques for doing things with Unicode text, we'll put them to work to do some useful things with Unicode text.

In this chapter, we'll look at performing various types of conversions and transformations on pieces of Unicode text. We'll look not only at converting among the various Unicode representations and between Unicode and other encoding standards, but we'll also look at two other important types of transformations on Unicode text: converting cased text between upper case and lower case, and the more general process known as transliteration, which is the process of converting text from one script to another and is often used as the basis for accepting Unicode text as input (especially when dealing with the Han and Hangul-based languages).

One of the things you'll notice is that most operations on Unicode text can be thought of as transformations of one kind or another on Unicode text. Searching and sorting, for example, involves mapping from Unicode to abstract sort-key values. Rendering text on the screen involves, among other things, mapping Unicode characters to glyph codes. We'll look at these more specialized transformations and the other processes that go along with them in the next chapters. In this chapter, we instead concentrate on Unicode-to-Unicode mappings and Unicode-to-legacy-encoding mapping.

The data structures and techniques we looked at in the last chapter form the basis of what we do here, and will form the basis of much of what we do in the next chapters. But each of these conversions has its own special quirks you need to be aware of. Each makes use of the data structures we looked at in its own unique way or supplements the more general mappings with additional, more specialized work.

Converting between Unicode encoding forms

We'll start by looking at Unicode-to-Unicode transformations. As we've seen, the Unicode standard comprises a single coded character set, but multiple encoding forms:

- UTF-32 represents each 21-bit code point value using a single 32-bit code unit.
- UTF-16 represents each 21-bit code point value using either a single 16-bit code unit (for code points in the BMP) or a pair of 16-bit code units (for code points in the supplementary planes).
- UTF-8 represents each 21-bit code point value with a single 8-bit code unit (for code points in the ASCII block), a sequence of two or three 8-bit code units (for code points in the rest of the BMP), or a sequence of four 8-bit code units (for code points in the supplementary planes).

In this section, we'll look at how you convert text between these three encoding forms. Unicode also comprises seven encoding *schemes* (or serialization formats) which take into account the differing byte-order properties of different machine architectures. Converting between byte orderings is trivial, so we won't cover that here, but we *will* look at implementing SCSU, a common encoding scheme that isn't officially part of the Unicode standard.

Usually, you pick one of the Unicode encoding forms for representing text in your application (or some API you've decided to use effectively makes the choice for you) and stick with it, but environments that make mixed use of multiple encoding forms are not all that uncommon. It's not unusual, for example, to have an application that uses one format for its internal representation (typically UTF-16, but UTF-32 is becoming more common) and another more compact (or, sometimes, more portable) UTF (typically UTF-8) in its data-file format. Or you have a client-server application that uses one format (e.g., UTF-16) on the client side and another (perhaps UTF-32) on the server side. Of you have a client-server application that uses the same format for internal processing on both sides (again, usually UTF-16 or UTF-32) but a different, more compact format (usually UTF-8) in the communication protocol between client and server. Mixed environments like this abound.

For one concrete example, consider the Java runtime environment. The `String` and `StringBuffer` classes in Java (and the `char` primitive type) use UTF-16 to store the characters, but the Java serialization facility (the format that gets written or read when you do `writeObject()` or `readObject()` on an object that implements the `Serializable` interface) uses UTF-8 for serializing strings. The conversion to UTF-8 and back when an object is serialized and deserialized in Java is transparent to the programmer, but it's still useful to be aware of this. In many other mixed environments, this conversion isn't transparent—you have to do it manually.

Pretty much all APIs that support Unicode at all provide the ability to convert a stream of text from one UTF to another, so this is code you rarely have to write yourself. Most systems just treat one UTF as the native character format and use their normal character-code-conversion facility to convert to the others (or to the non-native endian-ness). In Java, for example, the internal character format is UTF-16 in the native endian-ness of the host processor, and conversion to all other UTFs (or UTF-16 in the opposite endian-ness) is handled through the conversion APIs in the `java.io` package (typically accessed through the `Reader` and `Writer` classes).¹¹²

112 The one exception to this is the aforementioned serialization facility, which bypasses the `java.io` code and does the conversion itself for performance.

We talk about conversion to non-Unicode encodings later in the chapter, but even though most systems use the same API for conversion between UTFs that they use for conversion from Unicode to something else, the implementations are usually different. This is because while conversion between Unicode and most non-Unicode encodings has to be table-driven, conversion between UTFs is entirely algorithmic. You can take advantage of this to eliminate the mapping tables entirely or to cut them down greatly in size, which can also improve performance by eliminating time spent loading the table into memory. Table-loading time aside, the actual conversion can usually be made just as fast, or almost as fast, as a table-driven conversion.

Converting between UTF-16 and UTF-32

Conversion between UTF-16 and UTF-32 is fairly simple, since the vast majority of the characters you'll run into in a typical body of text will be BMP characters, which have the same representation in both formats. Only the supplementary-plane characters, which are generally quite rare, need special treatment. You can optimize for the BMP characters and save time.

Conversion from UTF-32 to UTF-16 looks something like this in Java (just to treat all formats equally, I avoid `char` in all these examples):

```
public static short[] utf32to16(int[] in) {
    int size = in.length;
    short[] out = new short[size * 2];

    int p = 0;
    for (int i = 0; i < size; i++) {
        if (in[i] > 0xffff) {
            out[p++] = (short)((in[i] >> 10) + 0xD7C0);
            out[p++] = (short)((in[i] & 0x03FF) + 0xDC00);
        }
        else
            out[p++] = (short)in[i];
    }
    short[] result = new short[p];
    System.arraycopy(out, 0, result, 0, p);
    return result;
}
```

Pretty obvious. You check the code point to see whether it's in the BMP or not and take the appropriate action based on the result. If it's in the BMP, you just truncate it to a `short`. Otherwise, the UTF-32 code unit is a supplementary-plane character and we have to turn it into a surrogate pair. The low surrogate is simple enough: just mask off all but the bottom 10 bits and add `0xDC00`. The high surrogate is a little trickier: You're actually cramming the top 11 bits from the UTF-32 code unit into the bottom 10 bits of the high surrogate. Out of these 11 bits, the bottom six come over unchanged. The top five bits get converted into four by treating them as an integer value and subtracting 1 (this works because 16 [or `0x10`] is the highest value these five bits can have). This gives you a 10-bit value to which you can add `0xD800` to get the high surrogate. In the example, we're doing this by shifting the eleven bits we care about over to the right and adding `0xD7C0`. The `0xD7C0` is the combination of knocking 1 off the top five bits (subtract `0x40`) and adding `0xD800` to the result.

Going back the other way is a little bit uglier because you have to check each UTF-16 code point to make sure it's not a surrogate. There's really no good way to do this except by doing two

comparisons, but if you compare against 0xD800 first, the vast majority of the time this comparison will be true and you won't have to do the second comparison. All that's above 0xDFFF are the private use and compatibility zones, and characters in these zones are relatively infrequent.

So in Java that gives you this:

```
public static int[] utf16to32(short[] in) {
    int size = in.length;
    int[] out = new int[size];

    int p = 0;
    int q = 0;
    while (q < size) {
        int temp = in[q++] & 0xffff;
        if (temp < 0xD800 || temp > 0xDFFF)
            out[p++] = temp;
        else {
            out[p] = (temp - 0xD7C0) << 10;
            out[p++] += in[q++] & 0x03FF;
        }
    }
    if (p == size)
        return out;
    else {
        int[] result = new int[p];
        System.arraycopy(out, 0, result, 0, p);
        return result;
    }
}
```

The actual conversion of the surrogate pair back into the UTF-32 code point value is simply the reverse of the calculations we did in the previous example. No magic here.

Of course, the above example is oversimplified a bit. If you can guarantee that the conversion routine will always be handed well-formed UTF-16 (that is, surrogate values only occur in high-low pairs), or if you're willing to accept slightly garbled results when you get malformed UTF-16, the above code will work fine. Otherwise, you have to do a little more work to account for unpaired surrogates. Traditionally, unpaired surrogates just get converted the same way non-surrogate values do, giving you something that looks like this:

```
if (temp >= 0xD800 && temp <= 0xDBFF && q < in.length
    && (in[q] & 0xffff) >= 0xDC00
    && (in[q] & 0xffff) <= 0xDFFF) {
    out[p] = (temp - 0xD7C0) << 10;
    out[p++] += in[q++] & 0x03FF;
}
else
    out[p++] = temp;
```

With slightly more complicated code, you could instead flag the error by converting unpaired surrogates into U+FFFD, the REPLACEMENT CHARACTER, which by convention is used to indicate the position of an un-convertible character in the source text; eliminate any trace of the unpaired surrogates from the output text; or throw an exception.

The extra code to ensure correctness is kind of ugly and takes a little bit more time, but you're sure you get the right answer. You have to make a tradeoff here between speed and absolute correctness.

Converting between UTF-8 and UTF-32

Going between UTF-32 and UTF-8 is more complicated because UTF-8 is a more complicated encoding standard, with anywhere from one to four code units per character.¹¹³ Here's one way of carrying out that conversion:

```
public static byte[] utf32to8(int[] in) {
    int size = in.length;
    byte[] out = new byte[size * 4];

    int p = 0;
    for (int i = 0; i < size; i++) {
        int c = in[i];
        if (c <= 0x007F)
            out[p++] = (byte)c;
        else {
            if (c <= 0x07FF) {
                out[p++] = (byte)((c >> 6) + 0xC0);
                out[p++] = (byte)((c & 0x3F) + 0x80);
            }
            else {
                if (c <= 0xFFFF) {
                    out[p++] = (byte)((c >> 12) + 0xE0);
                    out[p++] = (byte)(((c >> 6) & 0x3F) + 0x80);
                    out[p++] = (byte)((c & 0x3F) + 0x80);
                }
                else {
                    out[p++] = (byte)((c >> 18) + 0xF0);
                    out[p++] = (byte)(((c >> 12) & 0x3F) + 0x80);
                    out[p++] = (byte)(((c >> 6) & 0x3F) + 0x80);
                    out[p++] = (byte)((c & 0x3F) + 0x80);
                }
            }
        }
    }
    byte[] result = new byte[p];
    System.arraycopy(out, 0, result, 0, p);
    return result;
}
```

The code's a little complicated, but you can't really improve on it much in terms of performance. The one thing that kind of sucks about this approach is that you do three comparisons for each character for code point values above U+07FF. This includes the Han and Hangul areas and all the Indic

113 Earlier versions of ISO 10646 specified a version of UTF-8 that used up to six code units per character, but more recent decisions by the committees maintaining both Unicode and ISO 10646 have declared the encoding space above U+10FFFF off limits. Five- and six-byte sequences are only necessary to represent values in the area that's now off-limits for encoding; therefore, five- and six-byte UTF-8 are now illegal.

scripts, as well as a smattering of other scripts and all of the supplementary planes. For many scripts, including most of the stuff in the supplementary planes, they occur rarely enough not to be a problem. But it's a fairly heavy price to pay for, say, whole documents in Chinese, or Japanese, or Hindi, or Thai.

If you're willing to take up some space for a translation table, you can turn the series of `if` statements into a `switch` statement, which effectively gives you one comparison for all the code point ranges. You can do this with a 256-byte lookup table by using the leading byte of the UTF-16 code point value as an index into the table. This approach, for obvious reasons, works much better when you're going straight from UTF-16 to UTF-8, but since you have to do some extra work to get from a UTF-16 surrogate pair to a four-byte UTF-8 sequence, you lose a little speed with the supplementary-plane characters in exchange for gaining some on the BMP characters.

Going back from UTF-8 to UTF-32 is more interesting. The obvious straightforward approach is something like this:

```
public static int[] utf8to32(byte[] in) {
    int size = in.length;
    int[] out = new int[size];

    int p = 0;
    int q = 0;

    while (q < size) {
        int c = in[q++] & 0xff;

        if (c <= 0x7F)
            out[p++] = c;
        else if (c <= 0xDF) {
            out[p] = (int)(c & 0x1F) << 6;
            out[p++] += (int)(in[q++] & 0x3F);
        }
        else if (c <= 0xEF) {
            out[p] = (int)(c & 0x0F) << 12;
            out[p] += (int)(in[q++] & 0x3F) << 6;
            out[p++] += (int)(in[q++] & 0x3F);
        }
        else if (c <= 0xF7) {
            out[p] = (int)(c & 0x07) << 18;
            out[p] += (int)(in[q++] & 0x3F) << 12;
            out[p] += (int)(in[q++] & 0x3F) << 6;
            out[p++] += (int)(in[q++] & 0x3F);
        }
        else
            out[p] = 0xFFFFD;
    }

    int[] result = new int[p];
    System.arraycopy(out, 0, result, 0, p);
    return result;
}
```

Again, this is straightforward and gets the job done and performs reasonably well. Again, it doesn't check for malformed UTF-8 sequences (except for byte values from 0xF8 to 0xFF, which can't occur in legal UTF-8). In the presence of malformed UTF-8, this algorithm will produce garbled results, so it's only acceptable if that's acceptable, or if you can guarantee that the UTF-8 fed to this algorithm is well-formed.

Again, if you're willing to pay a little extra to keep some tables around, you can make the algorithm faster and more elegant and also have it handle malformed UTF-8 more gracefully. The approach shown below is based on a 128-byte lookup table.

```
private static final byte X = (byte)-1; // illegal lead byte
private static final byte Y = (byte)-2; // illegal trail byte
private static final byte[][] states = {
    //00 08 10 18 20 28 30 38 40 48 50 58 60 68 70 78
    // 80 88 90 98 A0 A8 B0 B8 C0 C8 D0 D8 E0 E8 F0 F8
    { 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      X, X, X, X, X, X, X, X, 1, 1, 1, 1, 2, 2, 3, X },
    { Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y,
      0, 0, 0, 0, 0, 0, 0, 0, Y, Y, Y, Y, Y, Y, Y, Y },
    { Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y,
      1, 1, 1, 1, 1, 1, 1, 1, Y, Y, Y, Y, Y, Y, Y, Y },
    { Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y, Y,
      2, 2, 2, 2, 2, 2, 2, 2, Y, Y, Y, Y, Y, Y, Y, Y }
};
private static final byte[] masks = { 0x7F, 0x1F, 0x0F, 0x07 };

public static int[] utf8to32b(byte[] in) {
    int size = in.length;
    int[] out = new int[size];

    int p = 0;
    int q = 0;
    int state = 0;
    byte mask = 0;

    while (q < size) {
        int c = in[q++] & 0xff;

        state = states[state][c >>> 3];

        switch (state) {
            case 0:
                out[p++] += c & 0x7F;
                mask = 0;
                break;

            case 1:
            case 2:
            case 3:
                if (mask == 0)
                    mask = masks[state];
                out[p] += c & mask;
                out[p] <<= 6;
        }
    }
}
```

```
        mask = (byte)0x3F;
        break;

    case Y:
        --q;
        // fall thru on purpose

    case X:
        out[p++] = 0xFFFFD;
        state = 0;
        mask = 0;
        break;
    }
}

int[] result = new int[p];
System.arraycopy(out, 0, result, 0, p);
return result;
}
```

We can eliminate the chain of `if` statements and some of the redundant code in the first example by using a state table to keep track of where you are in a multi-byte character. The state table also gives us a convenient way to detect malformed sequences of UTF-8. This example actually distinguishes between longer-than-normal and short-than-normal sequences. If it sees a leading byte when it's expecting a trailing byte, (a shorter-than-normal sequence, represented by `Y` in the state table) it throws out the incomplete sequence it's seen up until then, writes `U+FFFFD` to the output as a placeholder for the corrupted character, and starts over with the leading byte it encountered prematurely (i.e., it treats the premature leading byte as good and the truncated sequence that precedes it as bad). If, on the other hand, it sees a trailing byte where it expects a leading byte (a longer-than-normal sequence, denoted by `X` in the state table), or a byte which is illegal in UTF-8, it treats that byte as an illegal sequence all by itself.

Even here, though, we're not protecting against non-shortest-form UTF-8. An ASCII character represented with three bytes will come through just fine with either of these converters. You can actually guard against these illegal combinations without altering the basic approach by making the state table a little more complicated. I'll leave this as an exercise for the reader.

As always, many variations on this basic approach are possible. For example, you can eliminate the right shift in the state-table lookup by making the table bigger, and you can eliminate some of the manipulation of the `mask` variable (and the `if` statement based on it) by just looking up the mask every time. This also requires a larger table of masks. You can fetch both the mask value and the new-state value in a single array lookup by munging the values together mathematically in the table (for example, we know we only need six state values, and the bottom three bits of every mask value are 1, so we could use the bottom three bits for the state and the rest of the value for the mask), but then you need extra code to take the two values apart.¹¹⁴

114 One thing to keep in mind if (God forbid) you try to use my code examples verbatim in your code: This particular example won't compile as-is. Integer literals in Java have type `int`, and can't be used in an array literal of type `byte[]` without casts. I left these casts out to make the table more readable.

Converting between UTF-8 and UTF-16

Conversions between UTF-8 and UTF-16 are left as an exercise for the reader. Conceptually, you really can't get from UTF-8 to UTF-16 (or vice versa) without going through UTF-32. You can cut out some of the work that'd go on if you just, say, called the 8-to-32 function followed by the 32-to-16 function, but you can't cut out much of it.

Implementing Unicode compression

The Unicode compression scheme (aka the Standard Compression Scheme for Unicode, or SCSU, which is documented in Unicode Technical Standard #6), isn't officially a Unicode Transformation Format, but sort of informally qualifies, as it's yet another way of translating a sequence of abstract Unicode numeric values into patterns of bits. Unlike UTF-16 and UTF-32 (and, to a lesser extent, UTF-8), it isn't really suitable as an in-memory representation of Unicode. This is because it's a stateful encoding—particular byte values mean different things depending on the bytes that have come before. This is arguably true of trailing bytes in UTF-8 or low surrogates in UTF-16, as their exact meanings depend on the preceding lead bytes or high surrogates, but there are two important differences: 1) In both UTF-8 and UTF-16, leading code units, trailing code units, and single-unit characters (code units that represent a whole character by themselves) are represented by distinct numerical ranges—you can tell which category a particular code unit is in just by looking at it. In UTF-8, the leading bytes of different-length sequences also occupy different numerical ranges—you can tell how long a UTF-8 sequence is supposed to be just from looking at the leading byte. 2) When you do see a code unit that represents only part of a character, the number of units forward or backward you have to scan to get the whole character is strictly limited. (If you happen upon a non-surrogate code unit in UTF-16 or an ASCII value in UTF-8, the code unit stands on its own; if you happen upon a surrogate in UTF-16 or a leading byte in UTF-8, you can tell from it exactly how many code units you have to scan and in which direction to get the rest of the character; if you happen upon a trailing byte in UTF-8, you have to scan backwards no more than three bytes to find the beginning of the character, which will then tell you where the end of the character is.)

These two things aren't true in SCSU: A decompressor for SCSU maintains a “state” which can be changed by each byte as it's read: the current state determines the meaning of the next byte you see (which may, in turn, change the state again). You can't drop into the middle of an SCSU byte stream and tell what's going on (i.e., correctly interpret any arbitrary byte in an SCSU byte stream) without potentially scanning all the way back to the beginning of the file to figure out what state you're in. You don't *always* have to scan all the way back to the beginning of the file, and you can often make an educated guess as to what certain bytes mean, but there's nothing limiting your maximum look-back: you *may* have to scan all the way back to the beginning of the file.

Instead, SCSU is useful as a serialization format: as you read the text into memory (or otherwise deserialize it), you decompress it, converting it into one of the other representations for use in memory. You do the opposite when writing the text back out. UTF-8 is often used for this purpose because it's backward compatible with ASCII and because the representation is more compact than UTF-16 for Latin. But UTF-8 actually is *less* compact than UTF-16 for many scripts, especially the East Asian scripts. A Japanese document is 50% larger in UTF-8 than it is in UTF-16 (or, generally speaking, a legacy Japanese encoding). SCSU solve this problem: It's somewhat more expensive to convert to and from than UTF-8 and can't be used as an in-memory representation like UTF-8 can, but it's also backward compatible with ASCII and is more compact (or equally as compact as) both UTF-8 and UTF-16 for all scripts. In most cases, SCSU text is almost as compact as a legacy encoding for a given script. SCSU also lends itself well to further compression with a more-general compression scheme such as LZW. And on top of this, it's also backward compatible with Latin-1, and not just with ASCII.

The International Components for Unicode library includes conversion utilities for SCSU, but most other libraries don't. If you're not using ICU and want to use Unicode compression, you may well need to code it yourself. Here we'll look at how to do that.

SCSU allows for many different ways of encoding the same piece of Unicode text: unlike, say, UTF-8, it doesn't prescribe one particular alternative as the "canonical representation" for a particular piece of text. This means that even though SCSU is always decoded in the same way, there are many different ways in which Unicode text can be *encoded* in SCSU. A particular encoder can choose to optimize for speed or compactness.

The basic idea is this: SCSU has two "modes," single-byte mode and double-byte mode. In single-byte mode:

- The byte values from 0x00 to 0x1F, except for 0x00, 0x09, 0x0A, and 0x0D (the ASCII null, tab, line-feed, and carriage-return characters, respectively), are "tag" bytes which control the interpretation of the other byte values.
- The other byte values from 0x00 to 0x7F are used to represent the bottom seven bits of a Unicode code point, with the upper fourteen bits being determined by the decoder's state. There are eight "static windows," that is, eight ranges of 128 Unicode code point values that can be represented by these byte values. Normally they represent the characters in "static window 0," which is the Unicode code point values from U+0000 to U+007F, but they can be temporarily shifted to represent one of the other seven windows.
- The byte values from 0x80 to 0xFF are also used to represent the bottom seven bits of a Unicode code point, with the upper fourteen bits being determined by the decoder's state. There are eight "dynamic windows," that is, eight ranges of 128 Unicode code point values that can be represented at any one time. Unlike the static windows, however, these eight dynamic windows can be repositioned almost anywhere in the Unicode encoding range.

In double-byte mode, the bytes are interpreted in pairs as big-endian UTF-16, except that the byte values from 0xE0 to 0xF2, when in the lead-byte position (corresponding to the high-order byte of code point values in the private-use area), are generally interpreted as tag values.

The tag values do various things:

- Switch to single-byte or double-byte mode long enough to encode one Unicode code point.
- Switch to single-byte or double-byte mode indefinitely.
- Interpret just the next byte (in single-byte mode) according to some particular static or dynamic window (this changes the interpretation of all byte values in single-byte mode).
- From now until it's changed again, interpret the byte values from 0x80 to 0xFF according to some particular dynamic window (the values from 0x00 to 0x7F still get treated as either tag values or members of static window 0—you can't indefinitely switch the values from 0x00 to 0x7F to some other static window).
- Reposition a particular dynamic window to refer to a different range of values in the Unicode encoding space.

The initial positions of the windows mean text containing only characters from the ASCII and Latin-1 blocks has the same representation in SCSU as it does in Latin-1, unless it contains control characters other than the tab, line feed, and carriage return. (Note that this is different from UTF-8, where all the control characters come through without change, but the characters in the Latin-1 block turn into two-byte sequences—like the East Asian scripts, Latin-1 text is actually better-served by SCSU than by UTF-8.) Text in many other scripts can be represented in SCSU entirely by one-byte values as well, with a single extra byte at the beginning to shift the values from 0x80 to 0xFF to the appropriate

window. For rarer scripts that fit into a 128-code-point window, you need two or possibly three bytes at the front of the file, but can otherwise represent everything with a single byte per character.

Scripts that span more than 128 code point values can often still be represented fairly well with single-byte SCSU, with more tag bytes scattered throughout to shift the windows around. For those that occupy too wide a span of code point values for the shifting-window technique to work well, you can simply shift to Unicode mode—the representation is now simply big-endian UTF-16 with an extra byte on the front to tag it as such.

The code to do all this isn't too ridiculous, although it's more complicated than converting out of UTF-8. Here's one fairly straightforward interpretation (again, in Java, albeit somewhat abbreviated):

```
// single-byte-mode tag values:
//=====
// shift to specified window for one byte
private static final int SQ0 = 0x01;
// SQ1-SQ6 omitted for brevity
private static final int SQ7 = 0x08;

// define a new window using extended semantics and shift to it
private static final int SDX = 0x0B;

// shift to double-byte mode for one pair of bytes
private static final int SQU = 0x0E;

// shift to double-int mode
private static final int SCU = 0x0F;

// shift to specified window
private static final int SC0 = 0x10;
// SC1-SC6 omitted for brevity
private static final int SC7 = 0x17;

// define a new window and shift to it
private static final int SD0 = 0x18;
// SD1-SD6 omitted for brevity
private static final int SD7 = 0x1F;

// double-byte-mode tag values:
//=====
// shift to single-byte mode using specified window
private static final int UC0 = 0xE0;
// UC1-UC7 omitted for brevity

// define a new window and shift to it (and to single-byte mode)
private static final int UD0 = 0xE8;
// UD1-UD7 omitted for brevity

// quote a single Unicode character (can start with tag byte)
private static final int UQU = 0xF0;

// define a new window using extended semantics and shift to it
```

```
private static final int UDX = 0xF1;

private static final int[] staticWindows = {
    0x00, 0x00, 0x0100, 0x0300, 0x2000, 0x2080, 0x2100, 0x3000 };
private static final int[] initialDynamicWindows = {
    0x80, 0xC0, 0x0400, 0x0600, 0x0900, 0x3040, 0x30A0, 0xFF00 };
private static final int[] dynamicWindowOffsets = {
    0x0000, 0x0080, 0x0100, 0x0180, 0x0200, 0x0280, 0x0300, 0x0380,
    0x0400, 0x0480, 0x0500, 0x0580, 0x0600, 0x0680, 0x0700, 0x0780,
    // ...and so on through...
    0x3000, 0x3080, 0x3100, 0x3180, 0x3200, 0x3280, 0x3300, 0x3380,
    0xE000, 0xE080, 0xE100, 0xE180, 0xE200, 0xE280, 0xE300, 0xE380,
    // ...and so on through...
    0xFC00, 0xFC80, 0xFD00, 0xFD80, 0xFE00, 0xFE80, 0xFF00, 0xFF80,
    0x0000, 0x0000, 0x0000, 0x0000, 0x0000, 0x0000, 0x0000, 0x0000,
    // ...nine more rows of zeros...
    0x0000, 0x00C0, 0x0250, 0x0370, 0x0530, 0x3040, 0x30A0, 0xFF60
};

public static void decompressSCSU(
    InputStream in, OutputStream out)
    throws IOException {
    boolean singleByteMode = true;
    boolean nonLockingShift = false;
    boolean nonLockingUnicodeShift = false;
    int currentWindow = 0;
    int tempWindow = 0;
    int[] dynamicWindows = (int[])initialDynamicWindows.clone();

    int b = in.read();
    while (b != -1) {
        while (b != -1 && singleByteMode
            && !nonLockingUnicodeShift) {
            switch (b) {
                case SQU:
                    nonLockingUnicodeShift = true;
                    break;

                case SCU:
                    singleByteMode = false;
                    break;

                case SQ0: case SQ1: case SQ2: case SQ3:
                case SQ4: case SQ5: case SQ6: case SQ7:
                    nonLockingShift = true;
                    tempWindow = b - SQ0;
                    break;

                case SC0: case SC1: case SC2: case SC3:
                case SC4: case SC5: case SC6: case SC7:
                    currentWindow = b - SC0;
                    break;

                case SD0: case SD1: case SD2: case SD3:
```

```
case SD4: case SD5: case SD6: case SD7:
    currentWindow = b - SD0;
    b = in.read();
    dynamicWindows[currentWindow] =
        dynamicWindowOffsets[b];
    break;

case SDX:
    int hi = in.read();
    int lo = in.read();
    currentWindow = hi >>> 5;
    dynamicWindows[currentWindow] =
        0xD80DC00 + ((hi & 0x1F) << 21)
        + ((lo & 0xF8) << 16)
        + ((lo & 0x07) << 7);
    break;

case 0x80: case 0x81: case 0x82: case 0x83:
// ...plus the other values from 0x80 to 0xff
    if (nonLockingShift) {
        writeChar(out,
            dynamicWindows[tempWindow]
                + (b - 0x80));
        nonLockingShift = false;
    }
    else
        writeChar(out,
            dynamicWindows[currentWindow]
                + (b - 0x80));
    break;

default:
    if (nonLockingShift) {
        writeChar(out,
            staticWindows[tempWindow] + b);
        nonLockingShift = false;
    }
    else
        writeChar(out, b);
    break;
}
b = in.read();
}

while (b != -1 && (!singleByteMode
    || nonLockingUnicodeShift)) {
    if (nonLockingUnicodeShift || b < 0xE0 || b > 0xF0) {
        writeChar(out, (b << 8) + in.read());
        nonLockingUnicodeShift = false;
    }
    else {
        switch (b) {
            case UQU:
```

```
        nonLockingUnicodeShift = true;
        break;

    case UC0: case UC1: case UC2: case UC3:
    case UC4: case UC5: case UC6: case UC7:
        singleByteMode = true;
        currentWindow = b - UC0;
        break;

    case UD0: case UD1: case UD2: case UD3:
    case UD4: case UD5: case UD6: case UD7:
        singleByteMode = true;
        currentWindow = b - UD0;
        b = in.read();
        dynamicWindows[currentWindow] =
            dynamicWindowOffsets[b];

    case UDX:
        singleByteMode = true;
        int hi = in.read();
        int lo = in.read();
        currentWindow = hi >>> 5;
        dynamicWindows[currentWindow] =
            0xD800DC00 + ((hi & 0x1F) << 21)
            + ((lo & 0xF8) << 16)
            + ((lo & 0x07) << 7);
    }
    }
    b = in.read();
}
}
}

private static void writeChar(OutputStream out, int c)
    throws IOException {
    if (c < 0x10000) {
        out.write(c >> 8);
        out.write(c & 0xff);
    }
    else {
        out.write(c >> 24);
        out.write((c >> 16) & 0xff);
        out.write((c >> 8) & 0xff);
        out.write(c & 0xff);
    }
}
```

It's fairly long, but reasonably simple and efficient. This particular version writes UTF-16 to the output, since this is what you get when you're in double-byte mode. The window-offset arrays are arrays of `int`, but their elements aren't UTF-32 code points; they're either single UTF-16 code points padded out to 32 bits (pretty much the same thing), or they're UTF-16 surrogate pairs packed into a single `int`. The helper function `writeChar()` takes care of accounting for the second case. You could easily modify this routine to output UTF-32 instead: the window-offset arrays would then

be UTF-32 values, the code that handles the SDX and UDX characters would have to change, and the double-byte mode would have to convert UTF-16 surrogate pairs into UTF-32.

The example above also assumes it's looking at well-formed SCSU and will act funny or throw exceptions if the text it's reading isn't well formed.

There are a bunch of ways of producing SCSU in the first place. A very simple algorithm that optimizes for Latin-1 text would look like this:

- If the character is a Latin-1 non-control character (or if it's tab, CR, or LF) and we're in single-byte mode, just truncate it to a byte and write it to the output.
- If the character is a control character other than the ones listed above, write SQ0 to the output, followed by the character, truncated to a byte.
- If the character is not a Latin-1 character, but is a BMP character, and the character after it *is* a Latin-1 character, write SQU to the output, followed by the character's big-endian UTF-16 representation.
- If the character is not a BMP character, or it's not a Latin-1 character and neither is the following character, and we're still in single-byte mode, write SCU, shift to double-byte mode, and reprocess that character.
- If we're in double-byte mode and the character is a not a Latin-1 character or a private-use character, write its big-endian UTF-16 representation to the output.
- If we're in double-byte mode and the character is a private-use character, write UQU to the output, followed by the character's big-endian UTF-16 representation.
- If we're in double-byte mode and the character is a Latin-1 character, but the character following it is not, write its big-endian UTF-16 representation to the output.
- If we're in double-byte mode, the character is a Latin-1 character, and so is the following character, write UC0, shift back to single-byte mode, and reprocess the character.

A more complicated algorithm that tries to take into account all the Unicode characters might look something like this:

- Start out in single-byte mode, with the windows set up as they would be in the initial state of the decoder.
- If the current character is an ASCII printing character, or is tab, LF, or CR, truncate it to a byte and write it to the output.
- Otherwise, if the character fits into the current dynamic window (as defined by the last SCx we wrote, or the Latin-1 block if we haven't written any SCx codes), write the bottom seven bits of the code point value, plus 0x80, to the output.
- Otherwise, if the character fits into one of the static windows, or one of the currently-defined dynamic windows, but the next two characters don't fit into the same window, write the appropriate SQx code to the output, followed by the appropriate byte value (if the character fit into a static window, this is the bottom seven bits of the code point value; if it fit into a dynamic window, it's the bottom seven bits plus 0x80).
- If the current character doesn't fit into any of the currently-defined dynamic or static windows, and the next two characters don't fit into the same window, write SQU to the output, followed by the character's big-endian UTF-16 representation.
- If the next three characters in a row fit into the same dynamic window and it's not the current dynamic window but is one of the currently-defined dynamic windows, write the appropriate SCx code to the output, followed by the bottom seven bits of the current character plus 0x80.

- If the next three characters in a row fit into the same window but it's not one of the currently-defined dynamic windows, emit an appropriate SDx code, followed by the appropriate defining byte (or, if necessary, SDX, followed by the appropriate defining two bytes) and the bottom 7 bits of the current character plus 0x80. [It'd be best to use some sort of "least-recently-used" algorithm to determine which window to redefine—a more sophisticated algorithm would also ignore ASCII printing characters in determining whether "the next three characters" fit into the same window.]
- If the next three characters in a row don't fit into the same window, write SCU to the output, switch to double-byte mode, and reprocess the current character.
- If you're in double-byte mode, write characters out in their big-endian UTF-16 representation (preceding private-use characters with UQU) for as long as you don't encounter a sequence of three or more characters that either fit into the same window or are ASCII printing characters. If you see such a sequence, write the appropriate code to switch back to single-byte mode (UCx, UDx, or UDX, with the necessary defining bytes, determined according to the same criteria as we used in single-byte mode), switch back to single-byte mode, and process the sequence in single-byte mode.

Unicode normalization

The other category of Unicode-to-Unicode transformations we need to consider is Unicode normalization. You may or may not need to care about Unicode normalization as such; often it happens as a byproduct of some other process. For example, character code conversion is usually written to produce normalized Unicode text; a legacy encoding usually naturally converts to one normalized form or the other. (Of course, this may imply that not all converters available on a particular system produce the *same* normalized form, unless they've all been designed to do this.) Similarly, keyboard layouts and input methods can be designed to produce normalized text (and usually do so without any extra work, although again this might mean that different input methods and keyboard layouts on a given system produce *different* normalization forms, unless they were expressly designed to all produce the same one).

If your text is all coming from an outside process, you may be able to depend on the process that originally produced it to produce it in some normalized form. (For example, the World Wide Web Consortium's early-normalization rules require that every process on the Web that produces text produce it in Normalized Form C, so processes downstream of process that produced the text (XML parsers, for example) can assume the text is already in Form C.

Finally, your application might not care whether the text is normalized or not. Generally, Unicode normalization gives you two things: You can compare two strings for equality using a straight bitwise comparison, and rendering processes don't have to account for alternative representations of characters they can display. If you're already doing language-sensitive comparison rather than bitwise comparison of your strings, you don't have to do Unicode normalization, since language-sensitive comparison has to do normalization anyway (or take the alternate representations of a given string into account). If the rendering process you're using can handle the alternative representations of a given string (or if you don't care), you also don't have to do Unicode normalization.

If you actually *do* have to do Unicode normalization, there's again a fairly high probability you'll have to code it yourself. As with Unicode compression, the International Components for Unicode provide an API for normalization, but most other Unicode-support APIs don't.

There are four Unicode normalization forms, but only three processes that need to be implemented to support them:

- **Normalized Form D:** Canonical decomposition.
- **Normalized Form C:** Canonical decomposition, followed by canonical composition.
- **Normalized Form KD:** Compatibility decomposition.
- **Normalized Form KC:** Compatibility decomposition, followed by canonical composition.

The original normalized form for Unicode was Normalized Form D, which consists of canonical decompositions. Normalized Form KD, compatibility decomposition, is there to allow text to be mapped from characters that are only there for round-trip compatibility with some other standard to non-compatibility characters (although, since many of the compatibility composites involve some kind of stylistic transform on the original character, information can be lost going to Form KD).

Normalized Forms C and KC evolved more recently. In real life, a lot of implementations were preferring composed representations of characters over the decomposed representations because the composed representations were more compact and because many of the legacy encodings (most of the ISO 8859 family of encodings, for example) use composed representations.

Converting to Forms C and KC involves both a decomposition and a recomposition mainly because of issues related to the order of combining marks. Say, for example, you have in your text an “a” with a macron on top and a dot underneath. There’s no code point assigned to this character, so you have to represent it with a combining character sequence. There *are*, however, precomposed characters for both a-macron and a-underdot, so you could represent a-macron-underdot either with a-macron and a combining underdot or with a-underdot and a combining macron. As it turns out, the correct representation in Normalized Form C is a-underdot plus the combining macron (because of the canonical accent order). If you have text with the other representation (a-macron plus the combining underdot), you can only get to the correct representation by fully decomposing it, reordering the accents, and *then* composing it again.

There isn’t any such thing as “compatibility composition.” Normalized Form KC uses compatibility decomposition as the basis for canonical composition instead of using canonical decomposition as the basis for canonical composition.

Canonical decomposition

Canonical decomposition actually consists of two separate processes: *canonical decomposition* and *canonical reordering*. Canonical decomposition is the process of taking any canonical composites in the input stream and replacing them with their decomposed forms. Canonical reordering is the process of taking any combining character sequences and sorting them according to their combining classes. When you decompose a composite character, the decomposed representation will already be ordered by combining class; canonical reordering is oriented toward getting combining character sequences *already in* the input stream, not those produced by decomposing composite characters, into canonical order.

Canonical decomposition. There are exactly two types of canonical decompositions: those that consist of two characters and those that consist of only one. A character with a one-character (“singleton”) decomposition is a character whose use is effectively discouraged—it’s not supposed to show up in normalized Unicode text, regardless of normalization form. Two-character canonical

decompositions, on the other hand, represent “real” composite characters that can be represented with combining character sequences.

There are actually cases where a character should decompose to a *three*-character (or longer) combining sequence (for example, \bar{u} should decompose to small u, combining diaeresis, combining macron). The standard is structured so that characters like this decompose to a combining sequence that contains *another* composite character (in this case, to \ddot{u} , combining macron)—this means you may actually have to perform several mappings to fully decompose a character. The decompositions are set up so that the composite character is always the first character in the decomposition—each iteration of the mapping effectively peels one mark off the end:

$$\begin{aligned} &\bar{u} \\ &\ddot{u} + \bar{\cdot} \\ &u + \ddot{\cdot} + \bar{\cdot} \end{aligned}$$

You can take advantage of the two-character limit in constructing a lookup table to use in decomposition. Since there are supplementary-plane characters that are canonical composites, you have to use either a two-level compact array or some kind of more generalized trie structure for your lookup table. You can have the data array consist of 32-bit values: For the two-character decompositions, the 32-bit value is the UTF-16 representations of the two characters packed into a single 32-bit value. For the singleton decompositions, the upper 16 bits are zero and the lower 16 bits are the UTF-16 representation of the character to map to.

Canonical decomposition on supplementary-plane characters. There’s one hitch in doing canonical decomposition this way: some of the supplementary-plane characters have decompositions that include other supplementary-plane characters. In Unicode 3.1, these fall into two main categories:

- Musical symbols that decompose to sequences of other musical symbols (for example, the eighth note decomposes to a black notehead, followed by a combining stem, followed by a combining flag).
- Compatibility ideographs that map to non-compatibility ideographs in Plane 2 (there are also many compatibility ideographs in Plane 2 that map to ideographs in the BMP).

There are a couple different ways you can handle this:

- You can assume that there are no BMP characters that include supplementary-plane characters in their decompositions (this assumption holds now and is likely to hold in the future, but might be invalidated by some future version of Unicode). You could use a single-level compact array with the packing scheme described above to handle the BMP characters and a separate lookup table for the non-BMP characters (or a separate table for each plane). For the Plane 1 characters (currently the musical symbols), you could code things assuming all of the characters in the decomposition are also in Plane 1 and use the packing scheme above. For the Plane 2 characters, you could assume the decompositions are all one-character decompositions and store the decompositions using UTF-32.
- You could use a different lookup table for non-BMP characters and just have the entries be 64 bits wide, packing together two UTF-32 code units.
- You could use a unified lookup table for the whole Unicode range using the format described above, but use a range of sentinel values (say, 32-bit values where the top 16 bits, treated as a UTF-16 value, would be in the surrogate range) to refer to an exception table containing strings. This approach makes a lot of sense if you’re storing compatibility decompositions in the same table (because you have to do this anyway—see below). If you do it this way, you don’t want to do the recursive lookup, though: just have the exception table contain the full decompositions of

all the characters. Since many of the decompositions in the musical-symbols block are suffixes of other decompositions, you can actually get fairly decent compression.

One thing you really can't do is ignore the supplementary-plane characters. Unlike most other processes on Unicode text, where you get to pick which characters you're going to support, you can't really pick and choose with Unicode normalization. That is, you don't have the option of normalizing some characters and leaving other characters alone. If you're producing text that will be (or may be) read by another process, and it purports to be in some normalization form, you can't cut corners. You can only get away with blowing off some of the mappings if you're normalizing for internal use and it doesn't have to be in an official normalized form, if you're willing to output error characters or throw an exception if you encounter characters you can't handle (not really an option if you're operating on text from an outside source and purporting not to mess it up other than to normalize it, since this isn't legal according to the Unicode standard), or if you can guarantee that certain characters will never appear in the text you're trying to normalize. If you can't guarantee your input character repertoire, and you're sending the normalized text on to another process, you have to handle all the characters.

Canonical decomposition on Hangul characters. There's one important exception to the rule that canonical decompositions only have one or two characters, and that's decomposition of Hangul syllables. Hangul syllables all have a canonical decomposition to the equivalent sequence of conjoining Hangul jamo characters, and a Hangul syllable can decompose to as many as three jamo.

But the nice thing about Hangul is that you can decompose Hangul syllables algorithmically; you don't need a table lookup. This means all you have to store in the decomposition table for Hangul is a sentinel value indicating the character is a Hangul syllable and needs special treatment, or you check if the character is a Hangul syllable before doing the table lookup in the first place.

The code to decompose a Hangul syllable looks something like this¹¹⁵:

```
// beginning of Hangul-syllables block, number of syllables
private static final char SBase = '\uAC00';
private static int SCount = 11172;

// beginning of leading consonants, number of leading consonants
private static final char LBase = '\u1100';
private static final int LCount = 19;

// beginning of vowels, number of vowels
private static final char VBase = '\u1161';
private static final int VCount = 21;

// beginning of trailing consonants, number of trailing consonants
private static final char TBase = '\u11A7';
private static final int TCount = 28;

// number of syllables with a given initial consonant
private static final int NCount = VCount * TCount;

public static char[] decomposeHangulSyllable(char syl) {
    int sIndex = syl - SBase;
    int l = LBase + (sIndex / NCount);
```

115 This example is taken almost verbatim out of the Unicode standard, pp. 54-55.

```
int v = VBase + (sIndex % NCount) / TCount;
int t = TBase + (sIndex % TCount);

if (t == TBase)
    return new char[] { (char)l, (char)v };
else
    return new char[] { (char)l, (char)v, (char)t };
}
```

Canonical reordering. Converting to Normalized Form D involved two conceptual passes. The first is to go through the whole string and replace every canonical composite with its canonical decomposition. You repeat this until the string is made up entirely of non-composite characters.

Once you've done that, you make another pass through the string to perform canonical reordering. You go through the string and locate every sequence of two or more characters in a row that have a non-zero combining class. (This generally means "every sequence of combining marks," although there are some combining marks that have combining class 0 specifically so things won't reorder around them—the enclosing marks, for example, inhibit reordering. On the other hand, all characters that *aren't* combining marks *do* have a combining class of 0.)

For each sequence of combining marks, sort it in ascending order by combining class, being careful not to disturb the relative order of any characters in the sequence that have the *same* combining class. Since you're only extremely rarely dealing with sequences of more than two or three marks, anything more than a simple bubble or insertion sort would be overkill.

Of course, this means you have to have a lookup table somewhere that tells you each character's combining class. Often, the libraries you're working with will have an API you can call to get a character's combining class. If not, you'll have to create a lookup table to do this. The compact array or another trie-based structure lends itself well to this: most characters have a combining class of zero, there are a number of large blocks of characters have the same (non-zero) combining class, and the combining class is a byte value. Because of its size (and relative infrequency), the combining class can easily be doubled up with other Unicode properties in a single lookup table, if you also need to look up the other properties. See Chapter 13.

Putting it all together. You probably don't want to just go through the whole string twice to convert it into Normalized Form D. In fact, you're very often doing normalization as part of another process (such as string comparison) where you don't really want to store the normalized string at all.

Unfortunately, you can't just consider one character at a time when doing canonical decomposition. You'll often have to look at multiple characters in the input to figure out what the next output should be, and you'll usually end up generating several output characters at one. This generally means that a character-by-character interface will require both input and output buffering. Here's a pseudocode summary of how an object that returns the canonically-decomposed version of a string one character at a time (and, in turn, gets the original string to be normalized one character at a time) might go about it:

- If there are characters in the output buffer, return the first character in the output buffer, and delete it from the output buffer.
- Otherwise, do the following to replenish the output buffer:
- Read a character from the input. If it's not a canonical composite, write it to the output buffer. Otherwise, look up the character's canonical decomposition. If the first character of the canonical

decomposition is also a canonical composite, look *that* up, and so on until the original character is fully decomposed. Write the resultant characters to the output buffer. [One way to do this is to use an intermediate buffer. As you look up each decomposition, write the *second* character, if there is one, to the intermediate buffer. If the first character is decomposable as well, look up its decomposition and write the second character to the intermediate buffer. Repeat this process until you get a decomposition where the first character does not have a decomposition. When this has happened, write the *first* character to the buffer. At this point, you have the decomposition in the intermediate buffer in reverse order and you can reverse the order as you copy the characters into the output buffer.

- Peek at the input. If the next input character has a non-zero combining class, read it from the input, decompose it following the above rules, and write the decomposed version to the output buffer. Continue in this way until you reach a character with a combining class of zero. Do *not* read this character from the input yet.
- If the output buffer contains only one character, clear the buffer and return that character.
- Otherwise, if the output buffer has more than two characters, look for a pair of adjacent characters (starting with the second and third characters) where the first character in the pair has a higher combining class than the second. If you find such a pair, exchange the two characters. Compare the new first character in the pair with the character that now immediately precedes it and exchange *them* if their combining classes are out of order. Continue in this way until you've found the appropriate place for what was originally the second character in the pair, and then go back to looking for out-of-order pairs. (In other words, we do an insertion sort to get the combining marks in order by combining class without disturbing the relative order of characters with the same combining class.)
- Remove the first character from the output buffer and return it to the caller.

Compatibility decomposition

Compatibility decomposition (i.e., conversion to Normalized Form KD) is just like canonical decomposition. In fact, it's a superset of canonical decomposition. For every canonical composite in the string, you replace it with its decomposition. In addition, for every *compatibility* composite in the string, replace it with its decomposition (in fact you may have to do both on the same character: there are characters whose canonical decomposition includes a compatibility composite). Repeat until there are no characters left in the string that have either a canonical or a compatibility decomposition. Finish as before, by performing canonical reordering on the whole string.

The only difference is that you decompose compatibility composites. Unlike with canonical decompositions, a compatibility decomposition may have any number of characters. This is because with canonical decompositions, you could guarantee that the lead character in each partially-decomposed version of the character existed in the Unicode standard.

That is, for every character whose full decomposition consisted of more than two characters, you could guarantee that there would be appropriate characters in Unicode that would allow you to represent it as the combination of a single combining mark and another composite character. You don't have this guarantee with compatibility decompositions.

This means that you need a way to map from a single Unicode character to a variable-length string. You can use a two-level compact array (or generalized trie) for the main lookup table and have the results be pointers into an auxiliary data structure containing the result strings.

There are a couple ways of optimizing this: You could have the compact array map to 16-bit values. The value would either be a single UTF-16 code unit, or it'd be a pointer into the auxiliary string list (you'd probably reserve the surrogate and private-use values for use as these pointer values).

You could also have the compact array map to 32-bit values. The value could be a single UTF-16 code unit (zero-padded out to 32 bits), two UTF-16 code units, or the pointer into the auxiliary string list (again, using 32-bit values where the first 16 bits were surrogate or private-use code units).

If the compact array maps to 32-bit values, another way to do it would be to have the 32-bit value encode both an offset into the auxiliary string list *and a length*. The strings in the auxiliary list then wouldn't have to be null-terminated and you could be cleverer about how you packed things together in that list.

For example, consider the Roman numeral characters in the Number Forms block. These all have compatibility mapping to sequences of normal letters. You could store the compatibility mappings for the numbers from I to XII (U+2160 to U+216B) with the following string:

IXIIVIII

Twenty-six characters' worth of compatibility mappings stored using just eight characters' worth of string storage. Of course, if you're storing the one- and two-character mappings directly in the lookup table, you can lop the "I" off the front of the above example and you're stored thirteen characters' worth of compatibility mappings in seven characters' worth of storage. Still not bad.

Since there are no Unicode characters that have both canonical decompositions and compatibility decompositions (at the top level—a character can have both kinds of decomposition by virtue of having a canonical decomposition that contains a compatibility composite), you can also store both types of decompositions in the same table. This can save room, although it does mean you have to load all the compatibility mappings into memory even when you're only doing canonical decomposition, which could waste time and memory. The big problem with putting both kinds of mappings in the same table is that you need a way to tell which are the canonical mappings and which are the compatibility mappings. (Unless, of course, you only ever want to do compatibility decomposition.)

You could do this through lookup in a different table (perhaps combining it with the lookups necessary to get a character's combining class). If, on the other hand, you do it all in one table, you could do it by having separate auxiliary tables for canonical and compatibility mappings (and declaring that compatibility mappings always go to the auxiliary table, even if they fit in 32 bits), or by reserving some part of the table-lookup value as a flag to tell you whether a particular mapping is canonical or compatibility (this either makes each entry bigger, or it also forces you to go to the auxiliary table more than you'd otherwise have to).

Canonical composition

Normalized Forms C and KC are more recent normalized Unicode forms than Normalized Forms D and KD, which have been around as long as Unicode itself. They arose out of a need for a more compact canonical representation, and from the fact that many legacy encodings favor composite forms for character-mark combinations instead of favoring combining character sequences. These forms are specified in Unicode Standard Annex #15 (which also gives Normalized Forms D and KD their current names), which was officially incorporated into version 3.1 of the Unicode standard.

Part of the reason why the older versions of Unicode preferred decomposed forms over composite forms is that normalizing to the composed forms is inherently a more complicated business than normalizing to the decomposed forms, as we'll see.

To get to Normalized Form C or Normalized Form KC, you start by, respectively, converting the text to Normalized Form D or Normalized Form KD. In other words, you have to start by taking apart all the composite characters already in the text. Only then can you put them all back together. This is because of canonical ordering behavior. As we saw earlier in the chapter, there are two alternative ways of representing a-macron-underdot with precomposed characters: a-macron followed by a combining underdot and a-underdot followed by a combining macron. A-underdot followed by a combining macron is the canonical representation in Normalized Form C. If you start with a string that contains a-macron followed by a combining underdot...

ā + .

...the only way you can get from that to the preferred representation is to first decompose it...

a + ¨ + .

...perform canonical reordering on it...

a + . + ¨

...and *then* recompose what you can, yielding the preferred representation:

ḁ + ¨

So you start by doing canonical or compatibility decomposition and then you do canonical reordering. Only after you've done both of these, producing Normalized Form D or KD, do you perform the new process, *canonical composition*.

Canonical composition on a single pair of characters. The fun thing here is that instead of going from a single character to a pair of characters, you're going from a pair of characters to a single character. In effect, instead of having a lookup table with 1,114,112 entries, each consisting of either one or two characters, you need a lookup table with 1,114,112 *times* 1,114,112 entries. That's more than 1.24 *trillion* entries! It's a good thing each one of them needs to store only one character instead of a pair of characters.

Of course, only 13,038 of those entries actually have a character in them. (That's how many characters have canonical decompositions in Unicode 3.1—remember, this process never generates characters that have *compatibility* decompositions. Normalized Form C will leave such characters alone if they were in the text to begin with. Normalized Form KC won't, but neither form will produce new ones.) Actually, the real number is less than 13,038, for reasons we'll see shortly.

Because the lookup table is so sparsely populated, it lends itself well to trie-based schemes for reducing the memory required, although you need a more complicated scheme than a simple compact array because you're trying to model a compact array. You're also helped by the fact that no supplementary-plane characters participate in canonical composition (we'll see why in a moment), so the conceptual table you're trying to model is actually only 65,536 by 65,536, or 4,294,967,296 entries.

In fact, the conceptual table is smaller than that, because the character along one axis will always be a non-combining character and the character along the other axis will always be a combining character. Of course, this fact isn't necessarily of much practical significance, since the combining characters are scattered all over the BMP.

One approach to this problem would be to use a four-level trie. The first two levels are based on bits from the base character and the second two on bits from the combining character. Since in many cases you can discover just from looking at the top byte of the base character that it doesn't participate in any compositions, you can optimize for speed and space by pruning branches from the trie.

Also, since there's only a limited number of Unicode characters that participate in composition, you could take advantage of this fact and use a two-dimensional array with only the characters that participate in composition along the axes. This would be a *lot* smaller than a two-dimensional array with positions for all the characters. You could then use a simple compact-array lookup (or possibly even a hash function) to map from real code point values to the coordinates in this array.

Just as a single character can decompose into a sequence of three or more characters when performing canonical decomposition, so too can a sequence of three or more characters compose into a single character when performing canonical composition. Again, you do this by repeatedly applying the composition rules to pairs of characters. The details are a little trickier than for decomposition, however— we'll look at that later.

Canonical composition on a sequence of Hangul jamo. Just as it's possible to decompose a Hangul syllable algorithmically, it's also possible to recompose a Hangul syllable algorithmically. But there are some wrinkles.

When you canonically decompose a Hangul syllable, you end up with a very regular sequence of characters: exactly one initial consonant, exactly one vowel, and no more than one final consonant. These are the only sequences that can be recomposed.

It's legal to have multiple initial consonants or vowels or final consonants in a row in a sequence of conjoining Hangul jamo, and it's easy to see how this might happen: For example, the character

U+1115 HANGUL CHOSEONG NIEUN-TIKEUT (ㄴㅡ) could also be represented with two Unicode characters: U+1102 HANGUL CHOSEONG NIEUN followed by U+1103 HANGUL CHOSEONG TIKEUT. If you were to break a Hangul syllable down into individual marks this way, you could theoretically end up representing a single syllable with as many as eight code points (three initial consonants, two vowels, and three final consonants).

Unicode allows this, but does nothing to encourage it. None of the conjoining jamo characters have either canonical or compatibility decompositions: the jamo that consists of two or three marks (that can also be used on their own as jamo) exist on a completely equal footing with their brethren that only consist of a single mark: There's nothing in Unicode that will either break down U+1113 into U+1102 followed by U+1100 or combine U+1102 U+1100 into U+1113. In practice, pathological sequences consisting of more than one of a particular kind of jamo in a row (or missing an initial consonant or vowel) don't really happen: standard Korean keyboard layouts don't generate them, and rendering engines generally won't draw such sequences correctly.

(Of course, if you're actually *writing* a Korean input method, a Korean character renderer, or other software designed to operate on Korean text, you very well *may* care about these "pathological" sequences and their equivalences to more "well-formed" sequences. The appropriate equivalences aren't, however, specified by the Unicode standard and mapping based on them aren't part of Unicode normalization. Unicode 1.0 had some of these mappings, but they were removed from later versions and declared an application-specific mapping.)

So a canonical-composition process operating on Hangeul jamo has to start by checking to see whether it's looking at a leading-consonant/vowel or leading-consonant/vowel/trailing-consonant combination. If it isn't, it leaves the jamo characters alone and keeps searching for one of these combinations. (For the purposes of canonical composition, U+115F HANGUL CHOSEONG FILLER doesn't count as an initial consonant and U+1160 HANGUL JUNGSEONG FILLER doesn't count as a vowel.)

If it does see a lead-vowel-trail or lead-vowel sequence, then and only then does it perform the reverse of the transformation we looked at in the section on canonical decomposition. I'll leave the code to do that composition as an exercise for the reader.

One other wrinkle: If, God forbid, you're looking at a stretch of text that *mixes* conjoining jamo with precomposed syllables, you'll only get the right answer if you decompose the syllables before you compose the jamo. This is for the same reason as why you have to decompose any composite characters before composing everything: If you happen upon a sequence that consists of a syllable that only contains a leading consonant and a vowel, followed by a trailing-consonant jamo, the trailing consonant *can* combine with the syllable to form a different syllable. This is what you'll get if the original syllable had been spelled with jamo, so you have to make sure that also works if the first part of the syllable was spelled with an actual precomposed-syllable character.

Composition exclusion. Unicode Standard Annex #15, which defines the normalization forms, also specifies a list of characters with canonical decompositions that *aren't* allowed in Normalized Forms C and KC. This is called the *composition exclusion list*. This list can actually get longer over time, so it's supplied in the Unicode Character Database as `CompositionExclusions.txt`.

There are four classes of composite characters in the composition exclusion list:

- Characters with singleton decompositions. A character with a one-character canonical decomposition is a character that is either deprecated or discouraged in Unicode text, discouraged strongly enough so as to be prohibited in normalized Unicode text. These characters shouldn't happen in *any* normalization form, but need to be explicitly prohibited in Forms C and KC, since Forms C and KC are supposed to prefer canonical composites. (This is also why we know that canonical composition always works on pairs of characters: singleton decompositions are explicitly prohibited.)
- Composite characters with a non-zero combining class. There are a couple of combining characters that actually break down into pairs of combining characters in Normalization Forms D and KD. For example, there's a single code point value that represents both a dialytika (diaeresis) and a tonos (acute accent) on top of a Greek letter: it decomposes to separate diaeresis and acute-accent characters. That's the preferred representation in all normalization forms. Like the characters with singleton decompositions, composite characters with non-zero combining classes are effectively discouraged and shouldn't show up in normalized Unicode text.
- Script-specific exclusions. There are a few scripts where, for script-specific reasons, some composite characters are provided for compatibility with something, but their preferred representation is a combining character sequence. For example, the Alphabetic Presentation

Forms block includes a few pointed Hebrew letters. These particular letter-point combinations are fairly common in Yiddish or other non-Hebrew languages that use the Hebrew alphabet, but not in Hebrew, and not all that commonly in the other languages. If you normalize pointed Hebrew to Form C, you don't want this weird situation where 90% of the text uses the characters in the Hebrew block (with the points being represented with their own code points), but occasional letter-point combinations being represented with these compatibility composites from the Alphabetic Presentation Forms block. So these presentation forms are excluded.

- New characters. If, in version X of Unicode, you have to represent, say, x with a dot underneath using an x followed by a combining underdot character, but version X+1 of Unicode introduces a new precomposed x-underdot character, you get into trouble. Now if you take text in version X of Unicode and normalize it using a normalizer designed for version X+1, you wind up with text that's no longer in version X of Unicode and no longer readable by the program that originally produced it. To prevent this kind of thing, Normalized Forms C and KC were permanently nailed to version 3.0 of Unicode. Canonical composites added to Unicode after version 3.0 can never be produced by canonical composition. The canonical representation for these characters in *all four normalized forms* will always be the decomposed representation. This is why we know we don't have to worry about supplementary-plane characters in our canonical composition tables.¹¹⁶

You generally don't have to specifically worry about the characters in the composition exclusion list in your code at runtime: the table that maps base-mark pairs to composite characters just has to be set up so as not to map anything to characters in the composition exclusion list.

Canonical composition on a whole string. Applying canonical composition to a string of characters is a little more complicated than applying canonical decomposition to a string of characters. This is because you have to be careful not to do anything that would mess up the ordering of the combining marks in a combining character sequence. If you're starting with a string that's in Normalized Form D or KD, the code would look something like this:

```
public static char[] canonCombine(char[] in) {
    char[] result = new char[in.length];
    int p = 0;
    int p2 = 1;

    // start by comparing the first and second characters in the string
    int i = 0;
    int j = 1;

    while (i < in.length) {
        // if both characters are of combining class 0, we're done
        // combining things with the first character. Advance
        // past the combining character sequences in both
        // input and output
        if (combClass(in[i]) == 0
            && (j == in.length || combClass(in[j]) == 0)) {
            result[p] = in[i];
            p = p2;
        }
    }
}
```

¹¹⁶ You actually don't have a stability problem if the composite character *and all the characters in its full decomposition* are introduced in the same Unicode version, as is the case with the musical symbols in Unicode 3.1. You only have a stability problem if a newly-added composite character has characters in its full decomposition that were added in earlier version of Unicode. UAX #15 doesn't make that distinction, however: all composite characters added after Unicode 3.0 are disallowed in normalized text, even if their presence wouldn't cause a stability problem.


```
        ++p2;
        i = j;
        ++j;
    }

    // if the second character is a combining character and it
    // can combine with the first character, combine them,
    // storing the temporary result in the input buffer,
    // and advance past the combining mark without writing it
    // to the output
    else if (combClass(in[i]) == 0 && canCombine(in[i], in[j])) {
        in[i] = combine(in[i], in[j]);
        ++j;
    }

    // if the second character is a combining character and it
    // CAN'T combine with the first, this character (and any
    // others of the same combining class) will remain in the
    // output. Copy them to the output (keeping track of where
    // the base character is) and advance past them
    else if (combClass(in[i]) == 0) {
        int jCombClass = combClass(in[j]);
        while (combClass(in[j]) == jCombClass)
            result[p2++] = in[j++];
    }
}
}
```

This example assumes the existence of a function called `combClass()` that returns a character's combining class, a function called `canCombine()` that says whether a particular pair of characters can be combined into a single composite character, and a function called `combine()` that actually composes two characters into a single composite character. The basic idea is that you try to combine each base character (character with combining class 0) with as many combining characters after it as possible, with the exception that any characters that *can't* combine with the base character block any characters of the same combining class that follow them from combining with the base character.

You have to do this because two combining characters with the same combining class interact with each other typographically—their order tells you which one to draw closer to the base character. If the inner character can't combine with base, you can't combine the outer one with the base character without changing their visual presentation (i.e., making the outer character show up closer to the base character than the inner character does).

Putting it all together. Just as you can interleave canonical decomposition with canonical reordering, so too can you interleave both of these processes with canonical composition. If the string you're trying to convert to Normalized Form C is already in Form C, it's a big waste of time to convert it to Form D just so you can convert it back to Form D. You can avoid this by doing something like this:

- If the current character is a composite character, but it's not in the composition exclusion list and the next character is of combining class zero, leave it alone.
- If the current character is a composite character, perform canonical decomposition on it and, if necessary, perform canonical reordering on the result and any subsequent combining characters.

Then perform canonical composition on the resulting sequence. (If the character that got decomposed wasn't of combining class zero, you'll actually have to back up until you find a character that is and start both canonical reordering and canonical composition from there.)

- If the current character is *not* a composite character and the character that follows it is a combining character, just perform canonical composition on this character and the combining character(s) that follow it.
- Otherwise, leave the current character alone.

Optimizing Unicode normalization

The `DerivedNormalizationProperties.txt` file in the Unicode Character Database can be useful in optimizing your Unicode normalization code. For each Unicode character, this file tells you whether it can occur in each of the Unicode normalized forms. You can use this to do a quick check on each character to see whether anything at all needs to be done with it. If you're converting to Normalized Form C, for example, you don't actually have to drop into the full normalization logic until you see a character with the `CompEx`, `NFC_NO`, or `NFC_MAYBE` property. This can save you a lot of work.

The `DerivedNormalizationProperties.txt` file also includes properties that tell you whether a particular character will expand to more than one character when converted to a particular normalized form. You can use this to defer allocation of extra storage space for the result string until you really need it (or to find out before you start normalizing whether or not you can do it in place).

Testing Unicode normalization

One of the files in the Unicode Character Database is a file called `NormalizationTest.txt`. You use this file to determine whether your implementation of Unicode normalization is correct or not.

Each line consists of a set of Unicode code point values written out as hex numerals (each line also has a comment with the literal characters and their names to make everything more readable). The lines are divided into five parts, or "columns," separated with semicolons. The first column is some arbitrary Unicode string, which may or may not be normalized already. The other four columns represent that string converted into Normalization Forms C, D, KC, and KD respectively. If your implementation of normalization is correct, the following statements will be true for each line of the file:

- If you convert every sequence on the line into Normalized Form C, columns 2 and 4 shouldn't change, columns 1 and 3 should change to be the same as column 2, and column 5 should change to be the same as column 4.
- If you convert every sequence on the line into Normalized Form D, columns 3 and 5 shouldn't change, columns 1 and 2 should change to be the same as column 3, and column 4 should change to be the same as column 5.
- If you convert every sequence on the line into Normalized Form KC, column 4 shouldn't change, and the other four columns should change to be the same as it.
- If you convert every sequence on the line into Normalized Form KD, column 5 shouldn't change, and the other four columns should change to be the same as it.

The file is divided into three parts: Part 0 is a spot check for a number of interesting characters. Part 1 is a character-by-character test of every Unicode character that has a decomposition (in addition to the characters listed here, you should check all the other Unicode characters to make sure they don't change when any of the normalizations are applied to them). Part 2 is a fairly elaborate test of canonical reordering.

For more information, see the file itself. But the basic rule is that if your normalization code can't satisfy all of the conditions described above for every line in the file, it doesn't conform to the standard.

Converting between Unicode and other standards

It'd be nice if all the digitally-encoded text in the world was already in Unicode, but it isn't, and never will be. Not only are there innumerable pieces of encoded text and computer programs that process text out there that predate the introduction of Unicode and will never be updated, but there are and will always continue to be documents, computer systems, and user communities that don't and won't use Unicode. While one can hope that over time the computing public will move more and more toward using Unicode exclusively for the encoding of text, the reality is that there will always be text encoded using non-Unicode encodings, and that Unicode-based systems will frequently have to contend with text encoded in these encodings.

This means that some sort of conversion facility will be required for almost every system that exchanges text with the outside world. Of course, pretty much every Unicode-support library out there includes such a conversion facility—unlike compression and normalization, you'll rarely find yourself working with a library that doesn't include some sort of facility for character encoding conversion. The important question usually isn't "Am I working with an API that provides encoding conversion?" but "Am I working with an API that provides the conversions I need?" or "Am I working with an API that does conversion in the way I want it done?"

If the answer is yes, or the API you're using produces a close-enough answer that you can modify reasonably easily to get the answer you want, use it. There's no sense in going to the trouble to create all this yourself when you don't have to. But if you find you have to implement character conversion for whatever reason, read on...

Getting conversion information

The first issue is coming up with the information on just how you convert between such-and-such an encoding and Unicode. For a wide variety of encoding standards, you can find the answer on the Unicode Web site. The Unicode Consortium maintains an unofficial repository of mapping tables between Unicode and various other national, international, and vendor encoding standards. The tables are in a simple textual format where each line consists of a sequence of bytes (shown as two-digit hex values), a tab character, and a sequence of Unicode code point values (shown as four-or-more-digit hex values).

Interestingly, Unicode Technical Report #22 proposes a richer XML-based format for describing how to convert from some encoding to Unicode and back; as of this writing, this format wasn't being used for the mapping tables on the Unicode Web site, but certain third-party libraries, notably the open-source International Components for Unicode, are using it. There's also a lot of instructive information in this report on character encoding conversion, and it's worth reading.¹¹⁷

¹¹⁷ In fact, much of the information in the rest of this section, particularly the information on choosing converters and handling exceptional conditions, is drawn from UTR #22.

Converting between Unicode and single-byte encodings

The mechanics of converting from a single-byte encoding to Unicode are trivially simple: Just use an array of 256 code point values. Converting back the other way isn't much harder: use some form of trie structure (if the conversion doesn't involve any supplementary-plane characters, this can be a simple one-level compact array). Often, you'll get better compression if the trie stores numerical offsets rather than the actual values to map to, but this generally won't make a difference when you're converting to a one-byte encoding.

Converting between Unicode and multi-byte encodings

Converting between Unicode and a multi-byte encoding will generally require a trie structure of some type for both directions. If the non-Unicode encoding has both one- and two-byte code points, it often makes sense for the Unicode-to-non-Unicode table to behave as though all of the code points in the target encoding are two bytes long, representing one-byte code points with an invalid leading byte value (0x00 will work in most of them) and suppress that when emitting the resulting characters. This lets you always map to 16-bit values in the table and it also helps if you want to store numeric offsets in the table instead of actual code point values in the target encoding. With multi-byte encodings, there are often enough code points in the sequence being mapped to for this to help your compression in the trie.

Other types of conversion

There are a few interesting cases where you can't do simple table-based conversion and have to do something special.

Stateful encodings. A stateful encoding is one where bytes you're seen earlier in the conversion process affect the interpretation of bytes you see later in the process. We're not talking here about simple multi-byte encodings, where several bytes in a row combine into a single code point, but situations where whole code point values are overloaded with multiple semantics. Typically, you have a bank of code point values that each represent two (or more characters). Shift-in and shift-out codes in the text stream put the converter in a "mode" that specifies which of the two sets of characters that bank of code points is to represent. That mode persists until the next shift-in/shift-out character.

SCSU is an example of a stateful encoding. The old Baudot code, which used the same range of code points for letters and digits and used control codes to shift that bank back and forth, was as well. There are also a number of modern EBCDIC-based encodings that still use this technique.

For stateful encodings, your converter will have to keep track of its mode. You'll either need multiple to-Unicode conversion tables, one for each mode, or your tables will need to be able to store all of the overloaded values for those code points that are overloaded. The from-Unicode converter will also have to keep track of mode, with the table including information on which mode to use for the overloaded code points in the non-Unicode encodings. You'd emit a mode-shift character anytime you have to emit a code point with a mode you're not currently in.

ISO 2022. ISO 2022 is a special example of a stateful encoding. The ISO 2022 standard defines a set of common characteristics for coded character sets, but doesn't actually define a coded character set. Instead, it defines a character encoding scheme for representing characters from different coded character sets in the same text stream by making use of various escape sequences to change the meanings of characters in different parts of the encoding space. An ISO 2022 converter would

generally use other converters for the actual character code conversion; all it'd do is interpret (or emit) the escape sequences to shift coded character sets.

Algorithmic conversions. There are a few encodings where you can do all the work with a simple mathematical transformation. The classical example is, of course, ISO 8859-1 (Latin-1). You could convert from 8859-1 to Unicode simply by zero-padding the byte values out to 16 bits. You could convert from Unicode back to 8859-1 simply by truncating the value back down to 8 bits (although you probably want to do a little more work to handle exceptional conditions).

There are a number of other encodings where Unicode preserves the relative order of the characters in the code charts and you can convert to and from Unicode by just adding or subtracting a constant (or, in the case of some ASCII-based encodings, doing a range check and adding or subtracting one of two constants depending on the result).

Transformations on other conversions. If you're implementing a whole bunch of encodings, you may be able to save space for conversion tables by doubling up. For example, you might have two encodings, one of which is a superset of the other. You could use the same table for both of them and postprocess the conversion for the subset encoding. (This might work, for example with ISO 8859-1 and Microsoft code page 1252, which is a subset, although this is really only true if you ignore the C1 control-code space in ISO 8859-1.) As another example, there are a number of encoding schemes based on the JIS X 0208 character set, so you might use a table to convert to the abstract JIS 208 values and then algorithmically convert to the actual encoding scheme in use (for example, EUC-JP).

Another example is variants on the same conversion. For example, you may want 0x5C in a Japanese ASCII-based encoding to be treated as the yen sign, or you may want it to be treated as the backslash. Rather than having two different tables that differ only in their handling of this one character, you might just have a table for one and postprocess the result when you want the other variant.

Handling exceptional conditions

There are a number of exceptional conditions you need to handle when converting between any two encodings:

- The input text is malformed. You encounter a byte or sequence of bytes that isn't supposed to occur in that encoding. For example, U+FFFF is illegal in all Unicode transformations, 0xFF should never happen in UTF-8, a high surrogate is illegal in UTF-16 unless it's followed by a low surrogate, and so forth. A lot of encodings, particularly the variable-length ones, have restrictions like this.
- A byte or sequence of bytes is legal, but isn't actually assigned to a character. This may indicate malformed text, or it might indicate that the text is in a newer version of the encoding than your converter is designed to handle.
- A byte or sequence of bytes is legal and actually assigned to a character, but the encoding you're mapping to doesn't have a representation for that character.

You may want to treat all three conditions the same way, or you may want to treat them differently. It may also depend on the situation (if you're writing an API library, you may want the calling code to tell you what to do). You've got a number of choices for dealing with exceptional conditions:

- Halt the conversion and report an error. Generally, just choking on the input isn't a good idea, but if the conversion is restartable, this technique is often used to defer the job of handling the exception to the caller.
- Ignore the bytes from the source text stream that caused the exception. This is a pretty common strategy, but has the disadvantage of masking problems in the source text. The resulting text can look okay, with the user never realizing there was additional text there that didn't convert.
- Write out some kind of error character. If you're going to Unicode, this usually means writing out U+FFFD, the replacement character, when you encounter an illegal input sequence or one you can't map. (Most encodings have an analogous character—in ASCII, 0x1A, the SUB character, serves this purpose.) Sometimes, you might want to distinguish between unmappable and illegal input sequences. You could use U+FFFF in Unicode for the illegal sequences. You have to be careful with this technique, though, because U+FFFF should never happen in well-formed Unicode. It's okay for internal-processing purposes, but you'd have to be careful to convert U+FFFF to U+FFFD or strip it out before sending the text to some other process.
- If you encounter an unmappable character when converting to Unicode (that doesn't happen very often), you can map it to a code point in the private use space. The beauty of this is that if you later convert back to the source encoding from Unicode, the characters you couldn't handle can be restored. (If you know what you're converting to private-use code points, you can just treat them as ordinary characters inside your application as well.)
- Write out an escape sequence. If the text you're converting is in a format that supports numeric escape sequences, you can convert unmappable source code points into escape sequences. For example, if you're converting from Unicode to ASCII and you know the destination is a Java source file, you can convert the unmappable character to a sequence of the form “\u1234”. If you're going to an XML file, you can convert to an XML numeric character reference (e.g., “Ӓ”).
- Use a fallback mapping. If the character you're mapping can't be represented in the target character set, you might convert it to something reasonably similar. For example, you might convert the curly “” quotation marks to the straight quotation marks when going to ASCII.
You want to be careful when using fallback mappings, since if you convert from Unicode to something else using fallback mappings, then back to Unicode, you lose data. Sometimes this is okay, but if you need to preserve round-trip integrity, you can't use fallback mappings (or you need some way to flag them as such). In systems, such as HTML and XML, that can properly handle characters encoded as escape sequences, that's generally a better way to go.
For the same reason, you generally don't want multiple characters in one encoding to map to the same character in another encoding, because this also loses data in a round-trip conversion. Unless the mappings are truly for alternative representations (say, Unicode characters with singleton canonical mappings to other characters), you usually want to flag all but one of them as fallback mappings.

Dealing with differences in encoding philosophy

There are other types of transformations you might have to worry about when converting between Unicode and something else.

Normalization form. Usually, an encoding naturally converts to Normalized Form C or to Normalized Form D. Some will actually convert naturally to both (they don't include any composite characters and they don't include any combining characters—our old friend ASCII is an example). When converting from Unicode to one of these encodings, you generally want to convert to the appropriate Unicode normalized form first (or as part of the encoding conversion) to make sure you really catch everything representable in the target encoding.

If the non-Unicode encoding doesn't naturally convert to any Unicode normalized form (i.e., there are sequences in it that would convert to unnormalized Unicode text), it's generally a good policy to normalize the result after converting to Unicode.

Combining character sequences. Not all encodings that use combining character sequences encode them the same way Unicode does. If, say, you're converting to or from an encoding, such as ISO 6937, that puts the mark *before* the character it attaches to, you have to switch the base-mark pairs as part of the conversion.

Character order. Some older vendor encodings for Hebrew and Arabic use visual order instead of logical order to store the characters. In these cases, you actually have to apply the Unicode bi-di algorithm (or a reversed version of the bi-di algorithm) as part of the conversion.

Glyph encodings. There are some encodings (again, generally older vendor encodings) that are glyph-based. In these encodings, different glyphs for the same character get different code point values. Converting to an encoding like this actually involves the same kind of contextual glyph selection logic you need to render the text on the screen.

Choosing a converter

Choosing a converter is generally straightforward, but you can get into interesting situations if the converter you want doesn't exist but there's one there that's close. In situations like these, you can sometimes do a "best fit" mapping. You acknowledge you might lose data, but sometimes this is okay, as long as you know there will be errors and you don't purport to be doing the correct conversion.

The one case where things work fine is if the text you're reading is in encoding X, you have a converter for some encoding that's a superset of X, and you're going from X to Unicode. In this case, everything just works fine.

If you're in an environment that supports escape sequences (say, you're producing an XML file), you can also do the opposite: You want to output X, you have a converter for a subset of X, and you're going from Unicode to X. Some characters that are representable in X aren't representable in the subset you're actually converting to, but since you can convert them to escape sequences, you don't lose data.

Line-break conversion

Line-break conventions (i.e., which characters you use to signal the end of a line or paragraph) are generally specific to a particular platform or application, rather than a particular target encoding (of course, the target encoding defines which choices you have). If your converters have to be cross-platform, or used in some other situation where you can't nail down the line-breaking conventions, you may need to do line-break conversion in a separate pass, or separate out line-break conversion from the rest of the conversion logic.

Case mapping and case folding

English speakers are familiar with the concept of case. Every letter in the Latin alphabet comes in two “cases”: upper and lower (or capital and small). Like the Latin alphabet, the Greek, Cyrillic, and Armenian alphabets are also cased, or “bicameral.” (The letters in the Georgian alphabet also come in two distinct styles, which for a short time in history were treated as upper and lower case, but generally aren’t anymore.)

This means you will occasionally run into situations where you want to convert one letter (or a whole string) from one case to the other. Many word processors, for example, let you convert to upper or lower case with a single command or a character style.

On many systems, you convert everything to upper or lower case not for display, but to erase case distinctions when searching. In Unicode, this is done with a different process, known as case *folding*.

Unicode Technical Report #21 provides more information on performing case mappings on Unicode characters.

Case mapping on a single character

First of all, it’ll come as no surprise that you can’t just convert from one case to another by adding or subtracting some constant, as you can with ASCII. This is because not all case pairs are the same distance apart. For example, the case pairs in the ASCII block are 32 code points apart, but those in the Latin Extended A block are immediately adjacent to each other.

So you have to do a table lookup. The UnicodeData.txt and SpecialCasing.txt files in the Unicode Character Database include the case mappings for all the Unicode characters. You can use a simple one-level compact array as your case mapping table. As with code conversion, you get better compression if the table holds numeric offsets rather than actual code point values.

You do, however, need an exception table for case mappings. There are a number of reasons for this:

- **Titlecase.** Unicode includes a number of single code points that actually represent two characters. If one of these occurs at the beginning of a word, you need a special “titlecase” version that represents the capital version of the first letter and the lowercase version of the second letter and can be used at the beginning of a word. For the characters that have titlecase versions, you actually have to store *two* mappings: one to titlecase and one to the opposite case. (For all other characters, the “titlecase” form is the same as the uppercase form.)
- **Expansion.** Some characters actually turn into two characters when mapped to uppercase.
- **Context sensitivity.** Some characters map to different characters depending on where they occur in a word.
- **Language sensitivity.** Some characters have different mappings depending on the language.

The UnicodeData file includes the titlecase mappings for everything. The other special mappings, for historical reasons, are in a separate file, SpecialCasing.txt. For a complete treatment of SpecialCasing.txt, see Chapter 5, but to recap, some of the more important mappings in the SpecialCasing.txt file are:

- Many characters that expand when converted to uppercase. The most important and common example is the German letter ß, which uppercases to “SS,” but there are many others.
- The Greek letter sigma, which has two lowercase versions, one that occurs at the end of a word and another that occurs in the middle of a word (or when the letter sigma appears on its own as a “word”).
- Mappings to deal with Lithuanian. In Lithuanian, a lowercase i or j with an accent mark over it doesn’t lose the dot, like it does in other languages. (In other words, a lowercase i with an acute accent normally looks like this—*í*—but in Lithuanian it looks like this: *į* [fix baseline and spacing]. You can do this with a special font for Lithuanian, but if you want to explicitly represent it this way in Unicode, you have to use combining character sequences that use an explicit U+0307 COMBINING DOT ABOVE, and the case mapping code has to be smart enough to add and remove this character as necessary.
- Language-sensitive mappings for Turkish and Azeri that deal with the letter ı, whose uppercase form is İ. (The uppercase version of i becomes İ so you can tell them apart.)

One thing to keep in mind when doing case mappings is that they lose data and aren’t inherently reversible. If you start with “Lucy McGillicuddy” and convert it into uppercase, you get “LUCY MCGILLICUDDY.” If you now convert that to lowercase, you get “lucy mcgillicuddy.” Going to titlecase instead also doesn’t get you back your original string: you get “Lucy McGillicuddy.” In a word-processing program or other situation where you provide a case-mapping facility to the user, you need to make sure you save off the original string so you can restore it if necessary. (A simple Undo facility usually takes care of this requirement.)

Case mapping on a string

Mapping a string from one case to another is pretty simple if you’re going to uppercase or lowercase: Just apply the appropriate transformation to each character one by one. The only thing you need to keep in mind is that the string may get longer in the process (which may make converting case in place impossible or force you to fallback on the single-character mappings in UnicodeData.txt, even though they’re not always correct).

Going to titlecase is a little more complicated: Here, you have to iterate through the string and check context: If a character is preceded by another cased character (whatever its case), convert it to lowercase. Otherwise, convert it to titlecase (which, most of the time, will be the same as converting it to uppercase).

Of course, this doesn’t actually give you what you want all the time. “Ludwig van Beethoven” would turn into “Ludwig Van Beethoven” under this transformation, again possibly making it necessary to save off the original string before messing with it. More importantly, applying the transformation to “THE CATCHER IN THE RYE” gives you “The Catcher In The Rye,” when what you really want is “The Catcher in the Rye.” In other words, not every word in a real title is capitalized. Getting this right is language-specific, as the list of words you don’t capitalize varies from language to language. A simple algorithm won’t worry about this and will just capitalize everything.

Case folding

Now if what you’re really trying to do is erase case distinctions from a pair of strings you’re about to compare for equality, you don’t want to do it using the case-mapping stuff. Converting to any one

case will leave in some distinctions that converting the other way will erase, and doing both conversions is a little costly (and doing it in different orders will *also* produce different results).

The Unicode Character Database includes another file, `CaseFolding.txt`, specifically for this use. It maps all the cased letters to case-independent sequences of characters (generally, but not always, the mapping is equivalent to mapping first to uppercase, then mapping the result to lowercase).

`CaseFolding.txt` actually includes two separate sets of mappings: One set has a one-character mapping for everything. It'll generally give you the right answer, but will leave in some distinctions that you'd get rid of if you actually mapped to uppercase and then back to lowercase. The other set does have expanding mappings. It'll erase more differences, but is a little more complicated (and perhaps slower) to implement.

And again, the Turkish letter `ı` causes trouble. `CaseFolding.txt` includes two optional mappings for dealing with this: If you use them, it erases any case distinction, but also erases the distinction between `ı` and `i`. If you don't use them, it doesn't normalize out the case differences for `i` and `ı` (they and their uppercase counterparts will still compare different).

It's important to keep in mind that you only want to use the `CaseFolding.txt` mappings when you're doing *language-independent* comparisons, such as for building search trees or other internal uses. If you want truly language-sensitive behavior, you need a full-blown collation facility. We'll be looking at collation in the next chapter.

Transliteration

Finally, there's transliteration.¹¹⁸ The term generally refers to taking a piece of text written using some writing system and converting it to another writing system. This isn't the same thing as *translation*, which actually involves converting the text to a different language and changing the *words*. When text is transliterated, the words (and the language) stay the same; they're just written with different characters.

In computer text-processing circles, on the other hand, the word "transliteration" is used in a broader sense to refer to a broad range of algorithmic transformations on text. This is because a facility designed to convert text from one writing system to another will be flexible enough to also perform a bunch of other transformations as well.

So why is transliteration useful? Say, for example, you've got this customer list:¹¹⁹

김, 국삼

김, 명희

정, 명호

たけだ, まさゆき

ますだ, よしひこ

やまもと, のぼる

¹¹⁸ Most of the material in this section is drawn from Mark Davis and Alan Liu, "Transliteration in ICU," *Proceedings of the 19th International Unicode Conference*, session B14.

¹¹⁹ This example is ripped off directly from the Davis/Liu paper, *op. cit.*

Ρούτσι, Άννα
Καλούδης, Χρήστος
Θεοδωράτου, Ελένη

This will be a lot easier for an English-speaking user to deal with if you see it like this instead, especially if he doesn't have Japanese, Korean, or Greek fonts installed on his computer:

Gim, Gugsam
Gim, Myeonghyi
Jeong, Byeongho
Takeda, Masayuki
Masuda, Yoshihiko
Yamamoto, Noboru
Roútsē, Άnna
Kaloúds, Chrêstos
Theodōrátou, Elénē

Of course, a Russian-speaking user might be more comfortable with the list transliterated into Cyrillic letters instead.

In addition to converting text into some target writing system for display, a transliteration facility is also useful for a number of other types of mappings, including:

- Accepting keyboard input for a particular script using a keyboard designed for another script (e.g., entering Russian text using an English keyboard)
- Accepting keyboard input for Han and Hangul characters, where it takes multiple keystrokes on a conventional keyboard to enter a single character (building up a whole Hangul syllable out of several component jamo, for example, or entering Japanese or Chinese characters by entering a Romanized version of the text on an English keyboard)
- Converting text in a multi-scriptal language from one script to another (e.g., converting Serbo-Croatian from Cyrillic letters to Latin letters, Japanese from Hiragana to Katakana, Chinese from Simplified Chinese characters to Traditional Chinese characters, etc.)¹²⁰
- Case conversions
- Converting typewritten forms of certain characters to their typographically-correct equivalents (for example, converting -- to —, or converting " or `` to “)
- Allowing input of arbitrary Unicode characters by typing their hexadecimal code point values along with some special character
- Fixing various typing and spelling mistakes on the fly (think Microsoft Word's AutoCorrect feature, for example)
- Unicode normalization

Many of these transformations, especially the script-to-script conversions, are more involved, and require a more complicated approach, than the other conversions we've looked at. Not only do

¹²⁰ It's important to note that some of the more complicated conversions, such as from Han characters to any of the phonetic scripts, or between Simplified Han and Traditional Han, often require more-sophisticated logic, often involving linguistic analysis and dictionary lookup, that go beyond what a typical transliteration package can do. Conceptually, it's the same type of problem, but practically speaking, it requires different code.

you have to deal with the full generality of many-to-many mappings (where source and destination strings that are longer than a single character are the norm rather than the exception), but you also have to deal with:

- **Context.** A particular character or sequence of characters may only get mappings to some other sequence of characters when certain sequences of characters appear before or after it. For instance, " may only get mapped to “ when it’s preceded by a non-letter and followed by a letter.
- **Wildcard matches.** To specify the curly-quote mapping above, you aren’t going to want to spell out every single pair of surrounding characters that cause the mapping to happen; you’re going to want to specify character categories and let the engine expand this into actual pairs of characters.
- **Partial matches.** If the text you’re operating on is coming to you in chunks—say, from a buffered file API or from keyboard input—you may have to deal with situations where you have to wait to see more characters to know what mapping to perform. For example, if you’re mapping Latin letters to Cyrillic letters and have a rule that maps “s” to “с” and “sh” to “ш”, if you get an “s” in the input, you have to wait to see the next character (and see whether that character is an “h”) before you know what to do with the “s” (unless you know that the “s” is the last character in the input). This may involve either buffering characters you’re not done with yet or keeping track of which part of the input hasn’t been fully mapped yet.
- **Backup.** After you perform a mapping on a piece of text, you’ll generally want to start the next mapping with the first character you haven’t looked at yet, but you may want to back up into the result of the mapping you just performed instead. This can sometimes cut down drastically on the number of rules you need to achieve some effect.

For example, consider the palatalized Katakana combinations. You could just have rules that say that “kyu” turns into “キユ”, “kyo” turns into “キヨ”, “pyu” turns into “ピユ”, and so on, or you could have a rule that says that “ky” turns into “キ~y” and another rule that turns “~yu” into “ユ”. If there are x different consonants that can start a palatalized syllable and y different vowels that can end a palatalized syllable, you can reduce the number of mapping rules you need from x times y to x plus y using this technique. But it only works if you can back up after the first mapping and consider the “~y” again when determining what to do in the second mapping.

In essence, doing all this requires something similar to a generalized regular-expression-based search-and-replace facility. Without going into huge amounts of detail, one solution to the problem might involve the following data structures:

- A compact array or trie that maps single characters to character categories. This is important for keeping the main lookup table small when rules involve wildcard matches.
- A generalized trie structure that maps arbitrary-length sequences of character categories to result strings in an auxiliary data structure. You probably want the system to consider rules in some defined order (ensuring, for example, that you see the rule for “sh” before you see the rule for “s”), making a ternary-tree representation less than ideal. Most likely you’d base the main mapping table on a binary tree or a generalized tree of arrays. (That is, the list of matching rules or partially matching rules could be a simple linked list or an array.)
- A results data structure, probably a single memory block containing all the result strings. The strings could either be null-terminated or the main trie could contain pointers and lengths, possibly leading to better compaction.

The results data structure is interesting, because you have to store not only result strings there, but also information on which parts of the original matching sequence are context and don’t get

replaced and information on how far to back up into the result before starting the next matching operation. (Flags indicating which parts of the matching sequence are context could also be put into the main trie.) These are places where the new non-character code point values in Unicode 3.1 come in handy: you can use the non-character code points as tokens to indicate the backup and context positions, and then just “mark up” the result string with these tokens.

The algorithm would then look something like this:

1. Look up a category for the character at the starting position.
2. Use that category to look up the character in the main trie.
3. If the character leads you to another node in the trie, advance to the next character in the input and go back to step 1.
4. If the character leads you to an entry in the result table, replace the non-context part of the matching text in the input with the result text from the result table. The new starting position is either after the result text unless the result text specifies a backup position, in which case that’s the new starting position. Start over at the root level of the main trie and go back to step 1.
5. If you exhaust the input text while you’re pointing at something other than the root level of the main trie, you either have to pull in more text to figure out what to do, or you’re really done, in which case you use the “if the next character doesn’t match” rule at the current level of the main trie.

Again, this approach may not be the most effective for really complicated mappings that require more knowledge of the language, such as mapping Simplified Chinese to Traditional Chinese, but for most generalized mappings it works quite well.

Starting in the next chapter, we’ll look at more-complicated operations on Unicode text, many of which also involve mappings of the type we’ve looked at in this chapter.

CHAPTER 15 *Searching and Sorting*

Aside from drawing it on the screen (or some other output device), the most important processes that are usually performed on text are searching and sorting. And, of course, the main thing these two operations have in common is string comparison. As with the other processes we've looked at, there are a number of interesting complications involved in performing string comparison on Unicode strings, many of which stem from the larger repertoire of characters in Unicode. But unlike many of the other processes we've looked at, many of the issues that make good string comparison hard are not at all unique to Unicode, but are issues when dealing with text in any format.

Unlike many of the other processes on Unicode text that we've examined, language-sensitive string comparison is one of those features that virtually every Unicode support library provides. This means you'll rarely have to implement this yourself—you should usually be able to take advantage of some API that provides language-sensitive comparison. Of course, it might not do exactly what you want, or support the languages you're interested in. In this case you might have to write your own comparison routines, but it's often more of a matter of learning how to specify your desired results to the API you're using. Searching, on the other hand, is something you might well need to code on your own.

At any rate, it's still instructive to understand just what goes into doing Unicode-compatible searching and sorting. We'll start by looking at the basics of language-sensitive string comparison, then look at the unique considerations Unicode brings to the party, and finish by discussing some of the additional things to keep in mind when doing actual searching or sorting.

The basics of language-sensitive string comparison

The first thing to remember is that you can't simply rely on comparison of the numeric code point values when you're comparing two strings. Unless the strings that may be compared conform to a very tightly restricted grammar, this will always give you the wrong answer. (The one exception to

this is when the ordering and equivalences implied by the comparison routine will have no user-visible effects, but even then there are wrinkles to worry about—see “Language-insensitive string comparison” later in this chapter.)

This isn’t just a Unicode thing. Binary comparison will actually give you the wrong answers with most encodings (in fact, for every encoding standard, it’s probably possible to come up with a pair of strings that will compare wrong if you just use binary comparison, but it’s a bigger problem with some encoding standards). The numeric order of the code point (or code unit) values in most encodings is good enough to give you a reasonably acceptable ordering for many sets of strings (for example, the letters of the Latin alphabet are encoded contiguously and in alphabetic order in most encodings), but that’s usually about it.

Consider the following list of English words and names, sorted in ASCII order:

Co-Op of America
CoSelCo
Coates
Cohen
Cooper
Cooperatives United
Cosell
MacIntosh
MacPherson
Machinists Union
Macintosh
Macpherson
Van Damme
van der Waal
van den Hul
Vandenberg

If you squint hard enough, you can kind of convince yourself that this order makes sense, but it’s definitely not intuitive. The spaces, punctuation, and capital letters gum things up. The order should probably look more like this:

Coates
Cohen
Cooper
Cooperatives United
Co-Op of America
CoSelCo
Cosell
Machinists Union
MacIntosh
Macintosh
MacPherson
Macpherson
Van Damme
Vandenberg
van den Hul
van der Waal

In other words, you probably just want to consider the letters, with things like capital letters and spaces being secondary considerations. Sorting by the raw ASCII code point value doesn't give you that.

If you consider the Latin-1 encodings, things just get worse. Take the following list, for example:

archaeology
archway
archæology
co-operate
coequal
condescend
cooperate
coziness
coöperate
restore
resume
resurrect
royalty
rugby
résumé

Again, the accent marks gum things up. The alternate spellings of “cooperate” and “archaeology” wind up widely scattered, and “resume” and “résumé” wind up far apart. In general, the accented forms of letters end up widely separated from their unaccented counterparts, leading to weird results. Again, you probably expect the list to be sorted more like this:

archaeology
archæology
archway
coequal
condescend
cooperate
coöperate
co-operate
coziness
restore
resume
résumé
resurrect
royalty
rugby

Latin-1 has an additional problem. This is the fact that it's used to encode text in a lot of different languages, and the concept of “alphabetical order” isn't the same across all of them.

For example, consider the letter ö. In English, this is just the letter o with a diaeresis mark on it. The diaeresis is used to indicate that the o should be pronounced in a separate syllable from the vowel that precedes it, rather than forming a diphthong (such as in “coöperate” above), but the letter is still the letter o and should sort the same as if it didn't have the diaeresis.

In German, the diaeresis indicates a significant change in pronunciation, but the letter is still basically an o. Just the same, depending on who's doing the sorting, ö may sort either as a variant of o or as equivalent to "oe" ("oe" being an acceptable variant spelling of ö in German).

In Swedish, on the other hand, ö is a whole different letter, not just an o with a mark on it. It sorts with a number of other letters, such as å, ä and ø, at the end of the Swedish alphabet, after z.

It works the other way, too: In English, v and w are different letters. In Swedish, on the other hand, w is just a variant of v.

The upshot of this is not only that you can't trust the binary order of the code points, but that the actual order you use depends on the language the text is in. If you're sorting a group of Latin-1 strings, the order they end up in will be different depending on whether they're in English, Spanish, French, German, or Swedish (to name just a few). Sorting a list of Swedish names using the German or French sorting order is just as wrong as using the binary order.

This is important. The encoding isn't necessarily aligned with the language. Many encodings can be used to encode multiple languages. Some encodings tend to be used only with certain languages, but this isn't guaranteed and shouldn't be depended upon. In fact, sorting isn't dependent on the language the strings themselves are in; it's dependent on the language *of the person reading them*. An English-speaking user will want to see a list of Swedish names (or a mixed list of Swedish and English names) sorted in English alphabetical order.

All of this doesn't mean that you have to write a whole new string comparison routine for every language you might encounter. Instead, the normal approach is to use a rule-driven algorithm, where you have a generalized sorting routine and you pass it some sort of table of information explaining how to sort text in a particular language.

In its simplest form, you'd just map each character to an abstract weight value appropriate for the specified language and sort the strings according to the weights of their characters instead of the code point values. But the mapping of characters to weights is actually more complicated than this, even in non-Unicode encodings such as Latin-1.

Multi-level comparisons

For starters, you actually need more than one weight value. Consider this list:

car
CAT
Cat
cAT
cAt
caT
cat
CAVE

You want to sort according to the letters without taking the case differences into account, so that "car" sorts before all the "cat"s and "cave" sorts after them. But you don't want the order of the "cat"s to be determined by the whim of your sorting algorithm. The "cat"s aren't identical—you want

a defined ordering of some kind. For example, maybe the capital letter should come before its small counterpart.

Your first thought might be that you can get this effect by interleaving the weight values for the capital and small letters—for example, you put “C” before “c”, but both of them after “b” and before “D”. This doesn’t work. If you do that, you get this order:

```
CAT
CAVE
Cat
car
cat
cAT
cAt
caT
```

Note how “car” and “cave” wind up in the middle of the “cat”s instead of on either side of them.

What you have to do to get this effect is assign each letter *two* weight values: a *primary* weight and a *secondary* weight. The primary weights are the same for the capital and small versions of each letter. For example, both “A” and “a” have a weight of 1, “B” and “b” have a weight of 2, “C” and “c” have a weight of 3, and so on. The secondary weights, on the other hand, are different for the different cases. And since you already have the primary weights to differentiate things, you can actually reuse the same secondary-weight values for each letter. All of the capital letters have a secondary weight of 1, and all of the small letters have a secondary weight of 2.

You use these as follows: You compare two strings using their primary weights only. If they’re different, you’re done. **Only if they’re the same** do you go back and compare them again using their secondary weights. This way, a primary difference later in the string still beats a secondary difference earlier (e.g., “Cat” comes after “car” even though “c” comes after “C” because “r” comes before “t” and the difference between “r” and “t” is more important than the difference between “c” and “C”). This produces the ordering in the first example above.

To cover all languages, you generally need three or more weight values (or “levels”) per character. For most languages, those characters that are considered “different letters” are given different primary weights. Versions of the same letter with different accent marks attached to it (or other variants of the letter, such as “v” and “w” in Swedish) are given different secondary weights. Versions of the same variant of the same letter with different cases are given different tertiary weights.

Again, you compare the strings first using their *primary weights only*. If they compare equal, you break the tie by comparing them again using their secondary weights. If they’re still equal, you break the tie again using their tertiary weights. If they’re equal at this level, they really are equal and you either don’t care about their ordering or you differentiate based on a secondary key.

You don’t actually have to make three complete passes through the strings to compare them this way, of course. You can compare the weights at all levels in a single pass. It looks something like this:

```
public static int compare(char[] a, char[] b) {
    int firstSecondaryDifference = EQUAL;
```

```
int firstTertiaryDifference = EQUAL;
int i = 0;

while (i < a.length && i < b.length) {
    char cA = a[i];
    char cB = b[i];
    if (primaryWeight(cA) > primaryWeight(cB))
        return A_GREATER;
    if (primaryWeight(cB) > primaryWeight(cA))
        return B_GREATER;
    if (firstSecondaryDifference == EQUAL) {
        if (secondaryWeight(cA) > secondaryWeight(cB))
            firstSecondaryDifference = A_GREATER;
        else if (secondaryWeight(cB) > secondaryWeight(cA))
            firstSecondaryDifference = B_GREATER;
        else if (firstTertiaryDifference == EQUAL) {
            if (tertiaryWeight(cA) > tertiaryWeight(cB))
                firstTertiaryDifference = A_GREATER;
            else if (tertiaryWeight(cB) > tertiaryWeight(cA))
                firstTertiaryDifference = B_GREATER;
        }
    }
    ++i;
}
if (a.length > b.length)
    return A_GREATER;
if (b.length > a.length)
    return B_GREATER;
if (firstSecondaryDifference != EQUAL)
    return firstSecondaryDifference;
return firstTertiaryDifference;
}
```

Ignorable characters

Defining different levels of difference doesn't by itself give us the right answer for everything we've looked at so far. Consider the different variant spellings of "cooperate" (with a couple other words thrown in to clarify things):

concentrate
co-operate
cooperate
coöperate
copper

What weights do we give the hyphen to end up with this order? There's no primary weight you can give the hyphen that'll give you the right answer (short of making the hyphen sort between "n" and "o," which will produce goofy results when the hyphen appears in other words). What you want to do is give the hyphen *no* primary weight at all: At the primary level, this character should simply be transparent to the sorting algorithm. In other words, the hyphen should be *ignorable* at the primary level.

The hyphen is still there, of course, and still affects the string comparison. It's just that it doesn't become significant until we're looking for secondary differences. It's ignorable at the primary level, but not at the secondary level.

Generally, you deal with this by declaring some weight value as indicating "ignorable." If you encounter this in your comparison routine, this signals you to skip the character that gave you this weight value. This, of course, means you're no longer comparing characters in the same positions in the two strings. This makes it a lot harder to do all three levels in a single pass, but it's still possible. (Don't worry: it gets worse.)

Also, it doesn't make sense to treat a character as, say, ignorable at the secondary level but not at the primary level. Pretty much all Unicode comparison routines declare that a character that's ignorable at any particular level is also ignorable at all more-significant levels (e.g., a character that's ignorable at the secondary level is automatically also ignorable at the primary level; it might still be significant at the tertiary level, however).

French accent sorting

The unspoken assumption in everything we've looked at so far is that differences (at the same level) earlier in the string beat differences later in the string: "bat" comes before "car" because "b" comes before "c." The fact that "r" comes before "t" doesn't count because we found a difference earlier in the string.

Consider the French words "cote" (meaning "quote" or "quantity"), "coté" ("quoted" or "evaluated"), "côté" ("side"), and "côte" ("hillside").¹²¹ You'd normally expect them to sort like this:

```
cote
coté
côte
côté
```

The accented letters sort after their unaccented counterparts. But this is the English order for these words. In French, they sort like this instead:

```
cote
côte
coté
côté
```

Note the difference: In the first example, the two words with the unaccented "o" sort together, as do the two words with the accented "o." In the second example, the two words with the accented "e" sort together, as do the two words with the unaccented "e."

That's because in French, if words differ in spelling only in the placement of accents, accent differences *later* in the string are more important than differences *earlier* in the string—you start the comparison at the end and work your way back to the beginning.

¹²¹ Thanks to Alain LaBonté for suggesting the example.

In text-processing circles, this phenomenon is usually called “French accent sorting” or “French secondary order” even though it actually happens in a couple of other languages as well (such as Estonian and Albanian).

Contracting character sequences

In Spanish, the digraphs “ch” and “ll” actually count as separate letters. That is, “ch” isn’t a “c” followed by an “h”; it’s the letter “che.” “ch” sorts between “c” and “d”, and “ll” sorts between “l” and “m.” So this list of Spanish words is in alphabetical order:

casa
cero
clave
como
Cuba
chalupa
charro
chili
donde

In this case, you have a single set of weights assigned to a *pair* of letters rather than a single letter.¹²²

If you’re dealing with an encoding (such as Unicode) that uses combining character sequences, this would also be how you get an accented letter to be treated as a different letter from its unaccented counterpart. For example, in Swedish, “a” comes at the beginning of the alphabet and “ä” comes toward the end of the alphabet. If “ä” is represented as “a” followed by a combining diaeresis, you have the same phenomenon: two characters get mapped to a single set of weights that cause the combination to be ordered differently from either character individually. This is called a “contracting character sequence.” Contracting character sequences are also essentially what Unicode 3.2 is talking about when it talks about “language-specific grapheme clusters.” For the purposes of string comparison, a contracting character sequence is the same thing as a grapheme cluster: a sequence of multiple Unicode code points that behaves as a single character for the purposes of string comparison.

One interesting side effect of this is that the relative order of two strings can change if you append characters to the end of both of them or take substrings from the beginning. If a concatenation causes characters at the end of the original string to be involved in a contracting character sequence, or if taking a substring breaks up a contracting character sequence, order may change.

Expanding characters

The reverse also happens: you can have a single character map to a whole sequence of weight values. For example, if you want “archaeology” and “archæology” to sort together, then the “æ” needs to be treated as equivalent to “a” followed by “e.”

122 As an aside, because for a long time computer systems couldn’t handle contracting character sequences, Spanish speakers have gotten used to seeing “ch” and “ll” sorted the same way they sort in English. This is now called “modern Spanish sorting” and the version that treats “ch” and “ll” as single characters is called “traditional Spanish sorting.”

(Note here we say “equivalent” and not “the same.” You still want these strings to be considered different so that you get a defined ordering between them regardless of sort algorithm. Normally this is handled by having “æ” map not to “a”’s weights followed by “e”’s weights, but to a set of weights that has the same primary weight as “a” but a different secondary weight, followed by “e”’s weights [you don’t have to do anything special with the “e” weights, since you already made things different with the first set of weights].)

Expanding characters happen a lot in German. “ß” is equivalent to “ss.” You also often see it with the umlauted vowels: A vowel with an umlaut can also be spelled with the unadorned vowel followed by an “e,” so “ä,” “ö,” and “ü” sort as though they were “ae,” “oe,” and “ue.”

Context-sensitive weighting

Occasionally you’ll see a character behave differently for comparison purposes depending on the surrounding characters. The classic example of this is the long mark in Katakana. A vowel followed

by a long mark is equivalent to the vowel twice in a row. For example, アー is equivalent to

アア.

In other words, if a contracting character sequence is a many-to-one mapping and an expanding character is a one-to-many mapping, this is a many-to-many mapping. You treat each sequence involving the long mark independently: アー maps to two sets of weights that sort with a secondary difference from アア, イー maps to two sets of weights that sort with a secondary difference from イイ, and so on.

Putting it all together

So how do you implement all this? The way to think about it is that you do the comparison by creating a “sort key” from each of the strings being compared. A sort key is simply a sequence of weight values. To produce it, you go through the string and look up the primary weight for each of the characters, writing it to the sort key. Ignorable characters don’t result in anything being added to the sort key, expanding characters result in multiple weights being added to the sort key, and contracting sequences result in single weights corresponding to multiple characters in the string.

After you’ve derived the primary weights for all the characters, you write a separator value to the sort key, then you derive all the secondary weights, write another separator, and derive all the tertiary weights. (If you need to do French secondary sorting, you just write the secondary weights into the sort key backwards.) After you’ve done this with both strings, you’ve got a pair of sequences of values that can be compared with traditional bitwise comparison.

The idea behind the separator character is that you pick a value that’s lower than any of the actual weight values (usually you use zero and make sure your lowest weight value is 1). This way, if you’re comparing two strings where one is a proper substring of the other, you’ll wind up comparing the level separator of the shorter string, and not one of the secondary weights, to one of the primary weights for the longer string

An individual weight value in a sort key is called a “sort element” or “collation element.” In fact, a “collation element” is usually the concatenation of all the weight values for a particular character or sequence (i.e., the primary, secondary, and tertiary weight values). Instead of pulling them apart and arranging them in sequence so you get all the primary differences before the secondary and tertiary differences, a sort routine might build up a sequence of full collation keys and go through them multiple times, masking out all but a particular weight on each pass.

Generally, you don’t actually have to create the sort keys (although they can indeed be useful; we’ll look at that later). You just need to do the work to create them unit by unit until you find a difference.

Other processes and equivalences

There are a number of other things you might need to take into consideration when doing searching and sorting:

- When sorting titles, you usually don’t want to consider articles such as “a” or “the” at the beginning of a string. For instance, “A Farewell to Arms” usually sorts in the Fs, and “The Catcher in the Rye” usually sorts in the Cs. Of course, which words you disregard when they appear at the beginning of a string vary with language.
- Often, strings containing numbers are sorted as though the numbers were spelled out. For example, “9½ Weeks” would sort in the Ns, and “48 Hours” would sort in the Fs.
- Another choice with strings that contain numbers is to sort them number parts in numeric order so, for example, “10th St.” sorts after “1st St.” and “2nd St.” instead of between them.
- Sometimes you want abbreviations to be expanded so that, for example, “St. Paul” sorts before “San Diego” (or, for that matter, “Saint Tropez”), as though it had been spelled “Saint Paul.” (Of course, this can be a little tricky in many cases. In our example, “St.” is also an abbreviation for “Street,” so you’d need to be able to tell from context how to expand it.)
- When sorting lists of people’s names, you usually want to sort by the surname (which can be complicated, as the surname isn’t always the last word in a person’s name—this varies depending on language and, often, depending on the person).
- If a list includes words (people’s names, for example) in different scripts, you might want to sort everything by its spelling in one script in particular (for example, you might want “Михайл Горбачев” to sort in the Gs, just as if it had been spelled “Mikhail Gorbachev.”)
- When searching, you might want to catch variant spellings of the same word, especially the same word in different scripts (for example, Chinese using either simplified or traditional Han characters, or possibly even bopomofo or pinyin; Japanese in Hiragana, Katakana, Kanji, or possibly even Romaji; Korean in either hangul or hanja; Serbo-Croatian in either Latin or Cyrillic; Yiddish in either Latin or Hebrew letters; Mongolian in either Cyrillic or the traditional Mongolian script, etc.).

These kinds of things are typically not handled in the base string-comparison routine, but usually are done with some kind of preprocessing, such as passing strings through transliteration prior to putting them through the string-comparison process, or storing alternate representations that are used only for searching or sorting and not for display (a movie-title database might, for example, have an extra title field in each record that stores “Nine and a Half Weeks” as an alternate representation of “9½ Weeks,” or “Catcher in the Rye” for “The Catcher in the Rye,” and sort on this field instead of the regular title field).

Language-sensitive comparison on Unicode text

To the above considerations, which you have to deal with regardless of which encoding standard you're using to encode your characters, Unicode adds a couple of other interesting complications.

Unicode normalization

Unlike most other encoding schemes, many characters and sequences of characters have multiple legal representations. And one of the requirements of supporting Unicode is that (provided you support all of the characters involved) all representations of a character get treated the same. So, whether you represent “ä” with

```
U+00E4 LATIN SMALL LETTER A WITH DIAERESIS
```

or with

```
U+0061 LATIN SMALL LETTER A  
U+0308 COMBINING DIAERESIS
```

it should still look the same and behave the same everywhere. This means that a Unicode sorting facility must compare U+00E4 and U+0061 U+0308 as equal. It gets better with characters that have multiple combining marks on them: All five of the following sequences...

```
U+1ED9 LATIN SMALL LETTER O WITH CIRCUMFLEX AND DOT BELOW
```

```
U+1ECD LATIN SMALL LETTER O WITH DOT BELOW  
U+0302 COMBINING CIRCUMFLEX ACCENT
```

```
U+00F4 LATIN SMALL LETTER O WITH CIRCUMFLEX  
U+0323 COMBINING DOT BELOW
```

```
U+006F LATIN SMALL LETTER O  
U+0302 COMBINING CIRCUMFLEX ACCENT  
U+0323 COMBINING DOT BELOW
```

```
U+006F LATIN SMALL LETTER O  
U+0323 COMBINING DOT BELOW  
U+0302 COMBINING CIRCUMFLEX ACCENT
```

...should compare the same. (All are alternate representations for the Vietnamese letter ô.) Of course, the way you do this is to convert both strings being compared into one of the Unicode normalized forms, meaning all of the above sequences get converted either to the first one in the list (in Normalized Form C) or to the last one (in Normalized Form D). In practice, form D (fully decomposed) is the more common way to go, both because you have to go through Form D on your way to Form C (see the previous chapter) and because it's usually simpler to treat the accents independently (usually as characters that are ignorable at the primary level) than as part of their base letters (which might involve treating more characters as expanding characters).

Unicode normalization is usually done on the fly, a character (or combining character sequence) at a time, as part of the process of mapping characters to collation elements.

The requirement that equivalent representations be treated as identical applies only to sequences that are canonically equivalent (i.e., identical when mapped to Normalized Form D). Often, it makes sense to extend this requirement to compatibility equivalents as well (i.e., sequences that are identical when mapped to Normalized Form KD but not when mapped to Normalized Form D). Since mapping a character to its compatibility decomposition may lose data, compatibly-equivalent strings usually aren't treated as identical; the difference between a compatibility composite and its compatibility decomposition is usually considered a tertiary-level difference.

Korean is one interesting exception. Both sets of Hangul characters are arranged in the correct order such that you can use straight binary comparison to sort a list of strings in Hangul and get the right answer (of course, this breaks down if any of the strings contain non-Hangul characters). In fact, this is representation-independent: you get the same answer whether the strings all use the precomposed syllables or the conjoining jamo. This means you may not have to bother with converting Hangul syllables to Normalized Form D, which can be a useful performance optimization. Of course, if some of the strings use precomposed syllables and others use conjoining jamo, or if any of the strings use a combination of precomposed syllables and conjoining jamo, then you'll only get the right answer if you actually do convert everything to Normalized Form D first.

Because two Unicode strings can compare as identical without actually being bit-for-bit identical, Unicode-oriented comparison algorithms often include a fourth level of comparison. If two strings are identical all the way down to the tertiary level, their binary representations may be used to impose an ordering on the strings.

Reordering

The other interesting wrinkle Unicode brings to the party has to do with Unicode's representation of Thai and Lao. Most of the Indic scripts in Unicode are stored in logical order. This means that vowel marks always follow the consonants they attach to, regardless of which side of the consonant they attach to.

Thai and Lao, on the other hand, are stored in visual order. Top-, bottom-, and right-joining vowels follow the consonants they attach to in the backing store, but left-joining consonants precede the consonants. Split vowels with left- and right-joining components are represented with two components. One, representing the left-joining half, precedes the consonant, and the other, representing the right-joining half, follows it. The difference in encoding philosophy between Thai and Lao and the other Indic scripts has to do with the difference in encoding philosophy in the national standards the different Unicode blocks were based on; changing the encoding philosophy for Thai and Lao would have made interoperability between Unicode and the national standards more difficult and would have gone against standard practice, so Unicode has the inconsistency. (This is one of those situations where optimizing for truth and beauty would have led to fewer people using the standard.)

What makes this fun is that Thai and Lao still *sort* as though they were stored in logical order. (Strictly speaking, this issue isn't unique to Unicode, as the Thai national encoding standard shares it as well.) This means that if you encounter a Thai left-joining vowel followed by a Thai consonant, you have to exchange them before mapping them to collation elements. (The collation element for the consonant precedes the collation element for the vowel in the resulting sort key.)

[If someone can provide me with an appropriate list of Thai words, an example would probably be helpful here.]

Since your mapping tables need to support many-to-many mappings anyway (this usually comes for free once you're supporting one-to-many and many-to-one mappings, both of which are unavoidable), you could deal with this simply by including all of the pairs of characters that need to be exchanged as many-to-many mappings in your mapping table. But there are a lot of them and this would make the table very large, so most systems just handle this algorithmically.

A general implementation strategy

So putting together a good general-purpose string comparison algorithm for Unicode is actually quite complex. At the core, you need mapping tables capable of one-to-one, one-to-many, many-to-one, and many-to-many mappings. Usually, the table is optimized for one-to-one mappings. You either use a generalized trie structure or a simpler compact array augmented with exception tables to do the mapping. The main trie handles single characters and contracting sequences, and each maps either to a single collation element or to a spot in an auxiliary table of multiple-collation-element mappings (used for expanding characters and many-to-many mappings).

Generally speaking, all three weight values for a character can be crammed into 32 bits, so collation elements are usually 32-bit values. (In most real implementations, a fourth level, corresponding to the actual binary representation of the string is added, but since the fourth weight value in this case is the character's original hex value, it doesn't need to be stored).

There are two basic strategies for mapping raw Unicode text to collation elements. One is to perform decomposition and reordering in a separate preprocessing step and have the mapping table just map from normalized Unicode to collation elements. The other approach is to combine all of these operations into a single operation and use the mapping table for all of them.

Using a single table for everything can be faster, but results in a really big mapping table (the table would have to deal with canonical accent reordering, so it'd need to have separate mappings for all five representations of `ô`, whereas a table that depends on normalized Unicode might only need mappings for the `o`, the circumflex, and the underdot). Unless you're only handling a small subset of the Unicode characters, it typically makes more sense to handle decomposition and reordering in a separate step or pass.

Once you've mapped everything to collation elements, you then need to make another pass through everything to split apart the different weight values in the collation elements (so all the primary weights get compared before any of the secondary weights, etc.) and possibly to reverse the order of the secondary weights (if you have to do French accent sorting). Only then can you actually compare the sort keys to get an actual result.

So, to recap, language-sensitive comparison on Unicode strings proceeds in four conceptual phases:

- Decomposition and reordering
- Mapping normalized, reordered Unicode text to collation elements
- Splitting the weights in the collation elements and possibly reversing the secondary weights to produce collation keys
- Performing binary comparison on the collation keys to generate an actual result

Let's examine how this works on a pair of sample strings, say "Resume" and "résumé." The two strings start out like this:

Resume: 0052 0065 0073 0075 006D 0065
résumé: 0072 00E9 0073 0075 006D 00E9

The first thing you do is map them both to Normalized Form D:

Resume: 0052 0065 0073 0075 006D 0065
résumé: 0072 0065 0301 0073 0075 006D 0065 0301

Then you map to collation elements (the weights are separated by periods, and the values are illustrative):

Resume: [09CB.0020.0008] [08B1.0020.0002] [09F3.0020.0002] [0A23.0020.0002]
[0977.0020.0002] [08B1.0020.0002]

résumé: [09CB.0020.0002] [08B1.0020.0002] [0000.0032.0002] [09F3.0020.0002]
[0A23.0020.0002] [0977.0020.0002] [08B1.0020.0002] [0000.0032.0002]

Then you split out the weights (dropping the ignorable values) to get the sort keys (we use “0000” as a level separator):

Resume: 09CB 08B1 09F3 0A23 0977 08B1 0000 0020 0020 0020 0020 0020 0020
0000 0008 0002 0002 0002 0002 0002
résumé: 09CB 08B1 09F3 0A23 0977 08B1 0000 0020 0020 0032 0020 0020 0020
0020 0032 0000 0002 0002 0002 0002 0002 0002 0002 0002

Finally, you compare them:

Resume: 09CB 08B1 09F3 0A23 0977 08B1 0000 0020 0020 **0020** 0020 0020 0020
0000 0008 0002 0002 0002 0002 0002
résumé: 09CB 08B1 09F3 0A23 0977 08B1 0000 0020 0020 **0032** 0020 0020 0020
0020 0032 0000 0002 0002 0002 0002 0002 0002 0002 0002

Notice what has happened. The strings are the same at the primary level, and you end up with two sort keys that start with the same sequence of primary weights (the zero primary weights, representing the accents, which are ignorable at the primary level, get skipped when creating the sort keys). The first difference we see is the presence of the accent on the first “é” in “résumé,” which has a secondary weight of “0032.” This compares higher than the secondary weight of the “s” in “Resume” (“0020”), and so the strings compare different, with “résumé” coming second. The difference between the capital “R” in “Resume” and the small “r” in “résumé” winds up later in the sort key: the capital-ness of the “R” in “Resume” is represented by the “0008” in the first sort key. Even though this difference comes earlier in the strings, it comes later in the sort keys, since it’s a lower-priority difference.

Of course, for many strings some of these passes are no-ops, and the four passes can actually be interleaved so that a primary difference in the first character gets caught early rather than after you’ve jumped through all these hoops (or at least so that all four operations happen in a single pass with comparatively little use of temporary storage).

One nice thing about all this flexibility is you can base an implementation on UTF-16 (or even UTF-8) directly rather than having to convert things to UTF-32 before you can work on them. A supplementary-plane character can be represented with a UTF-16 surrogate pair and still get treated as a single character: It just gets treated as a contracting character sequence.

Another issue that's worth keeping in mind is that it usually doesn't make sense to write a string comparison routine whose behavior is hard-coded for just one particular language or group of languages. Typically string comparison is done as a rule-driven library, where you can provide rules to specify the exact sort order you're after. The Java `Collator` object, for example, takes rule strings of this form:

```
A , a < B , b < C , c < D , d < ...
```

The characters are listed in their desired order, with less-than signs used to indicate characters with primary differences and commas used to indicate tertiary differences. (Semicolons are used for secondary differences, and the ellipsis at the end is just there to indicate that the example would normally continue.) Most of the time, ordering rules are specified as deltas to a default ordering. The ampersand can be used to allow insertion of characters into the middle of the default rules. For example, the traditional Spanish sorting rules look like this:

```
& c < CH , Ch , cH , ch  
& l < LL , Ll , lL , ll  
& n < Ñ , ñ
```

The first line indicates that the “ch” digraph (in its different cased variations) sorts after `c` (and, by implication, before `D`, the next thing in the default order). The less-than sign indicates that “ch” has a primary difference from “c,” which makes the “ch” digraph behave as a completely different letter. The second line inserts the “ll” digraph between “l” and “M,” and the third line inserts `ñ` between `n` and `O`.

If you're doing this kind of generalized comparison routine, you've got basically two choices. You can either take a set of rules in a specialized metalanguage like the one above and build appropriate mapping tables on the fly, or you can simply choose from a set of pre-built mapping tables. The latter approach gives you better speed (because you don't have to take time to build the tables) and compression (because the developer can compress the tables at his leisure and just serialize out the result), but diminishes flexibility. Building the tables on the fly from rules in a metalanguage has the opposite set of tradeoffs. Some systems, such as the International Components for Unicode, actually allow for both approaches.

The Unicode Collation Algorithm

Anyway, all of this discussion of how you'd handle language-sensitive comparison in Unicode leads us nicely into our next topic. The Unicode standard actually defines a set of criteria for a Unicode-compatible string comparison routine. This is called the Unicode Collation Algorithm, or UCA for short.

I actually misspoke myself slightly in the above paragraph. Strictly speaking, the UCA isn't part of the Unicode standard. You can compare strings however you want and (as long as canonically equivalent strings compare equal) conform to the Unicode standard. The UCA is a companion standard, or a Unicode Technical Standard, UTS #10, to be exact. (The designation “10” doesn't mean there are nine other Unicode Technical Standards, by the way—the numbering of Unicode Technical Reports, Unicode Technical Standards, and Unicode Standard Annexes is consecutive

across all three types of documents, reflecting the fact that they were all once called Unicode Technical Reports.)

Actually, in the same way that Unicode is aligned with the ISO 10646 standard, the Unicode Collation Algorithm is aligned with ISO 14651, an international string-ordering standard. As with Unicode and 10646, the two standards are not identical, but implementations that conform to one generally conform to the other (or will with comparatively little work). One important difference: while Unicode is much more restrictive than 10646, 14651 is actually somewhat more restrictive than the UCA.

The Unicode Collation Algorithm calls for these features:

- It requires at least three distinct weighting levels, as described above, but provides for the possibility of more (most conforming implementations define a fourth level based on the original binary representations of the strings).
- It requires support for ignorable characters, contracting character sequences, expanding characters, and context-sensitive weightings.
- It strongly encourages, but doesn't require, French accent sorting. (The UCA calls this "backwards levels.")
- It provides for, but doesn't require, Thai vowel reordering.
- It provides for something called "alternate weighting," which we'll talk about in a minute.
- It specifies a default ordering for all Unicode characters. In the absence of tailorings, all UCA-conformant comparison implementations must sort any given set of strings into the same order.
- It provides for "tailorings," alterations to the default order normally used to produce the correct ordering for a particular language (the default ordering can generally get things right for only one language per script; this gives you the ability to get it right for other languages or customize the results for application-specific purposes). All UCA-conformant comparison implementations must, given the same set of tailorings, also sort any given set of strings into the same order.

The UCA follows the algorithm described in the previous section. You:

- convert everything to Normalized Form D,
- exchange Thai and Lao left-joining vowels with the characters that follow them,
- map everything to collation-element values (this may involve mapping single characters to multiple collation elements, mapping groups of characters to single collation elements, or mapping groups of characters to groups of collation elements in different ways than would happen if the characters were treated individually),
- create sort keys from the collation elements by reorganizing things so that all the primary weights come first (with ignorable characters omitted), followed by a separator value and all the secondary weights, followed by another separator value and all the tertiary weights, etc.,
- reverse the order of the secondary weights (or possibly some other level) if the tailoring you're using calls for it, and finally
- compare the resulting sort keys using straight binary comparison.

There is, however, one complication: Instead of just converting to Normalized Form D and mapping things sequentially to collation elements, you have to do one additional thing. You actually start by looking for the longest substring that has an entry in the mapping table, but before you map it, you check after it for extra combining marks. If you find any combining marks that:

- aren't separated from the substring you're considering by a non-combining character or a combining character with a combining class of 0,

- aren't separated from the substring you're considering by another combining mark of the same combining class, and
- when combined with the substring under consideration form a string that also has an entry in the mapping table,

you use that string instead and ignore the combining marks you were able to combine with the main character.

Why the extra work? Let's say your mapping table has a special mapping for the letter ô and you encounter the Vietnamese letter ơ in the string. In Normalized Form D, this decomposes to o, followed by the underdot, followed by the circumflex. You end up mapping each character to a collation element independently. The rule allows you to recognize ô as a sequence with its own mapping even with the intervening underdot. Instead of getting three collation elements, one for the o, one for the underdot, and one for the circumflex, you get two: one for ô, and one for the underdot.

The default UCA sort order

Other than specifying a good algorithm for comparing Unicode strings, the big thing the Unicode Collation Algorithm brings to the party is a default ordering for all the characters in Unicode. This is spelled out in a file called `allkeys.txt` which is available from the Unicode Web site. It's not part of the Unicode Character Database, since the UCA isn't officially part of the Unicode standard, but you can get to it via a link in UTS #10.

`allkeys.txt` consists of a list of mappings. Each mapping consists of a Unicode code point value (or, in some cases, a sequence of Unicode code point values), followed by a semicolon, followed by one or more collation key values. The Unicode code point value sequences are given (as in all the other Unicode data files) as space-delimited sequences of four-digit hex values. Each collation element is shown as a series of four period-delimited four-digit hex values in brackets. The four values represent the primary-level weight, the secondary-level weight, the tertiary-level weight, and the fourth-level weight respectively. So if you see an entry like this...

```
0061 ; [.0861.0020.0002.0061] # LATIN SMALL LETTER A
```

...it's telling you that the single Unicode code point value U+0061 (the small letter A, as the comment at the end of the line makes clear) maps to a single collation element. The primary weight is 0x0861, the secondary weight is 0x20, the tertiary weight is 0x02, and the fourth-level weight is 0x61 (which, not coincidentally, is the original character's raw code point value).

The contracting and expanding sequences in `allkeys.txt` are usually there to make sure the same thing happens with different normalization forms of the text. For example, here are two entries, one of which is a contracting character sequence:

```
04D1 ; [.0B05.0020.0002.04D1] # CYRILLIC SMALL LETTER A WITH BREVE
0430 0306 ; [.0B05.0020.0002.04D1] # CYRILLIC SMALL LETTER A WITH BREVE
```

The two entries make sure that the Cyrillic letter а sorts the same way in both its composed and decomposed forms (most of the time you can sort accented letters according to their decomposed forms—that is, you can just treat the letter and the accent mark as separate characters—but this is an example where the accent mark turns the character into a completely different letter in the languages that use it).

An expanding character's entry would look like this:

```
2474 ; [*0266.0020.0004.2474] [.0858.0020.0004.2474] [*0267.0020.0004.2474]
# PARENTHESESIZED DIGIT ONE; COMPATSEQ
```

Here the file is saying that the character U+2474, the parenthesized digit one, sorts as though it were a regular opening parenthesis, followed by a regular digit one, followed by a regular closing parenthesis (the weights are the same as for those characters, except for the tertiary-level weight on the first collation element, which preserves the fact that U+2474 isn't canonically equivalent to "(1)" written as three characters).

The entries in `allkeys.txt` are listed in order by the collation elements, although this isn't required. This makes it easy to see the order the characters sort into. The basic ordering boils down to this:

- Control and formatting characters, which are completely transparent to the sort algorithm (they're ignorable at all four levels).
- Various white-space, punctuation, non-combining diacritic, and math-operator characters, plus a bunch of other symbols.
- The combining diacritical marks, followed by a few more non-combining diacritical marks.
- Currency symbols, followed by a few more miscellaneous symbols.
- Digits and numbers, in approximate numerical order by digit (i.e., "10" comes between "1" and "2").
- Letters and syllables from the various scripts, in the following order: Latin, Greek, Cyrillic, Georgian, Armenian, Hebrew, Arabic, Syriac, Thaana, Ethiopic, Devanagari, Bengali, Gurmukhi, Gujarati, Oriya, Tamil, Telugu, Kannada, Malayalam, Sinhala, Thai, Lao, Tibetan, Myanmar, Khmer, Mongolian, Cherokee, Canadian aboriginal syllabics, Ogham, and Runic.
- Hangul (the file lists only the jamo, depending on canonical decomposition to get the precomposed syllables into the same order).
- Hiragana and Katakana (interleaved).
- Bopomofo.
- The Han ideographs and radicals. (They're not all listed; the regular Han ideographs have weights that can be derived algorithmically.)

The whole table can't be derived algorithmically, but many pieces of it can be, a fact that can be used to reduce the size of the table used in the implementation:

- Canonical composites (including the Hangul syllables) sort according to their decompositions. (There are exceptions to this in situations where it'll give the "wrong answer" from a linguistic standpoint. Generally in these cases, you'll see the decomposed version spelled out in the file as a contracting character sequence. See the Cyrillic example above.)
- Compatibility composites sort according to their decompositions, but the resulting collation elements have a different tertiary weight than the decomposed version would have. The tertiary weight is algorithmically derivable from the type of the decomposition as listed in the `UnicodeData.txt` file (e.g., "", "<small>", or "<compat>"). [This extra work imposes a definite ordering on, say, the different presentation forms for the Arabic letters.]
- The Han characters have their Unicode hex values as their primary weights, and 0x20 and 0x02 as their secondary and tertiary weights. The compatibility Han characters and the Han radicals are listed in the file and sort as equivalent to various regular Han characters.
- Code point values that shouldn't appear in the data stream at all, such as unpaired surrogate values or the value U+FFFF, should either cause an error or be treated as completely ignorable.

- Unassigned and private-use code points are given algorithmically-derived collation element values that cause them to sort after everything that's explicitly mentioned in the table (because there's no way to get this effect with a single code point value, they're all treated as expanding characters—the first collation element uses the first few bits of the Unicode hex value and adds it to a constant that puts the character at the end of the sequence; the second collation element uses the rest of the bits in the character's hex value).
- The fourth-level weight can be derived algorithmically. For characters with decompositions, the fourth-level weight is the original character value. For completely ignorable characters, it's zero. For Han and unassigned code points, it's usually 1. Generally speaking, for everything else, it's the original code point value.
- The tertiary weight can also usually be derived algorithmically. For compatibility composites, the weight is derived from the decomposition type. For characters with the general category “Lu” (uppercase letter), it's 0x08. For everything else, it's 0x02.

It's important to note that the weight values themselves are not normative; you can use any weight values you want, so long as they cause things to sort in the same order as the weight values given in `allkeys.txt`. The weight values in the table are specifically chosen so that they don't occupy overlapping numerical ranges (this allows you to dispense with separator characters in your sort keys) and so that you don't actually need sixteen bits for each weight value, as the table appears to imply. In fact, the secondary and tertiary weights can each be stored in eight bits, and the fourth-level weight can always be computed, so you can store a complete collation element in 32 bits.

In addition to the actual mapping entries in `allkeys.txt`, the file also has a heading at the top that specifies some global characteristics for the sort. This includes:

- Versioning info.
- A tag that specifies how to handle characters with alternate weightings (see below).
- A tag that specifies which levels need to be reversed in the final sort key (to allow French accent sorting—this is never done in the default sort order, but is here so that you can specify tailorings using the same file format as is used for the default order).
- A tag that lists which characters get swapped with the characters that follow them before being turned into collation elements (this should always be just the left-joining Thai and Lao vowels, but is included in the file in case future Unicode versions add more characters that need to be handled this way).

The UCA default order tends to lag a bit behind the actual Unicode version and is keyed to a particular version of Unicode. As I write this in December of 2001, the current version of the UCA is keyed to Unicode 3.1 and doesn't include mappings for the new Unicode 3.2 characters. It'll still produce a definite ordering for them, since there's a specific algorithm for giving weights to characters that aren't explicitly mentioned, but this does mean the order for those characters is likely to change when a future version of the UCA explicitly mentions them. (That will probably have happened by the time this book goes to press.)

Alternate weighting

You may have noticed that the collation-element values in `allkeys.txt` actually begin with an extra character. There'll be a period or an asterisk before the first weight value. For example, here's the entry for the space character:

```
0020 ; [*0209.0020.0002.0020] # SPACE
```

The UCA default ordering specifies a group of characters (the ones with the asterisks) that have “alternate weightings.” These are characters that, depending on the situation, you might want to treat as ignorable or you might want to treat as significant. Depending on the setting of a flag, a UCA-compliant comparison routine that implements alternate weighting (it’s an optional feature) can either use the weight values from the table as-is or perform various transformations on them. There are four choices for treatment of the characters with alternate weightings:

- Non-ignorable. The characters have the weights listed in `allkeys.txt` and have primary differences with each other and with other characters. Consider the following example:

e-mail
e-mail
effect
email
exercise

Punctuation comes before letters, so the strings that contain the hyphens sort before the strings that don’t. The variant spellings of “email” are separated. [For the purposes of this example, the longer dash is intended to represent the ASCII hyphen-minus character U+002D and the shorter hyphen the Unicode hyphen character U+2010.]

- Ignorable (or “blanked,” the term used in UTS #10). The character has zero for its primary, secondary, and tertiary weights and its original code point value as the fourth-level weight. This effectively makes the character transparent to the comparison algorithm unless it has to refer to the actual code point values to break a tie. The same set of words as our previous example comes out like this:

effect
e-mail
email
e-mail
exercise

The three variants of “email” come out together, but the punctuation marks are completely transparent, and the strings’ binary representations are used to break the tie. Because U+2010 (the hyphen) is greater than U+002D, the two punctuated versions of “email” come out on either side of the nonpunctuated version.

- “Shifted.” This is similar to “ignorable.” Characters with alternate weightings get treated in the same way as “ignorable,” but all other characters get a fourth-level weight of 0xFFFF. This gives slightly better-looking results when strings compare equal at the first three levels, basically ensuring that all otherwise-equal strings with an ignorable character in some position group together in the sorted order, rather than being interspersed with strings that don’t have an ignorable character in that position.

effect
e-mail
e-mail
email
exercise

The three variants of “email” come out together, again. Again, the hyphens are transparent at the first three levels, but here we have a slightly better algorithmically-derived fourth level to break the tie. The shifting algorithm make sure the two versions of “email” with the hyphens come out grouped together.

- “Shift-trimmed.” This is the same as “shifted,” but the characters that don’t have alternate mappings have a fourth-level weight of 1 (or no fourth-level weight at all). This has the same effect as “shifted,” except that the strings with ignorable characters sort *after* the strings that don’t have them, instead of before.

effect
email
e-mail
e-mail
exercise

The three variants of “email” come out together, but this time the two versions with the hyphen sort after the version without.

The default sort order is set up in such a way that the primary weights for all the characters that have alternate mappings come before the primary weights for all the characters that don’t (except, of course, those that are always ignorable at the primary level, which have a primary weight of 0 and don’t need to be treated specially). This means that you don’t actually have to take up any extra storage in your mapping tables to indicate which characters have alternate mappings. You can look up their primary weights and check the primary weight against a constant.

Optimizations and enhancements

There are a number of possible optimizations and enhancements worth thinking about:¹²³

Computing sort key values. When it’s possible to compute a sort key rather than looking it up in the table, you don’t have to store it in the table. As we saw, for the default Unicode ordering, you can get away without storing table entries for the Han characters, the Hangul syllables, unassigned and private-use code point values, canonical composites, and compatibility composites. In addition, for the default Unicode ordering, the and fourth-level weights can always be computed and don’t have to be stored. Except for a few rare exceptional characters, the tertiary weights can, too.¹²⁴

Computing things rather than storing them in the table sets up a classical performance-versus-memory tradeoff, but it’s frequently worth making this tradeoff in the direction of saving memory—there are other ways of improving performance.

Reducing the size of the weight values. The UCA data table actually gives secondary-weight values that are too large to be stored in eight bits. This is due to the algorithm that was used to generate these weights. However, there are never anywhere near 256 different secondary-weight values for any particular primary weight. You can take advantage of this to save space storing the secondary weights (by reusing values to mean different things depending on the primary weight—remember, if the primary weights are different, you never compare the secondary weights).

Increasing the repertoire of available weight values. The UCA assumes 16-bit values for all the weights, but what if you actually need more than 65,536 primary weight values? This can happen, for example, if you’re specifying an ordering for all of the Han characters in Unicode 3.1, both the ones in the BMP and the ones in Plane 2. You can accomplish this by treating some of the characters as expanding characters. You set aside some range of primary-key values as special. Some of the characters map to two collation elements. The first of these elements has one of the “special” primary keys (chosen so as to get the right ordering on its own when compared with a “non-special” primary

¹²³ Most of these optimizations come straight out of the “Implementation Notes” section of UTS #10, but some are taken from Mark Davis, “Collation in ICU 1.8,” *Proceedings of the 18th International Unicode Conference*, session B10.

¹²⁴ This is actually true only of the default UCA sort ordering; language-specific tailorings will often have tertiary-level weights that can’t be algorithmically derived.

key), and the second is used for disambiguation. Each “special” value thus gives you an extra 65,536 primary key values. The only trick here is making sure that all characters that have a particular “special” value for the first primary key sort the same way relative to all of the characters that have “non-special” primary keys, or different “special” primary keys.

That was a little vague. Let’s consider a very simple example. Let’s say you need 100,000 primary key values. The first 65,535 can be done with single collation elements with primary-key values from 0000 to FFFE. The remaining 34,465 characters map to two collation elements: the first one has a primary weight of FFFF, and the second one has a value from 0000 to 86A0. The secondary and tertiary weights for the second collation element can actually be zero—this is one situation where it’s okay to have a collation element that’s “ignorable” at the second and third levels but not at the first.

Of course, it generally makes sense to use a more sophisticated technique than this one. For example, you might designate all odd primary-key values to be single-element values and all even values to be the leading element of a two-element value. You could then fix it so that more common characters mapped to single collation elements and less-common characters mapped to pairs of elements.

This technique can also be used, of course, to give you more than 256 secondary- or tertiary-key values without using more than 8 bits as well.

Dealing with tailorings. Much of the time (probably most of the time) you’re going to be dealing not with the UCA default ordering, but with a language-specific tailoring of the UCA default ordering. It’s really painful to just have a language-specific mapping table for each language that includes the parts of the UCA mapping table that didn’t change. If you set up your UCA table to leave “holes” between all the characters (using the technique described above, most likely), the tailorings can put things in the “holes.” If you do this, the mapping table for a tailoring doesn’t have to repeat all the default UCA stuff: you can look up a collation-element mapping in the tailoring’s mapping table, look it up in the UCA table if you don’t find it in the tailoring table, and generate it algorithmically according to UCA rules if you don’t find it in the UCA table.

Optimizing the comparison. You actually don’t have to go through all the steps of the UCA right from the beginning of the string. You can use binary comparison to start with and only drop into the UCA when you get to a character that’s different between the strings. The thing that makes this hard is that you may hit the difference in a non-leading code point of a contracting character sequence or a code point that’s going to be affected by normalization (such as by accent reordering). In these cases, you have to back up to the beginning of the characters that need to be considered as a unit. For instance, if the difference happens on a combining character or a low surrogate, you’d back up in both strings to the first non-combining character or the high surrogate. You also may have to keep around a table of characters that occur in non-leading positions in contracting character sequences (it wouldn’t have to include combining characters, since you can deal with these algorithmically) and back up until you get to a character that isn’t in this table. Only then can you drop into the full UCA to compare the rest of the strings. Despite the complication, this optimization can actually buy you quite a bit in performance.

Avoiding normalization. In many cases, you actually don’t have to do any normalization. If both strings being compared are already normalized (in either Form C or Form D), you can compare them without any transformation. If you encounter a sequence of characters that isn’t normalized (or a mixture of the two forms), you can drop in and do normalization then, on just that piece. This requires an extra table to tell you when you need to do normalization.

Inverting cases. Whether capital letters sort before small letters or vice versa is often a matter of personal preference. The UCA default order sorts small letters before capital letters. You can flip this algorithmically for the untailed UCA default order easily if you're generating the tertiary weights algorithmically. The normal rule is the tertiary weight is 0008 for the characters in the "Lu" category and 0002 otherwise (if it's a compatibility composite, you do this check first and the result is one of your inputs for figuring out the tertiary weight). Just change this so that you get 0008 on "Ll" and 0002 the rest of the time.

For tailored orderings, this is trickier. The simplest approach is to simply map each character to its opposite-case counterpart when generating the tailoring table (or generate the tailoring table the normal way and then swap the resulting collation-element mappings of any case pairs after the fact).

Another approach is to leave the tables and collation-element generation alone and flip the cases of the input characters as part of the normalization/reordering pass.

Reordering scripts. Sometimes you may want the relative order of the characters in a particular script to stay the same, but change the relative order of the different scripts. You can do this in a similar manner to how you deal with switching the cases. The simplest method is to pass a "raw" primary-weight value through an extra lookup table that rearranges the primary weights according to the desired order of the scripts.

Language-insensitive string comparison

Of course, you don't always want language-sensitive string comparison. After all, it's complicated and no matter how much you optimize it, it'll never be anywhere near as fast as binary comparison. Of course, if binary comparison gives you the wrong answer, it doesn't matter how fast it is. But there are times when you don't care about linguistically-correct ordering. You just want *some* kind of ordering, and the user will never see the results of that ordering.

There are two classical versions of this. One is where you only care about whether two strings are equal. If they aren't, you're not interested in which one comes before which. "Equal" can also be a language-sensitive issue, but if you're dealing with non-natural-language strings (such as programming-language identifiers or XML tags), that doesn't matter.

You can use straight binary comparison in this case, but only if you know (or are in a position to demand) that the input strings are both in the same normalized form. (The W3C character model [see Chapter 17] requires all text to be in Normalized Form C, for example.) If you can't depend on that, you still have to do Unicode normalization on the strings but can blow off all the other work.

The other case is where you're doing some kind of ordered index to speed up searching. For example, you're building a binary tree or you need a sorted list that you can do a binary search on. The common case, of course, is indexed fields in a database. In these cases, you need a well-defined order, but the exact properties of the order don't matter because the user will never see it. Again, you can use straight binary comparison here, as long as you can guarantee the strings will all be in the same normalized form (or if your index can have multiple entries for strings that are supposed to be equal). In the worst case, you'll again only have to worry about Unicode normalization and can avoid all the other apparatus.

Of course, if you want the comparisons to be case insensitive but still don't care about linguistically-appropriate results (again, program identifiers are the classic example), you have to pass things through a case-folding table (see Chapter 14) and compare the results rather than doing straight binary comparison on the strings. (The `DerivedNormalizationProperties.txt` file in the Unicode Character Database provides a set of mappings that let you do both case folding and normalization at once.) But again, you can avoid all the other string-comparison baggage.

There's one important thing to keep in mind when performing binary comparisons on Unicode strings. The exact binary format makes a difference. Specifically, if you sort a list of strings encoded in UTF-8 or UTF-32, you won't get them in the same order you would have gotten if they had been in UTF-16. The supplementary-plane characters will sort after the BMP characters in UTF-8 and UTF-32, but between U+D7FF and U+E000 in UTF-16 (this is because the surrogate mechanism uses values between 0xD800 and 0xDFFF to encode the supplementary-plane characters). This is vitally important to keep in mind if you're depending on binary comparison in a mixed-UTF environment.

There are two main ways of dealing with this, short of going to using the same UTF throughout your system. The simplest is probably to modify the comparison routine for one UTF to replicate the ordering of the other UTF. For one UTF you do normal binary comparison. For the other, you use a modified algorithm that produces the same results as the other UTF (for example, you use a modified algorithm on UTF-16 that causes the values from U+D800 to U+DFFF to sort in a block after U+FFFF). Such a thing is trivially easy to write and can be made almost as fast as regular binary comparison.

Another way of dealing with the problem that's become fairly common is to use a modified form of UTF-8 (this is the "CESU-8" encoding form described in DUTR #26) instead of the regular form of UTF-8. CESU-8 represents supplementary-plane characters using six-byte sequences instead of the normal four-byte sequences. The six-byte sequence is effectively the result of mapping from UTF-16 to UTF-8 using a UTF-32-to-UTF-8 converter: each surrogate gets independently converted into a three-byte UTF-8 sequence, rather than the surrogate pair getting treated as a unit and being converted to a four-byte UTF-8 sequence. (Encoding supplementary-plane characters this way is prohibited in real UTF-8: it violates the shortest-sequence rule and can be a security violation.)

This can work, with some loss of storage efficiency, as long as the outside world never sees the bowdlerized UTF-8 (or at least, it's not advertised as real UTF-8). The system still needs to be able to handle real UTF-8 coming from outside the system and to send real UTF-8 to the outside world. But sending the bowdlerized form to the outside world and advertising it as "real" UTF-8 is a Very Bad Idea. It's almost always better to fix the sorting routines.

It's also worth mentioning compressed Unicode, which can be very helpful for saving space in database implementations. SCSU is great for this, but because of its stateful nature, binary ordering based on SCSU will be wildly and unpredictably different from the binary order of uncompressed Unicode. There's an interesting alternative compression scheme for Unicode called BOCU that preserves the binary ordering of the uncompressed text that might be useful for compressing text in fields that also need to be indexed. (See Chapter 6.)

Sorting

There's not a lot to say about sorting per se, since you can use any sorting algorithm you wish to sort a list of Unicode strings. There are a number of things worth keeping in mind, however.¹²⁵

Collation strength and secondary keys

First, it's worth thinking a little about your sort algorithm and what it does with keys that compare equal. Sort algorithms such as the bubble and insert sorts are *stable*: records with identical keys will remain in the order of their original insertion after sorting the list. Fast sort algorithms such as the Quicksort and the merge sort are generally *unstable*: records with identical keys are not preserved in insertion order during a sort. In effect, the sort algorithms will put records with identical keys in random order.

It's important to keep this in mind, since the Unicode Collation Algorithm can wind up treating many strings that aren't bit-for-bit identical as equal, and an unstable sort algorithm will arrange them randomly. Of course, if the whole record *is* the key, this probably doesn't matter, since even if the strings aren't bit-for-bit identical, they're the same for all intents and purposes (for example, they'll typically look the same when drawn on the screen).

In other cases, though, it's worth giving this some thought. If you don't want things scrambling gratuitously when, for example, you add a new item to the list, you'll either want to define your comparison algorithm to treat strings different if it possibly can (this means treating differences at all levels as significant and possibly doing a bitwise comparison to break the tie if the language-sensitive comparison says the strings are equal) or fall back on a secondary key.

Usually disambiguating things with a secondary key makes more sense. In fact, if you're breaking ties using a secondary key, you may only want to treat primary-level differences in your primary key as significant.

Be careful doing this though, because it won't always get you the results you're really looking for. Let's say that your primary key is someone's last name and your secondary key is their first name. If you treat fourth-level differences in the last name as significant before going to the first name, you might end up with something like this:

van Dusen, Paul
van Dusen, Tom
vanDusen, Paul
vanDusen, Tom
Van Dusen, Paul
Van Dusen, Tom

(Okay, so it's a contrived example.) The first names wind up interleaved: each variant spelling of "Van Dusen" gets treated as a group in itself. (The versions with the spaces wind up separated

¹²⁵ Much of the material in this section also comes from the Davis paper, *op. cit.*, and from UTS #10.

because if the space is ignorable, differences in spaces are a fourth-level difference and the case differences are a tertiary difference.)

So let's say you decide to fix this by only treating primary differences in the last name as significant before going to the first name. Now the first names sort together...

van Dusen, Paul
vanDusen, Paul
Van Dusen, Paul
Van Dusen, Tom
vanDusen, Tom
van Dusen, Tom

...but if you're using an unstable sort algorithm, the last names come out in random order. You group by first name, but within each group, the variant spellings of the last name come out scrambled.

To get around this, you actually have to use a composite key. You might store the first and last name separately, but for sorting purposes, you consider them as a single key, concatenating the first name onto the end of the last name. *Then* you treat all your differences as significant. This way, primary differences in the first name beat lower-level differences in the last name, but you still consider lower-level differences in both names if both names are equal at the primary level. This gives you a reasonable result:

van Dusen, Paul
vanDusen, Paul
Van Dusen, Paul
van Dusen, Tom
vanDusen, Tom
Van Dusen, Tom

The "Van Dusens" correctly group by first name, but they sort into the same order within each group.

One interesting thing to watch out for here, though: Be careful about what you use as a field delimiter when you're putting together your composite keys. In the default UCA order, the comma, like the space, has alternate mappings. If you treat characters with alternate mappings as ignorable, as we did in the preceding example, shorter last names won't necessarily sort before longer last names they're a substring of. You'll go to the first name prematurely, possibly giving you something like this:

Van, Bill
van Dusen, Paul
Van Dusen, Paul
van Dusen, Tom
Van Dusen, Tom
Van, Steve

Here, Bill Van and Steve Van sort on either side of the Van Dusens, rather than grouping together before them, because the first letters of their first names are getting compared against the D in "Van Dusen." To prevent this, you've either got to tailor your sort order or be sure to pick a field delimiter that's not ignorable (and sorts before the letters). (Since you're just using this for internal processing

and not displaying it, it doesn't much matter what you use for the field delimiter for your composite keys, so long as it produces the right sorting results. The dollar sign fills the bill, for example.)

Exposing sort keys

Up until now, we've talked about sort keys mainly as an internal construct you use to get things to sort in the right order. You calculate them on the fly when you compare two strings.

But it's also frequently worth it to keep sort keys around rather than just creating them (or, more often, partially creating them) when you do a comparison. This is because it's pretty expensive to create a sort key. If you're doing a lot of comparisons involving the same strings, it's a lot cheaper to put the sort keys together once and then do all the comparisons directly on the sort keys using a simple binary comparison.

Most of the time when you compare two strings on the fly, there's a difference in the first few characters, so you only have to go through all the work for a few characters rather than the whole string. So unless there's a lot of similarity in the strings you're comparing, on-the-fly comparison will be faster unless you're doing a *lot* of comparisons. But if you're sorting a long list or you do a lot of searches on a long list, or you have a lot of similar strings in the list, the balance shifts toward keeping the sort keys around. Most Unicode-aware string comparison APIs do give you a way to create a sort key from a string in addition to giving you a way to compare two strings on the fly.

There are a few things to keep in mind if you're going to persist sort keys (i.e., not just create them long enough to do a long sort, but store them in a database or file). The most important is that it's completely meaningless to compare sort keys generated from different sort tables. For example, if you compare a sort key generated from a Spanish sort order to one generated from a German sort order, the result is meaningless. This is because the same weight values may very well mean different things in different sort orders. Say, for example, that a German sort order assigns primary weights of 1, 2, 3, 4, and 5 to A, B, C, D, and E. A Spanish sort order might instead assign these weights to A, B, C, CH, and D, since "ch" is a separate letter in the Spanish alphabet. You can see how comparing sort keys using generated from these two sets of weights would produce goofy results.

A system that always uses the same values for the UCA default ordering and puts tailorings in the "cracks" between will produce less goofy results, but still doesn't really work. "Ch" will sort differently depending on whether it's in a string that had a Spanish sort key or one with a German sort key.

This sort of falls into the "Doctor, it hurt when I do this!" category ("Then don't do that!"), but it can also arise between different versions of the same sort order (for example, going from sort keys based on the Unicode 2.1.8 version of the UCA default order to sort keys based on the Unicode 3.0 version of the UCA default order). One way to avoid this is to append a version tag to the beginning of every sort key. You can use a hash-code generation algorithm to create a hash code based on the actual contents of the mapping table and use this as the version tag. It's then simple to detect a version mismatch. Or if you just compare the sort keys, you'll at least wind up with a bunch of sublists, one for each collation table used to create any of the sort keys.

Minimizing sort key length

The other thing that becomes important if you're persisting sort keys is minimizing their size. There are a number of strategies for minimizing sort key size:

Eliminating level separators. You don't have to waste space putting level separators in your sort keys if you use distinct ranges of numeric values for the primary, secondary, and tertiary weight values (this doesn't really work for the fourth-level weights, so you usually still need a separator here). However, this tends to make some of the other optimizations more difficult.

Shrinking the secondary and tertiary weights. It doesn't make sense to use a 16-bit value for the secondary and tertiary weights, since we know we can usually use 8-bit values for both. To make this work, you still need a level separator after the primary weights, since you may need all the 8-bit values for the secondary and tertiary weights (you might still be able to blow off the one between the secondary and tertiary weights, but since this only saves a byte, it's probably not worth worrying about). The level separator after the primary weights needs to be two bytes, since the primary weight values are two-byte values, unless you make sure that you don't use the separator-byte value as either byte in a primary weight value.

Avoiding zero. If you're using C-style zero-terminated strings to store your sort keys, you have to avoid using the byte value 0 anywhere in the sort key except at the end. You can use 1 or 0x0101 as a level separator (this is why the tertiary weights in the UCA table start at 2), and you have to make sure your primary-weight values don't include 00 as either byte. (This knocks out 512 possible primary-weight values.) Avoiding 00 as either byte of a fourth-level key value is tougher—you could add one to each byte of a fourth-level key value and use two bytes to express the byte values 0xFE and 0xFF.

Run-length encoding. Run-length encoding doesn't really buy you much at the primary level, where there doesn't tend to be much repetition, but it can buy you a lot at the secondary and tertiary levels, where there's often lots of repetition. In particular, the "neutral" tertiary weights—the ones that get used for uncased or unaccented characters—tend to show up a lot in most sort keys.

UTS #10 describes one way of going about this. At each level, you take all the weight values higher than the one that repeats (all values higher than 0x20 for the secondary weights, and all values higher than 0x02 for the tertiary weights) and move them to the opposite end of the byte. The highest weight value becomes 0xFF, the second-highest weight value becomes 0xFE, and so on. This opens up a "hole" in the middle above the value that repeats. Each value in this "hole" represents some number of repetitions of the value that repeats. For example, if we're talking about the tertiary weights, 0x03 represents 02 02, 0x04 represents 02 02 02, and so on.

Actually, it's a little more complicated than this. The range that represents repeating values actually gets divided in half. For each number of repetitions of the repeating value, there are actually two different byte values, a "high" value and a "low" value.

It works like this. Let's say we're talking about the tertiary weights. The value that repeats is 0x02, and the highest tertiary-weight value is 0x1F. The values 0x00 and 0x01 get represented the normal way (in the standard UCA algorithm, 0x01 would be the level separator and 0x00 wouldn't happen). The values 0xE3 to 0xFF represent 0x03 to 0x1F, and the values from 0x03 to 0xE2 represent varying repetitions of 0x02. 02 02 can be represented by either 0x03 or 0xE2, 02 02 02 by either 0x04 or 0xE1, 02 02 02 02 by either 0x04 or 0xE1, and so on. Whether you use 0x04 or 0xE1 to

represent 02 02 02 02 depends on the value that comes next. If it's higher than 0x02, you use 0xE1. If it's lower than 0x02, you use 0x04.

Of course, it's possible to have more repetitions of the repeating value in a row than can be represented with a single byte. You handle this in the obvious way: Use as many of the value representing the maximum number of repetitions as you need, followed by one more for the remainder. For example, if 0x80 is the highest "low" value, representing 127 repetitions of 0x02 (you can put the partition between "high" and "low" values wherever you want), you'd represent 300 repetitions (when followed by a value lower than 0x02) of 0x02 with 0x80 0x80 0x2D.

This technique will never give you longer sort keys than the uncompressed format, and can often result in tremendous amounts of compression. (And, of course, they sort the same way as the uncompressed versions, as long as everything follows the compressed format.)

Searching

The other main thing that string comparison is used for is searching.¹²⁶ By now, it should be apparent that there are many situations where sequences of Unicode code points that aren't bit-for-bit identical should still compare equal, and thus should be returned as a hit from a text searching routine. Not only do you have to deal with things like Unicode canonical and compatibility equivalents, but you also have to deal with linguistically equivalent strings (for example, if you're searching for "cooperate," you probably would want it to find "coöperate" and "co-operate" as well).

By now, you've probably figured out that the way to do this is not to match Unicode code points, but to convert both the search key and the string being searched to collation elements and match them instead.

In fact, searching is where different levels of equivalence really come in handy. If you want to ignore case differences, for example, you can declare that anything that's equivalent down to the secondary level matches and ignore tertiary-level differences. To ignore accent variations as well, you'd treat only primary-level differences as significant and ignore both secondary-level and tertiary-level differences.

Doing this requires catching the collation process at a slightly earlier point than you do for sorting. Instead of going the whole way and converting the search key and the string being searched into sort keys, you stop a step earlier when they're still sequences of collation elements. That is, you skip the step where you segregate all the weight values into separate parts of the sort key and (possibly) reverse the order of the weights at the secondary level. Instead, you catch things while all of the weights are still stored together in one numeric value. Typically, you convert the search key into a sequence of collation elements and just convert the string being searched into collation elements a character at a time while you're searching.

Ignoring secondary- or tertiary-level differences in this situation is simple: they occupy certain ranges of bits in the collation element, so you can just mask out the parts of the collation element you're not interested in. In our previous examples, a 32-bit collation element is divided into a 16-bit primary

126 Much of the material in this section is drawn from Laura Werner, "Unicode Text Searching in Java," *Proceedings of the Fifteenth International Unicode Conference*, session B1, September 1, 1999.

weight with 8-bit secondary and tertiary weights. If you're only interested in primary differences, you just mask off the bottom 16 bits of the collation elements before doing your comparison. You also have to be careful of ignorable characters: if, after masking off the parts you don't care about, you end up with zero (or whatever value is used to signal an ignorable character), you skip that collation element without using it in the comparison.

There is one additional complication: If you find a match, it doesn't count as a match unless it begins and ends on a grapheme cluster boundary: If accent differences are important, you don't want a match if you're searching for "cote" and find "coté" in the text, but if "coté" is in the text in its decomposed form, "cote" is a prefix of it, and you'll find it, even though you're not supposed to and wouldn't if "coté" was represented in memory in its precomposed form. So when you get a prospective match, you have to check whether it begins and ends on a grapheme boundary: "cote" wouldn't match "coté" because the matching sequence of code points ends in the middle of a grapheme cluster: it includes the "e" but not the accent. This is very similar to the work you have to do to implement a "Find whole words only" function, which we'll look at in a minute.

Depending on the API you're using, getting the collation elements in the form you need them may be difficult (unless, of course, you have a searching API). In Java, you can use the `CollationElementIterator` class to access the collation elements. The ICU libraries for both C and Java actually provide a searching API on top of the `CollationElementIterator` utility.

The Boyer-Moore algorithm

Sequential access to the collation elements in a string before they're split into the multiple zones of a sort key is all you need to do a simple brute-force string search in a Unicode-compatible linguistically-sensitive way. But there are, of course, better, more efficient ways to search a string for another string.

One of the classical efficient string-search algorithms was first published by Robert S. Boyer and J. Strother Moore of the University of Texas in a 1977 article in *Communications of the ACM*.¹²⁷ The Boyer-Moore algorithm takes advantage of the fact that each comparison actually tells you more than just that two characters don't match.

Let's say you're searching for the word "string" in the phrase "silly spring string."¹²⁸ In a traditional brute-force search, you'd start with the search key and the string being searched lined up at the first character and start comparing until you find a difference:

s	t	r	i	n	g														
s	i	l	l	y		s	p	r	i	n	g		s	t	r	i	n	g	

The "s"es compare equal, but then you fail comparing "t" to "i," so you slide "string" over one and try again:

¹²⁷ Robert S. Boyer and J. Strother Moore, "A Fast String Searching Algorithm," *Communications of the Association for Computing Machinery*, Vol. 20, No. 10, pp. 762-772, 1977.

¹²⁸ This example is lifted directly out of the Werner paper, *op. cit.*, because I couldn't think of a better example.

manageable. For Unicode, a 64K table is unwieldy, and a 1MB-plus table is even more unwieldy, which would seem to make the Boyer-Moore algorithm unworkable.

Using Boyer-Moore with Unicode

Actually, it's worse than 1 MB because you're comparing collation elements, not Unicode code points, and collation elements are 32 bits wide. A table with four billion elements is definitely unworkable.

Of course, the Boyer-Moore algorithm actually isn't unworkable with Unicode. The solution is one of those clever "Why didn't I think of that?" solutions. It was first discovered (and patented¹²⁹) by Mark Davis of IBM (actually Taligent at the time of the patent). The key is in how you handle a character that occurs in the sort key twice: you use the shorter shift distance.

When working with Unicode, you still use a 256-element table (actually, it can be any size you want: the principle is the same) and you use a hash function of some kind to map the collation element values down into this one-byte space. There are a number of different ways you can perform this hashing, but the simplest is to just isolate the low byte of the primary weight from the collation element and use that as your index into the table. More complicated hash functions can give you somewhat better distribution, but it's usually not worth the trouble.

If you have multiple collation elements with the same low byte in their primary weights, the entry in the table is the shortest shift distance for any of the collation elements that map to that table entry. As with characters that occur multiple times in the search key, this may mean you sometimes don't shift as far as you might have been able to, giving up a little performance, but you still get the right answer, and in real-world use you don't give up much performance.

There are two additional complications. One is ignorable characters. If they occur in the string being searched, they don't hurt. They may mean your shift distances aren't quite big enough, but you'll still get the right answer. If they occur in the search key, on the other hand, you have to filter them out before calculating your shift distances or you risk shifting too far.

The other fun thing is matching up the collation elements to the characters in the original string you're searching. If you treat the shift values as how many *code points* to move ahead, rather than how many *collation elements* to move ahead (which can be cheaper, since you can avoid mapping some characters to collation elements), you have to watch out for expanding characters and adjust your shift values to account for the possibility of skipping over expanding characters.

If, on the other hand, you go by collation elements, everything's already expanded and you don't have to worry about shifting too far. In this case, though, you have to watch out for false matches. You might find a match that starts on the second (or later) collation element of an expanding character. You'd need to do some extra bookkeeping to detect this condition.

¹²⁹ IBM has agreed to license the patented code freely. This is the approach used in ICU, which is an open-source project.

“Whole words” searches

Many times you want to search a document for a string, but you only want a match to count as a hit if the match is a whole word (that is, if you’re searching for “age,” you don’t want the system to tell you it found those letters in the middle of “cabbages”).

To avoid this, you can do one of two things. You can either check after you find a match to see whether both ends of the match fall on word boundaries, or you can adjust the Boyer-Moore algorithm so that when it shifts forward, it shift to the next word boundary following the position the shift table says to shift to. (In fact, you can check to see whether the end of the search key falls on a word boundary at the current position and know you don’t have a match without doing any comparisons at all, although doing one comparison usually helps more.)

But how do you tell whether a position in the text is a word boundary? There are basically two approaches to this problem: the pair-based approach and the state-machine-based approach. The pair-based approach looks at two characters and decides whether there’s a word boundary between them. The state-machine-based approach essentially parses the text (or some subset of it) to find word boundaries according to a state machine. It can consider more context than just two characters and can sometimes produce more accurate results (and, in fact, is required for some languages).

These are also the main approaches to deciding where to put line breaks when laying out a paragraph of text. Because of this, I’m going to defer a full treatment of these techniques to the next chapter when we talk about laying out lines of text.

Even when you’re not doing a “whole words” search, you have to be sure that a matching sequence begins and ends on a grapheme cluster boundary, as we saw earlier in the chapter. The same techniques that are useful for locating word boundaries are also useful for locating grapheme cluster boundaries.

Using Unicode with regular expressions

In addition to doing a normal search where you match a literal pattern to a document (with potentially varying degrees of looseness in what constitutes a “match”), text is often searched for text that matches a regular expression. A regular expression is an expression that describes a category of related strings. Any string in that category is said to “match” the regular expression. For example a regular expression like this...

```
p[a-z]*g
```

...might be used to specify a category consisting of all the words that start with p and end with g. “Pig,” “pug,” “plug,” and “piling” would all match this regular expression. If you were to search a string for this regular expression, the search would stop on any of those words (or any other word that began with p and ended with g).

I’m not going to explain how to write a regular-expression engine—there are enough books out there that do this. But as always, there are things to think about when implementing a regular-expression engine for Unicode.

Unicode Technical Report #18 gives a set of guidelines for writing a Unicode-compatible regular expression engine. It provides for three levels of Unicode compatibility, as follows:

- Level 1 is basic Unicode support. This is the minimal set of features necessary to make a regular expression engine usable with Unicode.
- Level 2 goes beyond Level 1 by adding in a deeper understanding of how Unicode text is structured.
- Level 3 provides complete locale sensitivity of the kind we've been looking at for regular searching and sorting.

There isn't the room to go into lots of detail on exactly what each of these levels mean, but here's a brief summary. Level 1 requires or recommends the following:

- The engine must (when operating on UTF-16 text) correctly work on UTF-16 code units instead of treating them as pairs of bytes.
- The regular-expression language must have some way to refer to a Unicode character by its code point value (for example, you could specify the Greek capital letter sigma with something like `\u03A3`). Since not all systems can display (or allow you to type) every character in the Unicode repertoire, this is essential.
- The language must have some way to specify Unicode categories. In an ASCII-based regular expression language, it's feasible to express "all letters" by saying `[a-zA-Z]`, but since there are thousands of "letters" in Unicode (tens of thousands if you count the Han characters), this isn't feasible. You need a special token that at least lets you specify groups of characters with the same Unicode general category.
- The language needs some way to remove characters from a set of characters. This would let you say things like "all Latin consonants" by saying something like `[a-z]-[aeiou]`. This goes hand in hand with the previous need for categories. Without it, there's no reasonable way to say things like "all currency symbols except the cent sign."
- If the engine supports case-insensitive matching, it's got to be smart enough to handle all the Unicode case equivalences. For Level 1 support, single characters that map to multiple character when uppercased don't have to be supported, but multiple mappings do. A case-insensitive match to the Greek letter sigma, for example, has to match both lowercase forms.
- If the regular-expression language supports matching of line boundaries, the support has to conform to the Unicode newline guidelines, recognizing as newlines not only all of the various ASCII and Latin-1 characters that can be used as newlines, but also LS and PS (U+2028 and U+2029) and the CRLF sequence.

Level 2 adds the following features to Level 1:

- It must allow specification of supplementary-plane characters by code point value and must treat them as first-class characters, even if it's actually operating on UTF-16 text where they're stored as surrogate pairs.
- It must correctly match canonically-equivalent strings.
- Its concept of a "character" should match a normal user's concept of a "character" (the Unicode documents tend to refer to this concept as a "grapheme"). This basically not counting as matches apparent matches that don't begin and end on grapheme cluster boundaries.
- It should be fairly intelligent about word boundaries. In particular, combining marks and formatting characters should be ignored when examining pairs of characters to determine whether they constitute a word boundary.
- Case-insensitive matches should work correctly with characters (such as the German letter β) that can turn into multiple characters when uppercased.

Level 3 goes whole hog and provides fully language-sensitive support. This basically builds on top of the collation-element stuff we've been looking at in most of the rest of the chapter. It's a lot

slower than the language-insensitive stuff we've been looking at so far, and should generally be an option rather than a requirement.

Level 3 support involves:

- Mapping the Pi and Pf categories to Ps and Pe in a way that's appropriate for the specified locale. (Basically this means taking the ambiguity out of the way Unicode handles quotation marks.)
- Treating as a “grapheme” anything that speakers of the specified language would consider to be one. Basically this means dealing not only with the generic Unicode graphemes of Level 2, but also treating as “graphemes” anything specified as a contracting character sequence in the locale's sort order. In Spanish, for example, “ch” would be considered a single grapheme.
- Respecting the locale's conception of what a “word” is. For languages such as Thai or Japanese, this may involve some fairly sophisticated linguistic analysis.
- Allowing for locale-dependent loose matching. This is essentially providing the ability to allow sequences with only second- or third-level differences according to the sort order to match if desired. For example, “v” and “w” are variants of the same letter in Swedish, so a pattern containing “v” should still match a “w” in the text (in Swedish).
- Range specifications should be according to the language's sort order, not the order of the Unicode code point values. This means, for example, that [a-ø] would have a different meaning in English (where it'd only match about half the alphabet) and Swedish (where it'd match the entire alphabet).

CHAPTER 16 *Rendering and Editing*

If the second-most-important process that gets performed on text is string comparison (the basis of sorting and searching), by far the most important process that gets performed on text is rendering: that is, drawing it on a computer screen or printing it on a piece of paper. A lot of the time this goes hand in hand with providing a user with some way of entering text into the computer in the first place and with some way of editing the text once it's been entered.

Rendering and editing text are both enormously complicated subjects, and we're not going to go into all the gory details here. Again, most of the time you'll be able to take advantage of the facilities provided by your operating environment to allow you to render and edit Unicode text. And much of the complexity isn't unique to Unicode, but is inherent in the various scripts Unicode can represent and there to be dealt with no matter what encoding scheme you're using to represent those scripts. But as with the other processes we've looked at, there are enough interesting things about rendering and editing that are unique to Unicode that it's worth taking a long look.

There are essentially three main categories of things that have to be done when rendering text, regardless of encoding, all of which become more interesting when Unicode is your encoding. They are:

- Splitting a large body of text up into lines and paragraphs. This can be as simple as looking for spaces in the text and breaking a line at the last space that comes before the desired right margin, but for many scripts it's much more complicated.
- Arranging the individual characters on a line of text. Again, this can be as simple as putting them in a single-file line running from left to right, but can get a lot more complicated, especially when the Unicode bi-di algorithm must be used.
- Figuring out which glyph to draw for each character. Often, this is a simple one-to-one mapping, but for many scripts the rules for deciding which glyphs to draw for a character or series of characters are much more complicated.

In addition to these considerations, we'll take a brief look at some of the considerations involved in editing Unicode text.

Line breaking

If you're drawing more than a few words, the first thing you have to do when drawing text on the screen is to break it up into lines. Usually a text-rendering process will maintain something called a *line starts array* that simply lists the offset in the backing store of the first character of each line of rendered text. In order to draw a block of text, the first thing you need to do is build that line starts array.

Dividing a block of text up into lines generally proceeds in three basic conceptual phases:

- Locate all the positions in the text where a line break *has* to go. This divides the text into a series of blocks. For the purposes of this discussion, we'll call these blocks "paragraphs," although it's possible in most word-processing programs to have forced line breaks within a paragraph.
- For each paragraph that can't fit entirely on one line, locate all the positions in that paragraph where it would be legal, according to the rules of the text's language, to break a line. (For example, it's not okay to break a line in the middle of an English syllable, and it's not okay to break a line in the middle of an English word unless you're capable of hyphenating it.) This divides the paragraph into a series of short units. Again, for the purposes of this discussion, we'll call these units "words," even though you can run both into situations where two or more real words form a single unit for the purposes of line breaking and situations where a single real word breaks into multiple units for the purposes of line breaking.
- Assemble the words into lines. The simplest and most common algorithm for this is to pack as many words as possible onto the first line, pack as many of the remaining words as possible onto the second line, and so on until you've assembled all of the words in the paragraph into lines. More complicated algorithms are also possible, and in fact professional typesetting software usually uses more sophisticated algorithms than the simple "maximal packing" algorithm we're all familiar with.

As with the other processes we've looked at, these three conceptual phases can be (and usually are) interleaved, and it's usually possible to abbreviate some of the phases (for instance, you can usually get away without finding every single word boundary in a paragraph you're dividing into lines).

For that matter, because it can be expensive to calculate the entire line starts array for a long document, you usually see line starts calculated on an as-needed basis: typically, the line starts are only calculated for the part of a document that's visible on the screen (often with some slop on either side to make scrolling faster) and more line starts are calculated as more text is scrolled into view.

In addition, most text-processing software will optimize the calculation of line starts so that only a minimal set of line starts are recalculated when there's a change to the text. (In fact, often line *lengths* are stored instead of actual line *starts*, since you can get away with updating less stuff this way.) For instance, a localized change such as the addition or deletion of a single character or word will only affect the line starts in a single paragraph. In fact, if you're using the maximal-packing algorithm, it only affects the line starts that follow the location of the change: you can start at the line of the change and recalculate line starts until you either get to a line whose starting position didn't change or to the end of the paragraph. (Algorithms other than the maximal-packing algorithm usually require re-wrapping the whole paragraph, or propagating in *both* directions from the change, since they consider more context.)

Line breaking properties

Most text-layout issues are language-specific (or transcend language) rather than being unique to Unicode, although by allowing a user to mix languages or scripts, certain issues may come up more often with Unicode than with other encoding schemes.

The most interesting piece of the line-breaking process from the standpoint of dealing with Unicode is the process of boundary analysis: dividing a document up into paragraphs and dividing a paragraph up into words. As with most other processes, the Unicode standard gives both rules and guidelines for doing this analysis on Unicode text. You'll find these rules in Unicode Standard Annex #14. (I say "rules and guidelines" because some of the material in UAX #14 is normative and some is just informative.)

UAX #14 sets forth a set of character properties that describe how a line breaking process should treat the characters. Some of these properties are normative; most are informative. The complete listing of which characters have which line-breaking properties is in the file `LineBreak.txt` in the Unicode Character Database.

UAX #14 basically divides the entire Unicode repertoire into several broad categories:

- Characters that are always followed by a line break ("breaking characters")
- Characters that may never be preceded or followed by a line break ("glue characters")
- Characters that stick to the characters that precede them and are otherwise transparent to the algorithm ("combining characters")
- Spaces, which affect the treatment of other characters
- Characters whose treatment is dependent on context
- A few special categories

The first category, breaking characters, is the most important, and is normative. These characters force a line break to happen afterwards, no matter what. This is how you divide a document up into paragraphs. The characters in this category mark the end of a paragraph. The characters in this category include LS, PS, and FF. LF, CR, and the CRLF combination also fall into this category, but because the CRLF combination gets treated as a single breaking "character," CR and LF get their own special categories: basically LF behaves the same as the other breaking characters, and CR also does, unless it's followed by LF. A break can never occur between CR and LF.

All of the other categories are concerned with how you divide up paragraphs into words for the purposes of line breaking. (Remember we're using "paragraph" and "word" to refer to units of text that are divided and assembled by a line-breaking algorithm—they usually, but not always, map onto the common-sense definitions of "paragraph" and "word".) For the rest of this section, we'll use the phrase "word break" to refer to a boundary between "words" for the purposes of line breaking—that is, a "word break" is a place in the text where it's legal for a line break to happen.

Glue characters are the next most important category. There is never a word break on either side of a glue character. Basically, it's these characters' job to suppress any word breaks that would otherwise occur (that is, they "glue" characters or words together). This category basically includes the WORD JOINER, all the characters with "NO-BREAK" in their name (NBSP, ZWNBS, etc.), plus a handful of special-purpose characters that also have "gluing" semantics.

The third normative category contains the combining characters. This includes not only the characters you'd expect (the combining and enclosing diacritical marks, the characters in the “Mc” and “Me” categories), but also a few other things. Specifically, the non-initial conjoining Hangul jamo (i.e., the conjoining vowels and final consonants) also fall into this category, as do all of the Indic vowel signs. The basic idea is that all of these characters are basically non-leading characters in combining character sequences. For the purposes of breaking, the whole combining sequence should be treated as a single character, with the first character in the sequence determining the sequence's behavior. Thus, all of the characters in the sequence are “glued” together, and the sequence behaves as a unit.

Interestingly, the invisible formatting characters also get treated as combining characters. Of course, they're not really combining characters, but they get treated the same way because they're each intended only to have an effect on one process and to be invisible to all the others. The easiest way to treat a character as invisible for line breaking is to treat it the same way as a combining character. (The exceptions, of course, are ZWSP and WORD JOINER, the invisible formatting characters whose whole job is to affect line breaking.)

UAX #14 also calls out a few special characters for normative purposes. The zero-width space (ZWSP) always precedes a word break. Other spaces can affect the behavior of surrounding characters in interesting ways. The object replacement character relies on out-of-band information to determine how it behaves for word-breaking purposes. Surrogate pairs stick together like other combining character sequences. (This rule shouldn't strictly be necessary, since surrogates aren't really characters. A UTF-16 surrogate pair should behave according to the rules for the character it represents.)

The other line-breaking categories are all informative and are designed, like the default Unicode Collation Algorithm ordering, to give a reasonable default line-breaking behavior for every character in Unicode. Without going into detail on all the categories, this is the approximate behavior they're designed to achieve (for full details, see UAX #14 itself):

- There's generally a word break after a series of spaces, although this is dependent on the characters on either side of the run of spaces.
Basically, the one exception to the rule about word breaks after spaces (other than the normative behaviors we already looked at) has to do with punctuation in some languages. Some punctuation marks in some languages (for example, the question mark in French) are separated from the words that surround them with spaces on both sides, like this:

Parlez-vous français ?

Even though there's a space, there shouldn't be a word break between most punctuation marks and the word that precedes them.
- There's generally a word break after a series of dashes or hyphens, except when they precede a numeral (in this case, a dash or hyphen might be being used as a minus sign). [The big exceptions are the em dash, which has a word break both before and after, and the Mongolian hyphen, which has the word break before instead of after.]
- Word breaks don't occur in the middle of a sequence of dot-leader characters.
- Certain punctuation marks, such as periods or commas, have a word break after them except when they occur in the middle of a numeral. Other punctuation marks, such as the question and exclamation marks, always have a word break after them.
- Opening punctuation (such as the opening parenthesis) never have a word break after them, and closing punctuation never have a word break before. Since quotation marks are often ambiguous, they don't have word breaks on either side.

- In the languages that use the Han characters, what constitutes a “word” is difficult to figure out and often ambiguous, so it’s actually okay to break a line in the middle of a word. So the Han characters and the Han-like characters (i.e., Hangul, Hiragana, Katakana, Bopomofo, etc.) generally have word breaks on both sides. However, there are various punctuation and diacritical marks (the ideographic period, small Kana, the Katakana long-vowel mark, etc.) that can’t occur at the beginning (or sometimes the end) of a line and “stick” to the preceding or following character as appropriate.
- Fullwidth presentation forms get treated the same as Han characters. A bunch of other characters that don’t have fullwidth variants (such as the Greek and Cyrillic letters) get treated either as Han characters or as regular Latin letters depending on their resolved East Asian width property.
- Thai and Lao are the other interesting case. Like Japanese and Chinese, they don’t use spaces (or any other mark) between words. But unlike them, you’re still only allowed to break lines between actual words. This can be accomplished by using the ZWSP, but better systems actually analyze the text to determine where the word boundaries are. UAX #14 just puts the affected characters into their own special category that indicates more information is needed to parse a run of these characters into words.

The important thing to remember here is that these category assignments are designed to give a reasonable default behavior. Sometimes (although not as often as with the UCA) tailorings are necessary to get good line-breaking behavior for a particular language. For example, the quotation-mark characters have language-dependent semantics (a particular mark might be an opening quote in one language, but a closing quote in another). Their behavior can be nailed down in a language-specific tailoring. Another example is Korean, which can be written either with or without spaces between words. If it’s written without spaces, Hangul syllables behave the same way as Han characters. If it’s written with spaces, they behave the same way as Latin letters.

Implementing boundary analysis with pair tables

So how do you actually implement code to do word boundary analysis?¹³⁰ The classical approach to this problem is to use a pair table. In its simplest incarnation, a pair table is a two-dimensional array of Boolean values. For each possible pair of adjacent character categories, the table tells you whether the position between them is a word boundary or not.

This, of course, is simple to implement: You use the techniques we’ve already looked at to map each character in the pair to a category and then use the categories to look up the answer in the pair table.

A straight pair-based algorithm can be pretty limiting, however, because it only considers two characters. For instance, you can’t correctly handle sequences of combining characters with a simple pair-based algorithm, nor can you correctly handle the case where a punctuation mark is separated from the preceding word by a space (without making simplifying assumptions, anyway). Usually, the pair-based algorithm is extended to cover situations like this: you purposely skip over a run of combining marks or spaces and look up the characters on *either side* of the run in the pair table. Instead of a Boolean, the table contains tags that tell you things like “there’s always a break here,” “there’s a break here only if there weren’t any intervening spaces,” “there’s a break here only if there *were* intervening spaces,” etc.

This is the approach suggested by UAX #14, which gives both sample code for implementing this algorithm and a sample pair table that implements the semantics it describes.

¹³⁰ Most of the material in this section is drawn from my own paper, “Text Boundary Analysis in Java,” *Proceedings of the Fourteenth International Unicode Conference*, session B2.

Even this can break down in some sufficiently complex situations, although those situations are relatively rare.¹³¹ You can continue to extend the pair-table algorithm to deal with the complex cases, but this requires more hard-coding and is less flexible.

To see how this works, let's say we're implementing a drastically simplified algorithm. Our algorithm has only three character categories: letters, whitespace, and sentence-ending punctuation. A "word" is a sequence of zero or more letters, followed by a sequence of zero or more whitespace characters, followed by a sequence of zero or more sentence-ending punctuation characters, followed by a sequence of zero or more whitespace characters (the idea is to implement something that will correctly handle French punctuation). If you implement this as a pair table, it looks like this:

	ltr	punct
ltr	*	
punct	X	*

Notice that we don't include the whitespace category in the pair table. It has to be handled algorithmically to get the effect we want. Also notice that there are three things that can appear in the cells: an asterisk, an X, or nothing. The asterisk means "put a break between the two characters only if there's intervening whitespace." The X means "always put a break between the two characters, intervening whitespace or not," and the empty cell means "never put a break between the two characters, intervening whitespace or not."

So to implement this, you need code that reads a character and takes note of its category, reads more characters until it reads a non-whitespace character and takes note of whether it saw any whitespace characters, and *then* looks up the categories of the two non-whitespace characters in the pair table, using the result (and whether or not there was whitespace between them) to decide whether to place the break there.

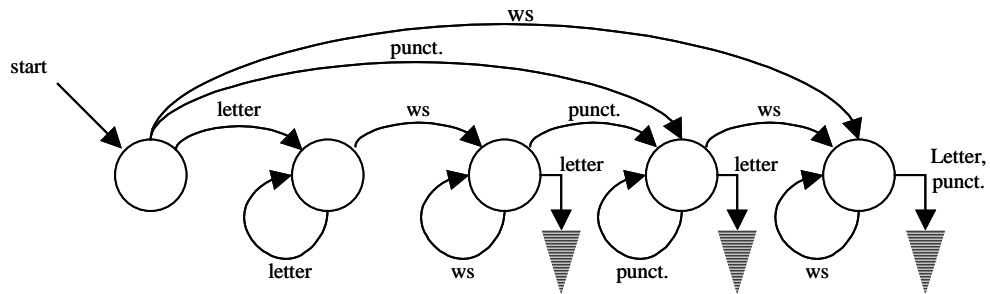
Implementing boundary analysis with state machines

An alternative approach is to recognize that locating word break positions is essentially a parsing problem and apply the traditional tools and techniques that are used for parsing text into tokens. Here, the traditional approach is to use a finite state machine.

In this approach, you still use a two-dimensional array, but this time only the columns represent character categories. Each row represents a state of the state machine. You start in some well-defined "start" state and then for each character, you look up its category and then use the category and the current state number to determine the next state. You continue to read characters, look up categories, and look up new states until you transition to a well-defined "end" state, which tells you the character you just read is the first character of the next "word" and you should put a word break between it and the preceding character. Then you start over in the start state with the first character in the next word (the character you just looked at) and do the whole thing again.

So if you use a state machine to implement our example algorithm, you get a set of state transitions that look like this:

¹³¹ In fact, the ones I can think of actually are related to other processes that use the same algorithms (we'll look at some of these later), not line breaking.



If you model this transition diagram as a two-dimensional array, you get this:

	ltr	ws	punct
start	1	4	3
1	1	2	3
2	STOP	2	3
3	STOP	4	3
4	STOP	STOP	4

Notice that we don't have to special-case the handling of the whitespace now; the state transition table takes care of that. This means you can consider however much context you need to determine a break position, and you don't have to write special-case code to handle the situations where you need to consider more than two characters.

Now the one thing that isn't clear from this example is what you do when you need to look at several characters after the break position to determine where the break position is going to go. The way we've got things set up above, when we transition to the "stop" state, the break position always goes before the character that caused us to go to the "stop" state. This means that you always have to be able to tell you're at a break position after looking only at the character immediately following it.

One simple way of doing this is to deviate a bit from having the parser be a pure deterministic state machine. You mark either each transition or each state (each state is almost always sufficient) with a flag. Instead of using your transition to the "stop" state as your indication of where to put the break, you use the flags. You start with a prospective break position at your starting point (or, in a lot of algorithms, one character after your starting point, to avoid endless loops). Each time you enter a flagged state (or follow a flagged transition, if you go that route), you update the prospective break position to be before the character you just looked at. The break goes wherever your prospective break position was at the time you transition to the "stop" state. This approach will let you consider multiple characters after the break position in the relatively rare circumstances where you have to do that.

Of course, the marking technique isn't perfect: it always moves the mark forward, and you can, in very rare situations, run into cases where you find, after looking at more characters, that you actually need to move the break position *backwards*. I have yet to run into a real-world situation where this has happened, but if you did, that'd probably mean it was time to switch to using a nondeterministic state machine.

One problem with using a state-machine-based approach is that you generally have to start at the beginning of the text: the state machine assumes you're in some ground state (generally a known word boundary) when you begin. If you have to parachute into the text in the middle rather than parsing the whole thing (we'll see why this is important in a minute), you actually need to back up from the random-access position until you find a position where you can safely start using the state machine to do the parsing without worrying about getting the wrong answer. (For pair-based approaches, you might need this too, but you only have to worry about things you've already had to special-case, such as the spaces in our earlier example.)

Typically you can locate the “safe place to turn around” with straight pair analysis: you identify the pairs of characters that *always*, no matter what, have a word boundary between them and seek backwards until you see one of these pairs. Occasionally, however, you might have to use a more sophisticated state-table-based approach to avoid seeking back to the beginning of the paragraph every time.

Performing boundary analysis using a dictionary

Correctly handling Thai and Lao is more difficult. These languages don't use spaces between words, but still require lines to be broken at word boundaries (as opposed Chinese and Japanese, for example, which also don't use spaces between words but permit you to break lines almost anywhere).

The classical way of handling this problem is to use a dictionary-based algorithm: you compare the text to a list of known words and use the matching sequences to parse the text into words.

This is fairly complicated. The simplest approach (find the longest word that matches the beginning of the text, put a word break after it, then repeat this process for the remainder of the text until you've used up everything) doesn't work most of the time. It might be, for example, that the first word in the text *isn't* the longest word in the dictionary that matches; it might be a shorter word that starts with the same letters. If you use the longest match, it may cause you to come to a sequence of letters that don't match any word in the dictionary.

If you were doing this on English without spaces, for example, and the sentence you were trying to parse was “themendinetonight”, you'd wind up trying “theme” as your first word and discovering this left you with something that began with “nd” as your second word.

What you actually have to do is find the set of words in the dictionary that enables you to use all the characters. This basically involves trying various sequences of words until you find one that fits—the approach is similar to a traditional maze-solving algorithm: try each sequence of turns one at a time until you find the sequence that gets you all the way through the maze. In our example, “theme” would lead you to “nd,” so you'd fall back on “them,” which would eventually lead you to “eto,” so you'd fall back on “the.” With a few more wrong turns and backtrackings, you'd eventually end up with “the men dine tonight.”

Of course, if there are *multiple* sequences of words that match the text you're parsing, you're stuck. You've either got to just pick one and hope for the best or use something like digram or trigram analysis (analyzing the relative frequencies of the various two- or three-word sequences in the text to figure out which sequence of words is the most likely, an approach that requires another table of digrams or trigrams and their frequencies).

You're also stuck if you hit a word that isn't in the dictionary, be it an uncommon word or just a typo. There are various fallback approaches that can be used here. You could assume you've seen a typo and use auxiliary information to guess what the word should have been. You can break the text up so that the longest series of letters you could parse successfully gets kept and you start over with the remainder of the text after skipping an offending character. Or you could fall back on algorithmic approaches: While it isn't possible to locate *word* boundaries in Thai purely algorithmically, you can locate *syllable* boundaries algorithmically.

Of course, you don't want to use the dictionary-based approach on everything. After all, punctuation and spaces do occur in Thai text and always mark word boundaries. Generally you use the techniques we looked at earlier in the chapter and then fall back on the dictionary only when you encounter a sequence of two or more Thai or Lao letters in a row.

A couple more thoughts about boundary analysis

It's probably evident that parsing Unicode text looking for logical boundaries is useful for things other than just line breaking. This is how you could tell, for example, whether a search hit falls on word boundaries or not (for a "find whole words" search). It's how you know what to select when a user double-clicks or triple-clicks at some arbitrary place in the document, and it's the approach you use to count words or sentences in a document.

One important thing to keep in mind, though, is that the exact definition of "word" varies depending on what you're doing with the information. For double-click selection or spell checking, you might need a dictionary-based algorithm where you don't for line breaking (such as for Japanese or Chinese). For "find whole words" searching, you want to exclude whitespace and punctuation from a "word," where you generally want to include these things for line breaking. For word counting, you might want to keep together sequences of characters (such as hyphenated phrases) that the line-break routine will divide up. And so on. A generalized algorithm that lets you substitute your own pair tables or state tables will let you do all these things without writing a whole new parser each time.

Performing line breaking

So once you've analyzed a paragraph into "words" (i.e., units of text that have to be kept together on one line), how do you use that information to actually divide the paragraph up into lines?

There are many different ways of doing this, but the most common (probably because it's the simplest) is to use the maximal-packing algorithm. The basic approach here is to fit as many words as possible on the first line, then fit as many of the remaining words as possible on the second line, and so on until you reach the end of the paragraph.

One obvious way to approach this is to interleave the process of parsing the text into words with the process of measuring it. You know how many pixels you have available, and you just parse into words until the total number of pixels exceeds that. Then you insert a line break and start over with the next word.

An alternative approach is to leap forward and then work your way back. If you're wrapping by number of characters instead of pixels (or you're working in a monospaced font where they correlate), you can skip forward the maximum number of characters that'll fit and then back up until you reach a word boundary. This approach can also work if you already know the widths of the text you're working with. If you've got a good idea of the average width of the characters in the font

you're using, you could jump forward to the spot in the text where the end of the line would be if every character was the average width, measure to see how close you got, and then work your way forward or backward from there, as appropriate, parsing the text into words.

One point to keep in mind is that you usually can't measure the text a character at a time to figure out how much text will fit on the line. This (generally) works for English, but for many other scripts, the characters interact typographically: the width of a character is dependent on the characters that surround it (because the glyph that will be displayed depends on the surrounding characters). This means you have to measure text at least a word at a time. Of course, this also means you have to perform glyph mapping on the text. In fact, since this may depend on rearrangement, you may also have to perform glyph rearrangement before you can break the text into lines accurately (this can be tricky, because glyph rearrangement depends on line breaking in some scripts—you may have to do speculative glyph rearrangement and then redo it if your assumptions turn out to be wrong).

Another point to keep in mind is that the widths of pieces of text may depend on where the line break actually goes: If the text includes the SOFT HYPHEN, for example, which is invisible unless it occurs at the end of a line, you have to make sure you take its width into account if the break is going to occur after it. (If your line breaking routine does hyphenation, you likewise have to take the width of the hyphen into account.) In fact, in some languages, such as Swedish, words actually have different spellings when they're hyphenated, so extra letters may be added (or deleted) before or after the hyphen in some languages.

Because of all these complications, the usual course is to measure *all* of the text that would go on the line if the line break were in some position (with all the necessary processes happening on the text to get an accurate measurement). You pick a spot for the line break, measure the entire resulting line, decide if it's too long or too short, skip an appropriate distance in the appropriate direction, try again, and so on until you've got an optimally-fitting line.

Hyphenation and justification are two techniques that often come into play when dividing text into lines. If you can't find a good fit, for example, you may have to hyphenate a word to get a better fit. Hyphenation in most languages involves the same kind of dictionary-based techniques used for regular line breaking in Thai. Other approaches may involve adding or removing inter-word or inter-character spacing to get the right margin to line up better (in languages like Thai when you don't normally have inter-word spacing, you may get a minimal amount of inter-word spacing adding by a justification routine to line up margins). Sometimes, alternate glyphs may actually be used for certain letters to get things to line up better (this works particularly well in Arabic, for example, where there are often a number of different ways pairs of characters may join together that yield different widths; Arabic justification also makes use of extender bars called *kashidas* to introduce extra space between certain pairs of letters when necessary to make a word wider).

Often, advanced techniques like these are used globally across all the lines in a paragraph to get optimal-looking results. For instance, a more-advanced line breaking algorithm might:

- Break in less-“optimal” places to avoid leaving a single word alone on the last line of the paragraph.
- Avoid hyphenation on certain lines to avoid “ladders” (i.e., several lines in a row that end with a hyphen).
- Deviate from maximal packing so as to keep the amount of inter-word or inter-character spacing more consistent from line to line of a justified paragraph.
- Deviate from maximal packing in a way that keeps the “color” (the approximate density of the ink) consistent across the paragraph.

- Allow the text to go slightly outside the margin on either side (this is called “optical justification”: round characters, for example, look lined up with the margin when they actually go a little beyond the margin and don’t look lined up when they actually are).

Line layout

Once you’ve decided which words are going to go together on a line, the next thing is to determine the order in which the characters will be arranged on that line.

For the majority of scripts, this is straightforward: The characters just march in succession, one by one, from the left-hand side to the right-hand side (or, in some languages, from the top of the page to the bottom). There are a couple of situations where it gets more interesting than that:

- The most common is combining character sequences. Usually, a combining character sequence turns into a single glyph, but sometimes (especially when there are a lot of base-diacritic combinations) the rendering process has to draw two (or more) separate glyphs and figure out how to position them relative to each other. Vowel and tone marks in Thai must be stacked, for example. Points in Hebrew and Arabic need to be positioned correctly relative to the base character (in Arabic in particular, the rules for how to do this can get pretty complicated in a high-quality font). And so on. This sort of thing is usually done in the font itself, and we’ll look at this more closely when we talk about glyph selection.
- Indic scripts may have complex positioning rules. Vowel marks don’t always position relative to the consonant that immediately follows them: if there are multiple leading consonants in the syllable, sometimes they all come into play (left-joining vowels and the left halves of split vowels, for example, usually go at the extreme left of their syllables rather than merely exchanging places with the consonant to their left). The exact placement of viramas, *rephas*, and diacritical marks often depends heavily on the surrounding letters. Again, most of this is usually taken care of in the font and will be looked at more closely in the next section. The one exception is left-joining vowels. On many systems, they’re just treated as separate characters and rearranged by the line-layout algorithm before glyph mapping.
- The big challenge in line layout, and the one we’ll spend the rest of the section looking at, is bi-di. If you’re dealing with one of the right-to-left scripts (currently Hebrew, Arabic, Syriac, and Thaana), there’s complicated work that goes on when these languages mix with left-to-right languages (or, much of the time, with numerals) on the same line. The Unicode bidirectional layout algorithm was designed to make sure you always get a well-defined ordering when scripts of different directionality are mixed on a line.

We already looked at the Unicode bi-di algorithm in Chapter 8, but we focused on the results it’s designed to produce. Here we’ll look at little more closely at how to implement this behavior.

Unicode Standard Annex #9 describes the bi-di algorithm in excruciating detail (it’s also explained in Section 3.12 of the Unicode standard itself, but UAX #9 is more up to date). The bi-di algorithm works by dividing up the line into *directional runs*, groups of characters that all run in the same direction (left to right or right to left). Within a directional run, the characters are laid out in the normal straightforward manner: the characters run one by one either from right to left or from left to right. The big challenge comes in figuring out how to order the runs relative to each other.

To do this, the bi-di algorithm organizes the characters into *embedding levels*. Each level consists of a sequence of directional runs and embedded levels. The items in each level are arranged relative to each other in a consistent direction (again, either left to right or right to left).

To see how this works, consider the following sentence:

Avram said **מזל טוב** and smiled.

It consists of three directional runs at two levels. The outermost level runs left to right and contains the two outer directional runs, “Avram said” and “and smiled,” and the inner level. The inner level, which runs right to left, contains the inner directional run, “מזל טוב.”

Without doing explicit embeddings, the deepest you can get is three runs: You can have a left-to-right outer level with right-to-left levels embedded within it. If a right-to-left level contains a numeral (numerals always run left to right), that’s an inner level that runs left to right. Combining character sequences in right-to-left text generally also get treated as inner levels that run left to right (this is so that after the character codes have been reordered for display, the combining marks still follow the characters they combine with). Of course, explicit embeddings can be used to get deeper levels of embedding, up to a maximum of 62 (which should be *way* more than enough for any rational piece of text).

To lay out a line, you need three buffers (actually, you can generally optimize some of them out, but bear with me for a few minutes). The first buffer contains either the characters, the glyph codes they map to, or sequence numbers mapping positions in the display back to positions in the backing store (or vice versa). We’ll call this the “character buffer.” The second buffer contains the resolved bi-di category for each character; we’ll call it the “category buffer.” The third buffer contains the resolved embedding level for each character. We’ll call it the “level buffer.”

The character buffer starts out containing all the characters in the line in their normal (logical) Unicode order. The category buffer starts out containing every character’s directional category as specified by the Unicode Character Database. The level buffer starts out as either all zeros or all ones, depending on the paragraph’s directionality. (The bi-di algorithm uses a convention that even-numbered levels are left-to-right levels and odd-numbered levels are right-to-left levels—if the paragraph’s directionality is left-to-right, everything starts out at level 0; if right-to-left, everything starts out at 1.) The paragraph’s directionality is usually specified with a property on the document or with styling information; in the absence of either, it defaults to the directionality of the first strong-directionality character in the paragraph.

The next thing you do is take care of any explicit embeddings or overrides. These can be specified either by out-of-band styling information or by the explicit directional characters in the General Punctuation block (that is, the LRE (“left-to-right embedding”), LRO (“left-to-right override”), RLE (“right-to-left embedding”), and RLO (“right-to-left override”) characters). When you see an RLE or RLO, you increment the embedding level of all characters between it and the next matching PDF (“pop directional formatting”) character to the next-higher odd value. For LRE and LRO, you increment the embedding level to the next-higher even value. (The sequences can nest; the incrementing happens at each nesting level.)

For the override characters, you also change the entries in the category buffer for the characters between them and the next PDF (except for any additional embedding or override characters you encounter). For LRO, everything between it and the next PDF (except characters embedded in a nested pair of explicit embedding/override characters) becomes L; for RLO, everything becomes R.

To see how this works, let’s look at another example. Let’s say you’re trying to get this as your displayed text:

Avram said מְזֹל טוֹב and smiled.

...an you split it into two lines so that only one of the Hebrew words can fit on the first line, you get this:

Avram said מְזֹל
טוֹב and smiled.

The other thing to keep in mind is that if you're doing styled text, you can't apply the bi-di algorithm separately to each style run—the bi-di algorithm may reorder the style runs. In other words, if you underline “said” and “מְזֹל,” you should get this...

Avram said טוֹב מְזֹל and smiled.

...and not this:

Avram said מְזֹל טוֹב and smiled.

Glyph selection and positioning

Arguably the most important part of drawing Unicode text is glyph selection and positioning. The Unicode character-glyph model (encode the semantics, not the appearance), combined with the inherent complexities of certain scripts, make rendering a complicated affair, at least for some scripts. For others, combining character sequences make things interesting.

Dealing with this sort of thing is usually the province of font designers, but it's important to have some knowledge of how various font technologies work, as some support for certain scripts may also require help from the rendering engine.

Font technologies

Let's start by talking for a few moments about font technology and how things are normally done in modern computerized fonts. This will help the following discussions of particular technologies and techniques make more sense.

In a modern text-rendering system, a font is a piece of software. It works in concert with the rendering engine itself to draw text on an output device, such as the computer screen or a piece of paper. A typical font may supply:

- Glyph images (bitmap fonts) or glyph descriptions (outline fonts). A glyph description may stand alone or may rely on other glyph descriptions to draw all or part of its glyph. (An accented-e glyph description might just call though to separate descriptions for the e and the accent, for example.)
- One-to-one mappings between code points in one or more character encoding standards and glyphs.

- Tables that allow for more-complex many-to-many mappings between code points and glyphs.
- Tables that allow the rendering engine to make alternations in the default position of a glyph depending on the surrounding glyphs.
- Tables that alter the default glyph mappings based on style settings chosen by the user.

There are two basic categories of font technologies: “bitmap” and “outline” fonts. In a bitmap font, each glyph has a bit-oriented image associated with it; this bit image is simply copied to the screen to draw that character. Because the character has essentially already been “rasterized” by the font designer, it’s designed for a particular combination of point size and device resolution. Bitmap fonts are essentially obsolete in modern personal computers and printers, but you still see them on PDAs and in various embedded applications (the displays on cell phones or copy machines, for example).

Outline fonts, on the other hand, are resolution independent. An outline font stores a geometric description of an ideal shape for each glyph. Each glyph description consists of zero or more ordered lists of points in an ideal coordinate space. Each ordered list of points gets connected together to form a closed shape called a “loop” or “contour” (this isn’t quite like a child’s connect-the-dots puzzle: some points in the loop are “off-curve” points—instead of actually being connected into the loop, they define the shape of a curve that connects together the preceding and following points). The loops get colored in to form the glyph (when loops intersect, there are rules that describe which parts get colored in, which is what allows you to have characters with “holes” in them). The rendering engine uses the idealized description to create a bitmapped image on the fly that is used to represent the character. An outline font can thus represent the same font at any size.

At small point sizes or coarse resolutions, however, the ideal coordinate space that is used to describe the glyphs begins to interact uncomfortably with the actual coordinate space where the rasterized image will be drawn, so outline fonts also include “hints” or “instructions” that are used to optimize display at low resolutions. The instructions for each glyph define an algorithm for deforming the idealized shape in a way that fits onto the output device’s coordinate system well while still honoring (as much as possible) the type designer’s intent.

The first major outline-font technology in the personal-computer arena was Adobe Systems’ PostScript, which first appeared on the Apple LaserWriter printer in 1986. This was followed a number of years later by Apple’s TrueType font technology, which first appeared in Apple’s System 7 **[or did it actually come out sooner?]** in **[when?]**. For a while after that, these two incompatible font technologies coexisted. When Microsoft adopted TrueType for Windows in **[when?]**, a situation arose where computers were generally using TrueType to render their characters on screen and printers were using PostScript to render the same characters on paper. Since it was possible for a low-end printer to rely on the host computer to rasterize the whole page and just send the bit image down to the printer, TrueType started to become more popular.

This led to a merger of the two technologies in **[when?]**, forming OpenType. OpenType fonts use the TrueType font file format, but can contain glyph outlines and hints in either the old TrueType format or in the old Postscript format, along with various tables of supplementary information. This format is currently supported by Microsoft and Adobe, as well as many other vendors.

Meanwhile, Apple went its own way, extending the basic TrueType format in other ways to create TrueType GX, which first appeared in **[when?]**. The basic TrueType GX format is still compatible with OpenType, but includes supplementary tables for various fancy typographical effects that aren’t

recognized by OpenType renderers. TrueType GX has since evolved into Apple Advanced Typography (AAT).¹³³

Old-style TrueType fonts are compatible with both OpenType and AAT, and newer OpenType or AAT fonts are generally (though not always) minimally compatible with each other, although (unless they were designed as dual-purpose fonts) they behave as old-style TrueType fonts on the system they weren't designed for.

[I'm reconstructing all of this from memory. I suspect I have details of the story wrong, in addition to missing all the dates when things came out. I suspect I'll also have to clean up the story a bit to avoid making enemies on one side or the other.]

OpenType and AAT fonts contain one or more *character mapping tables* (or “cmap” tables) that map code points in some encoding standard to internal glyph indices. A font can have multiple cmaps, allowing it to be used with more than one character encoding standard without the system having to perform code conversion before it can render. These days, fonts from all the major font vendors include cmap tables for Unicode (specifically UTF-16). The glyph indices aren't standardized—these are internal values used only by the renderer to access the glyph outlines. Glyph indices are two-byte values, meaning a single font can never contain more than 65,536 glyphs. The cmap table maps single characters to single glyphs; more complex mappings are handled differently in OpenType and AAT.

Poor man's glyph selection

Most Unicode-compatible systems begin by doing the stuff they can do *without* any cooperation from an underlying font engine (or from the fonts themselves). There's actually quite a bit of Unicode that can be covered using a rendering engine (or a font) that doesn't support complex typography. With suitable fonts but no special rendering technology, you can handle:

- Latin (precomposed characters only)
- Greek (again, precomposed characters only)
- Cyrillic
- Armenian
- Georgian
- Hebrew (unpointed, assuming you apply the bi-di algorithm to the characters before you pass the characters to the renderer)
- Ethiopic
- Cherokee
- Han, Hiragana, Katakana, and Bopomofo (assuming the fonts can hold sufficient numbers of characters)
- Hangul (precomposed syllables only)
- Yi
- Canadian aboriginal syllabics

¹³³ Apple's rendering engine for AAT fonts is called Apple Type Services for Unicode Imaging, or ATSUI (pronounced “at-soeey”), which gets my vote for funniest-sounding computer-industry acronym—I always want to say “gesundheit” when someone talks about it.

But you can actually go a bit beyond this by performing mappings on characters before sending them to the renderer. If you map to Normalized Form C before sending characters to the renderer, you pick up support for all combining character sequences that have an equivalent precomposed form.

But you can actually go a little beyond just this, especially if you do compositions that are excluded in Normalized Form C. For example, you can't do fully-pointed Hebrew this way, but you can do some other languages that only use a few combinations of Hebrew letters and points, such as Yiddish.

Many proportional-font technologies provide some ability for characters to overlap: you can actually have a character draw outside its designated bounding box. You can use this technique to get some support for combining character sequences that can't be normalized to Form C: you just have the glyphs for the combining characters overhang whatever happens to precede them. Of course, if the width of the preceding character can vary, you won't get very good-looking results, but it's better than nothing (this can also give you rather crude support for a few more scripts, such as Thai, pointed Hebrew, or Thaana, as long as you don't encounter characters that have more than one mark on them [such as Thai consonants with a top-joining vowel *and* a tone mark]).

If the font technology allows you to specify kerning pairs (pairs of characters that can be moved closer together than the default) you can use this to fine-tune the placement of combining marks for the characters they're being applied to, although this doesn't help you when marks have to be moved vertically or change shape to avoid colliding.

And you can also handle one script that requires complex typography: Arabic. The minimally-required contextual forms of the Arabic letters (and the minimal combinations of these forms with points) are in the Arabic Presentation Forms B block as compatibility composites. Using the rules in the ArabicShaping.txt file from the Unicode Character Database, you can perform minimal contextual shaping.

The algorithm for this is fairly simple:

- Map the text to Normalized Form C.
- Split the text into lines and apply the Unicode bi-di algorithm.
- If you encounter U+0627 ARABIC LETTER ALEF to the immediate left of U+0644 ARABIC LETTER LAM, check the joining class of the character to the right. If it's R or N, replace the lam and alef with U+FEFB ARABIC LIGATURE LAM WITH ALEF ISOLATED FORM. If it's L or D, replace the lam and alef with U+FEFC ARABIC LIGATURE LAM WITH ALEF FINAL FORM. (In real life, this step actually has to be extended to take into account pointed forms of alef as well as the unadorned version.)
- Go through each character on the line in visual order. Look up its joining category in ArabicShaping.txt, along with the categories of the characters on either side. Use this table to match the abstract character the Arabic block to an appropriate presentation form from the Arabic Presentation Forms B block:

class to left	current class	class to right	presentation form to use
any	N	any	independent
R or D	L	any	initial
L or N	L	any	independent
any	R	L or D	final

any	R	R or N	independent
L or N	D	R or N	independent
R or D	D	R or N	initial
L or N	D	L or D	final
R or D	D	L or D	medial

The abbreviations in the table should be fairly self-evident: L means the character connects to the character on its left but not the character on its right; R means the character connects to the character on its right but not the one on its left; D means the character can connect to characters on both sides (D stands for “dual-joining”). N means the character doesn’t connect to anything (it stands for “non-joining”). Anything not listed in `ArabicShaping.txt` is, by definition, in category N.

It’s theoretically possible to extend this technique to get somewhat better-looking results by also using the ligatures in the Arabic Presentation Forms A block. But this generally doesn’t work very well in practice (points and other marks still tend to get misplaced, for example) and very few fonts actually provide glyphs for these characters.

Arabic is the only script for which Unicode actually encodes presentation-form characters for all the contextual forms for each letter, but it’s possible to extend this technique to other scripts using special fonts and code points in the Private Use Area.

Glyph selection and placement in AAT

Apple Advanced Typography fonts go way beyond the basics to provide lots of advanced-typography features. Perhaps most important among the extra features in AAT fonts is the *glyph metamorphosis table*, or “mort” table, which provides the ability to do much more complex character-to-glyph mapping than is possible with the standard cmap table.¹³⁴

The transformations in the mort table are applied after the mappings in the cmap table have been applied. One advantage of this design is that it makes the font independent of the character encoding of the original characters: The original characters are mapped first into internal glyph indices using the cmap table, and then the resulting sequence of glyph indices can be further transformed into a *different* sequence of glyph indices (using the mort table) before being rendered.

AAT provides for two different kinds of glyph transformations. Non-contextual transformations are one-to-one mappings that aren’t based on surrounding characters. These mappings are instead based on global settings such as the internal state of the renderer. This lets you do things like choose between regular and old-style numerals, or between regular lowercase letters and small caps, depending on a font-feature setting. It also lets you pick different glyphs depending on whether text is being rendered vertically or horizontally (parentheses and hyphens, for example, rotate ninety degrees when used with vertical text).

¹³⁴ My main sources for this information were the TrueType Reference manual, available online at <http://fonts.apple.com/TTRefMan/index.html>, and Dave Opstad, “Comparing GX Line Layout and OpenType Layout,” <http://fonts.apple.com/WhitePapers/GXvsOTLayout.html>.

More interesting are contextual transformations. AAT makes use of finite state machines to allow glyph selection for a particular character to be based on arbitrary amounts of context on either side. In fact, the state-machine implementation actually allows you to modify the text stream as it goes, which permits some fairly complex mappings. All of this not only allows you to do the kind of contextual glyph selection necessary for cursive joining in Arabic, but also lets you do things like reorder glyphs (to deal with left-joining Indic vowel signs, for example), add glyphs to the text stream (to deal with split Indic vowel signs represented by single code points in the original text), and remove glyphs from the text stream (to make Indic viramas invisible, for example).

Contextual glyph transformations are also how you do ligature formation and accent stacking in AAT. In fact, an accented letter is just a special kind of ligature in AAT. AAT includes the concept of a “compound glyph.” A compound glyph doesn’t have its own outlines, but instead uses the outlines from one or more other glyphs, possibly applying transformations to them as part of the composition process. Essentially, it calls the other glyphs as subroutines. The letter é, for example, would generally have a compound glyph: It’d use the glyph for the regular letter e, plus the glyph for the acute accent, possibly repositioning the accent to appear in the right place above the e.

The one problem with this approach is that it doesn’t lend itself well to arbitrary combinations of accents and base characters, especially if a single character will have a bunch of accents applied to it. If the font designer didn’t think of the particular base-accent combination you want to use, figure out how the pieces should go together to draw it, and assign it a glyph code, an AAT font can’t draw it.

AAT fonts can include a lot of other interesting goodies, including:

Kerning. Kerning is the process of fine-tuning the positioning of two characters to look better, usually by moving them closer together. Normal font spacing usually leaves a little bit of white space between each pair of characters, but this can make certain pairs of characters appear to be too far apart. For example, the word “To” looks better if the crossbar of the T hangs a little over the “o,” and “AVANTI” looks better if the slanted strokes of the As and V are moved a little closer together so the letters overlap. AAT kerning tables not only allow for this kind of pair-based kerning, but also allow state-table-based algorithms that let you consider more than two characters at a time when deciding how to position the characters. It also supports “cross-stream kerning,” where a character is moved up or down based on the surrounding characters, which can be useful for things like hyphens and dashes (a hyphen should be positioned higher when it’s between two capital letters or two digits than when it’s between two small letters, for example).

Justification and tracking. Tracking is a global adjustment that moves all the characters on a line closer together or farther apart. Justification is, as we saw earlier, the process of systematically widening or narrowing a line of text using various techniques to take up a given amount of horizontal space. Both of these processes interact in interesting ways with the design of the type and with the script you’re using, and AAT allows you to add special tables to a font that specify just how these processes should take place on the text rather than leaving it to the rendering process. This lets the font designer choose, for example, when to justify by inserting *kashidas*, when to justify by increasing interword spacing, when to justify by adding intercharacter spacing, when to do a combination of these things, and in what proportions, and when to do things like break up ligatures.

Baseline adjustment. In English text, we’re used to the idea of text sitting on a “baseline.” With the exception of letters like the lowercase “p” that have “descenders,” letters all sit on an imaginary line. If a particular line of text mixes text of different point sizes, the differently-sized pieces of text line up along the baseline.

The baseline is in a different place in some other scripts. Han characters and their allied scripts (such as Hangul and Kana) commonly use a baseline that runs through the center of the line: differently-sized characters are centered vertically with respect to each other. Devanagari and other Indic scripts use a “hanging baseline”: the horizontal stroke at the top of most characters is the thing that stays in the same place as the point size changes.

If you mix characters that use different baselines on the same line of text, you need some way of relating the different baselines to one another so that things appear to line up correctly. AAT fonts can include a special table that includes information to facilitate this.

Optical alignment. Characters typically include a little bit of white space on either side so that they don’t touch when drawn in a row; this is called the “left-side bearing” and “right-side bearing.” This extra white space can lead to a ragged appearance when the text is lined up along a margin, especially if type size varies from line to line. In addition, straight and curved strokes don’t appear to line up right if the leftmost pixels of the curved stroke line up with the leftmost pixels of the straight stroke, requiring curved characters to actually be positioned slightly outside the margins. Certain punctuation marks in some scripts are also allowed to appear outside the margins. AAT fonts can include tables to facilitate this kind of positioning.

Caret positioning. An AAT font can also include tables that indicate where to draw the insertion point when it’s positioned between two characters that are represented by a ligature, and that let you do things like specify a slanted caret that’ll draw between the characters in an italic or oblique typeface.

These are just some of the special features that can be included in an AAT font. There’s a lot of other stuff that you can do as well. In addition, Apple’s ATSUI renderer handles a number of things for you automatically (the Unicode bi-di algorithm being the big one). A lot more work goes into designing a good AAT font, but it leaves less work in the end for the application developer.

Glyph selection and placement in OpenType

OpenType has some features in common with AAT, owing to their common lineage, but some important differences as well.¹³⁵ Both share the same overall file format (but contain different tables), and both can make use of TrueType outlines. In addition, OpenType fonts can use Adobe Type I Compact Font Format (CFF) outlines. Both AAT and OpenType also support bitmap glyphs, but have different formats for them.

OpenType has a lot of the same features as AAT, but takes a completely different approach. For example, its approach to complex glyph selection is totally different. It starts the same way—characters are mapped to glyphs in a one-to-one manner via the cmap table, which is the same as the cmap table in AAT, and then additional transformations are performed on the resulting glyph indices.

In OpenType, the additional transformations are defined in something called the *glyph-substitution table*, or GSUB table. The GSUB table uses a different operating principle from AAT’s mort table. Instead of a state-machine-based implementation, the GSUB table uses a string-matching

¹³⁵ For the section on OpenType, I relied on the OpenType font specification, as found at <http://partners.adobe.com/asn/developer/opentype/>, as well as on the Opstad paper cited earlier.

implementation. You specify a list of sequences of glyph indices and the glyph indices to map them to. There's a lot of flexibility in this. You can do the same kinds of non-contextual mappings that are possible in AAT (although OpenType fonts don't include a built-in concept of "font feature"; choices of glyph based on font features are left to the rendering engine).

There's also a lot of flexibility in the contextual mappings: you can specify literal strings to be replaced by other literal strings, or you can specify strings based on various "glyph classes" or user-defined collections of glyphs. You can also specify sequences that must appear before or after the sequence to be mapped in order for the mapping to take place. In essence, you have something very similar to a regular-expression-matching system at your disposal for specifying glyph substitutions. In effect, this system also uses a state machine for doing glyph substitution. The big difference between this and the AAT state-machine implementation is that the AAT state machine implementation allows the glyph stream to be modified on the fly while the state machine is passing over the text. This makes reordering and insertion of new glyphs, both necessary for proper rendering of many Indic scripts, more difficult in OpenType. It's still possible **[as far as I can tell]**, but most OpenType fonts for Indic scripts don't bother, leaving the work of dealing properly with split vowels to the underlying rendering engine, which must preprocess the text in ways analogous to the handling of the Unicode bi-di algorithm.

[Do I have this right? I know the IBM team had to do a bunch of extra work to handle Indic scripts using OpenType in Java, but I don't know exactly what the issues were—it does look like the font format does support the necessary mappings. What am I missing here?]

One really neat thing about OpenType fonts is their handling of combining marks and accents. OpenType has the compound-glyph idea that AAT uses, because it inherits it from the TrueType outline format (CFF outlines don't have this capability), but it provides a more powerful approach to accent application: the GDEF and GPOS tables. The GPOS table functions kind of like a fancier version of the TrueType kerning table (some early formats of which are still also supported by OpenType), but it allows adjustment of character positions in both directions, where the AAT kerning tables only allow adjustment in one direction. Better yet, the GPOS table lets you position glyphs relative to each other by mating *attachment points*, which are defined (along with some other stuff) in the GDEF table. In this way, for example, instead of saying that when you see an e and an acute accent next to each other, the accent should be moved to the left so many pixels and down so many pixels so it draws in the right place over the e, you can simply say that for all Latin-letter/top-joining-accent combinations, you should mate the top-center attachment point on the letter with the bottom-center attachment point on the accent. You then leave it to the GDEF table to specify where those points are on the various letters and accent marks. This kind of thing can drastically shrink the number of entries you need in the GPOS table (and, by eliminating the need for compound glyphs, the number of glyph indices you have to waste for letter-accent combinations).

These tables can be set up in such a way as to work right with multiple combining marks on a single base character (including arbitrary Unicode combining character sequences) and can do the same kinds of contextual matching that's possible in the GSUB table.

OpenType font files also include tables for justification, baseline adjustment, and caret positioning that do roughly the same things as their AAT counterparts, but they have a different format from the analogous tables in AAT fonts.

Special-purpose rendering technology

There are specialized applications that go outside the bounds of what normal rendering engines can do and require specialized rendering engines. For example, Unicode provides math and musical symbols, but mathematical formulas and music both require much more complex shaping and layout than normal written language.

Another interesting case is Arabic. Because of its cursive nature, there's an almost endless variety of combinations of characters that might have special forms (not to mention the rather complex placing of vowel points and other marks around the letters) and trying to get calligraphic-quality text with contextual forms and ligatures in a conventional font technology just becomes too complicated. To do really first-class Arabic typography requires more of an algorithmic approach. There exist specialized rendering engines just for Arabic that use special fonts and take this kind of algorithmic approach.

Compound and virtual fonts

In any language, you'll get the best-looking results if you use styled text and actually specify things like the font to use and the point size and style to use with it. But Unicode, combined with the Internet, raise the likelihood of running into text in some language and script other than the ones you read and write. And while if you don't read it, it might not matter so much what you see, you can generally make some sense out of it (or take it to someone who *does* read it) if it still shows up the way it's supposed to.

For this reason, you're seeing operating systems come with special Unicode fonts that include glyphs for all the Unicode characters (or at least a sizable subset). This technique can get rather unwieldy, however.

Another approach is to use a *compound font*. Some systems give you the ability to specify an ordered list of fonts to try using to display text if the font the document asks for isn't available or the document uses characters that aren't in that font. In this way, instead of having a huge "Unicode" font, you can have separate Greek, Hebrew, Arabic, Japanese, etc. fonts. This fallback list can be thought of as a "font" on its own—it just relies on various other fonts to supply the actual glyphs. This technique can be used to avoid having to carry huge multi-megabyte font files around to see all of Unicode and it allows users to fine-tune things to their liking (picking a preferred default font for their native language, for example, or preferring Chinese glyphs over Japanese glyphs for the Han characters).

Usually, at the root of a hierarchy like this is a "last resort" font. You can use this to get something more meaningful than the standard "missing" character for characters that aren't in any of the fonts in the fallback list. The "last resort" font usually has one glyph for each Unicode character block, and that glyph is used for all the characters in that block. In this way, if you don't have a Devanagari font, for example, you can at least still tell that a piece of text you can't read was in Hindi (or some other language that uses the Devanagari script).

[Is it possible to show a few example glyphs from Apple's last-resort font?]

Java and the MacOS are two examples of environments that use a virtual-font technique to handle Unicode's huge repertoire of characters.

Special text-editing considerations

Finally, it's worth taking a little time to look at a few things you need to keep in mind if you're writing a Unicode-compatible text editing utility. There are a number of Unicode-specific things to keep in mind when allowing a user to edit text.

Optimizing for editing performance

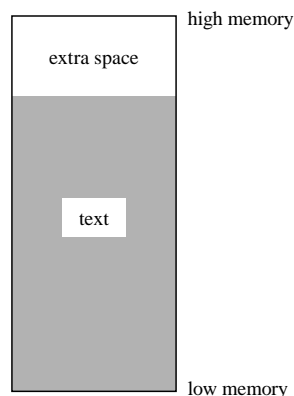
The first isn't actually Unicode-specific at all, but if you don't already know about it, it's worth taking a few minutes to discuss. That's a useful technique for speeding performance during editing.

Let's say you have a one-million-character document and you start typing at the beginning. If you've stored your text in one contiguous block (probably with some slop space at the end so you don't have to reallocate your character storage on every keystroke), you have to slide a million characters up two bytes in memory with every keystroke. This, obviously, is generally unacceptably slow.

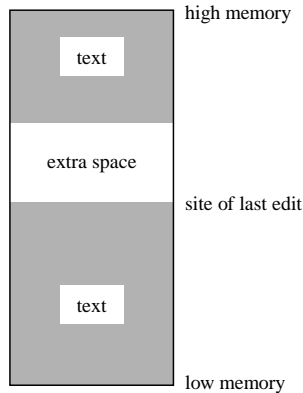
It generally makes more sense to break the document up into pieces. This way, you don't have to shove massive numbers of characters around whenever you're doing edits at the beginning of the document, and you don't have to allocate a new block and copy all the characters into it every time you overflow the available storage space. This technique also lets you read a document from disk in pieces rather than all at once and, if you're memory-constrained, lets you page parts of the document out to disk as you're working.

But full-blown block storage can be complicated to maintain; there's a lot of bookkeeping associated with keeping track of which blocks are where, which text is in which block, when to split or coalesce blocks, and so on. A simpler approach with many of the same performance benefits involves splitting the document up into only *two* blocks. It's called gap storage.

To understand how gap storage works, let's start by imagining your typical memory buffer holding a piece of text. Normally the buffer will be bigger than the actual text so you have room for it to grow as the user edits. It looks something like this:



But you don't have to keep the extra space at the end of your backing store. What the gap technique does instead is *move it around* as the text is edited. So after an edit in the middle of the document, the buffer looks like this:



The idea is very simple. Instead of moving all the text above the edit position up or down in response to an edit, you instead reposition the gap so that it occurs immediately after the position of the edit. The edit itself, whether an insertion or a deletion, doesn't cause any more characters to move at all: it behaves just as an edit at the very end of the document would with ordinary storage.

On average, you move a lot fewer characters around with this technique than you would the old-fashioned way. Better yet, if you're doing the edits in response to keyboard input, where large edits happen one character at a time, you get an even bigger win. If I type a ten-letter word at the beginning of the document, the gap is moved to the beginning of the document in response to the first keystrokes. The other nine characters all get inserted at the position of the gap and require no more moves at all! This approach can generally give you lightning-fast performance.

The bookkeeping associated with gap storage is minimal. Basically you have to keep track of where in the backing store the gap is and how big it is. Operations that only retrieve text from the buffer without changing it don't generally move the gap, so you need to keep track of whether a retrieval is directed at text before or after the gap. If a retrieval of text *crosses* the gap, then you may have to do a little more work (or move the gap), but the bookkeeping is still much simpler than doing the same thing with a full-fledged block storage mechanism.

Of course, if the text grows beyond the bounds of the storage buffer, you still have to reallocate the storage and copy the text into the new buffer, but since moving the gap can be done as part of the process of copying the text from the old to the new buffers, the incremental cost of the reallocation can be kept down. If you're memory-constrained, you have to rely on a virtual memory subsystem to do the paging in and out for you.

The gap technique can be even more fruitful for the other data structures associated with the text than for the text storage itself, such as the data structure that stores styles and other out-of-band information, or the line starts array.

Let's consider style storage for a moment. Typically, style information is stored in a *run array*, a data structure that associates records of style information with positions in the text. Each entry in a style run array refers to a *style run*: a contiguous sequence of characters that have the same styling

information. The entry specifies the position and length of the style run in the text and contains a pointer to another record somewhere that describes how that run of text is to be drawn (or possibly contains some other information about the text).

Now think about what happens when you add a character to a style run near the beginning of the document. All of the style runs that come later in the document have to have their positions updated so that they refer to the same characters they did before your addition. You can get around this by storing lengths instead of positions—this way only the run you’re actually changing needs to be updated—but now it becomes more difficult to find out which styles are in effect at an arbitrary position in the text. To do that, you have to start at the beginning of the run array and walk it, accumulating run lengths until you find the run containing the position you’re interested in.

It’s really best if you store both the positions and the lengths of the runs in the run array, but you don’t want to waste time updating every run in the array every time an edit happens to the document. Gap storage can be a big help here.

With run arrays, you’re not so much interested in the performance you save from not moving entries up and down every time there’s an edit. Run arrays are usually much smaller than character arrays, and whole runs are added and deleted a lot less frequently than they’re updated. Here, gap storage helps you with the task of updating run positions in response to edits.

You proceed the same way as before: Each time there’s an edit, you reposition the gap in the run array so that it occurs immediately after the run that’s affected by the edit. Then you only have to update that run’s length in response to the edit.

The magic is happening in the repositioning. The trick is the way you store the run *positions*. For the runs before the gap, you store the positions in the normal way: relative to the beginning of the text. For the runs *after* the gap, you store the positions relative to the *end* of the document instead.

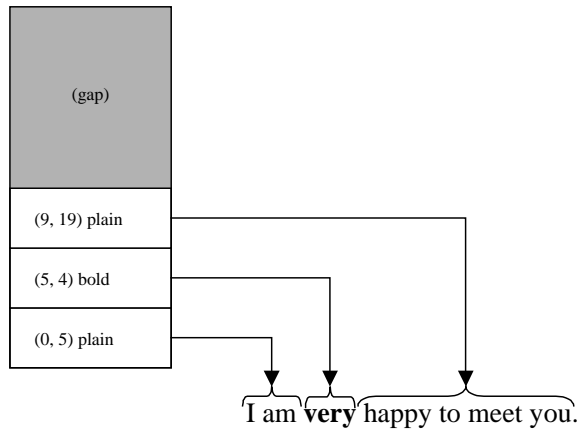
Now the position of a given run only gets updated when it switches polarity. If you have a bunch of localized edits in the same style run, the gap is immediately after it and the runs after it don’t have to have their positions recalculated (they’re relative to the end of the document, and while the end of the document will move, their positions relative to it don’t). When, later on, you make a change later in the document, the positions of the runs after your first edit are recalculated as part of moving the gap. The recalculation is simple: you just add up the lengths as you switch things.

Let’s say you have this sentence:

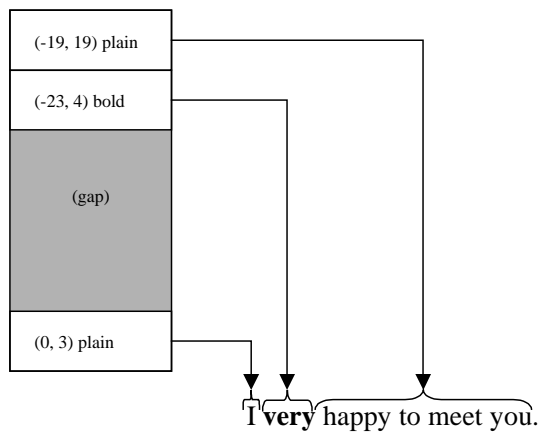
I am **very** happy to meet you.

It has three style runs: “I am,” “**very**,” and “happy to meet you.” The style run array starts out looking like this:

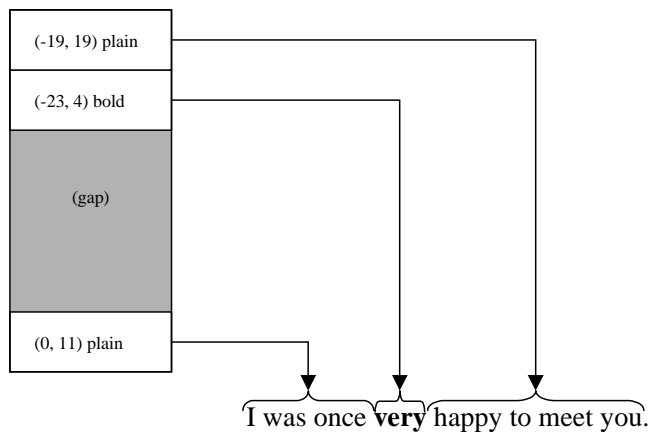
Special text-editing considerations



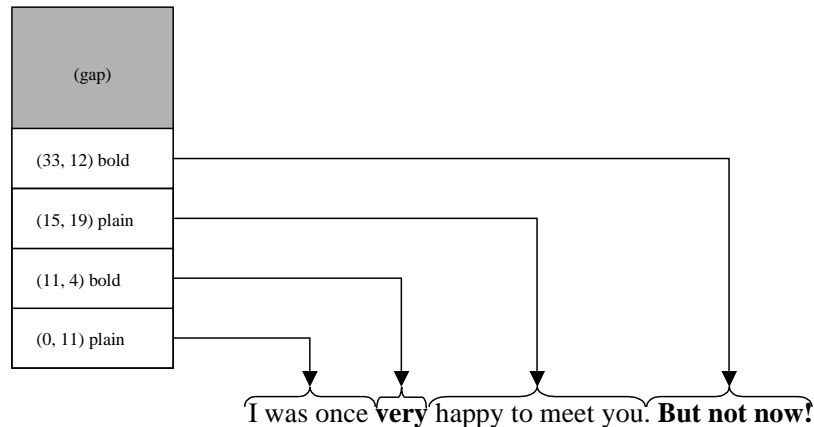
Now let's say you select "am" and delete it. Now the run array looks like this:



You replace "am" with "was once". The gap doesn't move and no positions need to be recalculated. Only the length of the run you're changing changes:



If you now go to the end of the document, switch to boldface, and add "But not now!", the gap moved back to the end, and all the positions get recalculated:



This technique can be very useful in keeping style information and other metadata up to date without sacrificing performance.¹³⁶ It's also useful for doing the same kind of thing with internal tables, the line-starts array in particular.

Accepting text input

I don't want to spend a lot of time dwelling on text input here, but it's worth emphasizing that Unicode doesn't assume a one-to-one mapping between keystrokes and characters any more than it assumes a one-to-one mapping between characters and glyphs. Code that accepts input from the user should keep this in mind.

For instance, combining character sequences can cause interesting things to happen. A Vietnamese keyboard might generate the sequence U+0061 U+0323 U+0302 (representing the letter â) with a single keystroke. Conversely, it might produce U+1EAD (the precomposed form of the same letter) with two keystrokes (one for the basic letter and another for the tone mark, potentially in either order).

A lot of European keyboards have "dead keys," keys that cause an accent to be applied to the next character you type, so the two-keystrokes-to-one-character mapping happens even in European languages. Most encodings tend to favor precomposed forms, so one keystroke producing multiple characters doesn't happen as often, but it definitely does happen (in Burmese, some split vowels can be represented only with combinations of two vowel signs, but the vowel might still be produced with a single keystroke, for example). It's important to keep this in mind when designing an input engine.

The more complicated case happens with the Han characters, which generally require something called an *input method*, a whole process that maps keystrokes (or possibly input gestures on some other device, such as handwriting or voice input) to characters. Because there are so many Han characters, it's completely impractical to have a keyboard that includes every character, so smaller keyboards (often conventional Latin keyboards) are used, and multiple keystrokes are used to enter each character. In Japanese, for example, it's fairly common to use a normal English keyboard and

136 The polarity-switching technique used with gap storage and applied to style runs is actually patented by IBM (the inventors are Doug Felt, John Raley, and me). **[What is the status of this patent? Can non-IBMers use this technique without owing IBM a license fee? Should I even be mentioning it if not?]**

enter text in romaji (the standard Latin-letter representation of Japanese). As letters are typed, they appear first as Latin letters, and then as whole syllables are completed, they turn into Hiragana. When a sequence of letters corresponding to the Romanization of some Kanji character is completed, the text turns into that character. (Actually, since many Kanji characters have the same romanization, there'll usually be an intermediate step where all possible Kanji characters corresponding to a given romanization are shown to the user in a menu and he picks one.)

Similar techniques are used for the other languages that use Han characters. The keystrokes might represent Latin letters in some recognized romanization, they might represent characters in a phonetic system such as kana, bopomofo, or hangul, or they might represent more abstract concepts such as groups of strokes. In all cases, not only are there complex mappings from series of keystrokes to input characters, but extra work must be done to manage the interaction between user and computer as keystrokes are assembled into characters.

The main point is that if you want to support East Asian languages, you may have to deal with the input-method issue, either by providing input methods or ways of writing them, or by tying into any input-method support provided by the underlying operating system.

Even on an American keyboard, you sometimes input methods coming into play, such as in implementing the “smart quotes” feature in most word processors that turns " into either “ or ” depending on context. A transliteration engine (see Chapter 14) generally plays a big role in input-method support.

Handling arrow keys

Finally, there's the issue of handling selection feedback and arrow keys. Again, this isn't always straightforward in Unicode.

Consider the handling of arrow keys, for example. What should happen when you press the right-arrow key? The right answer is *not* necessarily to move the insertion point forward one code point in the backing store. If, for example, the backing store contains U+0065 U+0301, the insertion point is before U+0065, and the screen shows the letter é, you don't want the insertion point to move forward one code point. If you do, you're now in the middle of the é, rather than past it. You either see no movement in the insertion point, it moves past the é and gives a misleading picture, or it gets drawn through the middle of the é, which is ugly and potentially confusing. If the user types “a” now, he ends up with “eá” instead of what he really meant, which was probably “éa”.

But it isn't just regular combining character sequences that can cause trouble. What should happen, for example, as you arrow through “कि ”? [spacing] There are two marks here with a fairly clear distinction between them, but the mark on the left (ि) actually comes *after* the mark on the right: the mark on the left is a left-joining vowel sign. Should you be able to arrow into the middle of a syllable like this?

What about something more complex like this:

क्रि

This syllable is actually made up of four code points. The main letter, *k*, is represented by क̄. The second letter is *r*, which is represented by a boomerang-shaped mark that gets subsumed into the क̄, appearing here only as the spike sticking out of the bottom of the loop. In fact, to get the *k* and *r* to join together in this way, there's actually a *virama* (a vowel-killer sign) after the क̄ to cancel its inherent vowel and make it combining with the *r*. The *virama* is invisible in this particular syllable: its presence is represented by the fact that the *k* and *r* have combined into a single glyph. Finally, you have the vowel, *i*, which is represented by the ि mark on the left.

Here the relationship between visual positions and positions in the backing store is very complicated. The fourth character in the backing store comes first, followed by the first character stacked on top of the third. The second character isn't visible at all.

[put together a picture that shows this clearly]

You could actually have the insertion caret appear inside the syllable cluster in cluster-specific positions (AAT fonts actually let you do this kind of thing), but it's very complicated and difficult and may actually be confusing to the user. Most of the time, the only reasonable solution to this is for the arrow keys to move syllable by syllable when dealing with Indic syllable clusters in the scripts that have complicated reordering or shaping rules. So you may not only have to keep regular combining character sequences together, but also Indic syllables.

You also have to watch out for this issue when handling mouse clicks. You don't want a mouse click to put the insertion point in the middle of a combining character sequence any more than you want the user to arrow into it.

The bottom line is that your arrow-key algorithm (not to mention your hit-testing algorithm, your selection-drawing algorithm, and possibly your line-breaking algorithm) has to honor grapheme cluster boundaries. It may be able to deviate a bit from Unicode 3.2's default definition of a grapheme cluster, but it has to treat the text in logical units that make sense to the user, units that may be multiple code points in length.

The parsing techniques we looked at earlier in the chapter for line breaking can also be brought to bear to parse text into grapheme clusters.

On top of all this, you may also encounter text that contains invisible formatting characters such as the zero-width joiner and non-joiner. You've got to skip these as well when you're handling the arrows keys unless, of course, your editor has some sort of "show invisible characters" feature turned on. Otherwise it just looks like the keystroke didn't register.

An even uglier problem happens in bidirectional text. The right arrow should move you to the right whether you're in left-to-right or right-to-left text. So if you're in right-to-left text, the right arrow actually moves you *backwards* through the backing store.

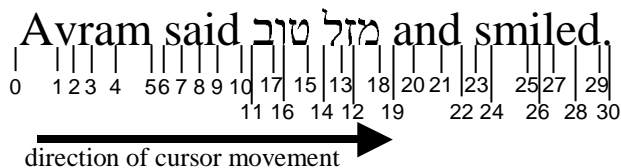
Furthermore, interesting things happen when you have text of mixed directionality on the same line. Any boundary between directional runs can be interpreted as two positions in the backing store. For example, in the following sentence...

Avram said מְזֵל טוֹב and smiled.

...the visual position between the first set of English letters and the Hebrew snippet could either represent position 11 in the backing store (after “Avram said” and before “מְזֵל טוֹב”) or position 18 (after “מְזֵל טוֹב” and before “and smiled”). The visual position between the Hebrew and the second English phrase can represent the same two positions. Furthermore, either logical position in the backing store can map to both visual positions on the screen. For example, if I’m at position 11 and I insert a Latin letter, it appears after (i.e., to the right of) “Avram said.” If I insert a Hebrew letter, it appears before (i.e., to the right of) “מְזֵל טוֹב”.

So you have to keep track of which you mean, you have to draw both insertion point positions, and you have to do some complex gymnastics as you go through the text with the arrow keys.

As you arrow through the text from the left, for example, the insertion point goes up one by one until it gets to 11 (the ambiguous visual position). Then it jumps to position 17 (between the last two Hebrew letters) and counts down until it gets to position 12. The next right arrow (the other ambiguous visual position) is position 18, after which the positions count up again. (If you left-arrow through the text, positions 11 and 18 switch places.)



I’ve seen some interesting algorithms to get this effect right. Basically they involve re-applying pieces of the Unicode bi-di algorithm on the fly as the user arrows through the text. Generally, though, you already need some kind of data structure that maps positions in the backing store to the order of the characters on the screen. You’ve generally had to produce this mapping to draw the characters. Instead of re-applying the bi-di algorithm on the fly, it’s generally easier to just keep this mapping around and use it to help determine what position in the backing store corresponds to the next visual position on screen.

Handling discontinuous selection

In addition to handling arrow keys, you have to deal with drawing selections. We talked about this in Chapter 8, but it’s worth repeating that if you’re dealing with bidirectional text, selection drawing becomes more complicated. If a selection crosses the boundary between directional runs, you have a discontinuity. The discontinuity can be either in the rendered feedback on screen (“logical selection”) or in the characters in the backing store (“visual selection”). For instance, if in our example sentence you click down between “Avram” and “said” and drag to the position between “מְזֵל” and “טוֹב”, you could either have meant to select a logically contiguous range consisting of “said” and “מְזֵל”...

Avram **said** **טוב** מזל and smiled.

...or you could have meant to select a visually contiguous range consisting of “said” and “טוב”:

Avram **said טוב** מזל and smiled.

Logical selection is generally easier to implement in spite of the fact that you may have several visually discontinuous highlight ranges to draw. After all, a selection highlight that spans a line boundary may also consist of two visually discontinuous highlight ranges.

Visual selection can be a little easier on the user, but can be trickier to implement. This isn’t just because you have to keep track of more than one selection range in memory. It also makes things like copy and paste harder. For example, it seems logical that if you select some text, choose “Cut” from the program’s “Edit” menu, and then change your mind and choose “Paste” from the “Edit” menu without doing anything in between, you should get back what you had originally. If you select “said” and “טוב” as we saw above, said “Cut,” and then immediately said “Paste,” you’d expect to get the original sentence back again. But remember that the two selected words weren’t originally contiguous in memory. This means the paste operation has to put them back in two logically discontinuous locations so they’ll wind up being drawn next to each other. “Said” goes between “Avram” and “מזל” in the backing store, but “טוב” goes between “מזל” and “and.”

If you’re not pasting at a directional boundary, the two pieces are stored contiguously in memory, but the order may be different because you want the stuff you copied to the Clipboard to come out in the same order as you originally saw it. So if you paste at the beginning of the sentence, you want to see this:

said טוב Avram מזל and smiled.

To get this effect, both “said” and “טוב” go before “Avram” in the backing store, and “said” precedes “טוב”. On the other hand, if you did a “Copy” instead of a “Cut” and then pasted the duplicate words between “מזל” and “טוב”, you’d want to see this:

Avram said טוב said טוב מזל and smiled.

The pasted text still comes out with “said” to the left of “טוב”, but because you’re pasting into right-to-left text now, “טוב” precedes “said” in the backing store.

Handling multiple-click selection

The last interesting thing you may have to worry about is double- and triple-click handling. Generally, double-clicking selects a word and triple-clicking selects a line, paragraph, or sentence, depending on the application. You can use the same parsing techniques that are used for locating

possible line-break positions to figure out where the selection range should begin and end in response to a multiple click.

Of course, the “word” you select with a double-click might or might not be the same thing as a “word” for the purposes of line breaking. You may or may not want to include whitespace and surrounding punctuation, for example. In languages that don’t use spaces between words, such as Japanese, you might want to honor word boundaries anyway using a dictionary-based parser. Or you may want to use a heuristic of some type that selects more text than you’ll find between possible line-break positions (which usually only include one character). In Japanese, for example, a common heuristic might be to select a contiguous range of Kanji followed by a contiguous range of Hiragana and any leading or trailing punctuation (if the user clicked in Kanji or Hiragana), or a continuous range of Katakana (if the user clicks in Katakana). This is usually a decent approximation of real word boundaries. (Of course, in Chinese this will select a whole sentence or paragraph—for Chinese it’s better to use a dictionary-based algorithm or just select a single character.)

Sentence parsing can also be tricky. In English, for example, you can’t tell for sure whether a period marks the end of a sentence or just the end of an abbreviation (you can tell when it’s being used as a decimal point or URL delimiter, though, because it isn’t followed by a space in these situations). One heuristic is to behave as though the period is at the end of a sentence if it’s followed by whitespace and the next letter after the whitespace (and any intervening punctuation such as parentheses or quotation marks) is a capital letter. If the next letter is a small letter, you’re not at the end of a sentence.

Of course, this is a pretty inexact heuristic. It’ll choke on a person’s name with a preceding title (such as “Mr. Jones”), for example. You can tune the algorithm by including a list of common abbreviations that are usually followed by capitalized words (such as “Mr.,” “Mrs.,” “Ms.,” “Dr.,” etc.), but of course this list is language-specific.

CHAPTER 17 *Unicode and Other Technologies*

Of course, most of the time you don't actually want to *implement* the Unicode standard; you just want to *use* it. In slightly over ten years of existence, Unicode has racked up an impressive array of supporters. A lot of other standards depend on it, and a lot of products use it in one way or another. In this last chapter, we'll take a look at some of the places where Unicode intersects with other technologies and standards.

Unicode and the Internet

Perhaps the area where Unicode will be most important over time is on the Internet. Here you have an area where people all over the world can communicate with each other and access information. If they just do this in English, use of the Internet will wind up limited mostly to English-speaking countries. Instead, there's a growing push for material on the Internet to be available in a wide variety of languages.

Of course, if you go beyond English, ASCII doesn't cut it any more as a method of representing characters, and if you go outside the Western European languages, neither does Latin-1. This means either designing the Internet protocols and the various pieces of software that speak them to handle many different character encoding standards, or it means Unicode. (In most cases, the actual answer is "both.")

Of course, a wide variety of standards and technologies is in use on the Internet. They were developed at different times by different people and are currently maintained by different standards bodies, so it's no wonder that they exist at different points along the adoption curve or take different approaches to integrating Unicode support. Below we take a look at a few of the most important standards.

The W3C character model

The World Wide Web Consortium, or “W3C” for short, is an industry group responsible for many of the more important Internet standards, including XML and HTML. It has a standing internationalization committee whose job it is to look over other W3C standards and make sure that they don’t contain biases in favor of particular languages or nationalities.

Since the W3C aims for all of their standards to be internationally usable, they’ve adopted Unicode as the standard character encoding for all of their standards that involve encoded characters. But just saying “the text will be in Unicode” doesn’t go far enough—you still have to worry about things like versions of Unicode, encoding formats, normalization forms, and so forth—so the W3C has issued a document, the *Character Model for the World Wide Web*, that attempts to nail down all these ancillary issues in the context of W3C standards.¹³⁷ All W3C standards published or updated after the character model was published have to follow its recommendations.

Here’s a quick summary of what the W3C character model specifies:

- It demands that W3C technologies never assume one-to-one mappings between characters (in all the myriad definitions of “character”), code points, glyphs, and keystrokes. (Unicode, of course, assumes this.)
- It proposes a list of more specific terms than “character” and “string” and provides definitions.
- It requires that all implementations of W3C specifications treat all text as though it were encoded in Unicode, although it allows other encodings to be used.
- It requires all W3C standards either to assume all text is in either UTF-8 or UTF-16 (or both, with an appropriately unambiguous way of distinguishing between the two), or to provide some mechanism for specifying the encoding used by various pieces of text. If the latter, it requires that the default be UTF-8 or UTF-16 if no encoding is explicitly specified, and it strongly recommends the use of the IANA charset identifiers for the purposes of identifying character encodings. (We’ll talk more about the IANA charset identifiers in a minute.)
- It places strong restrictions on the use of Unicode private-use characters. In particular, it requires that W3C specifications neither impose semantics on any private-use code point nor provide a mechanism for formally imposing semantics on private-use code points. The idea here is to preserve these code points’ identity as private-use code points, rather than allowing some of them to be usurped for some purpose by a particular protocol. W3C protocols have to use markup elements to identify characters that aren’t in Unicode (MathML, for example, has a special `mathlyph` element for specifying math symbols that aren’t encoded in Unicode).
- It requires the presence of some kind of escaping mechanism to allow the use of characters in text (such as the ampersand or less-than sign) that normally have special syntactic meaning, and to allow the representation of all Unicode characters even in documents that don’t actually use Unicode as their encoding. It doesn’t specify a particular format for an escape sequence (although it strongly recommends XML’s “`ሴ`” format), but it does require that escape sequences have an explicit end character (such as the semicolon in the XML escape sequence), allowing the code point value to have any number of digits.
- It imposes very specific normalization requirements: Processes that produce or modify text *must* produce it in Normalized Form C (escape sequences can’t be used to get around this: the text

¹³⁷ The document is available at <http://www.w3.org/TR/charmod>. At the time of this writing, January 2002, the most recent version was Working Draft 20, dated December 2001. A working draft doesn’t have normative force and may still be changing. But the draft was in the final stages of revision. By the time this book appears in print, the final version of this document should have been adopted and published.

must *still* be in Normalized Form C even after all escape sequences are turned into real characters). Processes that receive text *must* interpret it as Normalized Form C and must *not* normalize it themselves. Processes that receive text encoded using encodings other than Unicode *must* convert it to Unicode using converters that produce Normalized Form C.

The idea here is to place responsibility for normalizing the text as early in the process as possible: with the processes that actually produce the text in the first place. This lets processes that only receive text to be lighter-weight: they don't have to do anything special to make sure variant forms of the same strings compare equal. In fact, they're prohibited from normalizing it themselves because this would introduce a potential security hole: the text might be interpreted differently depending on the receiving process.

- It recommends discouraging or prohibiting the use of compatibility composites, but doesn't actually require (or even recommend) converting text to Normalized Form KC. (More on this below.)
- It defines two strings as equal only if, after converting both to the same Unicode encoding format and turning all escapes into real characters, they're code-point-for-code-point equal. (This assumes they were both in Normalized Form C to begin with.) One of the important implications of this is that string matching is case-sensitive.
- It recommends a couple of consistent methods of indexing particular characters in strings. In particular, it recommends indexing strings by abstract Unicode code point (regardless of actual storage format) or, failing that, by physical code unit, and it recommends a zero-based indexing format (actually, it recommends numbering the "cracks" between the code points or units, rather than the units themselves, which works out to the same thing).
- It recommends that standards allow the use of all Unicode characters in URL references, even though only ASCII characters are actually allowed in URLs, and specifies how the non-ASCII characters are to be converted into ASCII in real URLs. (More on this below.)
- It recommends that W3C standards specify a particular version of Unicode for their syntactic elements, but make an open (non-version-specific) reference to Unicode for their non-syntactic elements. It also specifies that open references always be to Unicode 3.0 or later.

Earlier drafts of the W3C character model actually restricted the allowable repertoire of Unicode characters to a subset of the whole Unicode repertoire (disallowing more than just the private-use characters), effectively making "W3C normalization" a different and more-restrictive beast than regular Unicode normalization. These requirements have been removed from the W3C character model, but appear instead in a separate document, "Unicode in XML and Other Markup Languages," which has been published jointly by the Unicode Consortium as Unicode Technical Report #20 and by the W3C as W3C Note 15.

The basic idea behind UTR #20 is that there are characters included in Unicode whose functions are normally done with markup tags when Unicode is used in markup languages such as XML. UTR #20 enumerates those characters that are redundant or conflict with markup tags and explains why they shouldn't be used and what should be done instead.

Basically, UTR #20 discourages the use of certain invisible formatting characters and certain compatibility composites, but not *all* of the characters in either category. In particular:

- It discourages the use of the LS and PS characters (you should use the <p> and
 tags instead), or specifies that they, like CR and LF, should be treated as plain whitespace rather than actual line and paragraph delimiters.
- It discourages the bi-di embedding and override characters in favor of markup tags that do the same things.
- It discourages the deprecated Unicode characters because they're deprecated.

- It discourages the interlinear-annotation characters in favor of markup tags that do the same thing.
- It discourages the object-replacement character, which is basically only for internal implementation use in the first place.
- It discourages the Plane 14 language-tag characters in favor of markup tags that do the same thing.
- It does *not* discourage the other invisible formatting characters, including the non-breaking characters, the zero-width joiner and non-joiner, the zero-width space, the left-to-right and right-to-left marks, and the various language-specific formatting characters, all of which either happen inside words where interposing markup tags would mess things up, or otherwise don't cause a problem when used with markup.
- It lists a fairly small subset of Unicode compatibility composites that it recommends converting to Normalized Form KC and a few others that it recommends replacing with markup. It recommends leaving most Unicode compatibility composites alone, however, either because they actually have semantic differences from their decomposed versions (in particular, the various styled versions of various letters, which are usually used as symbols) or because these aren't currently any markup that can produce the appropriate effect (for example, the various symbols intended for use in vertical CJK text).

XML

The Extensible Markup Language, or XML for short, is a document format for expressing structured data in a simple text-based manner. An XML file consists of arbitrary textual data interspersed with “markup,” text in specially-defined formats that imposes structure on the data and supplies extra information about it. XML doesn't specify a particular method for representing particular kinds of data, but merely provides a generalized mechanism for representing arbitrary structured data in a flat text file. Other standards, such as XSL (Extensible Stylesheet Language), MathML (Mathematical Markup Language), XHTML (Extensible Hypertext Markup Language), and SVG (Scalable Vector Graphics) specify how XML can be used to represent different specific kinds of data.

The XML standard, not surprisingly given its source, follows the recommendations in the W3C character model: It's based on Unicode as the character encoding for the text, but discourages the use of private-use characters, most compatibility composites, and most invisible formatting characters. Numeric character references (XML's escaping mechanism) are based on abstract Unicode code-point values. All of the XML syntax characters are defined in terms of their Unicode code point values. All XML parsers are required to understand both UTF-8 and UTF-16, but are permitted to understand non-Unicode encodings as well.

An XML document is assumed to be in Unicode unless it's explicitly tagged as being in some other encoding. If it's not tagged, it can be either UTF-8 or UTF-16; UTF-16 documents must begin with the byte-order mark, which is used both as a signal that the document is in UTF-16 and as a signal of the document's byte order. Text that doesn't begin with a byte-order mark and which doesn't begin with an explicit character-encoding declaration is assumed to be in UTF-8. If the document contains byte sequences that aren't legal in UTF-8, it's an illegal XML document.

Interestingly, the XML standard doesn't specify an encoding for the encoding declaration itself, but instead gives a method whereby the encoding of the declaration can be deduced (since the document either has to start with the declaration or with a byte-order mark, and the declaration has to start with “<?xml”, this isn't that hard). If the encoding specified in the declaration is found not to match the actual encoding used in the document, this is also an error.

In all the cases where strings have to be matched in XML (matching start tags and end tags, for example, or matching attribute names and values to their declarations), the matching happens using

simple binary comparison—in other words, an XML parser is required to assume the text has already been converted to Normalized Form C, and case differences are significant. One important implication of this is that identifier matching is case-sensitive: `</record>` is not a legal closing tag for `<RECORD>`.

XML includes an identifier syntax similar to the ones used in programming languages. It generally follows the Unicode identifier-syntax guidelines, with a few extra punctuation characters. But the XML standard lays out explicitly which characters are considered to be “letters,” “digits,” “combining marks,” and so on. These definitions are based on the Unicode Character Database general categories from Unicode 2.0, when the XML standard was adopted. This means more-recent “letters,” such as Sinhala, Khmer, Myanmar, and the various extensions to the Han repertoire, aren’t legal in XML identifiers (they are legal in XML character data, just not in identifiers). Most likely, the next version of the XML standard will update the repertoire of name characters to whatever version of Unicode is current at that time. Of course, documents that use the new name characters wouldn’t be backward compatible with XML 1.0 documents.

XML forms the basis for many other Internet document standards, such as XHTML, CSS, XSL, SVG, MathML, and so on, all of which therefore inherit XML’s Unicode compatibility.

HTML and HTTP

Users of the World Wide Web are, of course, familiar with Hypertext Markup Language (or “HTML”), the markup language used for most of the styled-text documents on the Web. HTML is similar to XML (both are based on SGML), but the syntax is a lot looser, and it’s been around a lot longer. The most recent version of HTML, HTML 4.1, uses Unicode as the base character set and generally follows the W3C character model. Earlier versions of HTML, however, were based on ISO 8859-1.

Documents in HTML 4.1 can actually be in any encoding, much as XML documents can, with the actual encoding specified in the document. HTML doesn’t have a designated place for this information like XML does, so it appears as an attribute in the `<META>` element in the document’s header.

There’s been a move to bring HTML into conformance with XML. That is, instead of HTML looking kind of like XML but having some important differences, it’d be nice for HTML to actually be based on XML, much as other document formats like XSL and SVG are. The result of this move is XHTML, an XML-based format that does the same things HTML does, but with a tighter syntax. XHTML, of course, inherits XML’s Unicode support.

The HTTP protocol that’s used by Web browsers has the ability to exchange information about the encodings of the documents being exchanged, but because of spotty application support, the encoding specified in a document’s HTTP header might not be the document’s actual encoding (and FTP and other file-sharing protocols don’t have any way at all of exchanging information about document encodings). Better to rely on the information in the document itself.

URLs and domain names

A Uniform Resource Locator, or URL, is a short text string that identifies and helps locate a particular resource on the Internet; it's sort of the Internet equivalent of a person's mailing address.¹³⁸

The syntax for URLs is very restrictive: basically only the uppercase and lowercase Latin letters, the Western digits, and a handful of special symbol and punctuation characters are allowed in URLs, which makes it possible to transmit a URL in almost any character encoding and through almost any protocol.

The problem, of course, is that the use of just the Latin letters introduces a bias toward languages that use the Latin alphabet. This wouldn't be a big problem if URLs were just internal identifiers, but they're not. These days, you see URLs everywhere, and they're intended to be mnemonic: You'll always see `http://www.ibm.com`, not `http://129.42.16.99`.

The URL specification explicitly provides a mechanism for escaping characters in URLs: You can use the percent sign followed by two hexadecimal digits to indicate an arbitrary byte value. The problem with this is that nothing specifically says what the byte values mean or how they're to be interpreted.

The industry is converging around always treating escape sequences in URLs as referring to UTF-8 code units. That is, the industry is converging around always interpreting `R%c3%a9sum%c3%a9.html` to mean `Résumé.html` (and always representing `Résumé.html` as `R%c3%a9sum%c3%a9.html`). If everybody agreed on this, then you could use illegal URL characters (such as the accented `é` in our example) in URL references in other kinds of documents (such as HTML or XML files) and know there's a universally-understood method of transforming that into a legal URL. Web browsers or other software could do the reverse and display URLs that include escape sequences using the characters the escape sequences represent (at least in the cases where they represent non-ASCII characters) and let you type them in that way.

The XML standard already allows URL references to be specified this way and specifies that they be converted to legal URLs by using the escape sequences to mean UTF-8 code units, and the W3C character model will extend this to the other W3C standards. Various other standards are also converging on this approach.¹³⁹

The %-escape approach, unfortunately, doesn't work for Internet domain names, which have a more restricted syntax: Each component of a domain name can consist only of Latin letters, Western digits, and hyphens, must begin with a letter, and must end with a letter or digit. Each individual

138 Technically speaking, a URL is a string that identifies a particular resource by specifying an access path—a way of reaching it, hence the term “locator.” There's also something called a URN, or Uniform Resource Name, which identifies a resource in a way that's independent of access path. The standards community uses the term URI, or Uniform Resource Identifier, as an umbrella term for both URL and URN. Informally, the term URL is used in a manner synonymous with “URI,” which is the convention used here. Also, in some documents, the “U” in all of these terms is said to mean “Universal” rather than “Uniform.”

139 My information here is mostly coming from the W3C Internationalization Working Group's paper on “URIs and Other Identifiers,” found at <http://www.w3.org/International/O-URL-and-ident>.

component of a domain name can't be longer than 63 characters, and the entire domain name (including the periods between components) can't be longer than 255 characters.

Some protocols and some software implementations let you go outside these restrictions, but doing so poses a big risk to interoperability. There have been a bunch of creative name-mangling schemes proposed for representing arbitrary Unicode strings using only legal domain-name characters, and the Internet Engineering Task Force, the other big body that maintains Internet standards, has an Internationalized Domain Names working group working on an official proposal. As of this writing, this work is still ongoing, and domain names are still restricted to Latin letters.¹⁴⁰

Mail and Usenet

Mail and Usenet (newsgroups) are the oldest services on the Internet, and their transfer protocols show that. While there are many proprietary email and bulletin-board systems out there with all kinds of features, the standard Internet protocols for mail and news are pretty simple and primitive.

The various mail and news protocols are all based on the Internet Engineering Task Force's Request for Comments #822 (or "RFC 822" for short).¹⁴¹ RFC 822, which dates way back to the dawn of personal computers, specifies the format for mail messages on the Internet. It defines a mail message as a series of CRLF-delimited lines of 7-bit ASCII text. RFC 822 goes on to give a detailed specification of the message header format and the address format, but does not go into any more detail as to the contents of the message body. The message body simply extends from the end of the header (delimited by a blank line) to the end of the message.

RFC 822 doesn't impose a line-length limit, but many mail protocols do. In particular, the SMTP protocol imposes a 1,000-byte maximum length on each line of a message (this counts the CR and LF at the end).

This format is fine for simple text messages (at least if you speak American English), but rather constraining, and it doesn't help you if you want to send things other than text in your email. Over the years, various methods have been devised for sending things other than text (or in addition to text, or text in encodings other than ASCII) in the body of an RFC 822 message. The one that eventually caught on is MIME, the Multipurpose Internet Mail Extensions, the most recent version of which is documented in RFCs 2045, 2046, 2047, 2048 and 2049.

MIME adds a few extra fields to the standard RFC 822 message header that are used to say more about what's in the message body. There's a `MIME-version` field that identifies the message as being in MIME format (and could be used to identify the version of MIME, except that there's currently only one version of MIME), there's a `Content-Type` field that identifies the type of the message body (e.g., text, image, multipart message body, undifferentiated byte stream, application-specific binary data, etc.), and there's a `Content-Transfer-Encoding` field that specifies how the

140 My sources here are RFC 1034, "Domain names: Concepts and Facilities," and RFC 2825, "A Tangled Web: Issues of I18N, Domain Names, and the Other Internet protocols", both found on the IETF Web site, www.ietf.org.

141 The IETF formally publishes a lot of documents as "requests for comments." Even though many of these do become formal standards, they're generally still referred to by their RFC numbers rather than the new numbers they get assigned when they become standards. The IETF actually encourages this, by setting things up so you search for all documents by their RFC number.

data type specified by the `Content-Type` field is encoded into lines of 7-bit ASCII characters. If the `Content-Type` starts with `multipart`, the message body is divided up into multiple parts, each of which gets its own `Content-Type` and `Content-Transfer-Encoding`. This lets you do things like send the same message in several alternative formats, attach binary files to a message, or include a separate message (for example, one you're forwarding) in the message body.

Many of the `Content-Types` also allow you to specify parameters. For text messages, the most important parameter is the `charset` parameter, which identifies the encoding scheme used to encode the text. (The term “charset” is a bit of a misnomer, harkening back to the days when there was a straightforward transformation between all coded character sets and character encoding schemes. Nowadays, the `charset` parameter specifies a character encoding scheme, not a coded character set.) A standard ASCII message is understood to have the following `Content-Type` declaration in its header:

```
Content-Type: text/plain; charset=US-ASCII
```

The Internet Assigned Numbers Authority (IANA), which is affiliated with the IETF (they're sister organizations under the umbrella of the Internet Society, or ISOC), is the registration authority for various Internet-related parameters, including things like domain names. The IANA registers `Content-Types` and `charsets`. There are a number of Unicode-related `charset` identifiers registered by the IANA:

```
utf-8      UTF-8
utf-16be   UTF-16 in big-endian byte order
utf-16le   UTF-16 in little-endian byte order
utf-16     UTF-16 with a byte-order mark to specify the byte order (if the BOM is
           missing, this is the same as utf-16be)
scsu       Standard Compression Scheme for Unicode
```

There are other `charset` identifiers for Unicode that are obsolete now. There aren't yet `charset` identifiers for the various flavors of UTF-32, but that'll probably happen eventually.

The `Content-Transfer-Encoding` field is important because the only character encoding that can safely be included in a message and expected to make it unscathed from source to destination without encoding is ASCII. For all other encodings, the only way to make sure the message body makes it where it's going without getting messed up is to use a special transfer encoding. For text, this takes the text in its native encoding (for us, Unicode) and applies an algorithmic transformation to convert it into 7-bit ASCII.

MIME specifies two basic ways of doing this: `quoted-printable` and `Base64`. `Quoted-printable` is best for situations where most of the encoded data actually is readable ASCII text and you want to preserve human readability. It represents the printable ASCII characters as themselves (except for the `=` sign). All other byte values are represented by an `=` sign followed by a two-digit hexadecimal number. Thus, the word “Résumé” in UTF-8 comes out like this in `quoted-printable`:

```
R=C3=A9sum=C3=A9
```

(The lower-case `é`, U+00E9, is 0xC3 0xA9 in UTF-8.)

Since `quoted-printable` doesn't mangle the original ASCII characters, it's useful for messages in Latin-based scripts. You might identify such a message with a header like this:

```
Content-Type: text/plain; charset=utf-8
Content-Transfer-Encoding: quoted-printable
```

Quoted-printable also requires that lines be only 80 characters long (counting the CRLF), so the = sign is also used at the end of a line to indicate that the CRLF following it was inserted by the protocol and not in the original message (CRLF sequences that aren't preceded by = are actually part of the message).

The other format, Base64, is better for pure binary data or for text where the majority of the text isn't readable as ASCII. It's also a little more robust in the face of various transformations than quoted-printable can be (if you only escape non-7-bit characters, control characters, and the = sign, a quoted-printable message can still be mangled on its way through an EBCDIC-based mail gateway and may mess up interpretation of a multipart MIME message).

In Base64, each sequence of three bytes (24 bits) is converted into four ASCII characters, with six bits of the original 24-bit sequence going to each character. This means you need 64 characters to represent each "nibble" (hence the name). The 64 characters chosen (the uppercase and lowercase letters, the digits, +, and /) were chosen because they're not syntax characters in anything that might appear in a MIME message and they exist in EBCDIC and will stay intact through an EBCDIC-based gateway. (A 65th character, =, is used at the end of messages as padding.) The word "Résumé" in UTF-8 would look like this in Base64:

```
UsOpc3Vtw2k=
```

This, of course, isn't human-readable, but it works well for cases where the majority of the text isn't ASCII in the first place. A Unicode-encoded Japanese message, for example, might have a heading like this:

```
Content-Type: text/plain; charset=utf-16be
Content-Transfer-Encoding: base64
```

Like quoted-printable, Base64 divides everything into 80-character lines (counting the terminating CRLFs), but since spaces and control characters aren't used in Base64, they can just be filtered out on the decoding side.

Of course, the normal MIME encoding only helps you with the body of a message; the header information still has to be in ASCII. One piece of the MIME specification, RFC 2047, addresses this issue with an escaping scheme that can be used in certain parts of an RFC 822 header to allow non-ASCII characters. Basically, it's a shorthand for the format above: you specify the encoding using a charset identifier, you specify the transfer encoding using B or Q, and then you have the encoded text. It looks like this:

```
Subject: =?utf-8?Q?Here's my r=C3=A9sum=C3=A9?=
```

None of these schemes is particularly beautiful, of course, and they impose some overhead on the size of messages, but they do let you use Unicode and other encodings in protocols that weren't originally designed for them. These days, most Internet mail clients understand MIME and do the work necessary so that you see that last example the way you expect to:

```
Subject: Here's my résumé
```

Unicode and programming languages

Most of the more recent programming languages now either use Unicode as their base internal character encoding or have a way to let you use it if you want. Many even allow Unicode characters in the syntax of the language in addition to using it in comments and literal strings.

The Unicode identifier guidelines

The Unicode standard actually gives guidelines for how programming languages (and other protocols such as XML) that want to use the full Unicode range in their identifiers can do it.¹⁴² The basic idea is to extend the common definition of identifier (a letter followed by zero or more letters or digits) to the full Unicode repertoire, and to do so in a way that allows for combining character sequences and invisible formatting characters.

The guidelines say that an identifier may start with anything Unicode considers a “letter” (all characters with general categories Lu, Ll, Lt, Lm, Lo, and Nl). Subsequent characters may be “letters,” combining marks (Mc and Mn, but not Me), “digits” (Nd but not No), connector punctuation (Pc), or formatting codes (Cf). Specs following these guidelines are encouraged to allow but ignore the formatting characters, but this isn’t required.

Since the actual characters in these categories change with each Unicode version, a specification making use of this guideline has to either list all the characters or nail itself to a specific Unicode version (XML nailed itself to Unicode 2.1, for example).

Individual specifications can do with these guidelines whatever they want, including ignore it, but following them ensures that users of non-Latin scripts will be able to construct identifier names in their native languages without any undue constraints being imposed by the systems they’re using. Most specs that allow Unicode in identifiers follow the Unicode identifier guidelines and add additional legal characters.

Java

Sun’s Java programming language has led the way in terms of Unicode support. The base `char` data type is specifically defined to be a UTF-16 code unit, and not just an unsigned 16-bit integer, and a Java `String` is a sequence of `char` and therefore inherits this characteristic.

One of the implications of this is that all the Java functions that operate on `String` and `char` know they’re operating on UTF-16 text and can do the right thing. The case-conversion functions in Java follow the mappings in the Unicode Character Database, for example. The various character-type queries on `Character` are also based on the Unicode categories.

Non-Latin characters are also legal in Java syntax. A Java source file can be in virtually any encoding, but is translated into Unicode internally at compile time. The specific syntax characters are all ASCII characters, but non-ASCII characters are allowed not just in literal strings and comments, but also in identifiers (see the above section). Java provides an escaping mechanism, the `\u1234`

142 See pp. 133-135 of the Unicode standard.

syntax, for putting Unicode characters into a Java source file when the source file isn't actually in Unicode (and provides the `native2ascii` tool to facilitate this).

The Java class libraries offer an extensive array of Unicode support facilities, including character encoding conversion, language-sensitive comparison, text boundary detection, rendering, and complex input-method support. There aren't specific APIs for text searching, transliteration, or Unicode normalization, but these are all available in open-source libraries. The `Unicode Collator` class doesn't automatically do the Unicode Collation Algorithm default order, but the implementation follows the UCA guidelines. The rendering and input facilities don't yet support every script in Unicode, but more scripts are added in every release (as of J2SE 1.4, the built-in text rendering software was OpenType compatible, implemented the Unicode bi-di algorithm, and supported Hebrew, Arabic, Devanagari, and Thai, in addition to the scripts that don't require contextual shaping or bidirectional reordering.)

C and C++

C and C++, being older, don't directly support Unicode. The `char` data type is simply defined as an integral type in C, so implementations can make it any size they want and the string-handling routines in the standard libraries can treat `char` data as any encoding they want. Most treat `char` as an 8-bit byte and use whatever the default platform encoding is (or sometimes just ASCII) as their encoding.

The `wchar_t` data type was added to allow for "wide" characters, such as Unicode, but doesn't necessarily help. Like `char`, the language spec doesn't impose any particular semantics on `wchar_ts`, so portable code can't depend on `wchar_t` data being Unicode. Worse, the C and C++ standards don't even require `wchar_t` to be big enough to hold a UTF-16 code unit; like the other integral types, it only has to be *as big as char*, not necessarily bigger.

Some compilers and runtime libraries treat `wchar_t` as Unicode, but if you're trying to write code that's portable to any C or C++ compiler, you can't depend on this. There are, of course, many Unicode support libraries available for C and C++, however.

Javascript and JScript

Javascript and JScript, the two Web-browser scripting languages from Netscape and Microsoft respectively, are both based on the ECMA 262 standard for scripting languages. ECMA 262, or ECMAScript, also uses Unicode as its base character type in strings. The third edition of ECMA 262 specifically nails this down to specify that string data is in UTF-16 and that the internal APIs assume Normalization Form C. It also stipulates that ECMAScript source code is Unicode and provides an escaping mechanism similar to Java's for including otherwise-untypable Unicode characters in source code. The third edition of ECMA 262 provides for a minimal set of Unicode-compatible string manipulation APIs; the idea was to keep things small and rely on whatever Unicode support is available from the host system. Most ECMAScript implementations also provide a method of accessing the Unicode support available from the host system.

Visual Basic

The current version of Microsoft Visual Basic is Unicode-compatible. Its internal character storage mechanism is Unicode, so all strings are automatically in Unicode. The `option compare text` statement enables locale-sensitive string comparison, which is Unicode-aware. The behavior here is

both case- and accent-insensitive (in other words, it ignores secondary and tertiary differences).

[I really feel like I should be saying more than this, but was having trouble finding information. Can someone help me out?]

Perl

[read Simon Cozens' paper and see whether it has anything to add to this]

Unicode support in Perl is a work in progress, and is phasing in gradually over the lifetime of Perl version 5.¹⁴³ The Perl community is moving in the direction of using UTF-8 as its internal encoding, but support is currently (as of Perl 5.6) incomplete, lacking, for example, the kind of transcoding I/O you need to use Unicode internally under all conditions.

The UTF-8 support in Perl has “character semantics.” This means that operations that index by or match individual bytes in other encodings match whole characters (actually, Unicode code points) in UTF-8, regardless of how many bytes they take up internally. Character-type queries and other related operations follow the specifications in the Unicode Character Database. The `\x` character escaping syntax has been extended to allow for values above 255: `\x{2028}` refers to the Unicode paragraph separator, for example. (For values between 128 and 255, the semantics are slightly different: `\xE9` is always the byte value 0xE9 (the Latin-1 é character), while `\x{E9}` is the Unicode é (internally 0xC3 0xA9, the UTF-8 representation of the letter é).

There's also a new `\p` token that matches all characters in a particular Unicode category. For example, `\p{Lu}` matches any Unicode uppercase letter. Actually, “category” here is more than just the general categories from the `UnicodeData.txt` file: you can use the `\p` syntax to match characters in a given script (`\p{InTibetan}`), or with a particular bi-di property (`\p{IsMirrored}`), or many other things as well as the general categories.

Currently, the UTF-8 support in Perl doesn't interact well with the locale support in Perl, which is based on the older POSIX locale standard, which is based on the old byte-oriented character encodings. This means things like locale-sensitive string comparison don't work well when operating on Unicode strings. **[Is this really true? Is it being fixed?]**

Finally, there's a `use utf8` pragma that enables the use of UTF-8 in Perl source files, where they can appear both in string literals and regular expressions and also in identifiers.

Unicode and operating systems

As with the other technologies we've looked at, the major operating systems all either have Unicode support now or are moving rapidly in the direction of adding it. Here's quick rundown.

Microsoft Windows

¹⁴³ My source for this information is <http://www.perldoc.com/perl5.6/pod/perlunicode.html>.

Microsoft has been gradually adding Unicode support to Windows for some time now.¹⁴⁴ Windows XP, NT and 2000 are built upon a Unicode base: all of the system APIs store strings as Unicode. (This is not true in Windows 95, 98, and Me—the two operating system lines have been united in Unicode support with the introduction of the home and professional versions of Windows XP.) The file system and networking protocols are Unicode-based too. The system uses the UTF-16 encoding internally and is based on Unicode 2.1.

For backward compatibility, all of the Win32 API functions that take strings as parameters come in two flavors: a version whose name ends in “A” (for ANSI) that takes the strings in one of the traditional byte-oriented Microsoft code pages, and a version whose name ends in “W” (for “wide”) that takes the strings in UTF-16.

Windows includes an extensive array of character encoding converters and provides a wide variety of Unicode-aware functions for such things as string comparison, boundary detection, case mapping, character-type queries, and so forth. **[I don't know the specifics here, and after some time rummaging around on the MS Web site, I'm convinced that I need an actual MSDN subscription to find out more. I can probably get access to that at work, but that's just not going to happen as part of putting this “first draft” together. I'll have to remember to come back and beef this section up. Any reviewer help here would also be greatly appreciated.]**

Finally, Windows 2000 includes a Unicode-aware text layout and display engine called Uniscribe that provides the ability to display multiscrypt text. It implements the Unicode bi-di algorithm and does the character shaping and accent stacking necessary to handle Arabic, Thai, Devanagari, Tamil, and Vietnamese (as well, of course, as those scripts that don't require complex shaping, such as Latin, modern Greek, Hebrew, and CJK). **[Does Uniscribe also handle vertical CJK text?]**

MacOS

MacOS 9 and earlier aren't Unicode based, but do have extensive Unicode support. Among the APIs are:

- the Text Encoding Conversion Manager, an extremely sophisticated and complete API for converting between various character encodings,
- Apple Type Services for Unicode Imaging (ATSUI), Apple's answer to Microsoft's Uniscribe, a Unicode-compatible text rendering engine that handles complex typography (it implements the Unicode bi-di algorithm and supports Apple Advanced Typography, where most other script-specific complex shaping behavior can be handled in individual fonts),
- the Multilingual Text Engine (MLTE), an ATSUI-based styled-text editing utility,
- the Text Services Manager, which provides Unicode-compatible input method support, and
- the Unicode Utilities, which provide Unicode-compatible language-sensitive comparison (including UCA support **[I'm guessing this; is it true?]**), language-sensitive boundary detection, and character property queries.

MacOS X includes all of the above services in its Classic and Carbon APIs and has counterparts to most of them in its Cocoa API. It includes a variable-length Unicode-based string object, CFString,

144 My sources are various documents found on the Microsoft Web site, particularly at <http://www.microsoft.com/globaldev>. **[I'm worried that some of the information I'm finding on the Microsoft site is out of date. Is there someone who can make sure I'm current on this stuff?]**

in its Core Foundation services. The underlying technologies are generally Unicode-based, although not all underlying technologies, owing to their diverse origins, use the same flavors of Unicode. In general, the application programmer doesn't have to worry about this, but there are issues at the boundaries of some of the environments (or in upgrading a MacOS 9 application to run under Carbon).¹⁴⁵

Varieties of Unix

It's difficult to generalize about Unicode support in Unix because there are so many different flavors of Unix and they all have different feature sets. All the big Unix vendors support Unicode to one degree or another.

[Can I, should I, go to the trouble of figuring out something more intelligent to say here about the individual Unix vendors? This is a difficult one because I was having trouble locating information and it'd be difficult to figure out which Unix versions merited discussion. I'm kind of hoping the comments on Linux below will suffice.]

One interesting development of the past several years has been the rise of Linux, the open-source version of Unix. In the past few years, Linux has not only gained popularity, but gained so much popularity that most of the major Unix hardware vendors have found themselves announcing support for it alongside their own proprietary versions of Unix. This might have the effect of bringing more consistency, if not outright standardization, to the Unix world.

In 1999, a bunch of Linux vendors and other supporters (including biggies like IBM, HP, and Sun) got together to form the Free Software Foundation, an industry standards body dedicated to coming up with a universal specification for Linux implementations—basically, a common feature set (or group of feature sets) for Linux implementations. Implementations conforming to the Linux Standards Base, or LSB, could be relied upon to have a certain set of features and APIs.

An important offshoot of the LSB effort was the Linux Internationalization Initiative, nicknamed LI18NUNIX, which was charged with standardizing the internationalization capabilities of Linux implementations. They issued the Linux Globalization Specification¹⁴⁶ for the first time in 2000. It includes specifications for every aspect of software internationalization, from low-level character manipulation APIs to tools support, rendering support, and input method support.

Among the things that the LI18NUNIX specification mandates is support for UTF-8 across a wide variety of locales. In addition to the standard Unix and C internationalization libraries, it also mandates (for the more-sophisticated implementations) the inclusion of ICU.

ICU, or International Components for Unicode, is an open-source project coordinated out of IBM. It consists of a broad array of utilities for Unicode support and related internationalization tasks including character encoding conversion, language-sensitive comparison and searching, boundary detection, character property queries, case mapping, transliteration, bidirectional reordering, and so on. It was originally based on the already-extensive internationalization libraries in Java, but goes beyond them (and has reimplemented them in C) to the point that there's now an ICU4J that contains

145 My sources are various documentation volumes found online at <http://developer.apple.com/techpubs>.

146 <http://www.li18nux.net/docs/html/LI18NUNIX-2000.htm>.

all the ICU features that aren't in the JDK (ICU4J is also mandated for LI18NUNIX-conforming Linuxes that support Java).

Unicode and databases

[There's supposed to be a section here that talks about how Unicode is handled in database systems, but I don't know anything about that and don't currently have access to any good research material. I thought I'd start with information on how the SQL standard handles Unicode, but SQL is an ANSI/ISO standard, which means it can't be found online. If I use that as my source, I'm going to have to go find a physical copy, probably in a library somewhere, or a contact who can give me the scoop. I'll do that in the next draft.]

[As for the individual database systems, the bottom line is that all the big database vendors support Unicode in their products. I don't know whether it's worth it to go into detail on just how they support it.]

[Should I even have this section, or should I leave it out?]

Conclusion

The bottom line of all this is that Unicode is firmly entrenched in the software community. The software that doesn't support it now (and is still being actively maintained) will. Interest in international markets and internationalized software continues to grow, and not only does Unicode solve a real problem, it's the only technology that solves that particular problem.

Unicode got where it is today through a carefully-maintained balance of careful, meticulous design and political know-how. Careful attention was paid to backward compatibility and interoperability, to upholding existing practice, to making sure everything was implementable in reasonable and reasonably-performing ways, to consistency, comprehensiveness, and elegance, and to making sure the needs of all the prospective user communities were taken seriously and addressed. That the Unicode standard is thriving (as evidenced by the increasing attendance at each successive Unicode Conference, in addition to the growing number of Unicode-compatible software products) is a testament to just how well these different and often conflicting objectives were balanced. It's doubly impressive when you consider just how many write-only standards there are out there. Unicode is not only an idea whose time has come, but an idea's *realization* whose time has come.

And which is here to stay, at least as long as we're still manipulating written language on computers.

Glossary

AAT. Abbreviation for **Apple Advanced Typography**.

Abstract character repertoire. A collection of characters with no reference to their encoded representation (for example, “the English alphabet”). See the discussion of character encoding terminology in Chapter 2.

Accent. A diacritical mark, usually modifying the sound of a vowel or indicating that a syllable has stress applied to it.

Accent stacking. The process of adjusting multiple diacritical marks applied to the same base character so that they don’t collide with each other. See the description of combining character sequences in Chapter 4.

Accented letter. A letter with one or more accented or other diacritical marks applied to it.

Accounting numerals. Chinese characters that are used in place of the “normal” Chinese numerals for accounting and commercial purposes, generally because they’re more difficult to alter. See the discussion of Han characters as numerals chapter 12.

Acute accent. A slanted line drawn over a vowel, usually to indicate primary stress or to change the sound of the vowel in some way.

Addak. A sign used in Gurmukhi to indicate consonant gemination.

Additive notation. A system for writing numerals in which the numeric value is derived by adding together the numeric values of all of the individual characters in the numeral. See the discussion of numerals in chapter 12.

Additive-multiplicative notation. A system for writing numerals in which the numeric values of some characters are added together and the numeric values of other characters are multiplied together (generally the order of the characters tells you whether to add or multiply). See the discussion of numerals in Chapter 12.

Alef maksura. A special form of the Arabic letter alef sometimes used at the end of a word. See Chapter 8.

Alifu. The name for the null-consonant character in Thaana. See Chapter 8.

Al-lakuna. The Shinala virama. See Chapter 9.

allkeys.txt. The data file giving the default weights of all Unicode characters for the Unicode Collation Algorithm. See Chapter 15.

Alphabet. A writing system in which the individual characters represent individual sounds (“phonemes”) in the language. A “pure” alphabet uses letters to represent both consonant and vowel sounds.

Alphabetic. The property applied to all letters in alphabetic writing systems.

Alphasyllabic. (or “alphasyllabary”) A writing system in with characteristics of both alphabets and syllabaries. An alphasyllabary typically consists of complex marks representing syllables, but which can be decomposed into individual marks (which over change order or shape) representing the individual sounds.

Alternate weighting. A property given to certain characters in the Unicode Collation Algorithm. These characters have several different possible sort weights that are user-selectable at the time a comparison is performed. See Chapter 15.

American National Standards Institute. One of the main standards-setting bodies in the United States; the United States national body in ISO.

American Standards Association. The old name of the **American National Standards Institute**.

Angkhankhu. A Thai punctuation mark used to mark the ends of long segments of text.

Ano teleia. A bullet-like character used as a semicolon in Greek. See the discussion of Greek in Chapter 7.

ANSI. Abbreviation for **American National Standards Institute**.

ANSI X3.4-1967. The official standard number of **ASCII**.

Anusvara. A character used in Devanagari and some other Indic scripts to indicate nasalization or an *n* or *m* sound at the end of a syllable. See Chapter 9.

Apple Advanced Typography. Apple’s advanced TrueType-derived font format. See Chapter 16.

Arabic form shaping. Another term for **contextual shaping** when applied to Arabic. This term is used with a couple of deprecated formatting characters that control whether contextual shaping is applied to the Arabic presentation-form characters. See the “deprecated characters” section of Chapter 12.

Arabic numerals. 1) The common name for the digit characters used in English and other Western languages. Also known as “European digits” or “Western digits.” These digits are *not* generally used with Arabic. 2) One of two sets of digits normally used with Arabic. These are *not* the same as the “Arabic numerals” used with English and the European languages.

Arabic-Indic numerals. The term used in Unicode to refer to the native digit characters used with Arabic and other languages that use the Arabic alphabet.

ArabicShaping.txt. A file in the Unicode Character Database that specifies a minimal set of rules for performing Arabic and Syriac contextual shaping. See Chapter 5.

Array. A data structure in which the individual elements are stored contiguously in a single block of memory.

Articulation mark. A mark applied to a musical note to indicate how it should be articulated.

ASA. Abbreviation for **American Standards Association**, the old name of **ANSI**.

Ascender. The part of a character that rises above the font’s x-height (such as the vertical stroke on the lowercase letter b).

ASCII. American Standard Code for Information Interchange. A character encoding developed from early telegraphy codes in the early 1960s that is still the most widespread character encoding standard used in data processing today. ASCII encodes the twenty-six letters of the English alphabet, plus the Western digits and a small selection of punctuation marks and symbols, and is a subset of a wide variety of other encodings.

ASCII_Hex_Digit. A property given to those Unicode characters that can be used as digits in hexadecimal numerals and are also in ASCII. (The upper- and lower-case letter A through F, plus the digits 0 through 9.)

Asomtavruli. A style of Georgian writing predominantly used in ecclesiastical texts and sometimes as capital letters. See Chapter 7.

Aspiration. Letting extra air escape the mouth in the pronunciation of a letter. The difference between the *t* at the beginning of “toy” (an aspirated *t*) and the *t* in “stamp” (an unaspirated *t*).

ATSUI. Apple Type Services for Unicode Imaging. Apple’s advanced text-rendering technology. See Chapter 16.

Attachment point. A location in a glyph shape in an outline font where a mark can be applied. An accent can be positioned correctly relative to different base characters by mating its attachment point to a particular attachment point on the base character. See Chapter 16.

Augmentation dot. A dot which, when applied to a musical note, lengthens the note by half its normal value.

Backing store. The area of memory containing the text that’s being displayed (or otherwise manipulated).

Backspace. 1) A control character in a telegraphy code that indicates that the printhead should be moved back one character position, or the direct descendants of a control code with this original use. 2) The key on a keyboard that deletes the character before the insertion point.

Backwards level. A level in a multi-level string comparison algorithm where differences later in the string are more significant than differences earlier in the string. See Chapter 15.

Bank. One of two or more alternate collections of characters that can be represented by the same numeric values in an encoding scheme (used particularly with older telegraphy codes such as the Baudot code). Special control characters are used to switch between banks. See Chapter 2.

Bantoc. A Khmer character that shortens the preceding vowel. See Chapter 9.

Bariyoosan. A Khmer punctuation mark used to mark the end of a section. See Chapter 9.

Base character. A character to which a diacritical mark can be applied.

Base64. A transfer encoding syntax that represents arbitrary binary data using only 65 printing characters that exist in almost all seven-bit encodings. See Chapter 17.

Baseline. An imaginary line against which all characters on a line of text are aligned. When text from different fonts or different type sizes is mixed on a line, the characters all line up along the baseline.

Baseline adjustment. The process of properly lining up text from scripts that use different baselines on a single line of text. (For instance, letters in Western scripts sit on the baseline, which letters in

most Indic scripts hang from the baseline. When Western and Indic writing is mixed on a single line of text, the baseline of one script has to be adjusted upward or downward to make the text appear to line up properly.) See Chapter 16.

Basic Multilingual Plane. The official name of **Plane 0** of Unicode and ISO 10646. This is the plane where the characters in all modern scripts are located. See Chapter 3.

Baudot. Emile Baudot, who invented a “printing telegraph,” the immediate predecessor of the teletype machine.

Baudot code. 1) The method of representing characters used with Baudot’s printing telegraph. 2) Also used to refer to a number of other telegraph codes that shared a similar five-bit, stateful structure. See Chapter 2.

BCD. 1) A method of representing numeric quantities using binary digits that preserves the ability to do decimal arithmetic on it. 2) Alternate name for **BCDIC**. See Chapter 2.

BCDIC. An in-memory character representation on early IBM computers based on the punched-card codes of the day. Four bits of the six-bit code were used to indicate the location of a punch in the rows 0 to 9, and the other two bits were used to indicate punches in the zone rows. See Chapter 2.

Beam. A line that connects together a group of musical notes (usually all the notes that occur in a single beat).

BEL. A control code that, in old telegraphy codes, caused a bell to ring (or some other audible signal to happen) on the receiving end of the communication link.

Bicameral script. A script with two sets of letters: upper case and lower case.

Bi-di. 1) Short for **bidirectional text layout**. 2) Used as a collective term for the scripts that require bidirectional text layout.

Bi-di algorithm. Short for **bidirectional text layout algorithm**.

Bi-di category. A property assigned to a character to indicate its directionality.

Bidi_Control. A character property applied to the invisible formatting characters that can be used to alter the default behavior of the bi-di algorithm. See Chapter 5.

BidiMirroring.txt. A file in the Unicode Character Database that links together characters with mirror-image glyphs. This file can be used to implement a rudimentary version of mirroring. See Chapter 5.

Bidirectional text layout. The process of arranging characters of mixed directionality on a single line of text.

Bidirectional text layout algorithm. The part of the Unicode standard that defines just how characters of mixed directionality should be arranged on a single line of text. For more on this algorithm, see Chapters 8 and 16.

Big5. An industry-developed character encoding standard for Traditional Chinese commonly used in Taiwan.

Big-endian. 1) A serialization format that sends and receives multiple-byte values most-significant-byte-first. 2) A processor architecture that stores the most significant byte of a multi-byte value lowest in memory.

Binary comparison. The process of comparing two strings by performing a simple numeric comparison of the individual code point values until a difference is discovered.

Binary tree. A linked data structure consisting of individual nodes that are arranged in a hierarchy and where each node points to no more than two other nodes.

Bindi. The Gurmukhi anusvara. See Chapter 9.

Bit map. 1) An array of Boolean values, usually packed together so that each Boolean value is represented by a single bit. 2) A series of bytes representing a visual image where each bit determines whether a particular pixel is lit or not.

Bitmap font. A font where the glyph images are represented as bit maps that can be copied directly into the display buffer. See Chapter 16.

Bit-mapped display. A display device that is controlled by a bit map in memory, where the bit map specifies which pixels are to be lit.

Bitwise comparison. Alternate term for **binary comparison**.

Blanked. An alternate-weighting mode in the Unicode Collation Algorithm where all character with alternate weights are treated as completely ignorable. See Chapter 15.

Block. A contiguous range of code point values in the Unicode encoding space set aside for a particular purpose (e.g., the Greek block contains characters used in Greek).

Block storage. A storage scheme for a long range of contiguous information (such as a text document) that stores it by dividing it up into a number of separate blocks in memory. See Chapter 16.

Blocks.txt. A file in the Unicode Character Database that gives the boundaries of the different “blocks” of the Unicode encoding space. See Chapter 5.

BMP. Abbreviation for **Basic Multilingual Plane**.

BOCU. Byte-Order Preserving Compression for Unicode. A Unicode compression scheme that will produce the same sort order as the raw UTF-8 or UTF-32 text. See Chapter 6.

BOM. Abbreviation for **Byte order mark**.

Bottom-joining vowel. An Indic vowel sign that is drawn below the consonant it modifies.

Boundary analysis. The process of analyzing a piece of text for boundaries between linguistic units such as words or sentences. Line breaking and double-click detection are two processes that perform boundary analysis.

Boustrophedon. Text where successive lines are written with alternating directionality (i.e., a left-to-right line is followed by a right-to-left line, which is followed by a left-to-right line, etc.). Used in ancient Greek and a few other ancient writing systems.

Box drawing character. A character representing a straight line, an angle, or some other piece of a box. A succession of box-drawing characters are used together to draw lines and boxes on character-oriented display devices.

Boyer-Moore search. A fast text-searching algorithm that takes advantage of the mix of characters in the search key to skip over positions in the text being searched where a match can't start. See Chapter 15.

Breathing mark. One of two marks used in polytonic Greek to indicate the presence and absence of an “h” sound at the beginning of a word or syllable.

Breve. A U-shaped diacritical mark, usually used to indicate the absence of stress on a particular syllable.

Byte order. The order in which the individual bytes of a multi-byte value are stored in memory or sent over a communication link.

Byte order mark. The Unicode code point value which, by convention, can be used at the beginning of a Unicode text file to indicate its byte order or to announce the Unicode character encoding scheme being used in that file. See Chapter 6.

C0. In an ISO 2022-based encoding, the range of code point values from 0x00 to 0x1F set aside for control characters, or the set of control characters allocated to this range. In most ISO 2022-based encodings, the C0 set consists of the ASCII control characters.

C1. In an ISO 2022-based encoding, the range of code point values from 0x80 to 0x9F set aside for control characters, or the set of control characters allocated to this range (and represented using escape sequences in a 7-bit encoding).

Candrabindu. A mark used in Devanagari and some other Indic scripts to indicate nasalization. See Chapter 9.

Canonical accent ordering. Alternate term for **canonical reordering**.

Canonical composite. A Unicode character whose canonical decomposition is different from the character itself (i.e., a Unicode character representing a character that can also be represented with a combining character sequence). See Chapter 4.

Canonical decomposition. The preferred method of representing a particular character in Normalized Form D. For most characters, this is the same as the character itself, but for characters that can also be represented with combining character sequences (i.e., canonical composites), the equivalent combining character sequence (or, if the same character can be represented with more than one combining character sequence, the equivalent sequence with the combining marks in canonical order). See Chapter 4.

Canonical order. For a given collection of Unicode combining marks applied to a base character, the order in which they are to be arranged in normalized Unicode text. See Chapter 4.

Canonical reordering. The process of putting a sequence of Unicode combining marks into canonical order. See Chapter 4.

Canonical representation. For a character that has multiple representations in Unicode, the representation that must be used in some Unicode normalized form. See Chapter 4.

Canonically equivalent. Two strings are said to be canonically equivalent if they're identical when converted to Normalized Form D. See Chapter 4.

Cantillation mark. A mark used in Hebrew and some other scripts to indicate accent and stress, especially for chanting in a liturgical context. See Chapter 8.

Capital. Alternate term for **upper case**.

Caret. A mark used to indicate a position in text where more text is to be inserted. Also used as an alternate term for **insertion point**.

Caron. A v-shaped diacritical mark drawn over some letters to alter their pronunciation.

Case folding. The process of converting a piece of text to a case-independent representation, usually done as part of comparing it to another piece of text when case differences between the two aren't significant. See Chapter 14.

Case mapping. The process of converting a cased letter (or series of cased letters) to a particular case. See Chapter 14.

Cased letter. A letter with different upper case and lower case forms.

CaseFolding.txt. A file in the Unicode Character Database containing mappings that can be used to perform case folding on a piece of Unicode text. See Chapter 5.

Case-insensitive comparison. The process of comparing two strings in such a way as to treat differences in case as insignificant (e.g., “hello”, “Hello”, and “HELLO” all compare equal).

Cc. The abbreviation used in the UnicodeData.txt file to identify characters as belonging to the “control character” general category. See Chapter 6.

CCII. Chinese Character Code for Information Interchange. A Taiwanese industry-developed encoding standard. See Chapter 2.

CCITT. Consultative Committee for International Telephone and Telegraph. A European standards body that sets standards for the telecommunications industry.

Cedilla. A hook- or comma-shaped mark that attaches to the bottom of some consonants to change their pronunciation.

CESU-8. Compatibility Encoding Scheme for UTF-16, 8-bit. A character encoding form for Unicode similar to UTF-8, but differing in its treatment of supplementary-plane characters. Supplementary-plane characters are basically represented as though they are converted first to UTF-16 and then the UTF-16 surrogates are converted directly to UTF-8. See Chapter 4.

Cf. The abbreviation used in the UnicodeData.txt file to denote character with the “formatting character” general category. See Chapter 5.

CGJ. Abbreviation for **combining grapheme joiner**.

Character. 1) The minimal unit of writing in some writing system (see **grapheme**). 2) Any of a collection of marks having the same semantic and generally considered to be the same by ordinary users of a writing system; the abstract concept of which a glyph is a concrete representation. (Compare **glyph**.) 3) The entity represented by a Unicode code point (or by a code point in some other encoding).

Character code. Alternate term for **code point**.

Character encoding form. A mapping from code points to code units. The way in which code points in a coded character set are represented in memory. The intermediate layer between a coded character set and a character encoding scheme in the Unicode character encoding model. See Chapter 2.

Character encoding scheme. A mapping from code points to bytes. The way in which code points in a coded character set are represented in a serialized medium such as a disk file or communication link. See Chapter 2.

Character generator chip. The device that maps from code point values to lit pixels on the screen in a character-oriented display device.

Character mapping table. The table in a TrueType-based font that maps from code point values to glyph indexes. See Chapter 16.

Character-oriented display. A display device that is controlled by a range of code points in memory, using a character generator chip to map from the code point values to the groups of lit pixels on the screen.

Character repertoire. See **abstract character repertoire**.

Character set. Usually used as a synonym for **coded character set** or **abstract character repertoire**.

CharMapML. Character Mapping Markup Language. An XML-based format for describing mappings between encoding standards.

charset. The parameter used in a number of IETF standards to specify the character representation used. A “charset” is equivalent to a **character encoding scheme**.

Chinese character. One of the characters originally used used for writing Classical Chinese and now used for a variety of languages, including Japanese and Korean.

Choon. The dash-like character used with Katakana characters used to indicate prolongation of a vowel sound.

Choseong. A Hangul jamo character representing the initial consonant sound of a syllable.

Chu Han. Vietnamese name for **Chinese character**.

Chu Nom. The characters used to write Vietnamese prior to the 20th century. The Chu Nom characters are constructed according to the same principles that guided the construction of the Chinese characters, but are specific to Vietnamese.

Circumflex accent. A mark similar to an upside-down V drawn above certain vowels to change their sounds.

CJK. Chinese, Japanese, and Korean.

CJK Joint Research Group. The group that originally developed the unified Han repertoire in Unicode. The forerunner of the **Ideographic Rapporteur Group**.

CJK Miscellaneous Area. The area of the Unicode encoding space running from U+2E80 to U+33FF and containing various characters used with the Han characters in various languages. This includes punctuation, symbols, and various phonetic writing systems used in conjunction with Han characters.

CJK Unified Ideographs area. The area of the Unicode encoding space containing the Han characters.

CJK-JRG. Abbreviation for **CJK Joint Research Group**.

CJKV. Chinese, Japanese, Korean, and Vietnamese.

CL area. The area in an ISO 2022-based encoding running from 0x00 to 0x1F and set aside for control characters. In most ISO 2022-based encodings, the ASCII control characters are assigned to the CL area.

Closing punctuation. Alternate name for **ending punctuation**.

cmap table. Abbreviation for **character mapping table**.

Cn. The abbreviation used in the UnicodeData.txt file for **noncharacter code point**.

CNS 11643. The official encoding standard of the Taiwanese government.

Co. The abbreviation used in the UnicodeData.txt file for **private-use character**.

Code chart. A chart that shows the mappings from characters to code point values in a character encoding standard.

Code page. Alternative name for **coded character set**.

Code point. 1) An abstract numeric value representing a character. 2) In Unicode 2.0, this term was used for what is now called a **code unit**. See Chapter 2.

Code point value. One of the values a code point can have.

Code set. In Extended UNIX Code, one of the four coded character sets it can switch between. See Chapter 2.

Code unit. A fixed-length numeric value representing a code point (or part of a code point) in memory. A code point may map to a single code unit, or to a sequence of code units. See Chapter 2.

Coded character set. A mapping of characters to code point values. See Chapter 2.

Code-switching scheme. A character encoding scheme that allows the representation of more than one coded character set, usually through the use of control characters or escape sequences to switch from one coded character set to another.

Coeng. The Khmer virama. See Chapter 9.

Collation. Another term for language-sensitive string comparison.

Collation element. A collection of weight values for a particular character or contracting character sequence according to a particular sort order.

Collation key. A sequence of weight values for a string according to a particular order. The weights are ordered in such a way that binary comparison of collation keys produces the same result as language-sensitive comparison of the original strings.

Collation strength. In language-sensitive comparison, the level of differences that are significant.

Combining character. A character which combines typographically in some way with the preceding character. Combining characters are generally accents, vowel signs, or other diacritical marks.

Combining character sequence. A sequence of Unicode code points that together represent a single character, usually a letter with one or more marks (such as vowels or diacritical marks) applied.

Combining class. A number from 0 to 255 assigned to every character in Unicode (non-combining characters all have a combining class of 0) and used to arrange the characters in a combining character sequence into canonical order.

Combining grapheme joiner. A character that can combine two grapheme clusters together into a single grapheme cluster.

Combining half mark. A combining character representing half of a double diacritic.

Combining mark. A mark, usually a diacritical mark, that attaches to the character that precedes it in the backing store.

Combining spacing mark. A combining mark that occupies space on the baseline.

Commercial numerals. Alternate term for **accounting numerals**.

Comp_Ex. Old name for **Full_Composition_Exclusion**.

Compact array. A data structure that can be used to represent an array containing repetitive patterns of values in a compact manner that preserves speed of lookup. See Chapter 13.

Compatibility area. The section of the Unicode encoding space dedicated primarily to compatibility characters.

Compatibility character. A character encoded in Unicode solely for round-trip compatibility with some source standard. Many, but not all, compatibility characters are also composite characters.

Compatibility composite. A character whose compatibility decomposition is different from its canonical decomposition; a compatibility character with a preferred representation in Unicode. See Chapter 4.

Compatibility decomposition. A mapping from a compatibility composite to a “preferred” representation consisting of one or more non-compatibility characters. The main difference between a compatibility decomposition and a canonical decomposition is that canonical decompositions are exactly equivalent to the original characters; compatibility decompositions, on the other hand, may lose data (usually formatting data).

Composite character. A character whose canonical or compatibility decomposition is different from the character itself. See Chapter 4.

Composite key. A key for sorting a record that is constructed from two or more other fields in the record. See Chapter 15.

Composition. The process of replacing a combining character sequence (or part of a combining character sequence) with a canonical composite character.

Composition exclusion list. The list of canonical composites that may not appear in any of the Unicode normalized forms. See Chapter 4.

CompositionExclusions.txt. The file in the Unicode Character Database that contains the composition exclusion list.

Compound font. A virtual font that draws on more than one other font for the actual glyph shapes.

Compound glyph. A glyph description in an outline font that doesn’t contain actual outlines, but instead combines together two or more other glyph descriptions. See Chapter 16.

Conform. To meet all the requirements set forth in a standard.

Conjoining hangul jamo. Unicode code points representing Hangul jamo that are meant to be used in sequences to represent Hangul syllables.

Conjunct consonant. The glyph or glyphs representing a series of consonants with no intervening vowel sounds in an Indic script. The characters may combine together into a ligature or may remain as separate glyphs, but all but the last character has a different glyph from the glyph it has in isolation, and the viramas that connect the characters together aren’t visible.

Connector punctuation. A punctuation mark that causes the words on either side to join together into a single word.

Consonant conjunct. See **conjunct consonant**.

Consonant gemination. The doubling or prolonging of a consonant sound, e.g., the “t” sound in “fourteen” as spoken by most Americans.

Consonant stack. A vertical stack of two or more glyphs representing consonants without intervening vowel sounds in Tibetan. The Tibetan analogue to the conjunct consonant in other Indic scripts.

Consonantal. Used to describe an alphabet in which the letters only represent consonants and vowel sounds are represented by applying diacritical marks to the consonants. (In most consonantal scripts, some vowel sounds actually are represented with letters, but these letters generally are also used as consonants.)

Content-transfer-encoding. The entry in a MIME message header specifying the transfer encoding syntax with which the message’s content is encoded. See Chapter 17.

Content-type. The entry in a MIME message header specifying the type of data the message (or message part) contains.

Context-sensitive weighting. A situation in language-sensitive comparison in which a particular character has a different set of weights depending on the surrounding characters. This is usually treated as a combination of contracting character sequence and expanding character. See Chapter 15.

Contextual shaping. The process by which the glyph to use to represent a given character depends on the surrounding characters.

Contour. Alternate term for **loop**.

Contracting character sequence. A sequence of code points that map to a single collation element. See Chapter 15.

Control character. 1) Alternate term for **control code**. 2) The Unicode general category consisting of the Unicode counterparts to the C0 and C1 code points in the ISO 2022 standard.

Control code. A single code point representing a control function.

Control function. A sequence of one or more code points that don’t represent characters, but instead affect the interpreting process in some way (managing the communication protocol, delimiting fields in structured data, telling the receiving equipment to do something, etc.) A control function consisting of a single code point is called a control code or control character. A control function consisting of a sequence of code points is often called an escape sequence.

Control picture. A visible glyph used to represent a control character.

Control point. Alternate term for **off-curve point**.

CR. Carriage return. A control character that tells the receiving equipment to return the print head to the home position. On some systems (notably the Macintosh), this is the new-line function.

CR area. The range of code point values from 0x80 to 0x9F in an ISO 2022-based encoding, reserved for control characters.

CRLF sequence. A code point representing the CR function followed by a code point representing the LF function. The combination of these two code points in a row is treated as the new-line function by some systems (notably DOS and Windows).

Cs. The abbreviation used in the UnicodeData.txt file for **surrogate character**.

Currency symbol. A symbol that precedes or follows a numeral to indicate that it represents a currency value.

Cursive. Written with characters that generally join together into a single unbroken line.

Cursive joining. Contextual shaping that causes a succession of characters to be drawn connected together into a single unbroken line.

Dagesh. A dot drawn inside a Hebrew letter to change its pronunciation, generally from a more liquid sound (such as *v*) to a stopped sound (such as *b*). See Chapter 8.

Dakuten. Two small strokes added to the shoulder of a Kana character to change the syllable's consonant sound to a voiced consonant. See Chapter 10.

Danda. A vertical stroke used as punctuation with many Indic scripts. It functions more or less like a period. See Chapter 9.

Dash. A property assigned to characters in the “dash punctuation” category that are dashes rather than hyphens. See Chapter 5.

Dash punctuation. A category of characters that includes hyphens and dashes. See Chapter 5.

Dasia. The Greek “rough breathing mark,” which was used to indicate the presence of an initial *h* sound at the beginning of a word that starts with a vowel. See Chapter 7.

Dead key. A key that generates a character (usually a diacritical mark) but doesn't advance the print head or carriage to the next character position.

Decimal digit value. A property assigned to characters in the “decimal digit number” category that specifies the numeric value the character represents when used as a decimal digit.

Decimal-digit number. A category of characters that includes all characters that are normally used as decimal digits in place-value notation.

Decomposition. 1) A property assigned to a character that specifies a sequence of one or more characters that is equivalent to that character. Usually the decomposition is either the character itself, a combining character sequence (when the character can also be represented using a combining character sequence), or some kind of preferred representation of the character. A canonical decomposition is strictly equivalent to the character; a compatibility decomposition may lose some data (usually formatting). 2) The process of replacing a character with its decomposition, or of converting a whole string to Normalized Form D or KD. See Chapter 4.

Decomposition type. 1) Whether a particular decomposition is canonical or compatibility. 2) A classification applied to compatibility decompositions that attempts to capture what information is lost by replacing the character with its decomposition.

Default ignorable code point. A character (or unassigned code point value) that should be ignored by default by a process that doesn’t know how to deal with it specifically. This property is generally there to specify code points that a font should treat as invisible (i.e., not draw) instead of handling by drawing a “missing” glyph.

DEL. A control function in ASCII and older telegraphy codes that either is ignored by the receiving equipment or tells the receiving equipment to ignore the preceding character. Typically represented by an all-bits-on signal. (The character derives from the practice of erasing an error in a paper-tape representation of a piece of text by punching out all the holes in the row with the error.) See Chapter 2.

Dependent vowel. In Indic scripts, a mark which, when applied to a consonant character, changes its inherent vowel to some other vowel sound. See Chapter 9.

Deprecated. 1) Strongly discouraged from use. Character assignments can’t be removed from Unicode, so if a character was encoded in error, it’s deprecated. 2) A character property assigned to characters that are deprecated.

Derived data file. A file in the Unicode Character Database whose contents can be derived from data in other files in the database (or from rules in the standard itself).

Derived property. A property whose members can be determined from other properties rather than by merely listing all the members.

DerivedAge.txt. A file in the Unicode Character Database that tells when each character was added to the standard.

DerivedBidiClass.txt. A file in the Unicode Character Database that lists the characters belonging to each of the bi-di categories.

DerivedBinaryProperties.txt. A file in the Unicode Character Database that lists the characters belonging to each of the binary properties (i.e., the properties that only have two values). Right now, the only property listed in this file is the “mirrored” property.

DerivedCombiningClass.txt. A file in the Unicode Character Database that lists the characters in each of the combining classes.

DerivedCoreProperties.txt. A file in the Unicode Character Database that lists characters with various properties. The properties can usually be derived from other properties in UnicodeData.txt and PropList.txt. See Chapter 5.

DerivedDecompositionType.txt. A file in the Unicode Character Database that lists the composite characters with each of the decomposition types.

DerivedJoiningGroup.txt. A file in the Unicode Character Database that classifies the Arabic and Syriac letters according to their basic shape.

DerivedJoiningType.txt. A file in the Unicode Character Database that classifies the Arabic and Syriac letters according to how they can join to their neighbors.

DerivedLineBreak.txt. A file in the Unicode Character Database that lists the characters with each of the line-break properties.

DerivedNormalizationProperties.txt. A file in the Unicode Character Database that lists groups of characters with similar behavior under normalization: for example, characters that can’t appear in normalization form, characters that can only appear in certain normalization forms, characters that turn into more than one character when converted to certain normalization forms, etc.

DerivedNumericType.txt. A file in the Unicode Character Database that groups the various Unicode numeric characters together according to how they can be used to represent numbers.

DerivedNumericValues.txt. A file in the Unicode Character Database that lists all the characters that can represent each numeric value.

DerivedProperties.html. A file in the Unicode Character Database that gives information on all the derived data files in the database.

Descender. The part of a character that extends below the baseline.

Diacritic. 1) A mark that is added to a letter or other character to changes its pronunciation or meaning in some way. 2) A property assigned to all Unicode characters that are used as diacritics.

Diaeresis. 1) Pronouncing a vowel as part of a new syllable rather than as part of a diphthong. 2) A mark (usually a double dot) applied to a letter to indicate diaeresis.

Dialytika. Greek name for **diaeresis**.

Dialytika-tonos. A single character representing both a dialytika and a tonos applied to the same letter.

Digit. A single character in a numeral written with place-value notation.

Digit value. A property applied to all characters that can be used as digits, specifying the numeric value they represent when used as digits.

Digraph. 1) A pair of characters used together to represent a single sound. 2) A Unicode code point representing two characters.

Dingbat. A special symbol, usually used decoratively, especially a character from the Zapf Dingbats font found on most laser printers.

Diphthong. 1) Two vowel sounds pronounced consecutively and sliding into one another. 2) A pair of vowels used together to represent a single sound, whether or not that sound is actually a diphthong.

Directional run. A sequence of consecutive characters with the same resolved directionality.

Directionality. A property that specifies the order in which consecutive characters are arranged on a line of text (and how consecutive lines are arranged on a page). For example, a series of characters with left-to-right directionality are arranged on a line such that the first character in memory appears farthest to the left, with the succeeding characters progressing toward the right. A character may have strong directionality, weak directionality, or neutral directionality. See Chapter 8.

Discontiguous selection. 1) A selection in a text-editing application that includes two or more ranges of characters that aren't contiguous in memory. 2) A selection in a text-editing application that two or more ranges of characters that aren't contiguous on the screen.

Display cell. 1) In a monospaced font, the space occupied by a single character. 2) In East Asian typography, the space occupied by a single Chinese character.

Disunify. To take two different meanings or glyph shapes represented by a single code point value and assign them different code point values. The different uses of the ASCII hyphen character (i.e., hyphen, dash, minus sign) are *disunified* in Unicode.

Dot leader. 1) A sequence of dots spaced at equal intervals and usually used to help lead the eye between two separated pieces of text (such as a chapter name and page number in a table of

contents). 2) A single character consisting of one or more dots intended to be used repeatedly in sequence to form a dot leader. See Chapter 12.

Double danda. A punctuation mark consisting of two vertical strokes used with many Indic scripts. A double danda usually marks the end of a larger run of text, such as a paragraph or stanza, than is ended with a single danda. See Chapter 9.

Double diacritic. A diacritical mark that applies to a sequence of two characters. See Chapter 4.

Draft Unicode Technical Report. A Unicode Technical Report that hasn't reached final approval yet, but whose content has generally been agreed on in principle. Unicode Technical Reports carry no normative force while still in draft form. See Chapter 3.

Dual-joining letter. An Arabic or Syriac letter that is capable of connecting cursorily to letters on either side.

DUTR. Abbreviation for **Draft Unicode Technical Report**.

DUTR #26. A Draft Unicode Technical Report defining CESU-8, a UTF-8-like encoding scheme used for compatibility on some systems. See Chapter 6.

Dynamic composition. The principle of being able to represent a character with a sequence of code points representing its constituent parts.

Early normalization. The principle set forth by the W3C that text must be produced in normalized form by the process that produces it (and kept that way by any process that modifies it). This allows processes that only receive and interpret text to assume it's in normalized form. See Chapter 17.

East Asian width. A property assigned to a character that specifies how many display cells it takes up in East Asian typography. See Chapter 10.

EastAsianWidth.txt. A file in the Unicode Character Database that gives the East Asian width property for each character.

EBCDIC. Extended Binary Coded Decimal Information Code. A family of encoding standards used on some IBM computers that evolved from BCDIC (which in turn evolved from earlier punched-card codes). The original version of EBCDIC added two bits to the six-bit BCDIC code, making it possible to represent a wider variety of punch patterns than was possible with BCDIC.

Ech-yiwn ligature. The ligature that is often formed when the Armenian letters ech and yiwn occur together.

ECMA. A pan-European standards-setting body specializing in information technology standards. ECMA is an ISO member body. The letters used to stand for European Computer Manufacturers' Association.

ECMA-6. The original international-standard version of ASCII. This became ISO 646. See Chapter 2.

ECMA-35. The standard that became ISO 2022. See Chapter 2.

ECMA-94. The standard that became ISO 8859-1. See Chapter 2.

ECMAScript. A standard specifying a core syntax and set of libraries for scripting languages. Netscape's Javascript and Microsoft's JScript, as well as some other scripting languages, are based on the ECMAScript model.

Editing. The process of using a piece of software to interactively change or compose a piece of text.

Elision character. A character that is used to indicate the absence of several other characters. The apostrophe is often used as an elision character in English.

Ellipsis character. 1) A character, such as the series of periods used in English, that is used to indicate the absence of one or more words. 2) **Elision character.**

Em. A unit of measurement along a line of text that is equivalent to the height of the line.

Em dash. A dash that is an em wide, usually used in English to indicate a break in continuity or to set off a parenthetical expression.

Em quad. A space that is an em wide.

Em space. Alternate term for **em quad.**

Embedding level. A number assigned to characters by the Unicode bi-di algorithm that governs how successive directional runs are to be arranged on a line relative to one another. See Chapter 16.

En. Half of an em.

En dash. A dash that is an en wide, usually used in English to separate ranges of numbers or as a minus sign.

En quad. A space that is an en wide.

En space. Alternate term for **en quad.**

Enclosing mark. A combining mark that surrounds the character (or grapheme cluster) it is applied to.

Encoding form. See **character encoding form**.

Encoding scheme. See **character encoding scheme**.

Encoding space. The range of possible code point values for a coded character set, often represented as a two- or three-dimensional array of cells containing characters. The size of an encoding space is often derived from the size in bits of the code point values. “8-bit encoding space” and “16 × 16 encoding space” are both synonymous with “256-character encoding space.” See Chapter 2.

Encoding standard. A standard that specifies a coded character set, character encoding form, character encoding scheme, or some combination of the above.

Endian-ness. Alternate term for **byte order**.

Ending punctuation. A category of characters that serve to mark the end of some range of text (such as a parenthetical expression or a quotation). See Chapter 5.

Erotimatiko. The semicolon-like mark used in Greek as a question mark. See Chapter 7.

Escape sequence. 1) A sequence of more than one code point representing a control function (traditionally, the ASCII ESC character was used as the first character of an escape sequence). 2) A sequence of more than one code point used to represent a single code point in a programming language or other structured text format (e.g., “\u1234” in Java).

Estrangelo. The oldest style of writing the Syriac letters. See Chapter 8.

EUC. Abbreviation for **Extended UNIX Code**.

EUC-CN. A version of Extended UNIX Code with code sets drawn from mainland Chinese encoding standards. See Chapter 2.

EUC-JP. A version of Extended UNIX Code with code sets drawn from Japanese encoding standards. See Chapter 2.

EUC-KR. A version of Extended UNIX Code with code sets drawn from Korean encoding standards. See Chapter 2.

EUC-TW. A version of Extended UNIX Code with code sets drawn from Taiwanese encoding standards. See Chapter 2.

Exception table. An auxiliary data structure used to store records that won’t fit in the main data structure, or cases where the rules specified in the main structure don’t hold. A data structure that needs to store mappings for character sequences might use a main table that holds single-character

mappings and points into one or more exception tables that hold the multiple-character mappings. A data structure that groups characters into categories for some process might consist of a main table that maps whole Unicode general categories to its categories and an exception table that contains category mappings for characters whose mappings don't match the mapping for their Unicode category. See Chapter 13.

Expanding character. A character that maps to more than one collation element. See Chapter 15.

Expands_On_NFC. A derived property containing all the characters that turn into multiple characters when converted to Normalized Form C.

Expands_On_NFD. A derived property containing all the characters that turn into multiple characters when converted to Normalized Form D.

Expands_On_NFKC. A derived property containing all the characters that turn into multiple characters when converted to Normalized Form KC.

Expands_On_NFKD. A derived property containing all the characters that turn into multiple characters when converted to Normalized Form KD.

Explicit embedding character. A character that forces the characters after it to be placed at a new embedding level. See Chapter 8.

Explicit override character. A character that overrides the inherent directionality of the characters after it. See Chapter 8.

Extended ASCII. A term used to refer to any of a large number of eight-bit encoding standards whose first 128 code point assignments are identical to those in ASCII, particularly ISO 8859-1.

Extended UNIX Code. A code switching scheme that allows for the representation of up to four different coded character sets: One is represented with code point values in the G0 space, one with code point values in the G1 space, and the other two with code point values from the G1 space preceded by special “introducer” characters. See Chapter 2.

Extender (property in PropList.txt).

Eyelash *ra*. A special half-form of the Devanagari letter *ra* used in certain conjunct consonants. See Chapter 9.

Fancy text. Alternate term for **rich text**.

Feather mark. Marks used in Ogham to mark the beginning and ending of a piece of text. See Chapter 11.

FF. Form Feed. A control character that tells the receiving equipment to advance the platen to the beginning of the next page.

FIELDATA. A telegraphy code developed by the U.S. Army and which is one of ASCII's most important antecedents. See Chapter 2.

FIGS. A control character that would cause the receiving equipment to interpret succeeding code points as digits or punctuation instead of letters.

Fili. The Thaana vowel signs. See Chapter 8.

Final form. The glyph shape a character takes on when it appears at the end of a word (or, in some cases, when it appears before a character that doesn't join to the preceding character).

Final-quote punctuation. A category assigned to half of the quotation-mark characters in Unicode. Whether the characters in this category are *actually* closing quotation marks or are instead *opening* quotation marks depends on the language of the text.

First series consonants. The Khmer consonants that carry an inherent *a* sound. See Chapter 9.

FNC. A derived property that gives a mapping from a character to a case-folded representation in Normalized Form C.

Fongman. A character used as a list bullet in Thai. See Chapter 9.

Font. 1) A collection of glyphs with a common typographical design. Several fonts may make up a typeface. 2) A computer program or set of tables that, in conjunction with a text rendering process, can draw the glyphs in a particular font.

Forfeda. Extra letters added to the Ogham alphabet after it began to be written on paper. See Chapter 11.

Format effector. Early term for **control character**.

Formatting character. The category assigned to Unicode characters which have no visual presentation of their own, but which exist to alter the way some process operating on the text treats the surrounding characters. See Chapters 5 and 12.

Fourth-level difference. In a multi-level comparison, a difference between fourth-level weights.

Fourth-level weight. In a multi-level comparison, the fourth weight value assigned to a character or contracting character sequence. In the Unicode Collation Algorithm, the fourth level weight is

usually either the character's actual code point value or is algorithmically derived from the character's code point value.

Fraction slash. A special slash character that causes surrounding sequences of digits to combine together into a fraction. If you put a normal slash between 1 and 2, you get “1/2”; if you put a fraction slash between them, you get “½”. See Chapter 12.

French accent sorting. Alternate name for **French secondary**.

French secondary. Treating the secondary level of a multi-level sort as a backwards level. So called because it's required to sort French properly.

FTP. File Transfer Protocol. A protocol used for sending files over the Internet.

Full_Composition_Exclusion. A derived property assigned to characters that can't occur in any of the Unicode normalized forms.

Fullwidth. Used to describe characters that take up a full display cell in East Asian typography. Unicode draws a distinction between characters that are explicitly fullwidth (i.e., have “FULWIDTH” in their names) and those that are implicitly fullwidth (such as the Han characters).

Furigana. Alternate term for **ruby**.

Futhark. The name for a Runic alphabet. See Chapter 11.

G source. The collective name for sources of ideographs submitted to the IRG by the Chinese national body. See Chapter 10.

G0. The range of the ISO 2022 encoding space running from 0x20 to 0x7F and set aside for printing characters. 0x20 is always the space and 0x7F is always the DEL character, leaving 94 code points available for other characters. The term also refers to the set of characters encoded in this space. In most ISO 2022-based encodings, the G0 set of characters is the ASCII printing characters. See Chapter 2.

G1. The range of the ISO 2022 encoding space running from 0xA0 to 0xFF and set aside for printing characters. The term also refers to one of three alternate sets of characters encoded in this space (G2 and G3 are the others). See Chapter 2.

Gakushu Kanji. A list of characters that all Japanese learn in elementary school. See Chapter 2.

Gap storage. A method of storing a large amount of contiguous data, such as the characters in a document being edited, that can greatly improve the performance of processes, such as interactive

text editors, that manipulate the data. The basic principle is that the memory block is larger than the data, to give it room to grow before a new block must be allocated, and the extra memory in the block is stored at the position of the last change rather than at the end of the block, minimizing the number of characters that have to be moved around when changes occur. See Chapter 16.

Garshuni. Arabic written using the Syriac alphabet. See Chapter 8.

GB 2312. The main Chinese national standard for character encoding.

Gemination. See **consonant gemination**.

General category. A property given to every Unicode character that classifies it according to its general purpose. See Chapter 5.

General scripts area. The area of the BMP running from U+0000 to U+1FFF and containing the characters for all of the modern non-CJK writing systems.

Geta mark. A mark sometimes used in Japanese typography to take the place of an ideograph for which a glyph can't be found.

GL area. In ISO 2022, the range running from 0x20 to 0x7F and set aside for the encoding of printing characters. In most ISO 2022-based encodings, the ASCII characters are placed into the GL range.

Glottal stop. The stoppage of sound caused by the closing of the glottis (the back of the throat). The hyphen in “uh-uh” is a glottal stop.

Glue character. Alternate term for **non-breaking character**.

Glyph. A concrete written or printed representation of a character or group of characters. See Chapter 3.

Glyph index. A number identifying a particular glyph description in a font file. The cmap table in a TrueType font maps from code points to glyph indexes.

Glyph metamorphosis table. In an AAT font, a table that performs complex mappings on glyph indexes. This is the table responsible for context-sensitive glyph selection, ligature formation, Indic vowel reordering, etc. See Chapter 16.

Glyph selection. The process of choosing the proper glyph to display for a character in a particular place.

Glyph substitution table. In an OpenType font, a table that performs complex mappings on glyph indexes. This is the table responsible for context-sensitive glyph selection, ligature formation, etc. See Chapter 16.

Glyphic variant. A different glyph representing the same underlying character. The term is usually used in opposition to “different character.” For example, \$, \$, and \$ are all glyphic variants of the same character—\$ and ¢ are different characters.

GR area. The part of the ISO 2022 encoding space running from 0xA0 to 0xFF and set aside for printing characters. See Chapter 2.

Grapheme. 1) A minimal unit of writing in some written language; a mark that is considered a single “character” by an average reader or writer of a particular written language. 2) In versions of Unicode prior to Unicode 3.2, this term was used for the concept that is now called a **grapheme cluster**.

Grapheme cluster. A sequence of one or more Unicode code points that should be treated as an indivisible unit by most processes operating on Unicode text, such as searching and sorting, hit-testing, arrow key movement, and so on. A grapheme cluster may or may not correspond to the user’s idea of a “character” (i.e., a single grapheme)—for instance, an Indic orthographic syllable is generally considered a grapheme cluster but may still be seen by an average reader or writer as several letters. The Unicode standard puts forth a default definition of grapheme clusters, but permits language-sensitive tailoring—for example, a contracting character sequence in most sort orders is a language-specific grapheme cluster. See Chapter 4.

Grapheme joiner. See **combining grapheme joiner**.

Grapheme_Base. A property given to all characters that are grapheme clusters unto themselves unless they occur next to a Grapheme_Extend or Grapheme_Link character. This is a derived property—all characters that aren’t Grapheme_Extend or Grapheme_Link characters are Grapheme_Base characters.

Grapheme_Extend. A property given to all characters that belong to the same grapheme cluster as the characters that precede them. Combining marks are Grapheme_Extend characters.

Grapheme_Link. A property given to all characters that cause the characters that precede and follow them to be considered part of the same grapheme cluster. The combining grapheme joiner and the viramas are Grapheme_Link characters.

Graphic character. Alternate term for **printing character**.

Grave accent.

GSUB table. Abbreviation for **glyph substitution table**.

Guillemet. A v-shaped mark used in a number of languages, notably French, as a quotation mark. See Chapter 12.

H source. Collective name for the various sources of ideographic characters submitted to the IRG by the Hong Kong Special Administrative Region.

Hacek. A v-shaped mark drawn above some consonants to indicate a change in sound. See Chapter 7.

Half mark. See **combining half mark**.

Half-form. The form that a consonant in Devanagari and some other Indic writing systems often takes when it's part of a conjunct consonant. The term comes from the fact that a half-form is usually an abbreviated version of the character's normal shape, usually omitting the vertical stem. Half-forms are glyphic variants of normal consonants and are not given a separate encoding in Unicode. See Chapter 9.

Halfwidth. The term given to characters that occupy half of a normal display cell in East Asian typography. Latin letters, for example, are generally considered halfwidth. The Unicode standard draws a distinction between characters that are explicitly halfwidth—these have “HALFWIDTH” in their names—and those that are implicitly halfwidth. Halfwidth characters don't always take up exactly half the space along the baseline of fullwidth characters; often they're proportionally spaced, and the term is used merely to indicate that they take up less space along the baseline than a normal ideograph. See Chapter 10.

Hamza. A mark used in Arabic to indicate a glottal stop. See Chapter 8.

Han character. Alternate term for **Chinese character**.

Han unification. The process of collecting Chinese characters from a variety of disparate sources and weeding out the characters from different sources that are really duplicates of each other. The characters that are duplicates are encoded once in Unicode, not once for every source encoding they appear in. The process of Han unification involves keeping careful track of just which characters from which sources are unified. See Chapter 10.

Handakuten. The small circle applied to certain Kana characters that turns the initial *h* sound into an initial *p* sound. See Chapter 10.

Hangul. The alphabetic writing system used for writing Korean.

Hangul syllable. 1) A cluster of Hangul jamo representing a single Korean syllable. 2) A single Unicode code point representing a Hangul syllable.

Hangzhou numerals. Alternate term for **Suzhou numerals**.

Hanja. Korean for **Chinese character**.

Hankaku. Japanese for **halfwidth**.

Hanzi. Mandarin for **Chinese character**.

Hard sign. A letter of the Cyrillic alphabet that indicates the absence of palatalization in a spot where it would normally occur.

Headstroke. The characteristic stroke at the top of all of the letters in most Indic writing systems. Depending on the script, the headstroke may be different shapes. See Chapter 9.

Hex_Digit. A property assigned to those characters that can be used as hexadecimal digits.

High surrogate. A code unit value that is used as the first code unit in a surrogate pair, or the code point value corresponding to it. The high surrogates run from U+D800 to U+DBFF.

Higher-level protocol. A protocol that is out of the scope of the Unicode standard but that uses the Unicode standard as its basis. XML is an example of a higher-level protocol.

Hindi numerals. Often used to refer to the native digits used with the Arabic alphabet. (*Not* to the digits used with Devanagari.)

Hint. Alternate term for **instruction**.

Hiragana. The more cursive form of Kana used for writing native Japanese words and grammatical endings. See Chapter 10.

Hit testing. The process of determining the position in the backing store that corresponds to a particular pixel position in a piece of rendered text.

HKSCS. Hong Kong Supplemental Character Set. A collection of characters used in the Hong Kong Special Administrative Region. Many of these characters are used specifically to write Cantonese.

Hollerith. Herman Hollerith, who developed a system of using punched cards for data processing.

Hollerith code. Any of a number of similar methods of representing alphanumeric data on punched cards.

HTML. Hypertext Markup Language. A method of representing styled text and hypertext links between documents. The standard method of representing Web page content.

HTTP. Hypertext Transfer Protocol. A protocol for requesting and transferring hypertext documents, and the standard communication protocol on the Web.

Hyphen (property in PropList.txt). A property assigned to Unicode dash-punctuation characters (and a couple of connector-punctuation characters) that are hyphens .

Hyphenation. The process of dividing words between lines of text, usually at syllable boundaries and with a hyphen added at the end of the first line, in an effort to improve the appearance of justified text.

HZ. A code-switching scheme used with Chinese text. See Chapter 2.

IAB. Abbreviation for **Internet Architecture Board**.

IANA. Abbreviation for **Internet Assigned Numbers Authority**.

ICU. Abbreviation for **International Components for Unicode**.

IDC. Abbreviation for **ideographic description character**.

IDS. Abbreviation for **ideographic description sequence**.

Identifier. A sequence of characters used to identify some entity. Variable and function names in programming languages are examples of identifiers.

Ideograph. 1) A character that represents an idea. 2) Alternate term for **Chinese character**.

Ideographic. A property assigned to ideographic characters in Unicode: specifically the Han characters and radicals.

Ideographic description character. A special character that indicates how the two or three ideographic characters or ideographic description sequences following it are to be combined to form a single character. Officially, ideographic description characters are printing characters, but a smart rendering engine can use them to approximate the appearance of an otherwise-unrepresentable ideograph. See Chapter 10.

Ideographic description sequence. A sequence of Unicode characters that describe the appearance of an otherwise-unrepresentable Chinese characters. An ideographic description sequence contains at least one ideographic description character.

Ideographic Rapporteur Group. A formal standing subcommittee of ISO/IEC JTC1/SC2/WG2 responsible for the maintenance and development of the unified Han repertoire in Unicode and ISO 10646. The IRG is made up of experts from wide variety of Asian countries that use the Han characters, as well as a few representatives of the UTC. See Chapter 10.

Ideographic variation indicator. A special character that can be used to indicate that the following character is only an approximation of the desired Chinese character, which can't be represented in Unicode. See Chapter 10.

IDS. Abbreviation for **ideographic description sequence**.

IDS_Binary_Operator. A property given to ideographic description characters that take two operands.

IDS_Tertiary_Operator. A property given to ideographic description characters that take three operands.

IEC. Abbreviation for **International Electrotechnical Commission**.

IETF. Abbreviation for **Internet Engineering Task Force**.

Ignorable character. A character that is to be ignored by a language-sensitive string comparison routine, or that is only to be considered significant at certain levels of comparison. The Unicode Collation Algorithm defines a large collection of characters (mostly spaces and punctuation) as ignorable by default; tailorings may specify more ignorable characters or make certain characters significant that are ignorable by default. See Chapter 15.

Independent form. The glyph shape taken by a character in the absence of any surrounding characters to which it can connect cursively. Usually used with Arabic and Syriac letters.

Independent vowel. A full-fledged letter (as opposed to a mark that gets applied to another letter) representing a vowel sound in an Indic script. Generally, independent vowels are used for word-initial vowel sounds and for vowel sounds that follow other vowel sounds. In some scripts, the “independent vowels” take the form of dependent vowels applied to a “null” consonant. See Chapter 9.

Index.txt. A file in the Unicode Character Database that lists all the Unicode characters by name; it’s basically a soft copy of the character names index in the Unicode book. See Chapter 5.

Indic script. One of a large collection of scripts used mostly in India, Southeast Asia, and certain islands in the Pacific and derived from the ancient Brahmi script. See Chapter 9.

Informative. A statement, property, or recommendation in a standard that need not be followed to conform to the standard. Informative elements are usually there to clarify or supplement the standard in some way.

Inherent directionality. The directionality a character has by virtue of its definition in the Unicode standard, as opposed to **resolved directionality**.

Inherent vowel sound. The vowel sound a consonant in an Indic script has when it has no dependent vowels attached to it. (Letters in Indic scripts actually represent syllables, not just consonants.) See Chapter 9.

Initial form. The glyph shape that a character takes when it appears at the beginning of a word, or when it’s preceded by a character that doesn’t connect to the following character.

Initial-quote punctuation. A property assigned to half of the quotation-marks in Unicode. Whether the characters with this property are *actually* opening quotation marks or are instead *closing* quotation marks depends on the actual language of the text.

Input method. A piece of software that translates user actions on some input device to characters that are passed as input to a text-editing utility. The term is usually used specifically to refer to code that permits input of Chinese characters with a series of keystrokes (and often interactive menu choices).

Insertion point. In an interactive text-editing application, the position, usually marked by a blinking vertical line, where any text the user types will be inserted.

Instruction. An element of a font that indicates how to deform the ideal shape of a glyph to better fit the coordinate system of the output device when the output device has low resolution. Instructions are used, for example, to make the characters in an outline font designed for a printer look good (or at least legible) on the screen.

Interact typographically. Affect the shapes or positions of characters around it. Two characters are said to interact typographically if the glyph shape(s) you get when the characters are juxtaposed in memory is different from the glyph shapes you get when the characters occur in isolation. A character *a* and the character *b* that follows it in memory interact typographically if 1) *b* is drawn somewhere other than in the next position after *a* in the dominant writing direction (e.g., in a left-to-right script, if *b* is drawn above, below, to the left of, inside, or overlaying *a*), 2) the shape of either *a* or *b* or both is different from the shape the character has in isolation, or 3) *a* and *b* combine together into a ligature.

Interlinear annotation. Annotations appearing between lines of text that help to clarify the pronunciation or meaning of unfamiliar characters or words. See Chapter 10.

International Components for Unicode. An open-source library of routines for performing different operations on Unicode text, including character encoding conversion, Unicode normalization, language-sensitive comparison and searching, boundary analysis, transliteration, case mapping and folding, and various other processes. See Chapter 17.

International Electrotechnical Commission. An international body that issues international standards for various technical and engineering disciplines.

International Organization for Standardization. An international body that issues international standards in a wide variety of disciplines and serves as an umbrella organization for the various national standards bodies.

International Phonetic Alphabet. A special set of characters used for phonetic transcription. See Chapter 7.

International Telegraphy Alphabet #2. An international encoding standard used for several decades in telegraphy. See Chapter 2.

Internationalization. The process of designing software so that it can be localized for various user communities without having to change or recompile the executable code.

Internet Architecture Board.

Internet Assigned Numbers Authority. A sister organization to the IETF that acts as a central registration authority for various internet-related identifiers and numbers, such as charsets, MIME types, and so forth. (The IANA used to be responsible for domain names, but isn't anymore.)

Internet Engineering Task Force. An informal organization of engineers who do Internet-related work and which issues various standards for the operation of different parts of the Internet. Most Internet-related standard, other than those relating to the World Wide Web, are IETF standards.

Internet-draft. A document issued by the IETF that either represents standards-track proposal in its very earliest stages of discussion or that contains only informative information. Internet-drafts carry no normative force.

Inversion list. A data structure that represents a collection of values with an array containing only the initial values in each contiguous range of values. See Chapter 13.

Inversion map. An extension of the inversion list that represents a collection of mappings from one set of values into another with an array that contains only the initial values for each contiguous range of source values that map to the same result value. See Chapter 13.

Invisible formatting character. A Unicode code point that has no visual presentation of its own, but affects the way some process operating on text treats the surrounding characters.

Invisible math operator. A Unicode code point representing a mathematical operation whose application is implicit (i.e., indicated by the simple juxtaposition of symbols in a mathematical formula). Invisible math operators have no visual presentation (although sometimes they affect the spacing of the surrounding characters), but exist to help software parsing a mathematical expression. See Chapter 12.

IPA. Abbreviation for **International Phonetic Alphabet.**

IRG. Abbreviation for **Ideographic Rapporteur Group.**

IRV. International Reference Version. For international standards such as ISO 646 that allow for variation between countries, the IRV is the version to use in the absence of a country-specific version. The International Reference Version of ISO 646 is American ASCII.

ISCII. Indian Script Code for Information Interchange. An encoding standard issued by the Indian government for representing various Indic scripts used in India.

ISO. Abbreviation for **International Organization for Standardization.**

ISO 10646. Unicode's sister standard. ISO 10646 is an official international standard, rather than an ad-hoc industry standard. Since the early 1990s, ISO 10646 has exactly the same assignments of

characters to code point values as Unicode and shares some other architectural characteristics with it as well. Unicode goes beyond ISO 10646 to specify more details of character semantics and implementation details, and there are some other subtle differences between the two standards. See Chapter 2 for details. The current formal designations for ISO 10646 are ISO/IEC 10646-1:2000 and ISO/IEC 10646-2:2001 standards; “ISO 10646” is used informally to collectively refer to both (or to earlier versions).

ISO 10646 comment. Informative text about a character that occurs in the ISO 10646 code charts and is included as a formal Unicode character property for compatibility with ISO 10646.

ISO 2022. An ISO standard which specifies a general encoding-space layout for various other encoding standards to follow, and which specifies an elaborate system of escape sequences for switching between different coded character sets.

ISO 6429. A standard that specifies semantics for control characters in various ISO 2022-derived encodings.

ISO 646. The international standard based on ASCII. The International Reference Version of ISO 646 is the same as ASCII, but the ISO 646 standard itself specifies a group of “national use” code points that can be assigned to characters other than the ones in the International Reference Version by various national bodies. With a few exceptions, such as JIS-Roman, the national-use variants of ISO 646 have fallen into disuse. See Chapter 2.

ISO 8859. A family of ISO-2022-compatible encoding standards for various European languages and languages used by large populations of people in Europe. The first 128 code point assignments in all of the ISO 8859 encodings are identical to the ISO 646 IRV, making all of them downward compatible with ASCII.

ISO 8859-1. An encoding standard for representing various Western European languages using the Latin alphabet. This the most common “extended ASCII” encoding standard and was the default character encoding for HTML up until recently. The first 256 code point assignments in Unicode are identical to ISO 8859-1. ISO 8859-1 is slowly being superseded by ISO 8859-15 (“Latin-9”), which is almost the same, but removes a few less-frequently-used character assignments and adds a few new characters, including the Euro symbol. See Chapter 2.

ISO-2022-CN. A code-switching scheme for Chinese. See Chapter 2.

ISO-2022-CN-EXT. A code switching scheme for Chinese similar to ISO-2022-CN, but including support for a greater number of characters. See Chapter 2.

ISO-2022-JP. A code-switching scheme for Japanese. See Chapter 2.

ISO-2022-KR. A code-switching scheme for Korean. See Chapter 2.

Isolated form. Alternate term for **independent form**.

Isshar. A mark used in Bengali to write the name of God. See Chapter 9.

ITA2. Abbreviation for **International Telegraphy Alphabet #2**.

J source. Collective name for sources of Han ideographs submitted to the IRG by the Japanese national body. See Chapter 10.

Jamo. The individual “letters” of the Hangul writing system. Two or more jamo, representing individual sounds, are arranged in a square display cell to form a Hangul syllable. Modern Hangul syllables can be represented in Unicode either with series of individual code points representing the jamo, or with a single code point representing the whole syllable. See Chapters 4 and 10.

Jamo short name. A property assigned to the Hangul jamo in Unicode that is used to derive names for the Hangul syllables. See Chapter 5.

Jamo.txt. A data file in the Unicode Character Database that specifies the Jamo short name property for the characters in Unicode that have this property.

Japanese Industrial Standards Commission. The main Japanese national standards-setting body, and the Japanese member of ISO.

Jinmei-yo Kanji. A list of characters officially sanctioned by the Japanese government for use in personal names. See Chapter 2.

JIS C 6226. Early term for **JIS X 0208**.

JIS X 0201. A Japanese encoding standard that includes code points for the ASCII characters and the Katakana characters.

JIS X 0208. The main Japanese encoding standard that encodes kanji characters.

JIS X 0212. A supplemental Japanese standard that adds more kanji to those already encoded by JIS X 0208.

JIS X 0213. A supplemental Japanese standard that adds more kanji to those already encoded by JIS X 0208 and JIS X 0212.

JISC. Abbreviation for **Japanese Industrial Standard Commission**.

JIS-Roman. Alternate term for **JIS X 0201**.

Johab. An encoding scheme for Korean that makes it possible to encode all of the modern Hangul syllables as single code points. See Chapter 2.

Join_Control. A property given to invisible formatting characters that control cursive joining and ligature formation.

Joiner. See **Zero-width joiner**.

Jongseong. A Hangul jamo representing a syllable-final consonant sound.

Joyo Kanji. A standardized list of Japanese characters used in government documents and newspapers.

JTC1. Joint Technical Committee #1 of the International Organization for Standardization and the International Electrotechnical Commission. This committee is responsible for developing and publishing information-technology standards.

Jungseong. A Hangul jamo representing a vowel sound.

Justification. The process of distributing extra white space on a line of text, usually to cause the text to be flush with both margins.

K source. Collective name for sources of ideographs submitted to the IRG by the Korean national bodies.

Kana. Collective name for Hiragana and Katakana, the two syllabic writing systems used in Japanese.

Kangxi radical. One of the radicals used to classify Chinese characters in the Kangxi dictionary.

Kangxi Zidian. A dictionary for Chinese completed in 1716 and widely considered definitive.

Kanji. Japanese for **Chinese characters**.

Kashida. A long connecting stroke used to elongate Arabic words to make them fit a specified space. See Chapter 8.

Kashida justification. The process of justifying a line of Arabic text by inserting kashidas.

Katakana. The more angular form of Kana used in Japanese to write foreign words and onomatopoeia. See Chapter 10.

Kerning. The process of adjusting the spacing between individual pairs of characters (usually by moving them closer together) to achieve a more pleasing appearance. See Chapter 16.

Key closure. The process of making sure that for every key in a trie-based data structure, all proper prefixes of that key are also included as keys. See Chapter 13.

Khan. A Khmer punctuation mark used at the ends of sentences. See Chapter 9.

Khomut. A Thai punctuation mark used at the end of a document. See Chapter 9.

Khutsuri. A style of Georgian writing that mixed *asomtavruli* and *nuskhuri* and was used primarily in ecclesiastical settings. See Chapter 7.

Killer. Alternate term for **virama**.

Koronis. A diacritical mark used in monotonic Greek to indicate elision of vowels. See Chapter 7.

KS C 5601. Old name for **KS X 1001**.

KS X 1001. One of the main Korean national character encoding standards.

Kunddaliya. A Sinhala punctuation mark that appears at the end of sentences. See Chapter 9.

L2. Committee L2 of the National Committee for Information Technology Standards, which is the American national body reporting to WG2.

Labialization. Beginning a vowel or ending a consonant with the lips puckered. English speakers tend to think of this as inserting a *w* sound between a consonant and a vowel.

Lam-alef ligature. The glyph that is formed when the Arabic letter lam is followed by the Arabic letter alef. This is the one mandatory ligature in Arabic. See Chapter 8.

Last resort font. A font used (usually by a compound font) to provide glyphs for characters that couldn't be displayed by any other font. A typical last-resort font includes different glyphs for each of the scripts covered by Unicode, allowing the user to at least get an idea of which language (or script) the missing characters were in and, thus, what kind of font would need to be provided to see them. See Chapter 18.

Latin-1. Alternate term for **ISO 8859-1**.

Left-joining letter. An Arabic or Syriac letter that can connect to a letter only on its left. See Chapter 8.

Left-joining vowel. An Indic dependent vowel that attaches to the left-hand side of the consonant it is applied to. See Chapter 9.

Left-to-right embedding. An explicit embedding character that starts a new embedding level and specifies its direction to be left-to-right. See Chapter 8.

Left-to-right mark. An invisible formatting character that has string left-to-right directionality and which can be used to cause characters of neutral directionality to resolve to left-to-right directionality. See Chapter 8.

Left-to-right override. An explicit override character that causes all following characters to be treated as having strong left-to-right directionality. See Chapter 8.

Legacy encoding. Generally used to mean “any encoding other than Unicode.” More specifically used to refer to whatever encoding standard or standards were in use before Unicode.

Letter. 1) A single base character in an alphabetic or alphasyllabic writing system. Usually represents a single consonant or vowel sound (although in many alphabetic writing systems, vowel sounds are represented with diacritical marks). 2) Unicode uses this word in many places (including the general categories) to refer to any character that is used to make up words and isn’t a diacritical mark, whether that character is a “letter” in the traditional sense, a syllable, or an ideograph.

Letter number. A category containing characters that represent numeric values but have the form or one or more letters. See Chapter 5.

Letter-mark combination. A letter with one or more diacritical marks applied to it.

Level separator. A sentinel value used in collation keys to delimit the weight values for one level from the weight values for the next level. The value is usually chosen to be lower than any of the weight values, so that shorter strings compare before longer strings. Level separators can be eliminated if the weight values for one level are disjoint from the weight values for the next level. See Chapter 15.

Lexical comparison. Alternate term for **binary comparison**.

LF. Line Feed. A control function that signals the receiving equipment to advance the platen to the next line (on some systems, this also implies moving the carriage back to the home position). The LF character is used as the new-line function on some systems, particularly most UNIX-based systems and the Java programming language.

LI18N/UNIX. The Linux Internationalization Initiative. An offshoot of the Linux Standards Base movement which published a specification for the internationalization capabilities that should be included in a Linux implementation. See Chapter 17.

Ligature. 1) A single glyph that represents more than one character and is formed when a certain sequence of characters occurs in succession. A ligature often has a shape that is very different from the shapes the characters that make it up have in isolation. 2) A glyph in a font file that is drawn

when a certain sequence of two or more code points occur together in the backing store. Note that these two definitions are not the same—a ligature in the normal sense may be drawn with two specially-shaped glyph descriptions in certain fonts, and an accented letter may be drawn with a “ligature” in a font, even though an accented letter is not normally considered a “ligature.”

Line break. The point in a piece of text where one line ends and another line begins.

Line break property. A property applied to all Unicode characters that describes how they are to be treated by a line-breaking routine. See Chapter 16.

Line breaking. The process of dividing a piece of text into lines.

Line layout. The process of arranging a series of characters on a line of text.

Line separator. 1) A control function that signals the beginning of a new line, but not the beginning of a new paragraph. 2) The Unicode character, U+2028, intended to unambiguously represent this function.

Line starts array. An array that stores the offset in the backing store of the first character on each line of a piece of displayed text (or the number of code points in each line). Line breaking can be thought of as the process of building the line starts array (or keeping it up to date as the text changes). See Chapter 16.

LineBreak.txt. The file in the Unicode Character Database that specifies each character’s line-break property.

Linked list. A data structure that stores a sequence of values using a separate memory block for each value, with pointers from one memory block to the next.

Linux Standards Base. A specification that specifies a standard feature set for Linux implementations. See Chapter 17.

Little-endian. 1) A serialization format that sends and receives multiple-byte values least-significant-byte-first. 2) A processor architecture that stores the least significant byte of a multi-byte value lowest in memory.

Ll. The abbreviation used in the UnicodeData.txt file for the “lowercase letter” property.

Lm. The abbreviation used in the UnicodeData.txt file for the “modifier letter” property.

Lo. The abbreviation used in the UnicodeData.txt file for the “other letter” property.

Locale-dependent. Dependent on the language and dialect of the text, and often on the nationality of the person writing or reading the text. Collation behavior is locale-dependent.

Localization. The process of converting an application for use by a new user community. Localization involves not just translating any user-visible text, but also altering things like pictures, color schemes, window layouts, and number, date, and time formats according to the cultural conventions of the new user community.

Logical character. Alternate term for **grapheme cluster**.

Logical order. The order in which the characters are pronounced and (usually) typed.

Logical selection. A selection that represents a contiguous range of code points in the backing store, but which may appear visually discontinuous.

Logical_Order_Exception. A property given to characters that in Unicode are stored in visual order rather than logical order. This property currently consists just of left-joining vowel signs in Thai and Lao. See Chapter 5.

Logograph. A character representing a single word. Sometimes used to describe Chinese characters.

Long s. An alternate glyph shape for the lowercase letter *s* that was used as a non-final form for the letter *s* in the eighteenth century. In handwritten text, it looks like an integral sign, and in printed text, it looks kind of like the letter *f*. See Chapter 7.

Loop. A series of points in an outline describing a closed shape consisting of line and curve segments. A glyph description usually consists of one or more loops. See Chapter 16.

Low surrogate. A code unit value used to represent the second code unit in a surrogate pair, or the code point values corresponding to these code unit values. The low-surrogate range is from U+DC00 to U+DFFF.

Lower case. In a bicameral script, the set of letters used to write most words.

Lowercase letter. The property given to lowercase letters in bicameral scripts.

LRE. Abbreviation for **Left-to-Right Embedding**.

LRM. Abbreviation for **Left-to-Right Mark**.

LRO. Abbreviation for **Left-to-Right Override**.

LS. Abbreviation for **line separator**.

LSB. 1) Abbreviation for **Linux Standards Base**. 2) Abbreviation for “least significant bit” or “least significant byte.”

Lt. Abbreviation used in the UnicodeData.txt file to represent the “titlecase letter” property.

LTRS. A control function in some old telegraphy codes that indicated to the receiving equipment that it should interpret succeeding code points as letters. See Chapter 2.

Lu. Abbreviation used in the UnicodeData.txt file to represent the “uppercase letter” property.

LZW. Lempel-Ziv-Welch. A general-purpose dictionary-based compression algorithm.

Macron. A line drawn above a vowel to indicate a long vowel sound.

Madda. A mark used with some Arabic letters to indicate a glottal stop. See Chapter 8.

Maiyamok. A mark used in Thai to indicate repetition of the syllable that precedes it. See Chapter 9.

Major version. The part of the Unicode version number before the first decimal point. The major version number is incremented any time a new version of Unicode is published in book form.

Mark. A collection of Unicode general categories containing the various combining characters.

Markup. A higher-level protocol that uses certain sequences of characters to impose structure on the text or supplement it with additional information.

Maru. Alternate term for **handakuten**.

Mathematical symbol. A category containing symbols that are normally used in mathematical expressions.

MathML. Mathematical Markup Language. An XML-based file format for describing mathematical expressions.

Medial form. The glyph shape used for a character when it occurs in the middle of a word, or for Arabic letters, when it can connect to letters on both sides.

MICR. Magnetic Ink Character Recognition.

MIME. Multipurpose Internet Mail Extensions. A method of transmitting different data types in a standard RFC 822 mail message. See Chapter 17.

Minor version. The second number in a Unicode version number. The minor version is updated anytime a new version of Unicode is issued that includes new characters or significant architectural changes, but that version isn’t issued in book form.

Mirrored. The property given to characters that have a different glyph shape when their resolved directionality is right-to-left than they have when their resolved directionality is left-to-right. The two glyph shapes are generally mirror images of each other, and often (but not always) mirrored characters come in pairs that can simply exchange glyphs in right-to-left text. See Chapter 8.

“Missing” glyph. A glyph that is drawn when the font doesn’t have a glyph for a particular code point value.

Mkhedruli. The modern form of the Georgian alphabet. See Chapter 7.

Modifier letter. A spacing character that has no pronunciation of its own but changes the pronunciation of the preceding letter. See Chapter 5.

Modifier symbol. Similar to a modifier letter, but used for spacing clones of combining marks and other characters that don’t necessarily get treated as letters. See Chapter 5.

Monospaced. Used to describe a font in which all the characters are the same width.

Monosyllabic language. A language in which all of the words have only one syllable. Many Chinese and related languages are generally considered monosyllabic, although most of them actually do have some polysyllabic words.

Monotonic Greek. A simplified system of Greek spelling that went into effect in the 1970s and greatly simplified the use of diacritical marks with Greek letters.

Morse. Samuel Morse, who invented the telegraph in 1837.

Morse code. The code that was used to send text with Morse’s telegraph (and is still used in certain applications today). It used variable-length sequences of long and short signals to represent characters.

mort table. Abbreviation for **glyph metamorphosis table**.

Multi-level comparison. The process of comparing a pair of strings in multiple passes, where certain types of differences are only significant if there are no higher-level differences between the strings (for example, case differences only count if the strings are otherwise the same).

Murray. Donald Murray, who invented the first teletype machine around 1900.

Muusikatoan. A Khmer mark that converts a second-series consonant to the first-series consonant. See Chapter 9.

Name. A property given to each character that identifies it and attempts to convey something about its meaning or use. See Chapter 5.

NamesList.html. A file in the Unicode Character Database that describes the NamesList.txt file.

NamesList.txt. A file in the Unicode Character Database that is used to produce the code charts in the Unicode standard. It lists all the characters by name, along with a variety of annotations. See Chapter 5.

Narrow. A property assigned to all characters that are implicitly halfwidth, and that characters of neutral or ambiguous East Asian width resolve to when used in Western typography.

Nasalization. Pronouncing a vowel sound with the nasal passage open. The first and last vowel sounds in the French phrase *en passant* are nasalized.

National Committee on Information Technology Standards. An American standards body specializing in information technology standards. NCITS acts as an umbrella organization for various information-technology-related ANSI working groups, including L2, the group responsible for representing the American position to WG2.

National digit shapes. A deprecated invisible formatting character that causes the code points in the ASCII range that represent digits to represent the native digits for the current language (usually Arabic) rather than the European digits they normally represent. See Chapter 12.

National use character. 1) A code point in ISO 646 (and some other encoding standards) that may represent a different character in different countries (i.e., different countries can have their own versions of ISO 646 that can be different from the International Reference Version and still conform to the standard). 2) A character represented by a national-use code point value.

NBSP. Abbreviation for **non-breaking space**.

NCITS. Abbreviation for **National Committee on Information Technology Standards**.

NCR. Abbreviation for **numeric character reference**.

Nd. The abbreviation used in the UnicodeData.txt file for the “decimal-digit number” general category.

NEL. New Line. A character in the C1 space (and in EBCDIC) used on some systems to represent the new-line function.

Nestorian. One of the modern writing styles of the Syriac alphabet. See Chapter 8.

Neutral directionality. No inherent directionality. A character with neutral directionality takes on the directionality of the surrounding characters.

New-line function. A control function that indicates the end of a line of text and the beginning of a new line. Depending on the application, the new-line function may also indicate the beginning of a new paragraph.

NFC. Abbreviation for **Normalized Form C**.

NFC_MAYBE. A derived property given to characters that may occur in Normalized Form C, but only in certain contexts. If a string contains NFC_MAYBE characters and no MFC_NO characters, extra analysis must be done to determine if it is, in fact, in Normalized Form C.

NFC_NO. A derived property given to characters that may not occur in Normalized Form C.

NFD. Abbreviation for **Normalized Form D**.

NFD_NO. A derived property given to characters that may not occur in Normalized Form D.

NFKC. Abbreviation for **Normalized Form KC**.

NFKC_MAYBE. A derived property given to characters that may occur in Normalized Form KC, but only in certain contexts. If a String contains NFKC_MAYBE characters and no NFKC_NO characters, extra analysis must be done to determine if it is, in fact, in Normalized Form KC.

NFKC_NO. A derived property given to characters that may not occur in Normalized Form KC.

NFKD. Abbreviation for **Normalized Form KD**.

NFKD_NO. A derived property given to characters that may not occur in Normalized Form KD.

Nigori. Alternate term for **dakuten**.

Nikahit. A Khmer sign used to indicate nasalization. See Chapter 9.

NI. Abbreviation used in the UnicodeData.txt file to represent the “letter number” general category.

NL. Alternate term for **NEL**.

NLF. Abbreviation for **new-line function**.

No. Abbreviation used in the UnicodeData.txt file for the “other number” general category.

Nominal digit shapes. A deprecated invisible formatting character used to cancel the effects of the National Digit Shapes character. See Chapter 12.

Non-breaking character. A character that prevents a line break from happening before or after it. It “glues” the characters before and after it together on one line. See Chapter 12.

Non-breaking hyphen. A hyphen with non-breaking semantics.

Non-breaking space. A space with non-breaking semantics.

Nonce form. A character coined for an ad-hoc use in one particular paper or situation. This term is usually applied to rare Chinese characters or mathematical symbols.

Noncharacter. A code point value that is specifically defined never to represent a character (this is different from private-use characters). Noncharacter code point values are set aside for internal application use (as sentinel values, for example) and should never appear in interchanged Unicode text.

Noncharacter_Code_Point. The property given to all noncharacter code points.

Noncognate rule. The rule that states that two similar-looking ideographs that would be unified based on glyph shape alone aren’t unified if they have different meanings and etymological histories. See Chapter 10.

Noninflecting language. A language where the forms of words don’t change depending on their grammatical context. Some Chinese languages are generally noninflecting.

Non-joiner. See **zero-width non-joiner**.

Non-joining letter. An Arabic or Syriac letter that does not connect cursively to the letters on either side.

Non-spacing mark. A combining character that takes up no space along the baseline. Non-spacing marks usually are drawn above or below the based character or overlay it.

Non-starter decomposition. A decomposition that consists entirely of combining marks. Non-starter decompositions are not allowed in Normalized Form C. See Chapter 4.

Normalization. The process of converting a piece of Unicode text to one of the Unicode Normalized Forms; the process of converting a sequence of Unicode code points that has other equivalent sequences to one particular equivalent sequence according to rules that dictate which sequences are preferred. Strings which are different but should be considered equivalent can be compared by normalizing them first (the normalized versions will be identical). See Chapter 4.

Normalization form. Alternate term for **normalized form**.

NormalizationTest.txt. A file in the Unicode Character Database that can be used to determine if an implementation of Unicode normalization actually conforms to the standard. See Chapter 5.

Normalized form. One of four different transformations on Unicode text, each of which converts all equivalent sequences of characters to a single preferred representation. The four Unicode normalized forms do this, however, according to different criteria. See Chapter 4.

Normalized Form C. A Unicode normalized form that prefers canonical composites over combining character sequences when possible. Combining character sequences are converted entirely to canonical composites or, if that isn't possible, to the shortest equivalent combining character sequence. Certain canonical composites, such as those with singleton decompositions, those with non-starter decompositions, and those that aren't generally used in certain scripts, are prohibited from appearing in Normalized Form C. See Chapter 4.

Normalized Form D. A Unicode normalized form that disallows canonical composites. Canonical composites are replaced with their canonical decompositions, and combining character sequences are arranged into canonical order. See Chapter 4.

Normalized Form KC. A Unicode normalized form that disallows compatibility composites, but prefers canonical composites over combining character sequences when possible. See Chapter 4.

Normalized Form KD. A Unicode normalized form that disallows both canonical and compatibility composites. Canonical composites are replaced with their canonical decompositions, compatibility composites are replaced with their compatibility decompositions, and combining character sequences are arranged into canonical order. See Chapter 4.

Normative. Used to describe specifications and descriptions in a standard that must be followed in order to conform to the standard.

Nukta. A dot applied to letters in some Indic scripts, usually to create extra letters to write languages other than the language the script was originally for.

NULL. A control character that is intended to have no effect, represented with the all-bits-off code point value. Some programming languages and systems, such as C, treat NULL as a noncharacter code point signaling the end of a string.

Null consonant. A letter representing the absence of a consonant sound. Usually used as a base to which vowel marks are applied to represent word-initial vowel sounds or the second vowel sounds in diphthongs.

Null vowel. A vowel mark representing the absence of a vowel sound after a particular consonant. A virama can be considered a special kind of null vowel.

Number. 1) A property given to Unicode characters that are generally used to represent numeric values. 2) A numeric value.

Numeral. The written representation of a numeric value, which may consist of one or more characters.

Numeric character reference. A sequence of characters in XML or some other markup language that together represent a single Unicode code point by its code point value. In XML, a numeric character reference takes the form “Ӓ”.

Nuskhuri. An older style of Georgian writing used in ecclesiastical texts. See Chapter 7.

Nyis shad. A Tibetan punctuation mark analogous to the Indic double danda and used to mark a change in topic. See Chapter 9.

Object replacement character. A special Unicode character that is used to mark the position of an embedded object in a stream of Unicode text. This character should normally not occur in Unicode text for interchange—the function is best left to a markup language such as HTML—but is useful as an implementation detail, where it can provide an implementation with a character in the text to attach extra information about the embedded object to. See Chapter 12.

OCR. Optical Character Recognition.

Off-curve point. A point in a loop in an outline font that is used to define the shape of a curve segment. See Chapter 16.

Ogonek. A hook-like mark that is attached to the bottom of some vowels in some languages to indicate a change in pronunciation.

On-curve point. A point in a loop in an outline font that defines the beginning of a curve or line segment and the end of another curve or line segment.

Onomatopoeia. A word whose pronunciation resembles the sound that the word describes: “Crash” is an example of onomatopoeia.

Opening punctuation. Alternate name for **starting punctuation**.

OpenType. An advanced font file format backed by Adobe and Microsoft.

Optical alignment. Adjusting the positions of characters along a margin so that they appear to line up. (If the leftmost pixels of a straight character like an H are aligned with the leftmost pixels of a round character like an O, they actually don’t appear to line up: the O must be shifted slightly to the left.) See Chapter 26.

Orthographic syllable. A group of characters in an Indic script that corresponds roughly to a spoken syllable. The letters in an orthographic syllable don't necessarily match the sounds in a spoken syllable: A consonant sound that is spoken at the end of one syllable is actually written at the beginning of the next orthographic syllable. For example, the "n" in "Hindi" is at the end of the first spoken syllable but at the beginning of the second orthographic syllable. An Indic orthographic syllable is one kind of grapheme cluster. See Chapter 9.

Other letter. The category assigned to the majority of "letter" characters in Unicode, including syllables, ideographs, and uncased letters. See Chapter 5.

Other number. The category assigned to all Unicode characters that represent numeric values but can't be used as digits in positional notation. See Chapter 5.

Other punctuation. The category assigned to all Unicode punctuation that doesn't fit into one of the more specific punctuation categories (generally all punctuation, other than dashes and hyphens, that doesn't come in pairs).

Other symbol. The category assigned to all Unicode symbols (the majority) that don't fit into one of the more specific symbol categories (generally, all symbols that aren't match or currency symbols and aren't used to modify another character).

Other_Alphabetic. A property assigned to the characters with the Alphabetic property that can't be determined algorithmically (certain categories of character, such as cased letters, automatically have the Alphabetic property). See Chapter 5.

Other_Default_Ignorable_Code_Point. A property assigned to the characters with the Default Ignorable Code Point property that can't be determined algorithmically (certain categories of characters, such as formatting characters, automatically have the Default Ignorable Code Point property). See Chapter 5.

Other_Grapheme_Extend. A property assigned to the characters with the Grapheme_Extend property that can't be determined algorithmically (certain categories of characters, such as non-spacing marks, have this property automatically). See Chapter 5.

Other_Lowercase. A property assigned to the characters with the Lowercase property (generally symbols made up of lower-case letters) that aren't in the Lowercase Letter general category. See Chapter 5.

Other_Math. A property assigned to all characters with the Mathematical property that don't belong to the Mathematical Symbols general category. See Chapter 5.

Other_Uppercase. A property assigned to the characters with the Uppercase property (generally symbols made up of upper-case letters) that aren't in the Uppercase Letter general category. See Chapter 5.

Outline font. Any font technology whose glyph descriptions are in the form of abstract descriptions of the ideal character shape rather than bitmapped images of the glyphs. See Chapter 16.

Out-of-band information. Information about a piece of text that is stored in a separate data structure (such as a style run array) from the characters themselves.

Oxia. Greek name for the acute accent. See Chapter 7.

Pair table. A two-dimensional table used in boundary analysis. Each axis is indexed by character or character category, and the cells of the table say whether (or under which conditions) there is a boundary between those two characters (or characters belonging to those two categories). See Chapter 16.

Paired punctuation. Punctuation marks that come in matched pairs and are generally used to mark the beginning and end of some range of text. Parentheses and quotation marks are examples of paired punctuation.

Paiyannoi. A Thai mark used to indicate elision of letters. See Chapter 9.

Palatalization. Modifying a consonant by ending it with the back of the tongue pressed toward the roof of the mouth, or beginning a vowel with the tongue in that position. English speakers tend to think of palatalization as inserting a *y* sound between a consonant and a vowel. The scraping sound at the beginning of “Houston” is a palatalized *h*. See Chapter 7.

Pamudpod. The Hanunóo virama.

Paragraph separator. A special Unicode character, U+2029, that is used to mark the boundary between two paragraphs. In many applications, the underlying system’s new-line function is treated as a paragraph separator.

Pc. The abbreviation used in the UnicodeData.txt file to represent the “connector punctuation” general category.

Pd. The abbreviation used in the UnicodeData.txt file to represent the “dash punctuation” general category.

PDF. 1) Abbreviation for **pop directional formatting**. 2) Portable Document Format.

PDUTR. Abbreviation for **Proposed Draft Unicode Technical Report**.

PDUTR #25. A proposed Unicode technical report describing how Unicode should be used to represent mathematical expressions.

PDUTR #28. The proposed Unicode technical report specifying Unicode 3.2.

Pe. The abbreviation used in the UnicodeData.txt file to represent the “ending punctuation” general category.

Perispomeni. The Greek circumflex accent (which actually may look like a circumflex or tilde depending on type style).

Pf. The abbreviation used in the UnicodeData.txt file to represent the “final-quote punctuation” general category.

Phonetic. A component of a Chinese character that supplies an approximate idea of the word’s pronunciation.

Pi. The abbreviation used in the UnicodeData.txt file to represent the “initial-quote punctuation” general category.

Pictograph. A character that represents the word for an object by depicting the object pictorially.

Pinyin. A system for writing Mandarin Chinese words using Latin letters and digits, designed and promulgated by the Chinese government.

Place-value notation. Alternate term for **positional notation**.

Plain text. Pure text with no accompanying metadata, such as formatting information, hyperlinks, or embedded non-text elements. Plain text is generally considered to be the minimal amount of information necessary to render a piece of text legibly (i.e., the written characters themselves, plus an extremely minimal amount of extra information necessary to render those characters legibly). See Chapter 1.

Plane. A particular range of code point values in a coded character set’s encoding space. Usually if an encoding space is divided into several different kinds of divisions (such as planes, rows, and columns), the plane is the largest such division. In Unicode, a plane is a range of 65,536 contiguous code point values whose first five bits (out of 21) are the same. These first five bits are considered to be the “plane number.”

Plane 0. In Unicode, the range of code point values from U+000 to U+FFFF. This is called the Basic Multilingual Plane, and is the range of code points assigned to the characters of all modern writing systems (except for some relatively rare Chinese characters).

Plane 1. In Unicode, the range of code point values from U+10000 to U+1FFFF. This is called the Supplementary Multilingual Plane, and is the range of code points assigned to various characters from historical (i.e., obsolete) writing systems and to various symbols that have fairly specialized uses.

Plane 2. In Unicode, the range of code point values from U+20000 to U+2FFFF. This is called the Supplementary Ideograph Plane, and is the range of code point values assigned to a wide variety of less-common Chinese characters.

Plane 14. In Unicode, the range of code point values from U+E0000 to U+EFFFE. This is called the Supplementary Special-Purpose Plane, and is the range of code point values assigned to a number of special-purpose characters.

Platen. The rubber roller that paper is wrapped around in a typewriter or teletype machine.

Po. The abbreviation used in the UnicodeData.txt file for the “other punctuation” general category.

Point. Any vowel mark or diacritical mark used in Hebrew or Arabic. See Chapter 8.

Pointed Hebrew. Hebrew written with most or all of the points (which are generally optional and omitted most of the time). See Chapter 8.

Polytonic Greek. The older system of Greek spelling, which includes a complicated system of diacritical marks. See Chapter 7.

Pop directional formatting. A special Unicode character with terminates the effect of the most recent explicit override character or explicit embedding character in the text. See Chapter 8.

Positional notation. The practice of writing a numeric value using a small number of characters whose value is determined by their position in the numeral. In a decimal integer, the rightmost digit is a multiple of 1, the next digit to the left is a multiple of 10, and the digits to the left represent multiples of progressively higher powers of 10. See Chapter 12.

Post-composition version characters. Canonical composites that are in the composition exclusion list because they were added to Unicode after Unicode 3.0, the version of Unicode on which Normalized Form C is based. See Chapter 4.

PostScript. A page description language and a collection of font formats developed by Adobe Corp. The PostScript outline format is one of two outline formats supported by OpenType. See Chapter 16.

Precomposed character. Alternate term for **composite character**.

Precomposed Hangul syllable. One of the Unicode code points representing a whole Hangul syllable. The precomposed Hangul syllables are all canonical composites and have canonical decompositions to sequences of conjoining Hangul jamo. See Chapter 4.

Presentation form. 1) The glyph shape used for a character in a particular situation. 2) A code point representing a particular glyph shape for a character that can have more than one glyph shape. Presentation forms are compatibility characters, and are generally compatibility composites with compatibility decompositions to the code points representing the character irrespective of glyph shape.

Primary key. The main field by which records are sorted. If two records' primary keys are equal, the secondary key determines the order. See Chapter 15.

Primary source standard. One of the main sources of characters for Unicode, and one of the encoding standards for which it was deemed important to maintain round-trip compability.

Primary weight. The weight value that is compared first in a multi-level comparison. Only if two strings' primary weights are all the same are secondary weights compared. In a typical sort order (and in the UCA default sort order), different primary weights are assigned to code points that are considered to represent different characters. For example, “a” and “b” have different primary weights. See Chapter 15.

Primary-level difference. A difference in the primary weights of two collation keys (or the corresponding difference in the original strings). See Chapter 15.

Printing character. A character with a visible glyph. The term is usually used to distinguish code points representing actual written characters from code points representing control functions (“control characters”) or which alter the appearance or treatment of surrounding characters (“formatting characters”).

Private Use Area. A range of code points whose meaning is purposely unstandardized. The code points in the private use area are available for ad-hoc use by users and applications, which are free to assign any meaning they want to them. In Unicode, the Private Use Area runs from U+E000 to U+F8FF. This is supplemented by two Private Use Planes.

Private-use character. Alternate term for **private-use code point**.

Private-use code point. A code point whose meaning is purposely unstandardized. In Unicode, the private use code point values are collected together in the Private Use Area and Private Use Planes.

Private-use high surrogate. A surrogate code unit value that serves as the high-surrogate value for code points in the Private Use Planes. The private-use high surrogates run from U+DB80 to U+DBFF.

Private-use planes. One of two planes of the Unicode encoding space set aside as extensions of the Private Use Area. The two Private Use Planes are Plane 15, which runs from U+F0000 to U+FFFFF, and Plane 16, which runs from U+100000 to U+10FFFF.

PropertyAliases.txt. A file in the Unicode Character Database that gives short and long names for all of the Unicode character properties, for use in things like regular-expression syntax. See Chapter 5.

PropertyValueAliases.txt. A file in the Unicode Character Database that gives short and long names for all the possible values for many of the Unicode character properties (numeric and Boolean values are left out, for obvious reasons), for use in things like regular-expression syntax. See Chapter 5.

PropList.html. A file in the Unicode Character Database that describes the various properties listed in PropList.txt. See Chapter 5.

PropList.txt. A file in the Unicode Character Database that lists a large number of character properties and the characters that have each of them. See Chapter 5.

Proportionally-spaced. Used to describe a font in which the glyphs have different widths, generally corresponding to the relative widths of the characters in handwritten text.

Proposed Draft Unicode Technical Report. A Unicode Technical Report that has been published while still in draft form. Unlike a Draft Unicode Technical Report, where the general content has been agreed upon in principle, everything about a Proposed Draft Technical Report is still in flux. Unicode Technical Reports carry no normative force while still in draft form.

Prosgegrammeni.

Ps. The abbreviation used in the UnicodeData.txt file for the “starting punctuation” general category.

PS. Abbreviation for **Paragraph Separator**.

Psili. The Greek “smooth-breathing mark,” which is used to indicate the absence of an *h* sound at the beginning of a word. See Chapter 7.

PUA. Abbreviation for **Private Use Area**.

Punctuation. The property assigned to all Unicode characters that are used as punctuation marks.

Quotation_Mark. The property assigned to all Unicode characters that are used as quotation marks.

Quoted-printable. A transfer encoding syntax that encodes generic binary data as ASCII text. Byte values corresponding to most ASCII printing characters are represented with those characters; other byte values are represented with a sequence of ASCII characters that contains the spelled-out hexadecimal byte value. This transfer encoding syntax is generally used to represent text in an ASCII-based encoding, such as UTF-8 or ISO 8859-1, in an RFC 822 message. See Chapter 17.

Radical. 1) The component of a Chinese character by which it is classified in dictionaries. Usually the radical is thought to convey an idea of the meaning of the character. 2) The main consonant in a Tibetan syllable. 3) A component of a Yi syllable by which it is classified in dictionaries. 4) The property given to Unicode code points that represent Chinese radicals.

Rafe. A diacritical mark used in Hebrew. See Chapter 8.

ReadMe.txt. A file in the Unicode Character Database that provides important background and overview information on the database. See Chapter 5.

Reahmuk. A Khmer mark that indicates a glottal stop. See Chapter 9.

Rebus. A method of writing words that uses pictures of objects to represent the sounds of those object's names. See Chapter 10.

Regular expression. An expression in some structured syntax (there are many regular-expression languages) that describes a pattern of characters. If a sequence of characters fits the pattern, it is said to “match” the regular expression. Regular-expression matching is used to do complicated searches on text. See Chapter 15.

Rendering. The process of converting a sequence of code unit values in memory to a visible sequence of glyphs on some output device, such as the computer's screen or a printer.

Rendering process. The software responsible for converting a series of code units in memory to a series of visible glyphs on an output device. The rendering process encompasses breaking the text into pages and lines, choosing appropriate glyphs for each character, arranging those glyphs on a line, and drawing the glyph images to the output device. The software usually encompasses both system software and logic in the font files themselves.

Repertoire. See **Abstract character repertoire**.

Repetition mark. A mark that indicates that the preceding character, syllable, or word should be repeated.

Repha. In Devanagari and a few other Indic scripts, the form that the letter *ra* takes when it's the first letter in a conjunct consonant. This is a hook-shaped mark that attaches to the top of the last consonant in the conjunct. See Chapter 9.

Replacement character. A special Unicode character intended to be used to stand in for an otherwise-unrepresentable character. It's usually emitted by character code converters when they encounter an illegal sequence of code units or a character that can't be represented in Unicode.

Representative glyph. The glyph shape shown for a character in the Unicode standard's code charts. As the name suggests, the glyph is chosen to clarify to the reader which character a code point is supposed to represent; it's generally not the only possible glyph shape for the character, nor is it necessarily even the *preferred* glyph shape. Representative glyphs have no normative force.

Reserved. Set aside for future standardization. For example, all unassigned code point values (except for those in the Private Use Area) are reserved—the Unicode Technical Committee is reserving for itself the right to make future character assignments to all these code point values, and they are not free for other uses in the meantime.

Resolved directionality. The directionality a character of neutral directionality takes on by virtue of its position in the text. For example, if a neutral character is preceded and followed by left-to-right characters, its resolved directionality is left-to-right.

RFC. Request for Comments. This is the term used for documents published by the IETF that are intended to become standards at some point. Because documents tend to stay in RFC form for a long time, many become de facto standards well before they become real standards, and those that do become actual IETF standards tend to continue to be referred to by their RFC number even after official publication as standards.

RFC 822. An IETF standard that specifies a standard format for text messages, such as email and Usenet messages. Mail messages on the Internet are almost always in RFC 822 format, although because RFC 822 is an ASCII-based format that requires explicit line breaks between the lines of a message, today most mail messages are sent in a format, such as MIME, that can encode much richer content within the framework of an RFC 822 text message. See Chapter 17.

Rhotacization. Pronouncing a vowel with the tongue curled upward. In American English, this is generally thought of as following the vowel with an *r* sound. The *a* in “car,” as spoken by most Americans, is a rhotacized *a*.

Rich text. Text with accompanying metadata, such as formatting information, language information, or embedded non-text objects.

Right-joining letter. An Arabic or Syriac letter that can connect to a character to its right, but not a character to its left. See Chapter 8.

Right-joining vowel. A dependent vowel in an Indic writing system that attaches to the right-hand side of the consonant it modifies. See Chapter 9.

Right-to-left embedding. An explicit embedding character that starts a new embedding level with left-to-right directionality. See Chapter 8.

Right-to-left mark. An invisible character with strong right-to-left directionality. This character can be used to alter the resolved directionality of neutral-directionality characters. See Chapter 8.

Right-to-left override. An explicit override character that causes all of the following characters to be treated as strong left-to-right characters.

RLE. Abbreviation for **right-to-left embedding**.

RLM. Abbreviation for **right-to-left mark**.

RLO. Abbreviation for **right-to-left override**.

Romaji. The Japanese term for Japanese written with Latin letters.

Romanization. Transliteration to Latin letters, particularly transliteration of languages that use Chinese characters to Latin letters. Pinyin is a system of Romanization.

Rough-breathing mark. Alternate term for **dasia**.

Round-trip compatibility. The ability to perform some transformation, followed by its inverse, and end with the original data in its original form. Two encodings are said to have round-trip compatibility if you can convert from one to the other and back to the first without losing data. Unicode includes a number of characters specifically to make this guarantee with respect to some source encoding.

Ruby. Japanese term for **interlinear annotation**.

Run array. A data structure that associates information with contiguous sequences of records in some other data structure (such as a sequence of characters in a character array).

Run-length encoding. A data compression technique that compresses sequences of repeating values by storing the repeating value or pattern once, accompanied by a repeat count.

SAM. Abbreviation for **Syriac abbreviation mark**.

Sc. The abbreviation used in the UnicodeData.txt file for the “currency symbol” general category.

SC2. Subcommittee #2 of ISO/IEC Joint Technical Committee #1, which is responsible for international coded-character-set standards. WG2, which is responsible for ISO 10646, is a working group of SC2.

Script. Alternate term for **writing system**.

Scripts.txt. A file in the Unicode Character Database that identifies the script to which each character in Unicode belongs. See Chapter 5.

SCSU. Standard Compression Scheme for Unicode. A technique for compressing Unicode text that is backward compatible with ISO 8859-1 and generally allows text encoded with Unicode to be stored in roughly the same amount of space as it would occupy if encoded with an appropriate legacy encoding. See Chapter 6.

Second series consonants. Khmer consonants that carry an inherent *o* sound.

Secondary key. A key which is used to determine the sorted order of two records only if their primary keys are the same.

Secondary weight. The weight value that is compared in the second pass of a multi-level comparison. Differences in secondary weights are only significant in a comparison if the primary weights are all

the same. In most sort orders, differences in secondary weight correspond to different variants of the same letter, such as the letter with different accents applied to it. For example, “a” and “ä” typically have the same primary weight but different secondary weights. Note that this is language-specific: In languages such as Swedish, where “ä” is a completely different letter, “a” and “ä” have different primary weights. See Chapter 15.

Secondary-level difference. A difference in the secondary weights of two collation keys, or the corresponding difference in the strings they were created from.

Separator. A property given to characters that mark divisions between units of text. There are three separator categories: the line and paragraph separators, which are each categories unto themselves, and a third category containing all of the spaces.

Serialization format. Alternate term for **character encoding scheme**.

Series-shifter. The Khmer *muusikatoan* and *tripsis* characters, which change a consonant from one series to the other. See Chapter 9.

Serto. One of the modern writing styles of the Syriac alphabet. See Chapter 8.

Shad. A Tibetan punctuation mark analogous to the Devanagari danda and used to indicate the end of an expression. See Chapter 9.

Shadda. A mark written over an Arabic letter to indicate consonant gemination.

Sheva. A Hebrew point used to indicate the absence of a vowel sound after a letter. See Chapter 8.

Shifted. An alternate-weighting setting in the Unicode Collation Algorithm that causes strings with ignorable characters in analogous positions to sort together. See Chapter 15.

Shift-JIS. A character encoding scheme that represents characters from the JIS X 0201 and JIS X 0208 standards with various one- and two-byte combinations. See Chapter 2.

Shift-trimmed. An alternate-weighting setting in the Unicode Collation Algorithm that causes strings with ignorable characters in analogous positions to sort together. See Chapter 15.

Shin dot. A dot drawn over the leftmost tine of the Hebrew letter *shin* to indicate it has an *sh* sound. See Chapter 8.

SHY. Abbreviation for **soft hyphen**.

Simplified Chinese. A method of writing Chinese that uses simplified (sometimes greatly simplified) versions of many of the Chinese characters. This method is used today in most of the People’s Republic of China and in Singapore.

Sin dot. A dot drawn over the rightmost time of the Hebrew letter *shin* to indicate it has an *s* sound. See Chapter 8.

Singleton decomposition. A canonical decomposition consisting of a single character, or a canonical composite with a one-character decomposition. Singleton decompositions are generally variants of characters that are included elsewhere in Unicode and disunified for round-trip compatibility with some other standard. These variants are effectively discouraged from normal use by prohibiting their appearance in any of the Unicode normalized forms. See Chapter 4.

SIP. Abbreviation for **Supplementary Ideographic Plane**.

Sk. The abbreviation used in the UnicodeData.txt file for the “modifier symbol” general category.

Slur. A musical marking that indicates that two or more notes should be played with no separation between them.

Sm. The abbreviation used in the UnicodeData.txt file for the “mathematical symbol” general category.

Small. When used to talk about a letter in a bicameral script, a synonym for **lower case**.

Smooth-breathing mark. Alternate term for **psili**.

SMP. Abbreviation for **Supplementary Multilingual Plane**.

So. The abbreviation used in the UnicodeData.txt file for the “other symbol” general category.

Soft hyphen. A code point that indicates to a text rendering process a location where a word can be hyphenated. When a soft hyphen appears immediately before a line break, it is rendered as a regular hyphen; otherwise, the soft hyphen is invisible.

Soft sign. A Cyrillic letter used to indicate palatalization in a location where it isn’t otherwise implied by the spelling. See Chapter 7.

Soft_Dotted. A property given to characters whose glyphs include a dot that usually disappears when a diacritical mark is applied. The lowercase *i* and *j* and their variants, whose dots go away when an accent is applied (for example, in “naïve”) have this property. (There are a few languages, such as Lithuanian, where the dot doesn’t disappear—this is usually represented with a specific combining-dot character following the *i*; case mappings then have to take care to add or remove the combining dot when appropriate.) See Chapter 5.

Sort key. Alternate term for **collation key**.

Source separation rule. The rule that dictates that two characters that are disunified in one of Unicode’s primary source standards will also be disunified in Unicode, even if they would otherwise be unified.

Space separator. The general category given to all of the space characters in Unicode, except for a few that explicitly have non-breaking semantics.

SpecialCasing.txt. A file in the Unicode Character Database that lists complex case mappings. For historical reasons, simple one-to-one case mappings are in `UnicodeData.txt`; case mappings that are specific to certain languages, only happen in certain contexts, or result in a single character turning into two or more are listed in `SpecialCasing.txt`.

Split caret. In a text-editing application, an insertion point marking a position in the backing store that corresponds to two positions in the rendered text. Generally, an insertion point is drawn at both positions in the rendered text, often with a tick mark or arrow that indicates where different kinds of text will go if the user enters more text. Split carets generally appear at the boundaries of directional runs. See Chapter 16.

Split vowel. An Indic dependent vowel with components that appear on two different sides (such as both the left- and right-hand sides, or above and to the right) of the consonant it modifies. Often, but not always, the two components of a split vowel are vowel sounds in their own right, used in combination to represent another vowel sound. See Chapter 9.

SSP. Abbreviation for **Supplementary Special-Purpose Plane**.

StandardizedVariants.html. A file in the Unicode Character Database that lists all combinations of a regular Unicode character and a variation selector that are legal and the exactly glyph shape each combination should produce. See Chapter 12.

Starting punctuation. A category of characters that mark the beginning of some range of text (such as a parenthetical expression or a quotation). See Chapter 5.

State machine. A function that remembers things from invocation to invocation. It's possible for a state machine to produce different outputs for the same input, as its behavior is dependent not only on the current input, but on some number of previous inputs as well.

State table. A state machine implemented as an array. The function looks up a new state and maybe an action based on the current state and the input value.

Stateful. When applied to an encoding scheme, indicates that the interpretation of particular code unit values can vary depending on which code unit values have been seen earlier in the text stream. The term typically isn't used to refer to variable-length encoding schemes such as UTF-8, where it takes several code units following a well-defined syntax to represent a single code point, but to schemes such as SCSU, ISO 2022, or the old Baudot code, where one may have to scan backwards an arbitrary distance (maybe all the way to the beginning of the text stream) to know for sure how a particular arbitrary byte in the middle of the stream is to be interpreted.

Stem. The vertical stroke in the middle of most Devanagari consonants that represents the core of the syllable and goes away when the consonant is part of a conjunct.

STIX. Scientific and Technical Information Exchange. A consortium of experts on mathematical typesetting that has derived a standard set of mathematical symbols. See Chapter 12.

Storage format. Alternate term for **character encoding form**.

Strong directionality. Directionality sufficient to terminate an embedding level and to affect the resolved directionality of neutral-directionality characters. A sequence of strong-directionality characters placed into a sequence of characters of the opposite directionality will cause that sequence to be split into separate directional runs at the same embedding level. Compare **weak directionality**.

Style run. A contiguous sequence of characters with the same styling information.

Styled text. Alternate term for **rich text**.

SUB. Substitute. The ASCII control character (0x1A) that is often used to mark the position of an otherwise-unrepresentable characters. The ASCII equivalent of Unicode's REPLACEMENT CHARACTER.

Subjoined consonant. In Tibtan, this refers to a consonant that occurs in a stack somewhere other than at the top. Since there's no virama in Tibetan, consonant stacks are represented as a regular consonant followed by one or more subjoined consonants. See Chapter 9.

Sukun. The Thaana null-vowel character. See Chapter 8.

Supervisory code. Alternate term for **control character** used in FIELDDATA and some other early encoding standards. See Chapter 2.

Supplementary Ideographic Plane. Alternate name for **Plane 2**.

Supplementary Multilingual Plane. Alternate name for **Plane 1**.

Supplementary planes. All of the planes in the Unicode encoding space except for the Basic Multilingual Plane.

Supplementary Special-Purpose Plane. Alternate name for **Plane 14**.

Surrogate. 1) A code unit value in UTF-16 that is half of the UTF-16 representation of a supplementary-plane character. In UTF-16, supplementary-plane characters are represented as a sequence of two code units: a high surrogate followed by a low surrogate. 2) A code point value in the range U+D800 to U+DFFF, corresponding to the surrogate code-unit values. 3) The general category assigned to code point values in the surrogate range. See Chapter 6.

Surrogate mechanism. The practice of representing supplementary-plane characters using pairs of code unit values from the surrogate range; now part of UTF-16. See Chapter 6.

Surrogate pair. A sequence of two UTF-16 code units from the surrogate range—a high surrogate followed by a low surrogate—used together to represent characters in the supplementary planes. See Chapter 6.

Suzhou numerals. A collection of commercial numerals used in parts of China and Japan. See Chapter 10.

SVG. Scalable Vector Graphics. An XML-based format for describing object-based graphics.

Syllabary. A writing system whose characters represent whole syllables. In a true syllabary, the characters representing syllables cannot be further decomposed into components representing individual sounds.

Syllable cluster. Alternate term for **orthographic syllable**.

Symbol. A group of general categories containing symbols and other miscellaneous characters.

Symbols area. The part of the Unicode encoding space running from U+2000 to U+2E7F and containing various punctuation marks and symbols.

Symmetric swapping. Alternate name for mirroring (i.e., the process of using mirror-image glyphs for characters with the “mirrored” property when they appear in right-to-left text). There are two deprecated characters that turn this behavior on and off. See Chapter 12.

Syriac abbreviation mark. A special character used with Syriac to indicate abbreviations. Syriac abbreviations are denoted with a bar drawn above the characters in the abbreviation: the Syriac abbreviation mark has no appearance of its own, but indicates to the rendering process that an abbreviation should be drawn over every character from the next character to the end of the word. See Chapter 8.

T source. Collective name for the various sources of ideographs submitted to the IRG by the Taiwanese national body.

Tag. Some sort of metainformation (such as a language tag) represented in Unicode plan text with tag characters.

Tag characters. A group of special characters that are invisible but can be used to place information about the text into the plain-text stream along with the text itself. Currently only language tagging (using the tag characters to mark certain ranges of text as being in a particular language) is supported by Unicode. Use of the tag characters is strongly discouraged in favor of higher-level protocols such as XML. See Chapter 12.

Tailoring. A set of modifications to the default behavior of a process, usually to adapt its behavior to the requirements of a particular language. The Unicode Collation Algorithm provides a tailoring syntax to let you change the default sort order to a language-specific one, and the Unicode standard also allows for tailoring of the rules for finding grapheme-cluster boundaries, although no formalized syntax exists for this yet.

Tate-chu-yoko. Japanese for “horizontal in vertical.” The practice of setting two or more characters horizontally in a single display cell in vertical Asian text. See Chapter 10.

Tatweel. Alternate term for **kashida**.

TCVN 5712. The Vietnamese national character encoding standard. See Chapter 2.

Technical report. See **Unicode Technical Report**.

Teh marbuta. A special form of the Arabic letter *the* used to write feminine endings. See Chapter 8.

Telegraph. An electrical device for sending and receiving written communication over long distances.

Telegraphy code. An encoding scheme used to represent text in the communication link between two telegraphs.

Teletype. A kind of telegraph where text being received is output as printed characters on a piece of paper and (usually) the text being sent is entered into the machine with a typewriter keyboard.

Terminal_Punctuation. A property given to punctuation marks used at the end of a sentence or clause.

Ternary tree. A linked data structure consisting of individual nodes that are arranged in a hierarchy and where each node points to no more than three other nodes.

Tertiary weight. The weight value that is compared in the third pass of a multi-level comparison. Differences in tertiary weights are only significant in a comparison if both the primary and secondary weights are all the same. In most sort orders, differences in tertiary weight correspond to case differences or to differences between characters whose compatibility decompositions are the same. See Chapter 15.

Tertiary-level difference. A difference in the tertiary weights of two collation keys, or the corresponding difference in the original strings. For example, the strings “HELLO” and “Hello” have a tertiary-level difference at the second character position.

TeX. A computer language for describing typeset text, used especially for mathematical and technical typesetting.

Text rendering process. See **rendering process**.

Thanthakhat. A mark sometimes used in Thai to mark a silent letter. See Chapter 9.

Tie. A marking that connects two musical notes together and indicates that they should be played as a single note with the combined value of the two tied notes.

Tilde. A diacritical mark that looks like a sideways *s* drawn above a letter.

Tippi. The Gurmukhi candrabindu.

TIS 620. The Thai national encoding standard.

Titlecase. The property given to single Unicode code points representing digraphs where the first letter in the digraph is upper case and the second is lower case (for all other characters, “titlecase” is equivalent to “uppercase”).

Tonal language. A spoken language where the pitch of the voice (or the change in pitch of the voice) is a significant component of pronunciation (significant enough that a change in voice pitch or pitch contour alone is enough to turn a word into a different word). Many Asian and African languages are tonal languages.

Tone letter. A spacing character that indicates the tone of the preceding letter in the written form of a tonal language.

Tone mark. A diacritical mark that indicates the vocal pitch contour that a letter or syllable should be spoken with.

Tonos. The accent mark used in monotonic Greek. See Chapter 7.

Top-joining vowel. An Indic dependent vowel that joins to the top of the consonant it modifies.

Tracking. Widening a line of text by adding extra space between all the characters.

Traditional Chinese. Chinese written using the more-complicated traditional forms of the characters. The term is usually used to draw a distinction with Simplified Chinese. Traditional Chinese is used in Taiwan and Hong Kong, among other places. See Chapter 10.

Transfer encoding syntax.

Translation. Converting a particular series of thoughts expressed in one language to the same series of thoughts expressed in another language.

Transliteration. 1) Converting text in a particular language from one writing system to another. The words stay the same; the characters used to write them change. 2) Performing an algorithmic transformation on a sequence of text according to some set of rules. Software that does this can be used for actual transliteration, but also for other types of transforms.

Trema. Alternate term for **dialytika**.

Trie. An n -way tree structure in which a node at any level of the tree represents the set of keys that start with a particular sequence of characters. See Chapter 13.

Triisap. A Khmer mark that converts a first-series consonant to a second-series consonant.

TrueType. A file format for describing an outline font. Both OpenType and AAT extend the TrueType file format with extra tables that allow for more complex typography than a traditional TrueType font, but in incompatible ways. See Chapter 16.

TrueType GX. Former name for **AAT**.

Tsheg. A dot used to mark word divisions in Tibetan. See Chapter 9.

Typeface. A collection of fonts with common design characteristics intended to be used together. (For example, Times is a typeface; Times Roman and Times Italic are fonts.)

Typographer's ellipsis. A single character representing the European ellipsis: a series of three periods in a row. Three normal periods in a row generally are spaced too close together; the typographer's ellipsis includes the optimal amount of spacing between the periods. See Chapter 10.

U source. Collective name given to sources of ideographs used by the IRG without their being submitted by a particular national body.

UAX. Abbreviation for **Unicode Standard Annex**.

UAX #9. The official specification of the Unicode bidirectional layout algorithm. See Chapter 8.

UAX #11. The official specification of the East Asian Width property. See Chapter 10.

UAX #13. A set of guidelines on how line and paragraph divisions should be represented in Unicode. See Chapter 12.

UAX #14. The official specification of the Unicode line breaking properties. See Chapter 16.

UAX #15. The official specification of the Unicode normalization forms. See Chapter 4.

UAX #19. The official definition of UTF-32. See Chapter 6.

UAX #27. The official specification of Unicode 3.1. See Chapter 3.

UCA. Abbreviation for **Unicode Collation Algorithm**.

UCAS. Abbreviation for Unified Canadian Aboriginal Syllabics. See Chapter 11.

UCD. Abbreviation for **Unicode Character Database**.

UCS. Abbreviation for Universal Character Set, an alternate name for ISO 10646.

UCS-2. One of the two character encoding forms specified by ISO 10646. UCS-2 represents characters in the Basic Multilingual Plane as 16-bit quantities. Characters in the supplemental planes can't be represented in UCS-2; the surrogate space is not used and is treated like a range of unassigned code point values. UCS-2 is now generally discouraged in favor of UTF-16. (Applications that can't handle surrogate pairs properly often bill themselves as "supporting UCS-2.")

UCS-4. One of the two character encoding forms specified by ISO 10646. This represents the abstract code point value as a 32-bit quantity, with the extra bits filled with zeros. With the most recent versions of ISO 10646 and Unicode, UCS-4 is equivalent to UTF-32 (in earlier versions, code point values above 0x10FFFF were theoretically allowed in UCS-4 but not in UTF-32, although this is of little practical significance, since nothing was encoded there).

Umlaut. A brightening of a vowel sound in languages such as German, or the mark (a double dot) placed on a vowel to indicate umlaut. (The diaeresis and the umlaut are essentially the same mark used to represent different things; Unicode unifies them and standardizes on the term "diaeresis" to refer to the mark itself.)

Unassigned. Used to describe a code point value that hasn't been assigned to a character. Unassigned code point values in Unicode are reserved for possible future encoding and not available for ad-hoc uses.

Unicode 1.0 name. A property that gives the name a character had in Unicode 1.0, in the cases where it differs from the character's current name. See Chapter 5.

Unicode Character Database. A collection of data files that specify fully the character properties given to each of the Unicode code point values, along with other useful information about the characters. For more information, see Chapter 5.

Unicode Collation Algorithm. A specification for a standard method of performing language-sensitive comparison on Unicode text. The UCA comprises an algorithm for comparing strings, a default order for all the Unicode code points in the absence of tailoring, a syntax for specifying tailorings, and some useful implementation hints. See Chapter 15.

Unicode Consortium. The organization responsible for developing, maintaining, and promoting the Unicode standard. The Unicode Consortium is a membership organization generally made up of corporations, nonprofit institutions, and governments with an interest in character encoding.

Unicode mailing list. A mailing list for general discussion of Unicode. The address is unicode@unicode.org. Membership is open to anyone; to join, go to the Unicode Web site.

Unicode normalization form. One of four methods of standardizing the Unicode representation of text that has multiple representations in Unicode. Normalized Forms D and KD favor decomposed forms, while Normalized Forms C and KC favor composed forms. Normalized Forms C and D normalize only canonically-equal strings to the same representation; Normalized Forms KC and KD extend this to characters whose compatibility decompositions are also the same. See Chapter 4.

Unicode scalar value. The term used in Unicode 2.0 for what is now called a **code point**.

Unicode standard. 1) *The Unicode Standard, Version 3.0* (Reading, MA: Addison-Wesley, 2000). 2) The full Unicode standard, comprising not only the book listed above, but also the complete suite of current Unicode Technical Reports and the Unicode Character Database.

Unicode Standard Annex. A technical report that carries normative force and is considered an addition to the standard itself. Unicode 3.1 and 3.2 are published as Unicode Standard Annexes, and certain sets of character properties added after Unicode 2.0, such as line breaking and East Asian width properties, are also published as Unicode Standard Annexes.

Unicode Technical Committee. The group within the Unicode Consortium that actually makes decisions about changing and extending the standard. This group is made up of experts from the Unicode Consortium's member organizations.

Unicode Technical Report. 1) A document issued by the Unicode Consortium that amends or supplements the standard in some way. The term, used this way, refers collectively to Unicode Standard Annexes, Unicode Technical Standards, and Unicode Technical Reports (in the more restrictive sense). This usage, and the fact that all three types of technical reports are numbered from the same series, stems from the fact that “Unicode Standard Annex” and “Unicode Technical Standard” are relatively new terms—everything used to be called a Unicode Technical Report. 2) Those documents that are still called Unicode Technical Reports are informative supplements to the Unicode standard, providing clarifying information useful implementation information, or techniques for using Unicode in certain situations.

Unicode Technical Standard. A supplementary standard issued by the Unicode Consortium and based on Unicode. Unicode Technical Standards are not part of the Unicode standard itself, but are separate, independent standards with their own conformance rules. The Unicode Collation Algorithm is one example of a Unicode Technical Standard.

Unicode Web site. The Web site maintained by the Unicode Consortium. This Web site serves as the definitive source for all things Unicode, in particular the most recent versions of the Unicode Character Database and Unicode Technical Reports. The URL is <http://www.unicode.org>.

UnicodeCharacterDatabase.html. A file in the Unicode Character Database that gives an overview of all the files in the database. See Chapter 5.

UnicodeData.html. A file in the Unicode Character Database that describes the contents and format of UnicodeData.txt. See Chapter 5.

UnicodeData.txt. A file in the Unicode Character Database that lists all the characters in the current version of Unicode, plus their most important properties: name, general category, decomposition, numeric value, simple case mappings, and some other things. See Chapter 5.

Unicore mailing list. A mailing list for internal discussion by members of the Unicode Consortium. Membership is open only to designated representatives of Unicode Consortium members and certain invited experts.

Unification. 1) Declaring two “characters” from different sources (or with different appearances) to be the same character and assigning them a single code point. 2) Declaring two scripts to be variants of one another and unifying the characters they have in common. See Chapter 3.

Unified Han Repertoire. The complete set of Han characters in Unicode (excluding the characters in the Han Compatibility Ideographs blocks), which is the result of unifying a wide variety of different collections of Han characters from different sources. See Chapter 10.

Unified Repertoire and Ordering. The name of the original set of Han characters in Unicode, developed by the CJK Joint Research Group.

Unified_Ideograph. The property given to the Chinese characters in Unicode. The difference between this property and the Ideographic property is that this property excludes the Han radicals.

Uniform early normalization. See **early normalization**.

Unify. See **Unification**.

Unihan. Nickname for **Unicode Han Repertoire**.

Unihan.txt. A data file in the Unicode Character Database that lists all of the Han characters in Unicode, along with their counterparts in a wide variety of sources and a variety of other information about them. See Chapter 5.

Uniscribe. The advanced text rendering engine in more recent versions of Microsoft Windows.

Update version. The final component of a Unicode version number, after the second decimal point. The update version number is incremented whenever changes are made to the Unicode Character Database without changes being made to other parts of the standard.

Upper case. In a bicameral script, the set of letters used at the beginnings of sentences and proper names and for emphasis.

Uppercase letter. The property given to upper-case letters in bicameral scripts. See Chapter 5.

URI. Uniform Resource Identifier. A string of characters that identifies a resource on the World Wide Web. URLs and URNs are types of URI. See Chapter 17.

URL. Uniform Resource Locator. A URI that identifies a resource on the Web by specifying a location and a path to follow to reach the resource. See Chapter 17.

URN. Uniform Resource Name. A URI that identifies a resource on the Web using a location-independent name.

URO. Abbreviation for **Unified Repertoire and Ordering**.

User character. Alternate term for **grapheme cluster**.

UTC. Abbreviation for **Unicode Technical Committee**.

UTF-16. 1) A Unicode encoding form that represents Unicode code points in the BMP with 16-bit code units and Unicode code points in the supplementary planes with pairs of 16-bit code points. 2) A Unicode encoding scheme based on UTF-16 that uses the byte order mark to indicate the serialization order of bytes in each code unit. See Chapter 6.

UTF-16BE. A Unicode encoding scheme based on UTF-16 that serializes the bytes in each code unit in big-endian order. See Chapter 6.

UTF-16LE. A Unicode encoding scheme based on UTF-16 that serializes the bytes in each code unit in little-endian order. See Chapter 6.

UTF-32. 1) A Unicode encoding form that represents Unicode code points as 32-bit values. 2) A Unicode encoding scheme based on UTF-32 that uses the byte order mark to indicate the serialization order of bytes in each code unit. See Chapter 6.

UTF-32BE. A Unicode encoding scheme based on UTF-32 that serializes the bytes in each code unit in big-endian order. See Chapter 6.

UTF-32LE. A Unicode encoding scheme based on UTF-32 that serializes the bytes in each code unit in little-endian order. See Chapter 6.

UTF-7. A transfer encoding syntax that makes it possible to use Unicode text in environments, such as RFC 822 messages, that can only handle 7-bit values. Rarely used anymore. See Chapter 6.

UTF-8. A Unicode encoding form that represents Unicode code points as sequences of eight-bit code units. 7-bit ASCII characters are represented using single bytes, characters from the rest of the BMP are represented using sequences of two or three bytes, and characters from the supplementary planes are represented using sequences of four bytes. See Chapter 6.

UTF-8-EBCDIC. Alternate name for **UTF-EBCDIC**.

UTF-EBCDIC. A UTF-8-like encoding form that maintains backward compatibility with EBCDIC rather than with ASCII. See Chapter 6.

UTR. Abbreviation for **Unicode Technical Report**.

UTR #16. The description of UTF-EBCDIC. See Chapter 6.

UTR #17. A set of definitions of terms for discussion character encodings. See Chapter 2.

UTR #18. A set of guidelines for supporting Unicode in a regular-expression engine. See Chapter 15.

UTR #20. A set of guidelines for using Unicode in markup languages such as XML. See Chapter 17.

UTR #21. A more detailed discussion of case mapping in Unicode. See Chapter 14.

UTR #22. A specification for an XML-based file format for describing mappings between Unicode and other encoding standards. See Chapter 14.

UTR #24. Defines the script-name property. See Chapter 5.

UTS. Abbreviation for **Unicode Technical Standard**.

UTS #6. The official definition of the Standard Compression Scheme for Unicode.

UTS #10. The official definition of the Unicode Collation Algorithm.

V source. Collective term for sources of Han ideographs submitted to the IRG by the Vietnamese national body.

Varia. The Greek name for the grave-accent character. See Chapter 7.

Variation selector. A special formatting character that has no appearance of its own, but tells the rendering engine to use a particular alternate glyph for the preceding character (if possible). See Chapter 3.

Vertical Extension A. Alternate name for the CJK Unified Ideographs Extension A block in Unicode. (The “vertical,” by the way, refers not to the fact that Han is often written vertically, but to the fact that the extension adds characters. A “horizontal” extension merely adds new mappings from Unicode to a source standard.)

Vertical Extension B. Alternate name for the CJK Unicode Ideographs Extension B area in Unicode.

Virama. 1) A mark that cancels the inherent vowel sound of an Indic consonant, causing it to be counted as part of the same orthographic syllable as the consonant that follows it. 2) In Unicode, a code point representing a virama may not always result in a visible virama in the rendered output; the code point is also used to signal the formation of a conjunct consonant. (Conceptually, this is the same thing: the canceling of the inherent vowel on an Indic consonant; all that's different is how this concept gets drawn.) See Chapter 9.

Virtual font. A font that doesn't include any glyph images of its own but merely refers to other fonts for the glyph images.

Visarga. A mark in various Indic scripts that indicates an *h* sound. See Chapter 9.

VISCII. Vietnamese Standard Code for Information Interchange. An informal character encoding standard for Vietnamese.

Visual order. The order in which the characters are displayed on the screen.

Visual selection. A selection that encompasses a group of characters that are continuous in the rendered text, but may be split into two or more separate ranges in the backing store.

Vowel bearer. One of three marks used in Gurmukhi with the dependent vowels to form independent vowels. See Chapter 9.

Vowel mark. A mark applied to a consonant to indicate the vowel sound that follows it.

Vowel point. Alternate term for **vowel mark**.

Vowel sign. Alternate term for **dependent vowel**.

VT. Vertical Tab. An ASCII control character sometimes used as a line-separator character in applications that use the system new-line function as a paragraph separator.

Vulgar fraction. A single character representing an entire fraction: the numerator, fractional bar, and denominator. See Chapter 12.

W3C. Abbreviation for **World Wide Web Consortium**.

W3C character model. A W3C recommendation that specifies how text is to be represented in all W3C data file formats. It mandates early uniform normalization, specifies UTF-16 or UTF-8 as the default character encoding, and specifies that numeric character references also have to be in normalized form. See Chapter 17.

W3C normalization. The normalized form that Unicode should take on the World Wide Web. This is essentially Normalized Form C, but with the extra provision that numeric character references, when

expanded, can't break normalization. Earlier versions of the W3C character model also included prohibitions on characters that made "W3C normalization" more restrictive than Normalized Form C, but these have been moved out into a separate document and are no longer a requirement.

wchar_t. A C and C++ data type representing a "wide character." Like `char`, which merely specifies an integral type large enough to hold a character and no larger than any of the other integral data types, `wchar_t` doesn't impose any semantics on the contents—it merely defines an integral type large enough to hold a "wide character." `wchar_t` must be at least as large as `char`, but isn't required to be any larger.

Weak directionality. Directionality that isn't strong enough to start a new embedding level. A sequence of characters with weak directionality is drawn with the correct directionality, but doesn't break up the enclosing directional run. See Chapter 8.

Weight. An arbitrary value assigned to a character that controls how that character sorts. In other words, if B should sort after A, B is assigned a higher weight value. See Chapter 16.

WG2. Working Group #2 of Subcommittee #2 of ISO/IEC Joint Technical Committee #1; the group responsible for maintaining the ISO 10646 standard.

White_space. A property assigned to all characters that should be treated as "white space" by programming-language parsers and similar applications. This generally consists of all space characters, the TAB character, the line and paragraph separators, and all of the control characters that are used as line and paragraph separators on various systems.

Wide. A property assigned to all characters that are implicitly fullwidth, and that neutral and ambiguous characters resolve to when used in vertical text.

WJ. Abbreviation for **word joiner**.

Word joiner. A special character that prevents word breaks from being placed between it and the preceding and following characters: It essentially "glues" those two characters together on one line. See Chapter 12.

Word wrapping. Breaking text up into lines in such a way as to make sure that words aren't divided across two lines (or are hyphenated). The term is generally functionally equivalent to "line breaking" now.

World Wide Web Consortium. An industry body that issues standards (called "recommendations") for various operations and data file formats on the World Wide Web.

Writing system. A collection of characters and rules for writing them that is used to write one or more written languages. The characters and rules may vary slightly between languages (for example, English and German have different letters in their alphabets, but can still both be considered to be written with the Latin alphabet), and the same language may be written with more than one writing system (Kanji and Hiragana for Japanese, or the Latin and Cyrillic alphabets for Serbo-Croatian).

WRU. “Who are you?” A control function used in early teletype machines as part of a handshaking protocol. In response to this signal, the receiving equipment would send back a short identifying string. See Chapter 2.

www.unicode.org. The URL of the **Unicode Web site**.

WYSIWYG. What You See Is What You Get. A term applied to applications where the document being operated on appears on the screen exactly the same way (or close enough) to the way it will look when printed.

X3. ANSI Committee X3, the committee responsible for information technology standards at the time ASCII was published. X3 has since morphed into NCITS.

X3.4. See **ANSI X3.4-1967**.

XCCS. Xerox Coded Character Set. One of the important predecessors of Unicode.

x-height. The height of the lowercase letter x in a Latin font (and, by extension, the height of the other lowercase letters that don’t have ascenders and descenders, and the height of the parts of the other lowercase letters other than the ascenders and descenders).

XHTML. Extensible Hypertext Markup Language. An XML-based page description language with functionality equivalent to HTML.

XML. Extensible Markup Language. A W3C standard for describing structured data in a text file. Forms the basis of many W3C data-file standards. See Chapter 17.

XSL. Extensible Stylesheet Language. An XML-based language for describing the look of a page of text.

Xu Shen. The compiler of the first known Chinese dictionary in AD 100.

Ypogegrameni. The iota subscript used with some vowels in ancient Greek to write some diphthongs. See Chapter 7.

Zenkaku. Japanese for **fullwidth**.

Zero width space. A special invisible character that indicates to a line-breaking algorithm a location where it is okay to put a line break. Generally used in situations (such as Thai) where this can't be determined from looking at the real characters. See Chapter 12.

Zero-width joiner. A special formatting character that has no appearance of its own, but indicates to a text rendering process that the characters on either side should be joined together in some way, if possible. This may involve cursive joining or forming of a ligature (if both are possible, it requests that both happen). See Chapter 12.

Zero-width non-breaking space. The character that, in versions of Unicode prior to Unicode 3.2, served the same purpose as the **word joiner**. The same code point also serves as the **byte order mark**.

Zero-width non-joiner. A special formatting character that has no appearance of its own, but indicates to a text rendering process that the character on either side should *not* be joined together. It is used to break a cursive connection that would otherwise happen, or to prevent the formation of a ligature. See Chapter 12.

Zl. The abbreviation used in the UnicodeData.txt file for the “line separator” general category, which consists solely of the line-separator character.

Zone row. A row of a punched card that is used to change the meanings of punches in other rows. See Chapter 2.

Zp. The abbreviation used in the UnicodeData.txt file for the “paragraph separator” general category, which consists solely of the paragraph-separator character.

Zs. The abbreviation used in the UnicodeData.txt file for the “space separator” general category.

ZWJ. Abbreviation for **zero-width joiner**.

ZWNBSP. Abbreviation for **zero-width non-breaking space**.

ZWNJ. Abbreviation for **zero-width non-joiner**.

ZWSP. Abbreviation for **zero-width space**.

ZWWJ. Zero-width word joiner. An old term for **word joiner**.

Bibliography

The Unicode Standard

The Unicode Consortium, *The Unicode Standard, Version 2.0*, Reading, MA: Addison-Wesley, 1996.

The Unicode Consortium, *The Unicode Standard, Version 3.0*, Reading, MA: Addison-Wesley, 2000.

Davis, Mark, “Unicode Standard Annex #9: The Bidirectional Algorithm,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr9/>.

Freytag, Asmus, “Unicode Standard Annex #11: East Asian Width,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr11/>.

Davis, Mark, “Unicode Standard Annex #13: Unicode Newline Guidelines,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr13/>.

Freytag, Asmus, “Unicode Standard Annex #14: Line Breaking Properties,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr14/>.

Davis, Mark and Dürst, Martin, “Unicode Standard Annex #15: Unicode Normalization Forms,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr15/>.

Bibliography

Davis, Mark, “Unicode Standard Annex #19: UTF-32,” version 3.1.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr19/>.

Davis, Mark; Everson, Michael; Freytag, Asmus; Jenkins, John H.; *et. al.*, “Unicode Standard Annex #27: Unicode 3.1,” version 3.1.0, May 16, 2001, <http://www.unicode.org/unicode/reports/tr27/>.

Wolf, Misha; Whistler, Ken; Wicksteed, Charles; Davis, Mark; and Freytag, Asmus, “Unicode Technical Standard #6: A Standard Compression Scheme for Unicode,” version 3.2, August 31, 2000, <http://www.unicode.org/unicode/reports/tr6/>.

Davis, Mark, and Whistler, Ken, “Unicode Technical Standard #10: Unicode Collation Algorithm,” version 8.0, March 23, 2001, <http://www.unicode.org/unicode/reports/tr10/>.

Umamaheswaran, V. S., “Unicode Technical Report #16: UTF-EBCDIC,” version 7.2, April 29, 2001, <http://www.unicode.org/unicode/reports/tr16/>.

Whistler, Ken and Davis, Mark, “Unicode Technical Report #17: Character Encoding Model,” version 3.2, August 31, 2000, <http://www.unicode.org/unicode/reports/tr17/>.

Davis, Mark, “Unicode Technical Report #18: Unicode Regular Expression Guidelines,” version 5.1, August 31, 2000, <http://www.unicode.org/unicode/reports/tr18/>.

Dürst, Martin and Freytag, Asmus, “Unicode Technical Report #20: Unicode in XML and Other Markup Languages,” version 5, December 15, 2000, <http://www.unicode.org/unicode/reports/tr20/>.

Davis, Mark, “Unicode Technical Report #21: Case Mappings,” version 4.3, February 23, 2001, <http://www.unicode.org/unicode/reports/tr21/>.

Davis, Mark, “Unicode Technical Report #22: Character Mapping Markup Language,” version 2.2., December 1, 2000, <http://www.unicode.org/unicode/reports/tr22/>.

Davis, Mark, “Unicode Technical Report #24: Script Names,” version 3, September 27, 2001, <http://www.unicode.org/unicode/reports/tr24/>.

Phipps, Toby, “Draft Unicode Technical Report #26: Copatibility Encoding Scheme for UTF-16: 8-Bit (CESU-8),” version 2.0, December 12, 2001, <http://www.unicode.org/unicode/reports/tr26/>.

Bibliography

Beeton, Barbara; Freytag, Asmus; and Sargent III, Murray, “Proposed Draft Unicode Technical Report #25: Unicode Support for Mathematics,” version 1.0, January 3, 2002, <http://www.unicode.org/unicode/reports/tr25/>.

The Unicode Consortium, “Proposed Draft Unicode Technical Report #28: Unicode 3.2,” version 3.2.0, January 21, 2002, <http://www.unicode.org/unicode/reports/tr28/>.

The Unicode Character Database, <http://www.unicode.org/Public/UNIDATA/>.

Other Standards Documents

ECMA-35, “Character Code Structure and Extension Techniques,” 6th edition, December 1994

ECMA-94, “8-Bit Single Byte Coded Graphic Character Sets - Latin Alphabets No. 1 to No. 4,” 2nd edition, June 1986.

Linux Internationalization Initiative, “LI18NUNIX 2000 Globalization Specification,” version 1.0, amendment 2. <http://www.li18nux.net/docs/html/LI18NUNIX-2000.htm>.

IETF RFC 1034, “Domain names: Concepts and Facilities.”

IETF RFC 2825, “A Tangled Web: Issues of I18N, Domain Names, and the Other Internet protocols.”

W3C Working Draft, “Character Model for the World Wide Web,” December 20, 2001, <http://www.w3.org/TR/charmod/>.

Books and Magazine Articles

Boyer, Robert S. and Moore, J. Strother, “A Fast String Searching Algorithm,” *Communications of the Association for Computing Machinery*, Vol. 20, No. 10, pp. 762-772, 1977.

Daniels, Peter T. and Bright, William, eds., *The World’s Writing Systems*, Oxford: Oxford University Press, 1995

Ibrah, Georges (trans. David Bellos, E. F. Harding, Sophie Wood, and Ian Monk), *The Universal History of Numbers*, New York: Wiley, 2000.

Jenkins, John H., “Adding Historical Scripts to Unicode,” *Multilingual Computing & Technology*, September 2000.

Bibliography

Kernighan, Brian W. and Ritchie, Dennis M., *The C Programming Language*, first edition
Englewood Cliffs, NJ: Prentice-Hall, 1978.

Lunde, Ken, *CJKV Information Processing*, Cambridge, MA: O'Reilly, 1999.

Nakanishi, Akira, *Writing Systems of the World*, Tokyo: Charles E. Tuttle Co., 1980.

Unicode Conference papers

Atkin, Steven and Stansifer, Ryan, "Implementations of Bidirectional Reordering Algorithms,"
Proceedings of the Eighteenth International Unicode Conference, session C13, April 27, 2001.

Davis, Mark, "Bits of Unicode," *Proceedings of the Eighteenth International Unicode Conference*,
session B11, April 27, 2001.

Davis, Mark, "Collation in ICU 1.8," *Proceedings of the Eighteenth International Unicode
Conference*, session B10, April 27, 2001.

Davis, Mark and Liu, Alan, "Transliteration in ICU," *Proceedings of the Eighteenth International
Unicode Conference*, session B14, April 27, 2001.

Edberg, Peter K., "Survey of Character Encodings," *Proceedings of the Thirteenth International
Unicode Conference*, session TA4, September 9, 1998

Gillam, Richard, "Text Boundary Analysis in Java," *Proceedings of the Fourteenth International
Unicode Conference*, session B2, March 24, 1999.

Jenkins, John H., "New Ideographs in Unicode 3.0 and Beyond," *Proceedings of the Fifteenth
International Unicode Conference*, session C15, September 1, 1999.

Ksar, Mike, "Unicode and ISO 10646: Achievements and Directions," *Proceedings of the Thirteenth
International Unicode Conference*, session B11, September 11, 1998.

Milo, Thomas, "Creating Solutions for Arabic: A Case Study," *Proceedings of the Seventeenth
International Unicode Conference*, session TB3, September 6, 2000.

Werner, Laura, "Unicode Text Searching in Java," *Proceedings of the Fifteenth International
Unicode Conference*, session B1, September 1, 1999.

Other papers

Emerson, Thomas, "On the 'Hangzhou-Style' Numerals in ISO/IEC 10646-1," a white paper
published by Basis Technology.

Halpern, Jack, "The Pitfalls and Complexities of Chinese to Chinese Conversion," an undated white paper published by the CJK Dictionary Publishing Society.

Online resources

Adobe Systems, "OpenType Specification v1.3,"
<http://partners.adobe.com/asn/developer/opentype/>.

Apple Computer, "Developer Documentation," [http:// developer.apple.com/techpubs](http://developer.apple.com/techpubs).

Apple Computer, "TrueType Reference Manual,"
<http://developer.apple.com/fonts/TTRefMan/index.html>.

Czyborra, Roman, "Good Old ASCII,"
<http://www.czyborra.com/charsets/iso646.html>.

Czyborra, Roman, "ISO 8859 Alphabet Soup",
<http://www.czyborra.com/charsets/iso8859.html>

Davis, Mark and Scherer, Markus, "Binary-Ordered Compression for Unicode," IBM DeveloperWorks, <http://www-106.ibm.com/developerworks/unicode/library/ubinary.html>.

Dürst, Martin, "URIs and Other Identifiers," [http://www.w3.org/ International/O-URL-and-ident](http://www.w3.org/International/O-URL-and-ident)

Jennings, Tom, "Annotated History of Character Codes,"
<http://www.wps.com/texts/codes>.

Jones, Douglas W., "Punched Cards: An Illustrated Technical History" and "Doug Jones' Punched Card Codes," <http://www.cs.uiowa.edu/~jones/cards>.

McGowan, Rick, "About Unicode Consortium Procedures, Policies, Stability, and Public Access," published as an IETF Internet-Draft, <http://search.ietf.org/internet-drafts/draft-rmcgowan-unicode-procs-00.txt>.

Microsoft, "Global Software Development," <http://www.microsoft.com/globaldev>.

Opstad, Dave, "Comparing GX Line Layout and OpenType Layout,"
<http://fonts.apple.com/WhitePapers/GXvsOTLayout.html>.

Perldoc.com, "Perl 5.6 Documentation: perlunicode,"
<http://www.perldoc.com/perl5.6/pod/perlunicode.html>.

Bibliography

Searle, Steven J., "A Brief History of Character Codes," found at
<http://www.tronweb.super-nova-co-jp/characcodehist.html>.

SIL International, "Ethnologue," <http://www.ethnologue.com>.