

LGN5809 - Genética Molecular

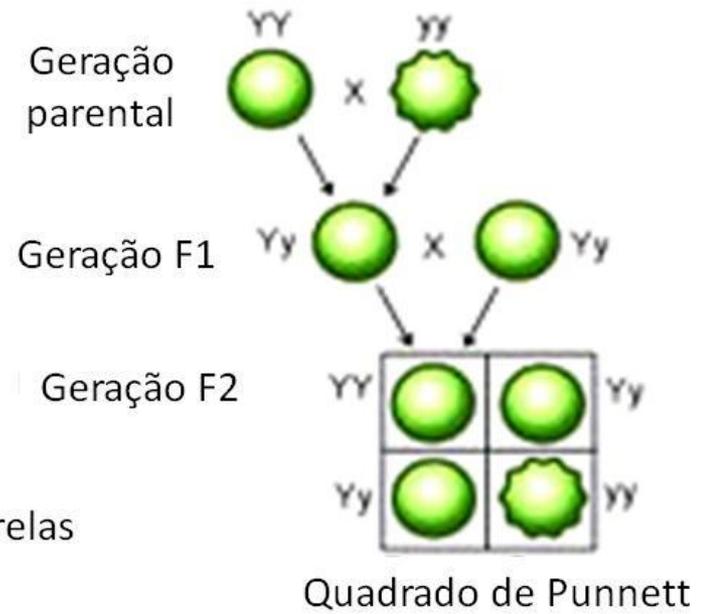
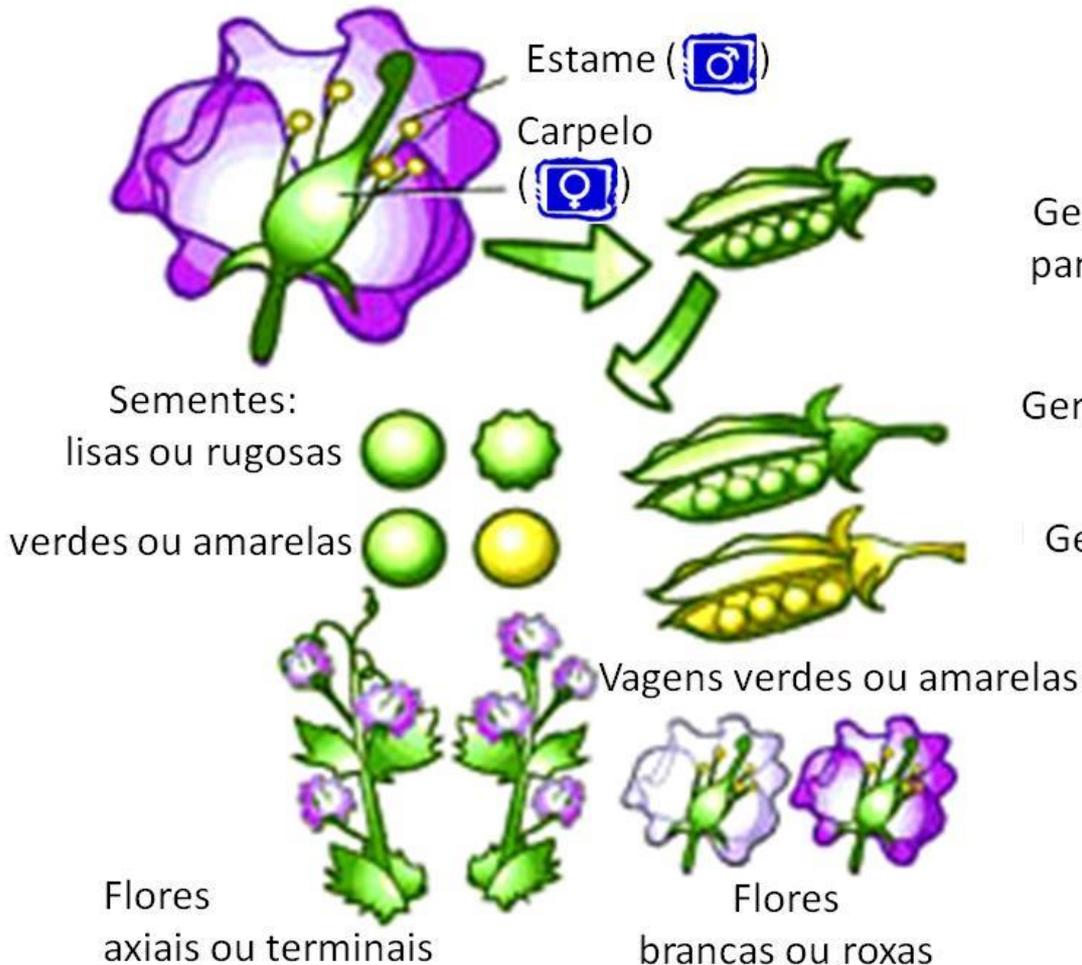
ARQUITETURA GENÔMICA: DOS GENES AO GENOMA

Maria Carolina Quecine
Departamento de Genética
mquecine@usp.br

SUMÁRIO

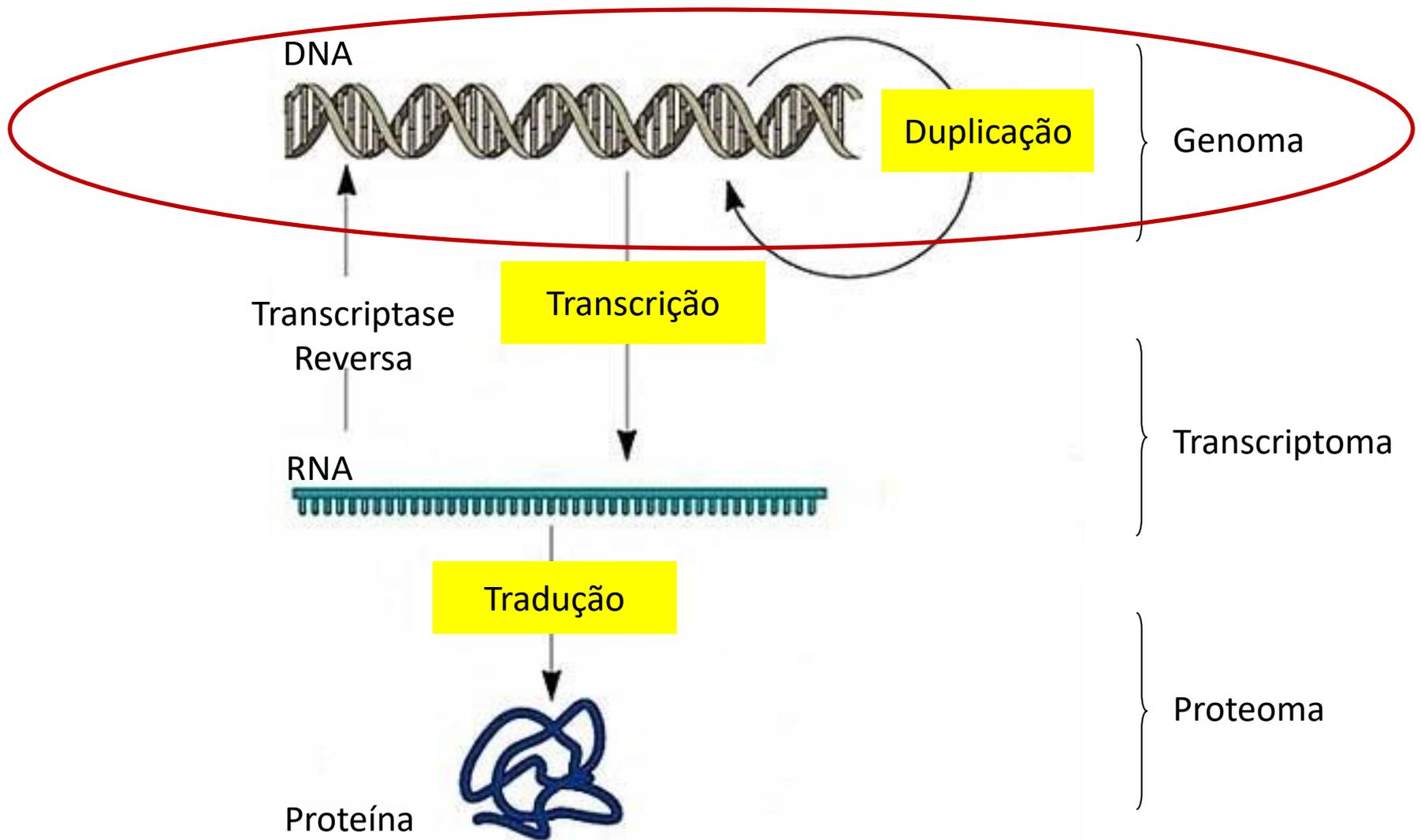
- Fluxo da informação genética;
- Conceito de genoma;
- Anotação de genomas;
- Conceito de genes;
- DNA lixo;
- Genoma de plantas;
- Desafios da genômica.
- Próxima aula

MENDEL: **FATORES CONSTANTES** QUE CONTROLAM CARACTERÍSTICAS FENOTÍPICAS



Leis de Mendel (1866)

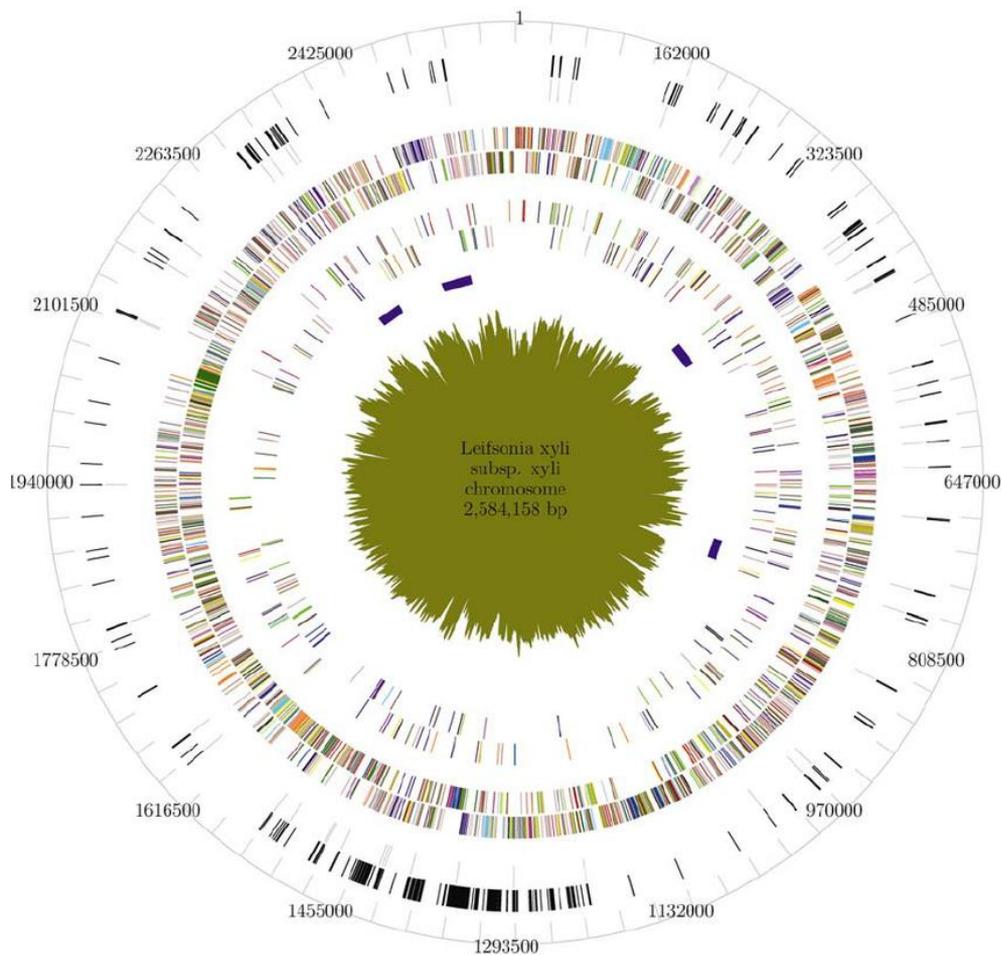
DOGMA CENTRAL DA BIOLOGIA



o **genoma** é toda a informação hereditária de um organismo que está codificada em seu **DNA**. Isto inclui tanto os **genes como as sequências não-codificadoras**.

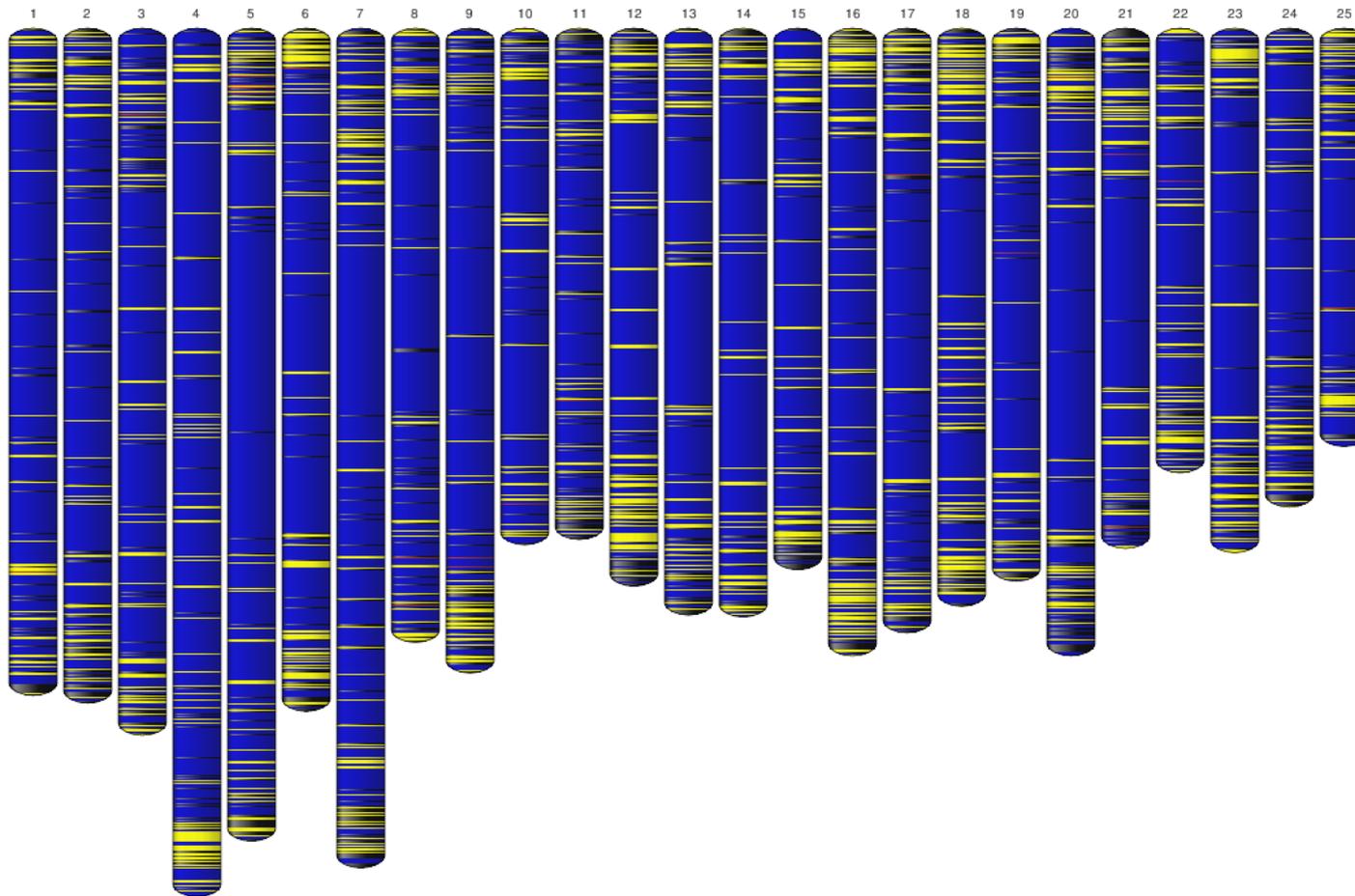
????





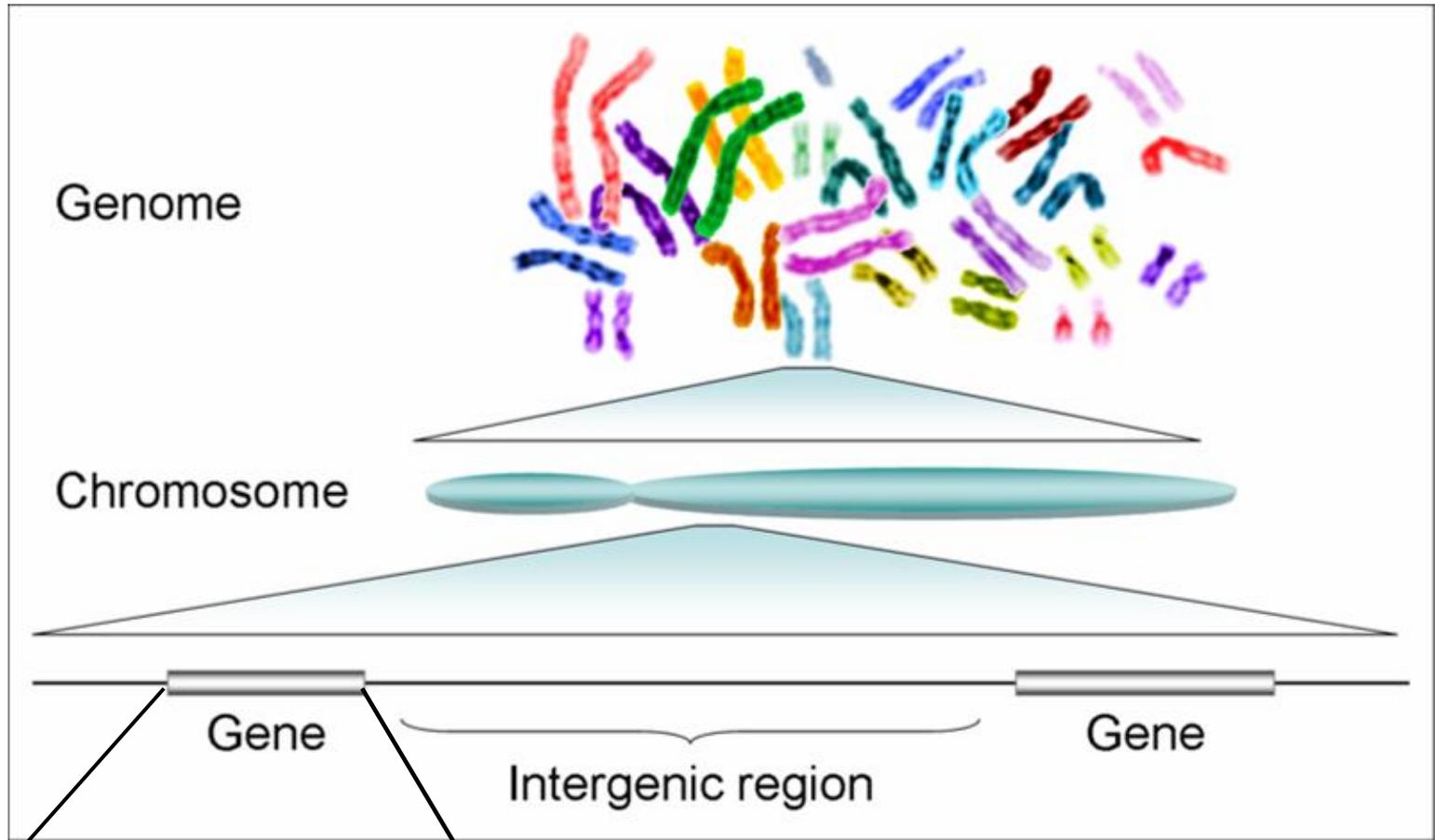
- | | |
|-----------------------------------|--|
| ■ Intermediary metabolism | ■ Cell structure |
| ■ Energy metabolism, carbon | ■ Cellular processes |
| ■ Regulatory functions | ■ Transport |
| ■ Biosynthesis of small molecules | ■ Mobile genetic elements |
| ■ Amino acids biosynthesis | ■ Pathogenicity, virulence, and adaptation |
| ■ Nucleotides biosynthesis | ■ Conserved hypothetical proteins |
| ■ Macromolecule metabolism | ■ Hypothetical ORFs |
| ■ DNA metabolism | ■ ORFs with undefined category |
| ■ RNA metabolism | |

Monteiro-Vitorello et al. 2004



- Finished sequence
- Whole Genome Shotgun sequence
- Scaffold gaps

<http://www.sanger.ac.uk/science/data/zebrafish-genome-project>



attgcgcgataccccggctaa

TTCATACTTGGTTAAGACCTTTACAAGCCGACCAACGTGGTGACAGTGTCGTCCTTTA
CGCACCGAATCCCTTTATCATTGAATTAGTAGAAGAGCGATACTTAGGACGTCTTCGG
ATGGAATCTTGGTCCCGTTGCCTGGAACGTCTTGAAACTGAATTCCCGCCAGAAGATG
TTCATACTTGGTTAAGACCTTTACAAGCCGACCAACGTGGTGACAGTGTCGTCCTTTA
CGCACCGAATCCCTTTATCATAATGAATTAGTAGAAGAGCGATACTTAGGACGTCTTC
GGGAATTGTTATCCTATTTCTCAGGAATACGTGAAGTAGTCCTTGCAATTGGCTCACG
ACCTAAAACAACAGAACTACCCGTACCAGTAGACACTACAGGACGTTTGTCTTCAACA
GTCCCATTTAACGGAAATCTCGACACACACTATAACTTTGATAATTTTGTGAGGGAC
GAAGCAATCAACTCGCTCGTGCTGCAGCTTGGCAAGCGGCACAGAAACCGGGAGACCG
TACTCACAACCCTCTATTGCTCTATGGTGGGACTGGTTTGGGTAAAACCCATTTAATG
TTTGCTGCAGGTAACGTAATGCGGGCAAGTAAACCCAACTTATAAAGTAATGTATCTTC
GTTTCGGAACAGTTTTTTCAGCGCCATGATAAGAGCGTACAAGATAAAAAGTATGGATCAT
AAGGGTAAAACCCATTTAATGTTTGTCTGCAGGTAACGTAATGCGGGCAAGTAAACCCAA
CTTATAAAGTAATGTATCTTCGTTTCGGAACAGTTTTTTCAGCGCCATGATAAGAGCGTA
CAAGATAAAAAGTATGGATCATAAGGGTAAAACCCATTTAATGTTTGTCTGCAGGTAACG
TAATGCGGGCAAGTAAACCCAACTTATAAAGTAATGTATCTTCGTTTCGGAACAGTTTTT
CAGCGCCATGATAAGAGCGTACAAGATAAAAAGTATGGATCATAAAAACGTAATGCGGGCA
AGTAAACCCAACTTATAAAGTAATGTATCTTCGTTTCGGAACAGGGTAAAACCCATTTA
ATGTTTGTCTGCAGGTAACGTAATGCGGGCAAGTAAACCCAACTTATAAAGTAATGTATC
TTCGTTTCGGAACAGTTTTTTCAGCGCCATGATAAGAGCGTACAAGATAAAAAGTATGGAT
CATAAAAACGTAATGCGGGCAAGTAAACCCAACTTATAAAGTAATGTATCTTCGTTTCGGA
ACAAAAACGTAATGCGGGCAAGTAAACCCAACTTATAAAGTAATGTATCTTCGTTTCGGA

GENOME ANNOTATION: FROM SEQUENCE TO BIOLOGY

Lincoln Stein

**O que buscamos são
sequências com alguma
função biológica!!!**

**No mínimo - entre 5 - 15% dos genes são
negligenciados!!**

Doi: 10.1038/350805296

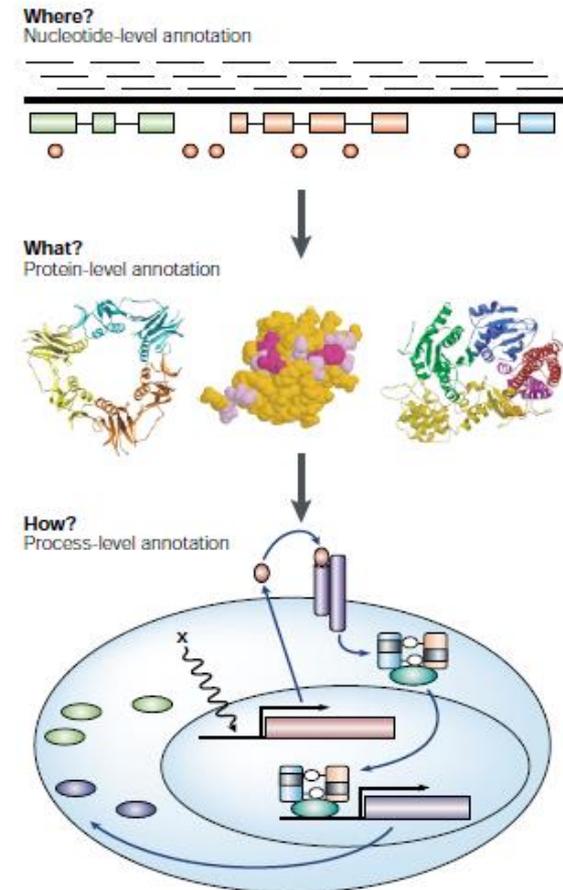


Figure 1 | The three layers of genome annotation: where, what and how?

Haemophilus influenza - 85% dos 1.8-Mb são regiões codantes.
Leveduras – menos que 70%.

Aves e vermes – menos que 25%.

Humanos – menos que 1%

Vários programas de predição

Box 1 | Resources and tools |

The following list provides brief descriptions of some of the software tools and resources mentioned in the article. Online, links are provided to these resources from this box, and from the text.

BLASTN, BLASTX,
BLASTP, PSI-BLAST <http://www.ncbi.nlm.nih.gov/BLAST/>
This family of sequence-similarity search tools allows you to rapidly search a query protein or nucleotide sequence against a large database of sequences, to identify sequences that are similar to the search sequence.

Ensembl <http://www.ensembl.org>
This web site, a joint project of the European Bioinformatics Institute and the Sanger Centre, seeks to make available a high-quality, consistent set of annotations on the human and mouse genomes.

e-PCR <http://www.ncbi.nlm.nih.gov/genome/sts/ePCR.cgi>
BLAST does not work well with short sequences, such as PCR primers. The ePCR program was developed to search rapidly for a primer pair (using their sequences and known physical separation) in a large sequence, such as a genome.

FlyBase <http://www.flybase.org/>
A fully realized model organism system database, presenting curated annotations on the *Drosophila melanogaster* genome, as well as rich information on the genetics of the organism, mutant strains and molecular resources.

GeneMark.hmm <http://genemark.biology.gatech.edu/GeneMark/>
Another gene-prediction program that uses hidden Markov models (HMMs). The online version supports numerous eukaryotic and prokaryotic genomes.

Genie http://www.fruitfly.org/seq_tools/genie.html
The gene-prediction program used to annotate genes in *Drosophila melanogaster*. The online version has been trained for human and *Drosophila* sequence.

GENSCAN <http://genes.mit.edu/GENSCAN.html>
This is probably the most widely used gene-prediction program. It uses HMMs to predict the presence of a gene given the raw DNA sequence. The online version provides prediction services for vertebrates, *Arabidopsis thaliana* and maize.

Grail <http://compbio.ornl.gov/Grail-1.3/>
One of the oldest gene-prediction programs still in use, this software uses a neural network to predict genes. The online version provides gene-prediction services for human, mouse, *Arabidopsis*, *Drosophila* and *Escherichia coli* sequences.



Auxilio com cDNA
Mas tem problemas:
Perdas de transcritos (raros)
Contaminação

Doi: 10.1038/350805296

Review

Review on the Computational Genome Annotation of Sequences Obtained by Next-Generation Sequencing

Girum Fitihamlak Ejigu and Jaehee Jung *

Department of Information and Communication Engineering, Myongji University,
Yongin-si 17058, Gyeonggi-do, Korea; girumfitex@gmail.com

* Correspondence: jhjung@mju.ac.kr

Received: 21 August 2020; Accepted: 16 September 2020; Published: 18 September 2020

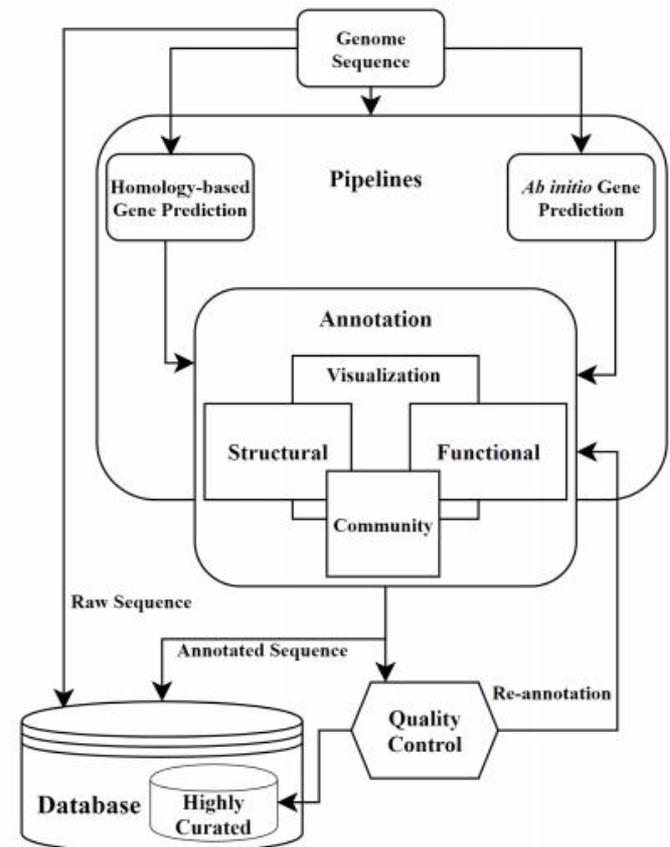
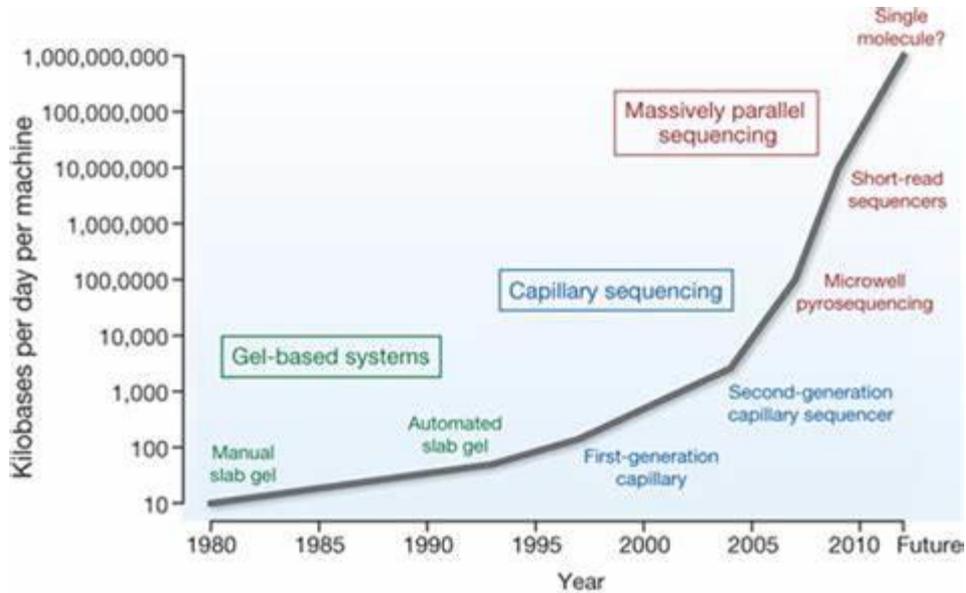
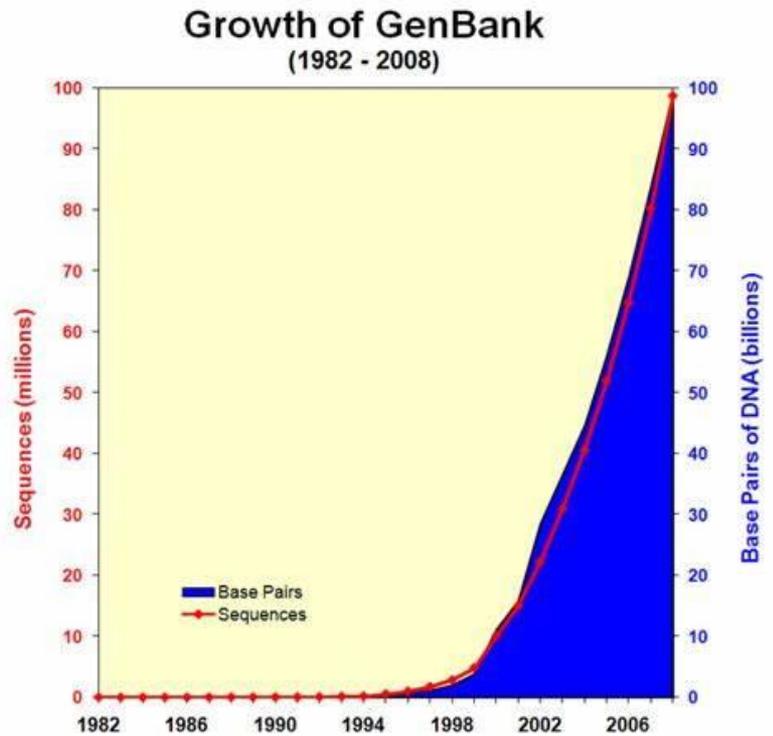


Figure 1. Genome annotation workflow.



**Desatualizados...mas,
é muita sequencia!!**



Muito dado – 2017!

The screenshot shows the UniProt website interface. The browser's address bar displays www.uniprot.org. The UniProt logo is in the top left, and a search bar with a dropdown menu set to 'UniProtKB' is in the top right. Below the search bar, navigation links for 'BLAST', 'Align', 'Retrieve/ID mapping', 'Help', and 'Contact' are visible. A mission statement reads: 'The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.'

The main content area features several highlighted sections:

- UniProtKB (UniProt Knowledgebase)**: This section is circled in red and contains two sub-sections:
 - Swiss-Prot (550,740)**: Manually annotated and reviewed.
 - TrEMBL (63,039,659)**: Automatically annotated and not reviewed.
- UniRef (Sequence clusters)**
- UniParc (Sequence archive)**
- Proteomes**
- News**: Includes social media icons for Blog, Twitter, Facebook, and RSS, and a link for 'Forthcoming changes'.
- Other resources**: Cross-ref. databases, Diseases (with 'XXX' text), and Keywords.

At the bottom, there are sections for 'Getting started' (with a 'Text search' link), 'UniProt data' (with a 'Download latest release' link), and 'Protein spotlight' (with a link to 'The Art Of Biocuration' from March 2016). The Windows taskbar at the bottom shows several open applications and the system clock indicating 11:06 on 29/03/2016.

SWISS-PROT TrEMBL, anotação automática de sequência de DNA Codante (CDS)

...2021...

Ejigu and Jung 2020.pdf x UniProt x +

uniprot.org

UniProt

UniProtKB Advanced Search

BLAST Align Retrieve/ID mapping Peptide search SPARQL Help Contact

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

UniProtKB
UniProt Knowledgebase
Swiss-Prot (564,638)
Manually annotated and reviewed.
Records with information extracted from literature and curator-evaluated computational analysis.
TrEMBL (214,406,399)
Automatically annotated and not reviewed.
Records that await full manual annotation.

UniRef
UniParc
Proteomes

New UniProt portal for the latest SARS-CoV-2 coronavirus protein entries and receptors, updated independent of the general UniProt release cycle.

SWISS-PROT TrEMBL, anotação automática de sequência de DNA Codante (CDS)

Literature citations
Cross-ref. databases
Taxonomy
Diseases
Subcellular locations
Keywords

UniProt release 2021_02
With a little help from my friend | SwissBioPics subcellular location visualization | Change of evidence codes for combinatorial evidence

UniProt release 2021_01
(Almost) all about that CBASS | Cross-references to VEuPathDB | Changes to humsavar.txt and related keywords | Reference proteomes downlo...

News archive

00:39
14/04/2021

Muito dado - 2023!

The image shows a screenshot of the UniProt website homepage. The browser's address bar shows 'uniprot.org'. The navigation menu includes 'UniProt', 'BLAST', 'Align', 'Peptide search', 'ID mapping', and 'SPARQL'. On the right, there are links for 'Release 2023_01', 'Statistics', and 'Help'. The main content area features four large tiles: 'Proteins UniProt Knowledgebase' (circled in red), 'Species Proteomes', 'Protein Clusters UniRef', and 'Sequence Archive UniParc'. Below these are sections for 'ProtNLM Predictions' and 'UniProt COVID-19 portal'. A black banner at the bottom contains a privacy notice and an 'Accept' button. The Windows taskbar at the very bottom shows the date as 11/04/2023 and the time as 20:50.

uniprot - Resultados da busca Ye x UniProt x +

uniprot.org

UniProt BLAST Align Peptide search ID mapping SPARQL

Release 2023_01 | Statistics Help

Proteins
UniProt Knowledgebase

Reviewed (Swiss-Prot) 569,213
Unreviewed (TrEMBL) 245,871,724

Species
Proteomes

Protein sets for species with sequenced genomes from across the tree of life

Protein Clusters
UniRef

Clusters of protein sequences at 100%, 90% & 50% identity

Sequence Archive
UniParc

Non-redundant archive of publicly available protein sequences seen across different databases

ProtNLM Predictions
Browse all the entries annotated with Google's ProtNLM predictions
What is ProtNLM?

UniProt COVID-19 portal
UniProt portal for the latest SARS-CoV-2 coronavirus protein entries and receptors, updated independent of the general UniProt release cycle

We'd like to inform you that we have updated our **Privacy Notice** to comply with Europe's new General Data Protection Regulation (GDPR) that applies since 25 May 2018.

Pesquisar

24°C Limpo 20:50 11/04/2023

Grande desafio...

RNAs não codantes??

Regiões regulatórias??

374-378 *Nucleic Acids Research*, 2003, Vol. 31, No. 1
DOI: 10.1093/nar/gkg108

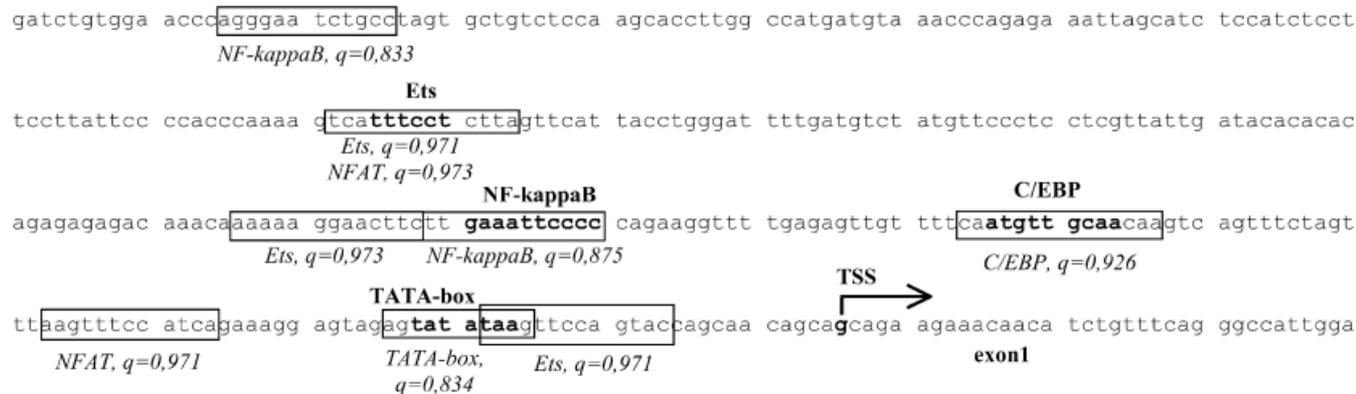
© 2003 Oxford University Press

TRANSFAC[®]: transcriptional regulation, from patterns to profiles

V. Matys^{1,*}, E. Fricke¹, R. Geffers¹, E. Gößling¹, M. Haubrock¹, R. Hehl², K. Hornischer¹, D. Karas¹, A. E. Kel¹, O. V. Kel-Margoulis¹, D.-U. Kloos¹, S. Land¹, B. Lewicki-Potapov¹, H. Michael², R. Münch¹, I. Reuter¹, S. Rotert¹, H. Saxel¹, M. Scheer¹, S. Thiele¹ and E. Wingender^{1,3}

¹BIOBASE GmbH, Halchtersche Strasse 33, D-38304 Wolfenbüttel, Germany, ²Institut für Genetik-Biozentrum, Technische Universität Braunschweig, Spielmannstrasse. 7, D-38106 Braunschweig, Germany and ³Gesellschaft für Biotechnologische Forschung mbH, Mascheroder Weg 1, D-38124 Braunschweig, Germany

Received September 16, 2002; Revised October 11, 2002; Accepted October 27, 2002



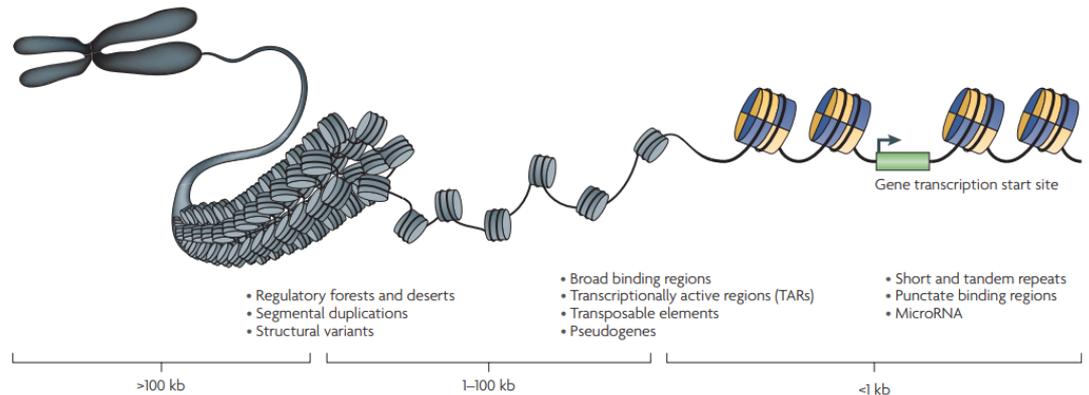
Annotating non-coding regions of the genome

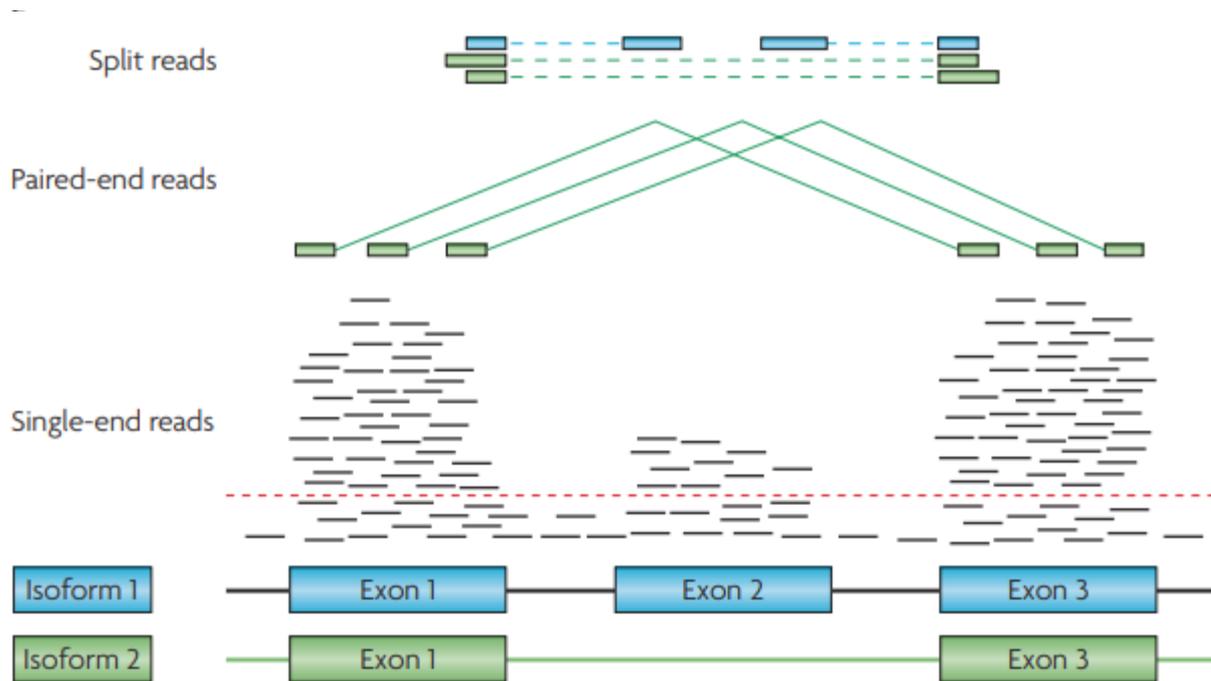
Roger P. Alexander^{*†}, Gang Fang^{*†}, Joel Rozowsky[†], Michael Snyder[§] and Mark B. Gerstein^{*†||}

Table 1 | Length, number and genome coverage of a representative collection of non-coding features

Classification	Property	Length (nucleotides)		Number of items	Genome coverage (Mb)	Genome coverage (%)
		Average	Longest			
<i>From comparative analysis</i>						
Short and tandem repeats	Simple repeat	63	2,961	415,917	26.1	0.84
	Satellite	1,444	160,602	8,997	13.0	0.42
	Low complexity	46	2,023	370,102	17.0	0.55
DNA transposons		215	3,625	459,524	98.6	3.17
Retrotransposons	LINEs	426	8,505	1,490,241	634.6	20.4
	Alu SINE element	261	614	1,186,885	309.7	9.97
Pseudogenes	Duplicated	6,607	181,882	2413	15.9	0.51
	Processed	723	15,732	8303	6.0	0.19
Segmental duplications		5,740	630 kb	26,469	151.9	4.89
Structural variants		8,761	3.3 Mb	96,874	848.8	27.3
<i>From functional analysis</i>						
Punctate binding sites	STAT1	446	9,079	~2,300	1.0	0.03
	CTCF	1,181	79,200	~35,000	41.4	1.33
	H3K4me3	1,759	71,025	~62,000	110.2	
Broad binding sites	H3K36me3	4,518	380,076	~130,000	589	
MicroRNA		89	150	718	0.063	
TARs		72	1,854	644,200	46.7	
Regulatory forests		3,890	35,165	68,900	268	
Regulatory deserts		27,107	203,691	72,500	1,970	

Box 1 | Catalogue of non-coding elements



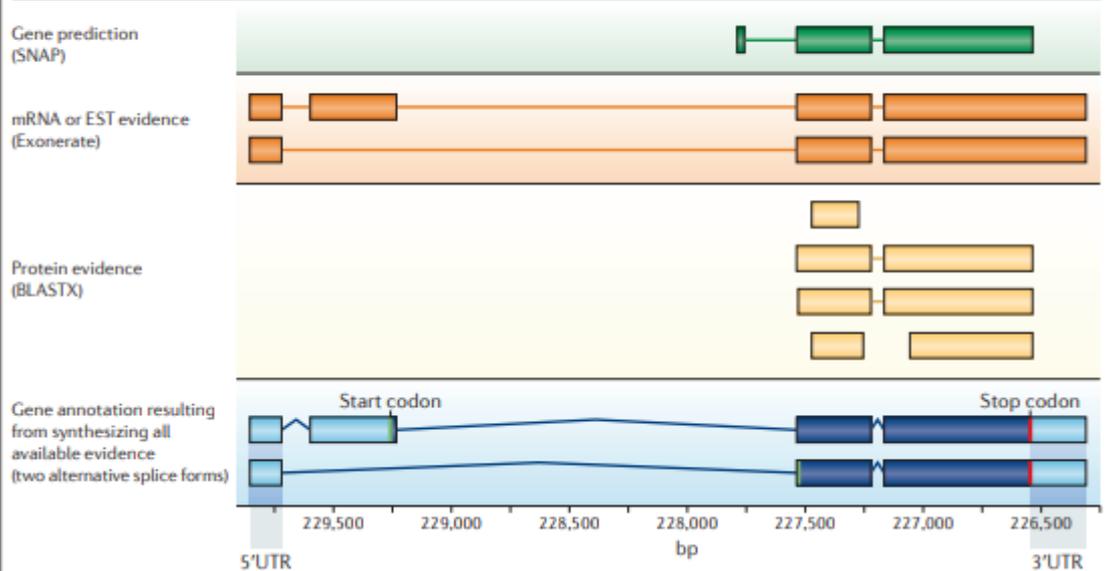


A beginner's guide to eukaryotic genome annotation

Mark Yandell and Daniel Ence

Abstract | The falling cost of genome sequencing is having a marked impact on the research community with respect to which genomes are sequenced and how and where they are annotated. Genome annotation projects have generally become small-scale affairs that are often carried out by an individual laboratory. Although annotating a eukaryotic genome assembly is now within the reach of non-experts, it remains a challenging task. Here we provide an overview of the genome annotation process and the available tools and describe some best-practice approaches.

Box 2 | Gene prediction versus gene annotation



Doi: 10.1038/nrg3174

MAS O QUE É UM GENE?

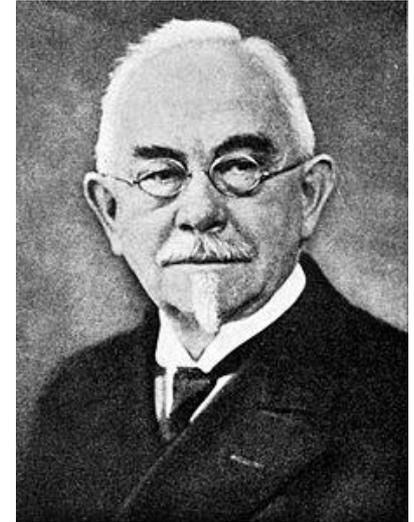


DEFINIÇÃO DE GENE

Pioneiro - Wilhelm Johannsen

1909 → gene

“The word gene is completely free of any hypothesis; it expresses only the evident fact that, in any case, many characteristics of the organism are specified in the germ cells by means of special conditions, foundations, and determiners which are present in unique, separate, and thereby independent way”



THE LAWS OF INHERITANCE.

Elemente der exakten Erblchkeitslehre. Deutsche wesentlich Erweiterte Ausgabe in Fünfundzwanzig Vorlesungen. By W. Johannsen. Pp. vi+516. (Jena: Gustav Fischer, 1909.) Price 9 marks.

WITHIN the last few years the output of exact experimental work upon phenomena of heredity has been very large, and the progress made, as compared with that of the previous forty years, has been astounding. In England it has chiefly been produced by investigators who have strictly segregated themselves either to the Mendelian or the biometrical schools, and who as a rule seem unable to realise the merits of the work of their rivals. One may pause in astonishment on reading, in a recent work issued by the head of the Mendelian school, that

GENE TÍPICO DE PROCARIOTOS

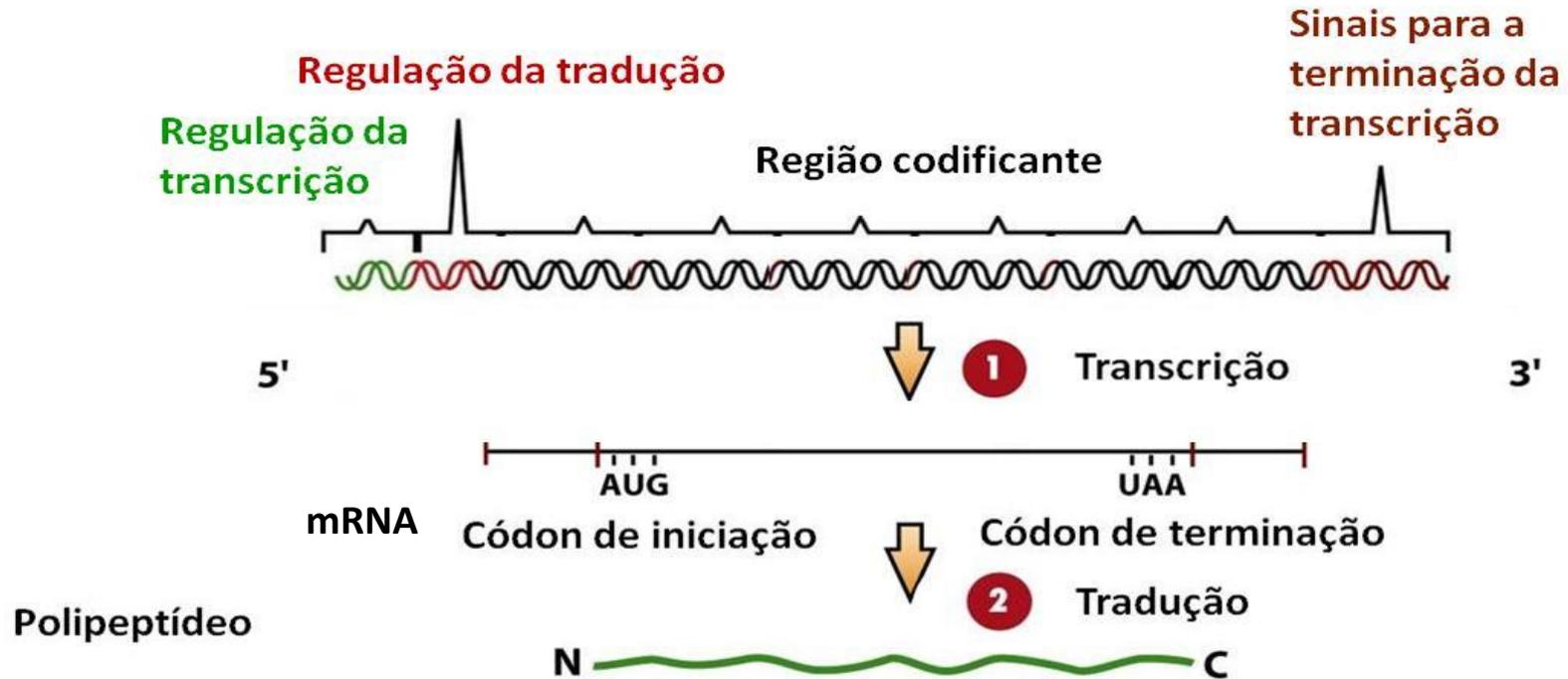
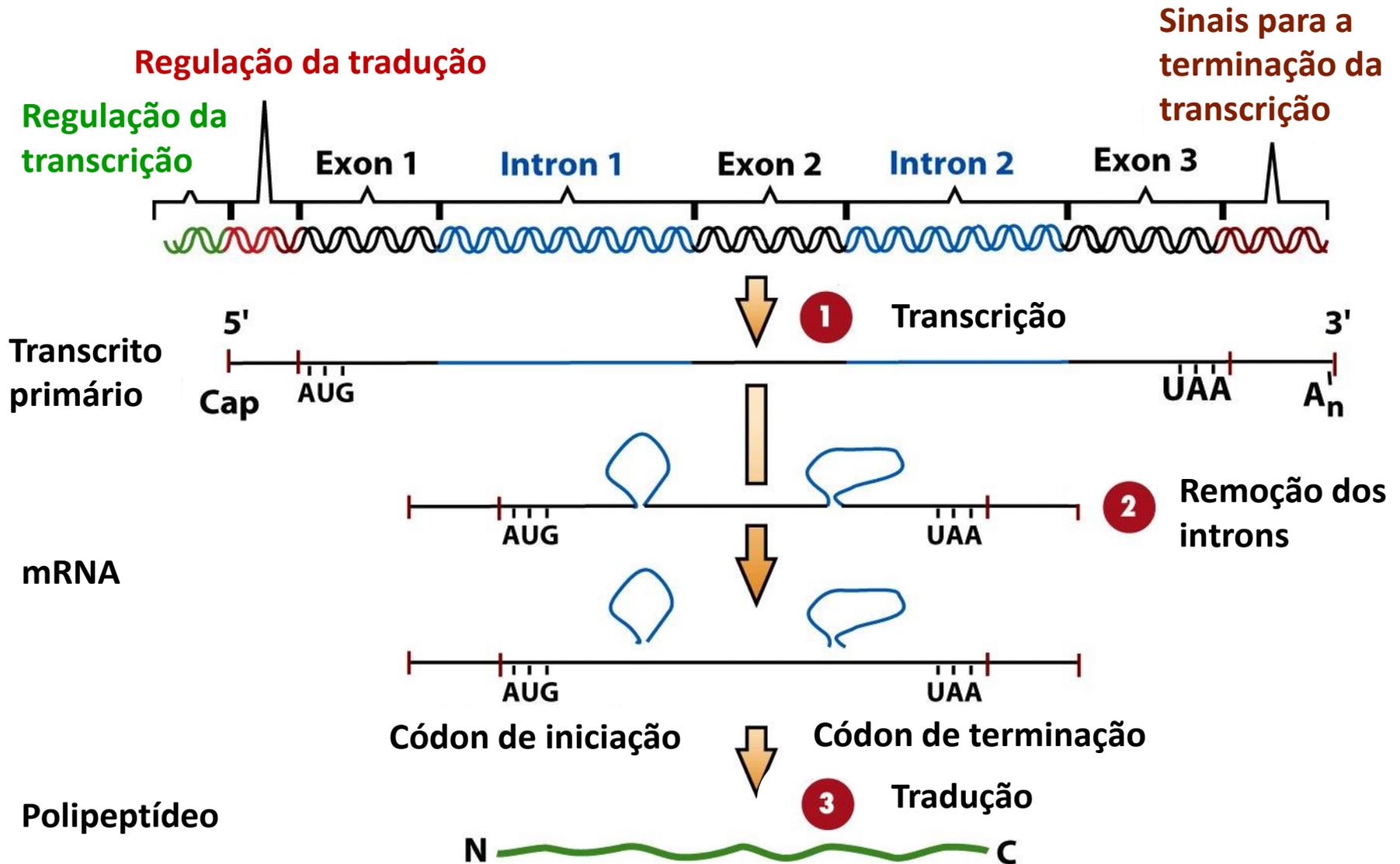


Figure 14-1b Principles of Genetics, 4/e
© 2006 John Wiley & Sons

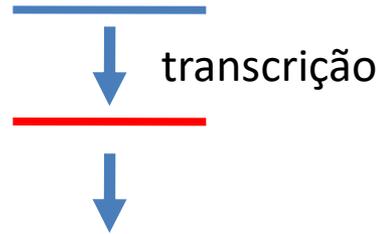
Unidade da informação genética que controla a síntese de polipeptídios!

GENE TÍPICO DE EUCARIOTOS

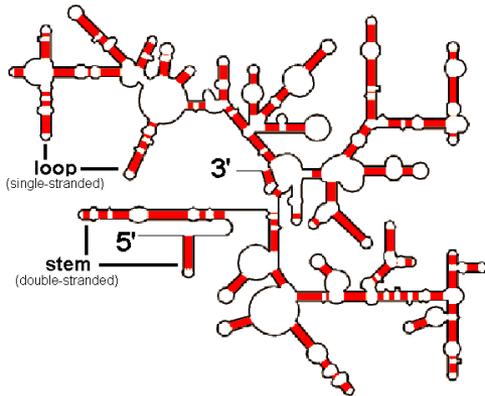


CONCEITO TAMBÉM SERVE PARA RNA ESTRUTURAL

DNA

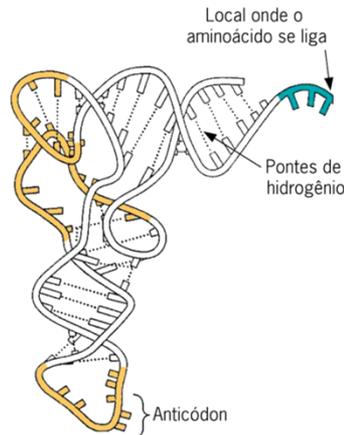


RNA ribossomal (rRNA)



RNA transportador (tRNA)
RNA ribossomal (rRNA)

RNA transportador (tRNA)



E os outros RNAs????

RNA em três dimensões.

What is a gene, post-ENCODE? History and updated definition

Mark B. Gerstein,^{1,2,3,9} Can Bruce,^{2,4} Joel S. Rozowsky,² Deyou Zheng,² Jiang Du,³ Jan O. Korbelt,^{2,5} Olof Emanuelsson,⁶ Zhengdong D. Zhang,² Sherman Weissman,⁷ and Michael Snyder^{2,8}

“A gene is a union of genomic sequences encoding a coherent set of potentially overlapping functional products”

HISTÓRICO DA DEFINIÇÃO DE GENE

Gene é uma unidade discreta de hereditariedade;

Gene é um locus distinto;

Gene codifica uma proteína;

Gene é uma molécula física;

Gene é uma unidade de transcrito;

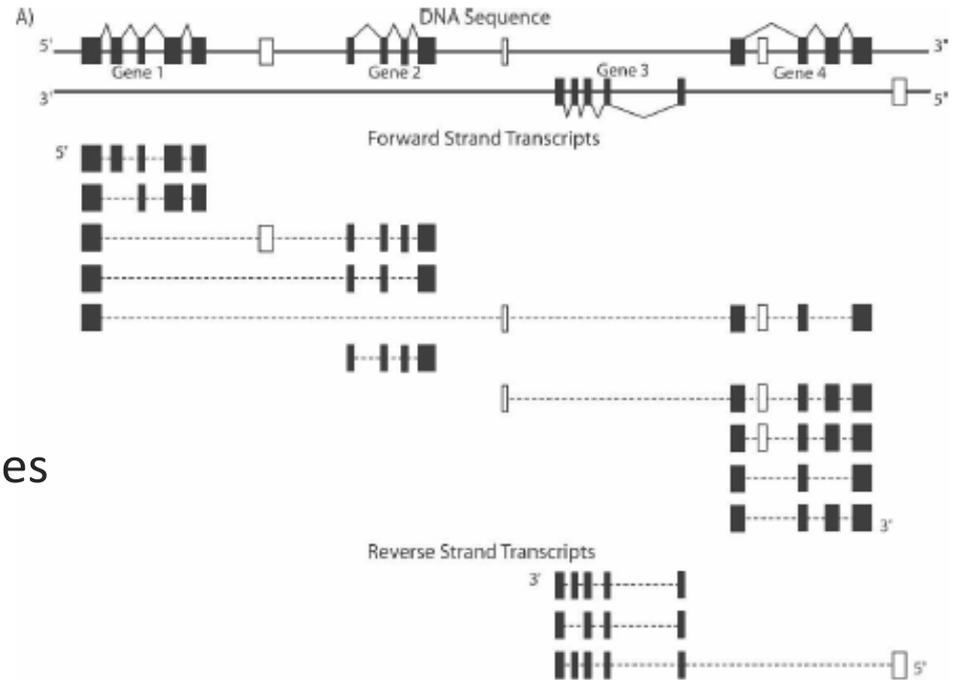
Gene é um quadro aberto de leitura (ORF);

SERÁ?

Table 1. Phenomena complicating the concept of the gene

Phenomenon	Description	Issue
<i>Gene location and structure</i> Intronic genes	A gene exists within an intron of another (Henikoff et al. 1986)	Two genes in the same locus
Genes with overlapping reading frames	A DNA region may code for two different protein products in different reading frames (Contreras et al. 1977)	No one-to-one correspondence between DNA and protein sequence
Enhancers, silencers	Distant regulatory elements (Spilianakis et al. 2005)	DNA sequences determining expression can be widely separated from one another in genome. Many-to-many relationship between genes and their enhancers.
<i>Structural variation</i> Mobile elements	Genetic element appears in new locations over generations (McClintock 1948)	A genetic element may be not constant in its location
Gene rearrangements/structural variants	DNA rearrangement or splicing in somatic cells results in many alternative gene products (Early et al. 1980)	Gene structure is not hereditary, or structure may differ across individuals or cells/tissues
Copy-number variants	Copy number of genes/regulatory elements may differ between individuals (Iafate et al. 2004; Sebat et al. 2004; Tuzun et al. 2005)	Genetic elements may differ in their number
<i>Epigenetics and chromosome structure</i> Epigenetic modifications, imprinting	Inherited information may not be DNA-sequence based (e.g., Dobrovic et al. 1988); a gene's expression depends on whether it is of paternal or maternal origin (Sager and Kitchin 1975)	Phenotype is not determined strictly by genotype
Effect of chromatin structure	Chromatin structure, which does influence gene expression, only loosely associated with particular DNA sequences (Paul 1972)	Gene expression depends on packing of DNA. DNA sequence is not enough to predict gene product.
<i>Post-transcriptional events</i> Alternative splicing of RNA	One transcript can generate multiple mRNAs, resulting in different protein products (Bergert et al. 1977; Gelinas and Roberts 1977)	Multiple products from one genetic locus; information in DNA not linearly related to that on protein
Alternatively spliced products with alternate reading frames	Alternative reading frames of the INK4a tumor suppressor gene encodes two unrelated proteins (Quelle et al. 1995)	Two alternative splicing products of a pre-mRNA produce protein products with no sequence in common
RNA trans-splicing, homotypic trans-splicing	Distant DNA sequences can code for transcripts ligated in various combinations (Borst 1986). Two identical transcripts of a gene can trans-splice to generate an mRNA where the same exon sequence is repeated (Takahara et al. 2000).	A protein can result from the combined information encoded in multiple transcripts
RNA editing	RNA is enzymatically modified (Eisen 1988)	The information on the DNA is not encoded directly into RNA sequence
<i>Post-translational events</i> Protein splicing, viral polyproteins	Protein product self-cleaves and can generate multiple functional products (Vila-Komaroff et al. 1975)	Start and end sites of protein not determined by genetic code
Protein trans-splicing	Distinct proteins can be spliced together in the absence of a trans-spliced transcript (Handa et al. 1996)	Start and end sites of protein not determined by genetic code
Protein modification	Protein is modified to alter structure and function of the final product (Wold 1981)	The information on the DNA is not encoded directly into protein sequence
<i>Pseudogenes and retrogenes</i> Retrogenes	A retrogene is formed from reverse transcription of its parent gene's mRNA (Vanin et al. 1980) and by insertion of the DNA product into a genome	RNA-to-DNA flow of information
Transcribed pseudogenes	A pseudogene is transcribed (Zheng et al. 2005, 2007)	Biochemical activity of supposedly dead elements

NOVA VISÃO DOS GENES

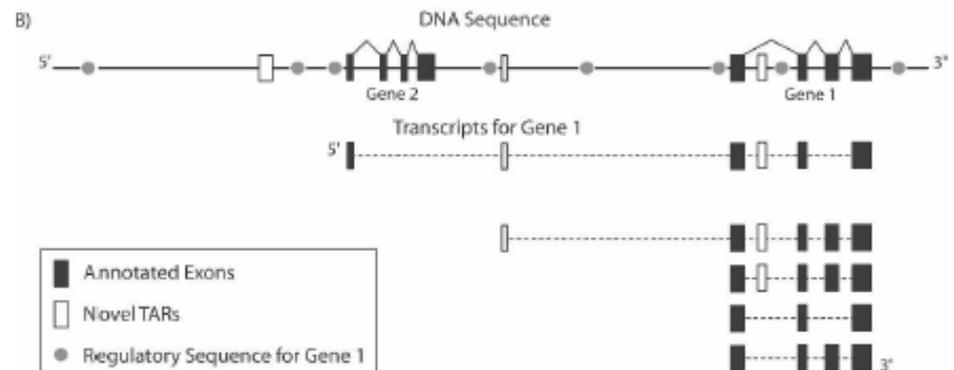


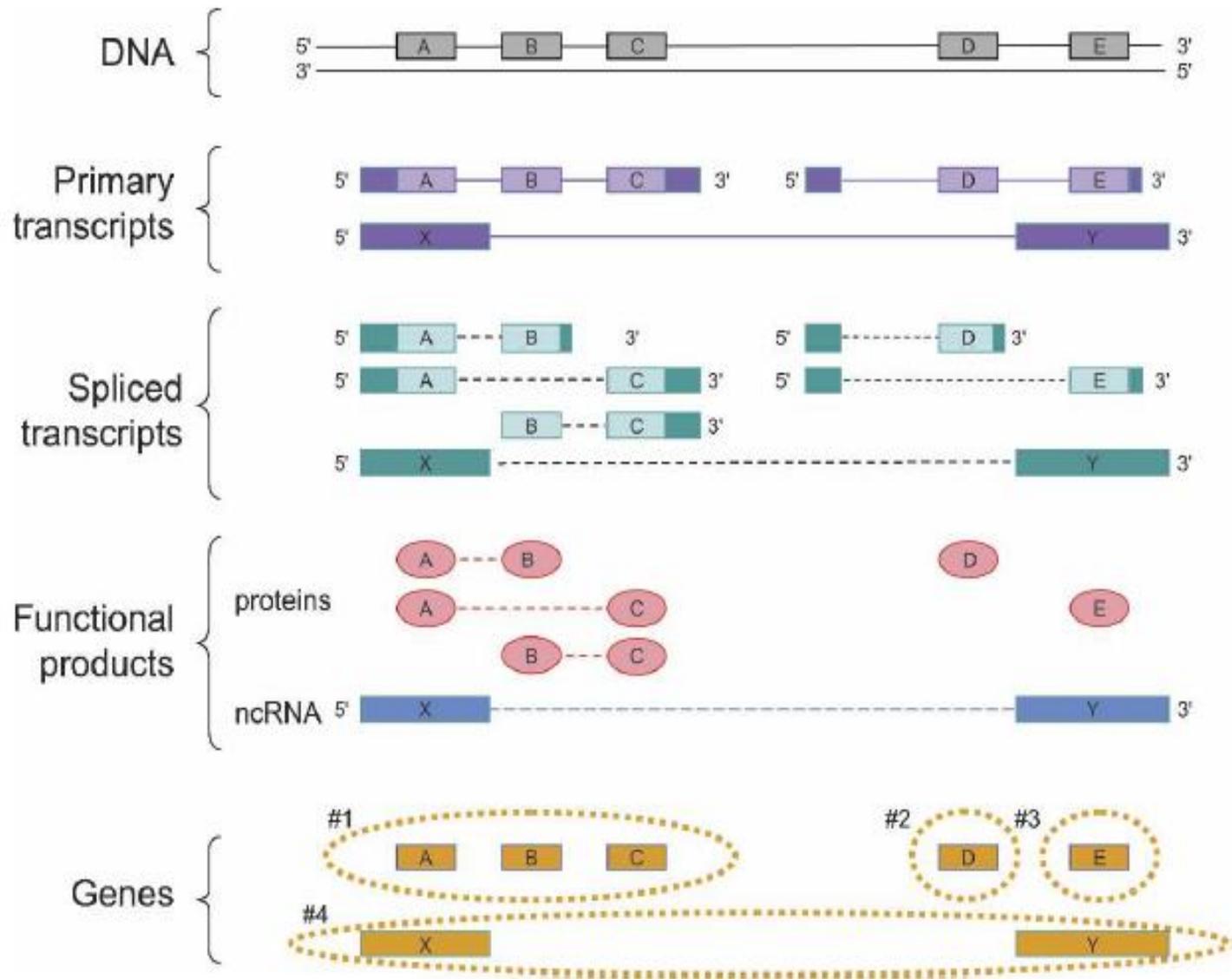
Sequências regulatória dentro dos genes

Muitas possibilidades...

TAR – regiões ativas de transcrição

TSS – sitio de inicio de transcrição





The Evolving Definition of the Term “Gene”

Petter Portin*¹ and Adam Wilkins[†]

*Laboratory of Genetics, Department of Biology, University of Turku, 20014, Finland and [†]Institute of Theoretical Biology, Humboldt Universität zu Berlin, 10115, Germany

Table 1 Abridged list of different propositions for a definition of the gene in the current era given by different authors

Essential Content or Character of the Proposition	Classification	Author(s)
These three first operational definitions give criteria, formal, experimental and computational, for identifying genes in the DNA sequences of genomes, annotation of genomes, and for specifying the function of genes	Operational	Snyder and Gerstein (2003)
	Operational	Pesole (2008)
	Operational	Stadler <i>et al.</i> (2009)
In these three following definitions, classified as molecular, the structural and the functional gene are conceptually distinguished and separated	Molecular	Scherrer and Jost (2007)
	Molecular	Keller and Harel (2007)
	Molecular	Burian (2004)
In this definition two gene concepts, “gene-P (preformationist)” and “gene-D (developmental)”, are distinguished	Complex	Moss (2003)
This definition presents three different concepts of the gene: instrumental, nominal and postgenomic	Complex	Griffiths and Stotz (2006)
This definition aims at to define the gene on the basis of its products and separates it from DNA	A new kind of redefinition	Waters (1994)

Received: 14 April 2021 | Revised: 5 October 2021 | Accepted: 8 December 2021

DOI: 10.1111/bioe.13006

ORIGINAL ARTICLE

bioethics  WILEY

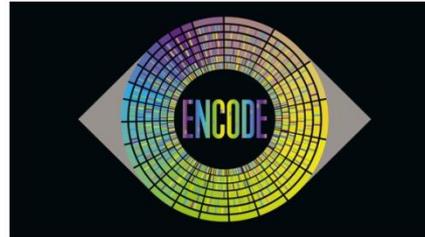
The ethical gene

Reuven Brandt 

COLLECTION | 05 SEPTEMBER 2012

Encode

Access the collected papers by exploring the thematic threads that run through them, with topics such as DNA methylation, RNA or machine learning.



<https://www.nature.com/collections/aghcdefffg/>

ENCODE

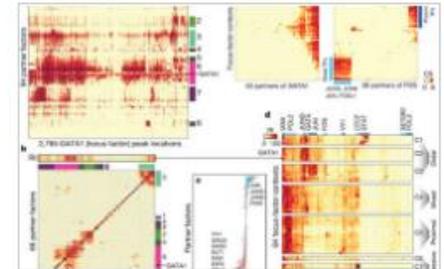
3 Characterization of intergenic regions and gene definition

The prevalence and analysis of ENCODE data are changing the definition and characterization of intergenic and genic regions

The cumulative coverage of transcribed regions in the 15 cell lines across the human genome is 62.1% and 74.7% for processed and primary transcripts, respectively (Supplementary Table 10 and Supplementary Fig. 22). On average, for each cell line, 39% of the genome is covered by primary transcripts and 22% by processed RNAs. No cell line showed transcription of more than 56.7% of the union of the expressed transcriptomes across all cell lines. When mapping the current RNA-seq data to the ENCODE pilot regions (Supplementary Table 10), we observed a similar, albeit higher, extent of transcriptional coverage of 73.3% for processed RNAs and 84.5% for primary transcripts. Previously reported estimates in these regions for processed and primary transcripts were 24% and 93%, respectively (Supplementary Table 2.4.3 and ref. 3). The increased genome coverage by processed RNAs stems largely from the inclusion of non-polyadenylated RNAs in the current study. Other than that, given the differences in the samples studied, the selection of pilot regions with high genic content, the increase of annotated genomic regions over time, and the different technologies used to interrogate transcription, both estimates are in reasonable agreement.

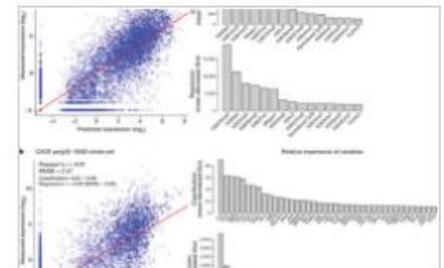
10 Characterization of network topology

ENCODE data analysis helps to describe the various types of regulatory "wiring" implicit in the genome



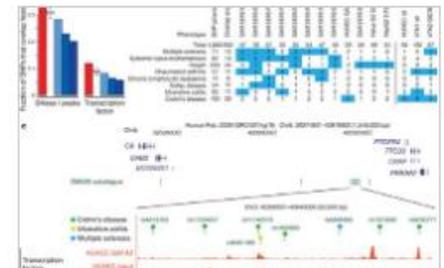
11 Machine learning approaches to genomics

ENCODE has applied machine learning approaches to enable integration and exploration of large and diverse data



12 Impact of functional information on understanding variation

ENCODE provides an initial interpretation of many human variants and plausible leads for the role of many variants identified in genome-wide association studies



CONSTRUÇÃO PRESENTE NA SOJA RR[®]

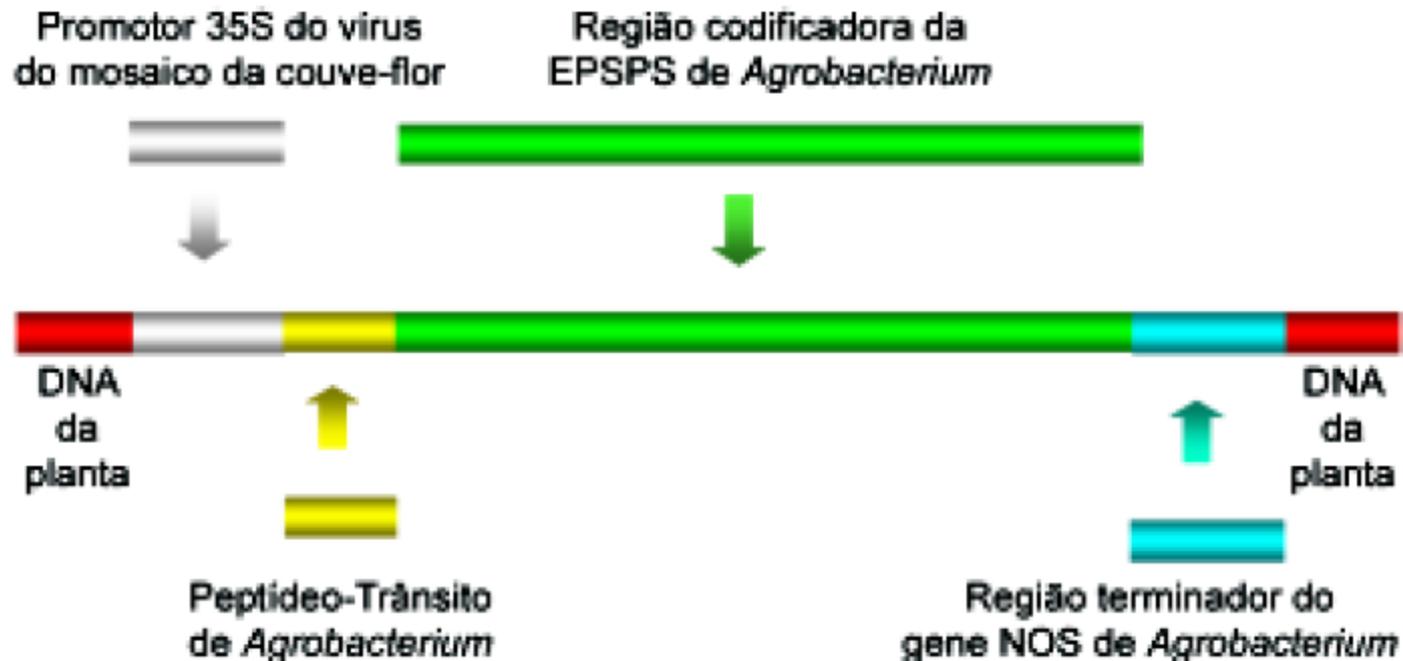
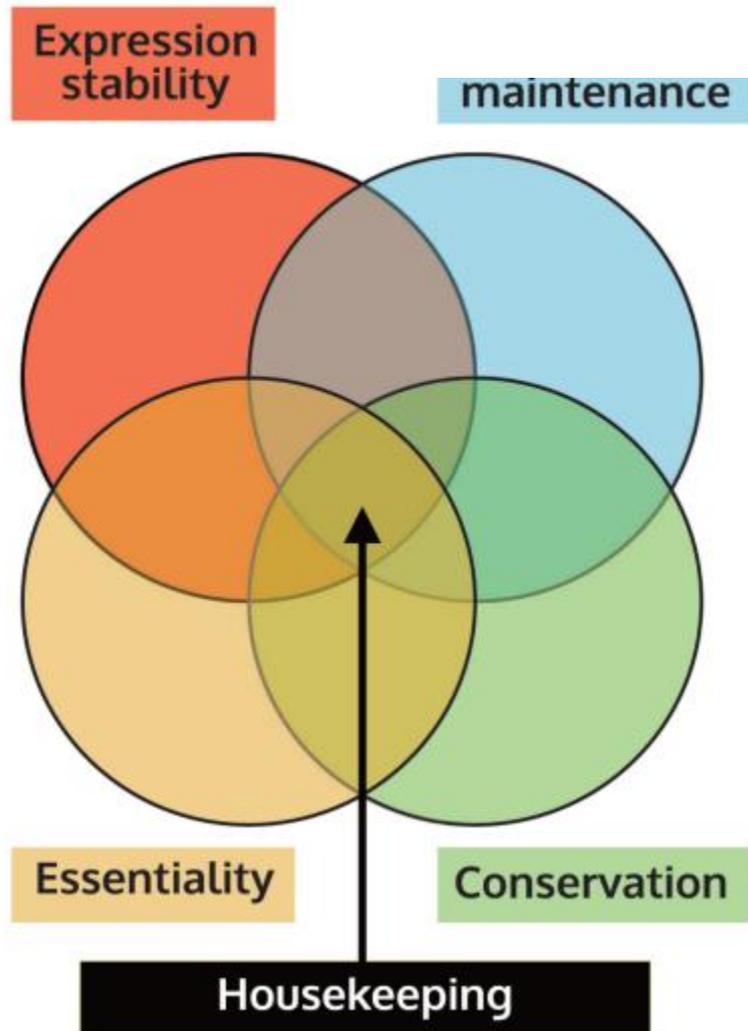


Figura 1 - Representação da construção presente na soja RR[®] (*Roundup Ready*). Região promotora 35S do vírus do mosaico da couve flor, peptídeo de trânsito de *Petúnia*, gene que codifica a proteína EPSPS, que confere a resistência ao herbicida, e o terminador do gene da nopalina sintase (NOS).

RESEARCH ARTICLE

What are housekeeping genes?

Chintan J. Joshi¹, Wenfan Ke², Anna Drangowska-Way², Eyleen J. O'Rourke^{2*},
Nathan E. Lewis^{1,3,4,5*}



Systems biology

Defining the extent of gene function using ROC curvature

Stephan Fischer ^{1,2} and Jesse Gillis ^{1,3,*}

¹Cold Spring Harbor Laboratory, Stanley Institute for Cognitive Genomics, Cold Spring Harbor, NY 11724, USA, ²Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris F-75015, France and ³Department of Physiology, University of Toronto, Toronto, ON, Canada

Systems biology

Transfer learning across ontologies for phenome–genome association prediction

Raphael Petegrosso¹, Sunho Park², Tae Hyun Hwang² and Rui Kuang^{1,*}

¹Department of Computer Science and Engineering, University of Minnesota Twin Cities, Minneapolis, MN 55455, USA and ²Department of Clinical Sciences, University of Texas Southwestern Medical Center, Dallas, TX 75390, USA

O MELHOR ENTENDIMENTO DOS GENES CONTRIBUI COM O MELHOR ENTENDIMENTO DA ARQUITETURA GENÔMICA!!!

HIV tipo I -19.750 b



Milho
2.5 Gb



Mamute
4.17 Gb



Escherichia coli
5 Mb



Humano
3 Gb

NÚMERO DE GENES EM EUCARIOTOS

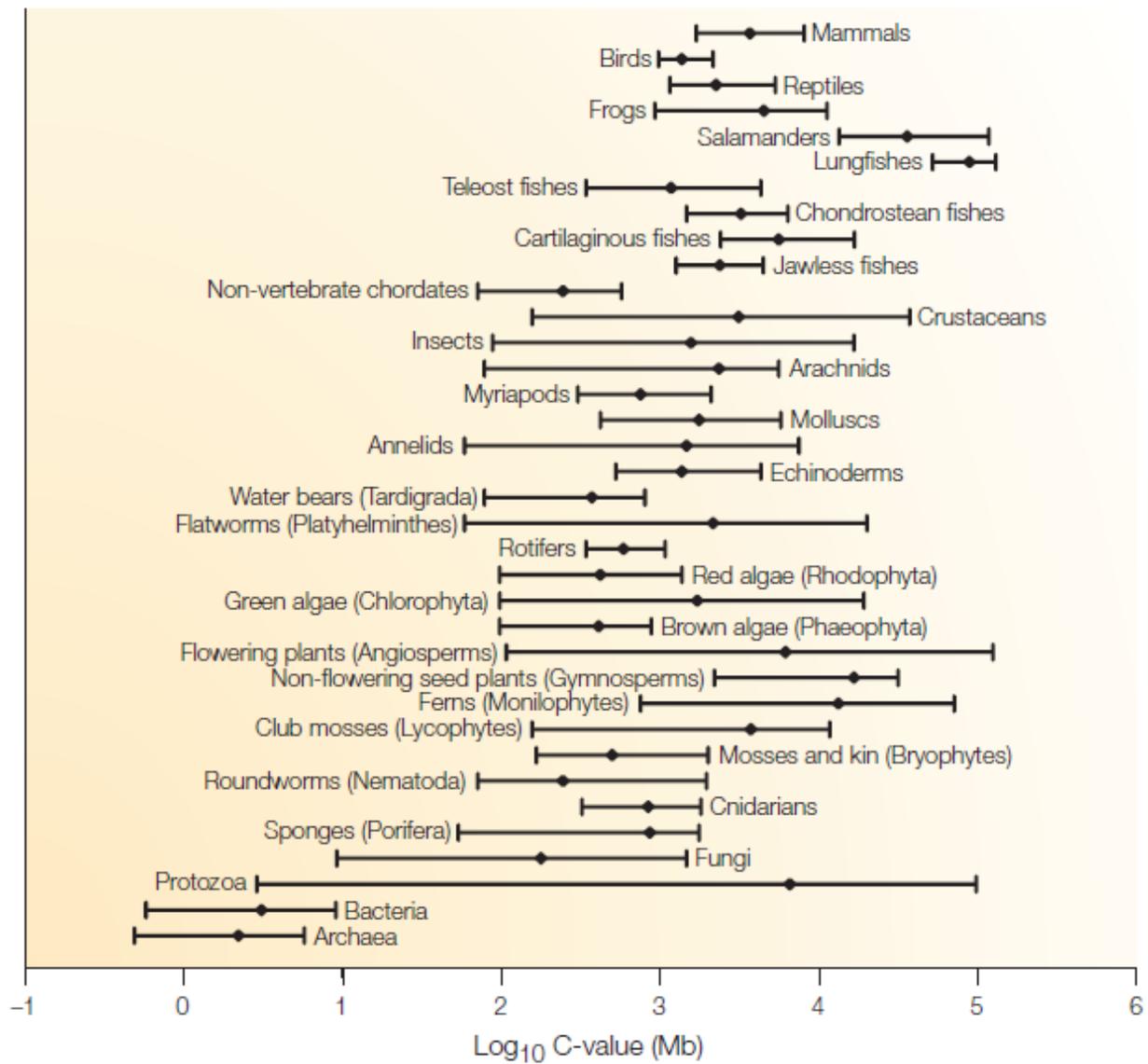
Espécies	Genoma (Mb)	Genes
<i>D. melanogaster</i>	165	~12.000
<i>S. cerevisiae</i>	13	~6.000
<i>C. elegans</i>	97	~20.000
<i>H. sapiens</i>	3.300	~25.000



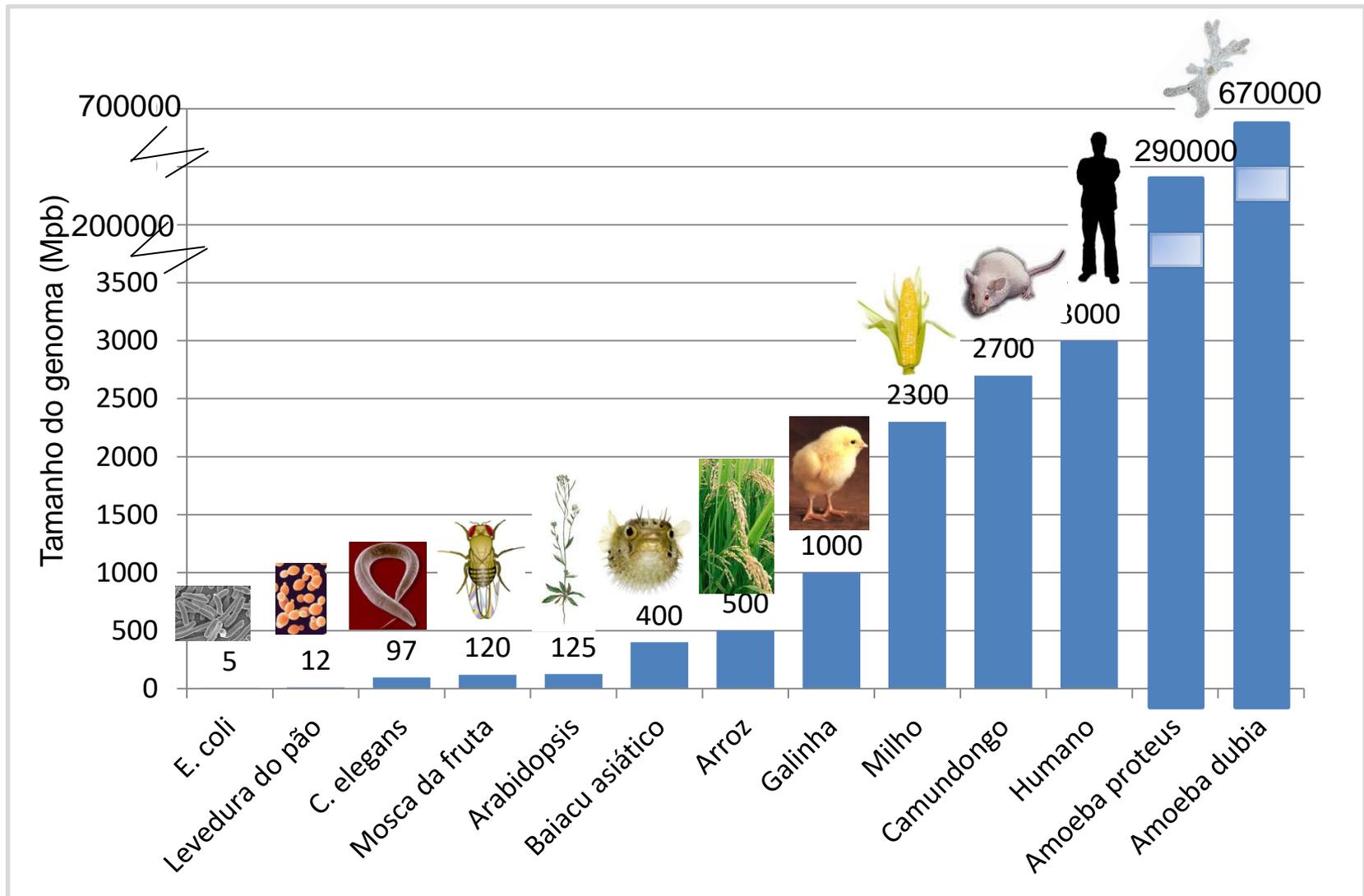
Table 2 | **Percentage of non-coding DNA in selected sequenced genomes**

Species name		Genome size (Mb)	Fraction of genome (%)				Source of gene annotations
Common	Scientific		Genic	Exonic	Non-coding		
				Intronic	Intergenic		
Yeast	<i>Saccharomyces cerevisiae</i>	12.2	73.5	72.9	0.6	26.6	Saccharomyces Genome Database (June 2008 build)
Nematode worm	<i>Caenorhabditis elegans</i>	100.3	59.2	28.1	31.2	40.8	WormBase (WS190)
Fruitfly	<i>Drosophila melanogaster</i>	168.7	48.2	18.3	30.0	51.8	FlyBase and Berkeley Drosophila Genome Project (BDGP; release no. 5)
Human	<i>Homo sapiens</i>	3,107	45.1	2.8	42.3	54.9	UCSC Genome Browser Known Genes table (hg18)

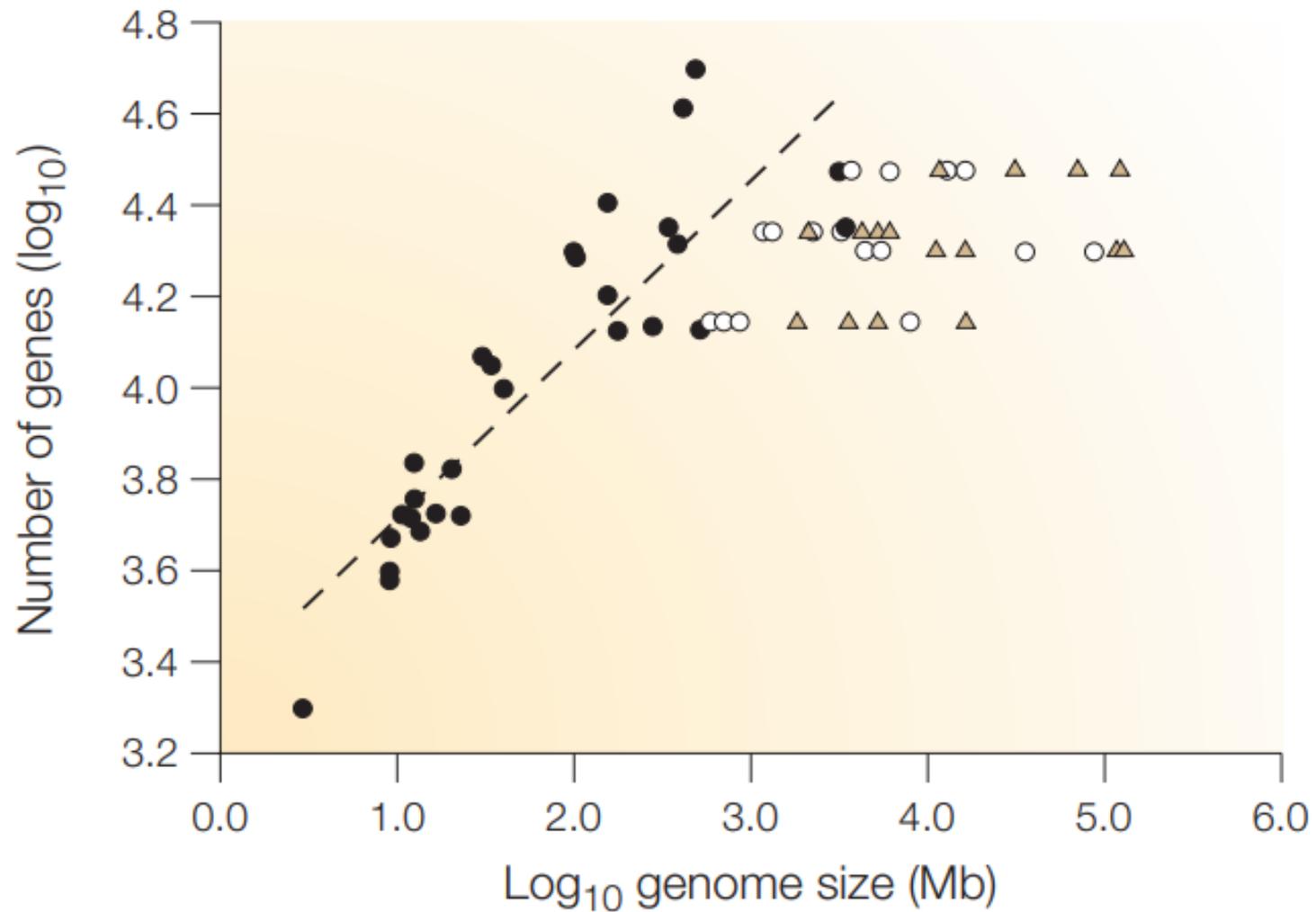
The genic fraction consists of both exonic and intronic sequence. The exonic fraction consists of both coding sequence (CDS) and 5' and 3' UTRs. Strictly speaking, UTRs are non-coding, so the exon fraction is a slight overestimate of the fraction of coding sequence in the genome.



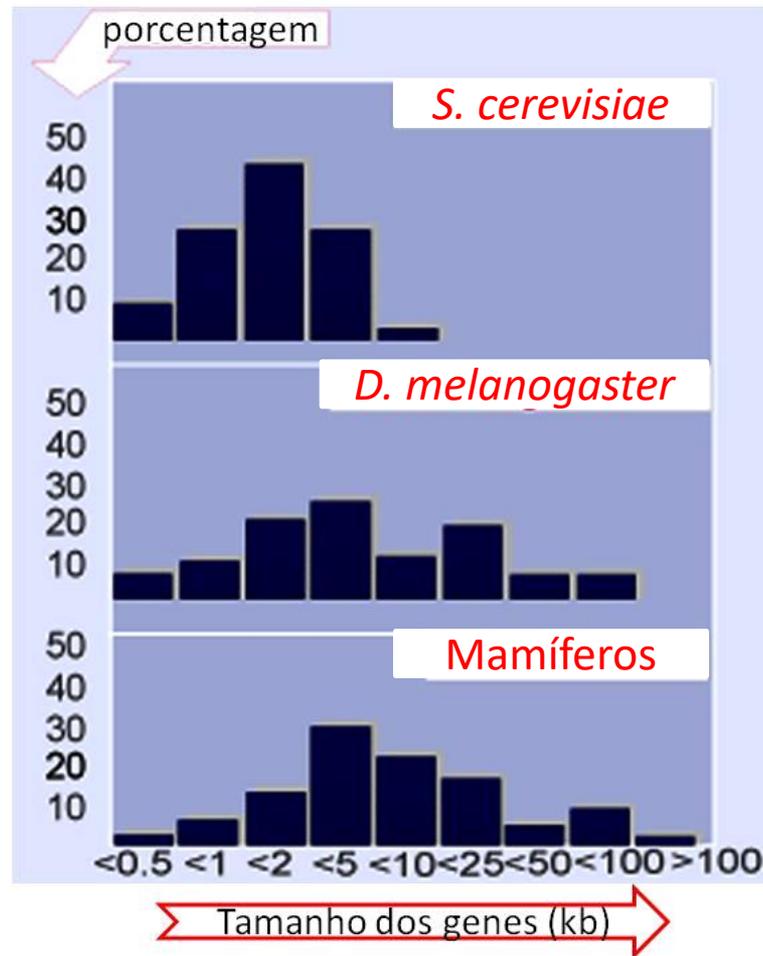
COMPARAÇÃO NO TAMANHO DE GENOMAS



A complexidade de um organismo não é diretamente proporcional ao tamanho do genoma; alguns organismos unicelulares possuem muito mais DNA que os humanos.

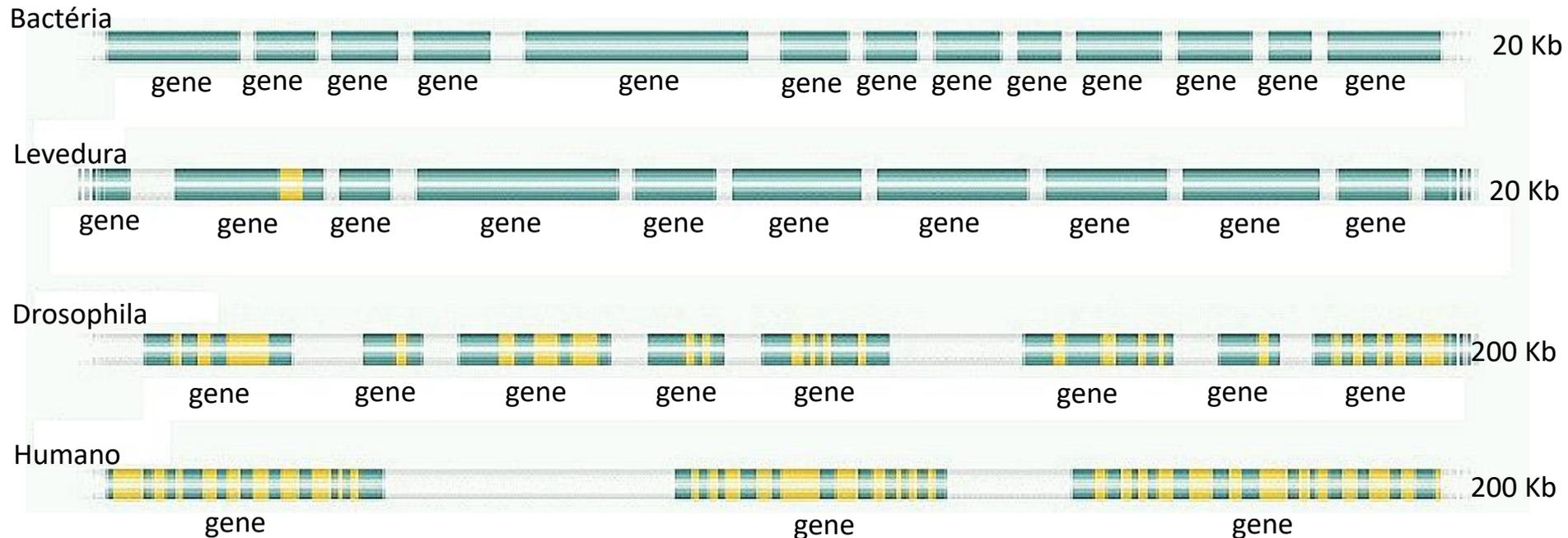


ORGANIZAÇÃO DOS GENES EM EUCARIOTOS



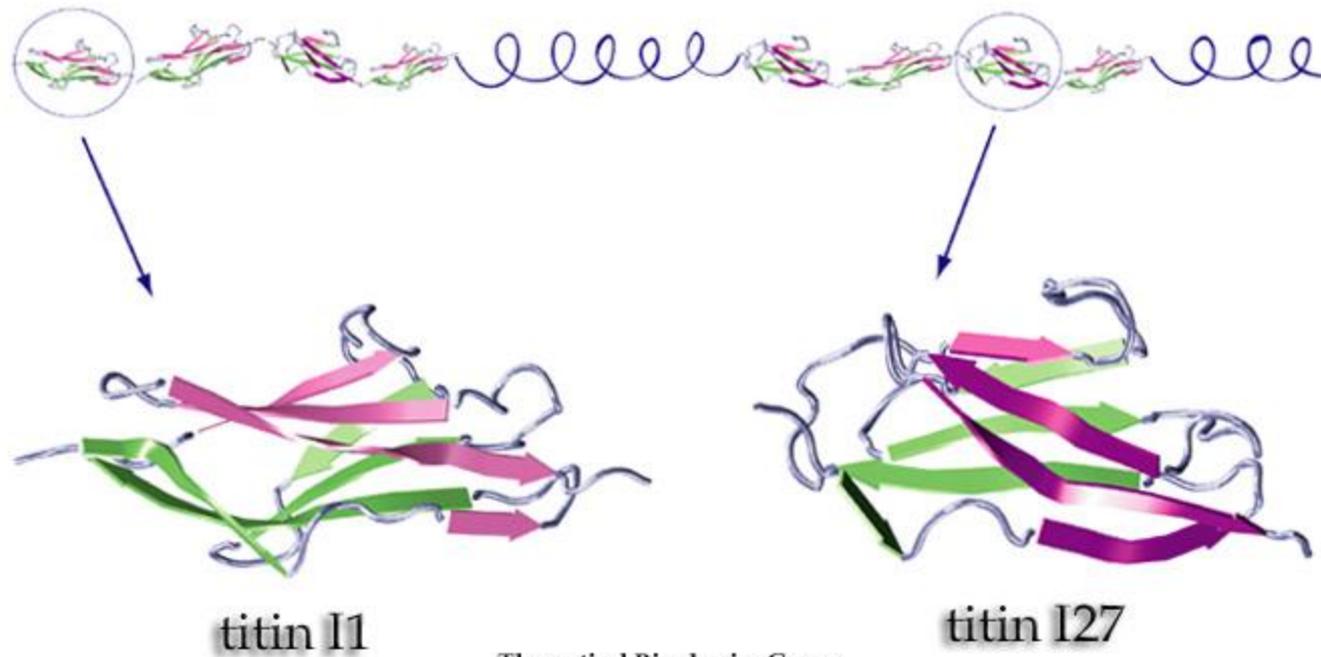
**Enorme variação no tamanho dos genes dos mamíferos!
De onde vem isso?**

GENES NA MOLÉCULA DE DNA

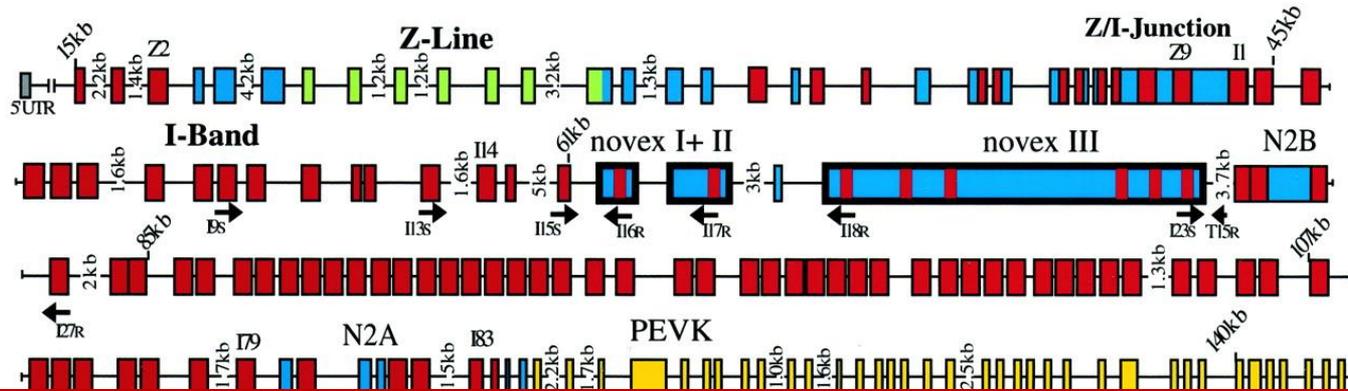


Grande variação nos tamanho dos genes geradas pela presença dos íntrons!

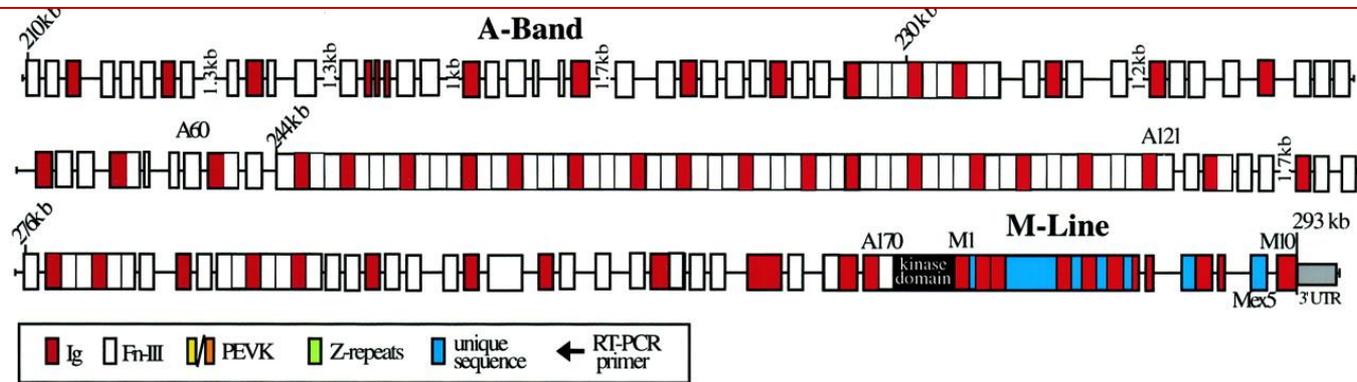
A proteína gigante do músculo titina contém 38 138 resíduos de aminoácidos (contém 363 exons) que tem um papel importante na contração e elasticidade dos músculos.



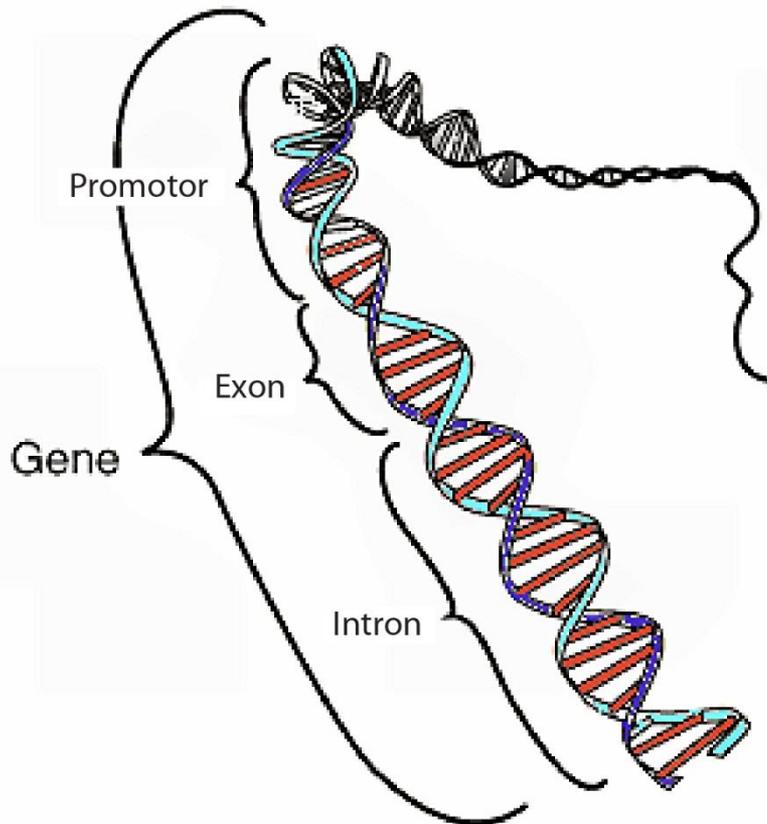
Estrutura de exon-intron do gene titin (293 kb)



O TAMANHO MÉDIO DE UM ÉXON HUMANO É DE 150 pb!!!!



Bang, M.-L. et al. Circ Res 2001;89:1065-1072

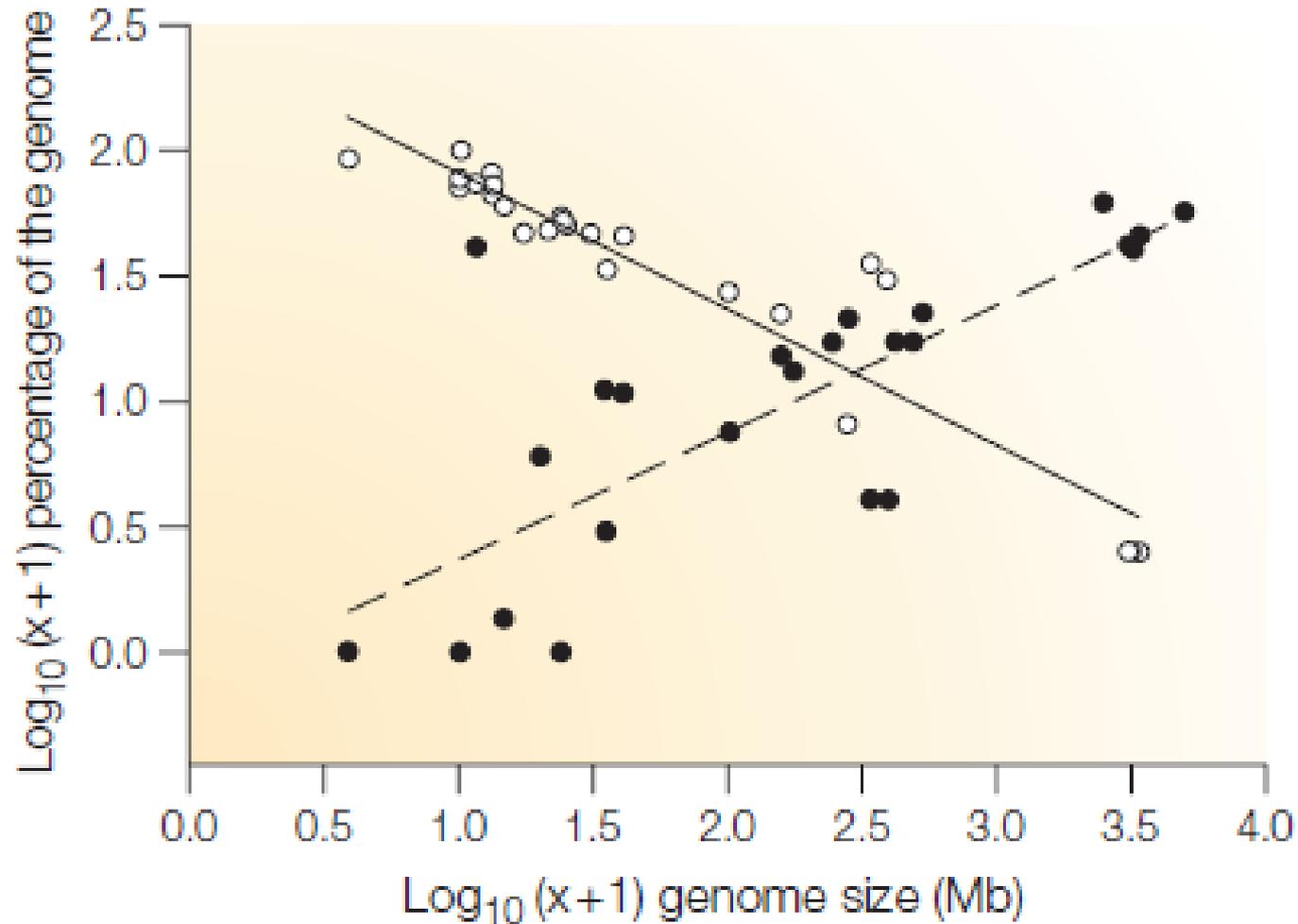


DNA lixo???

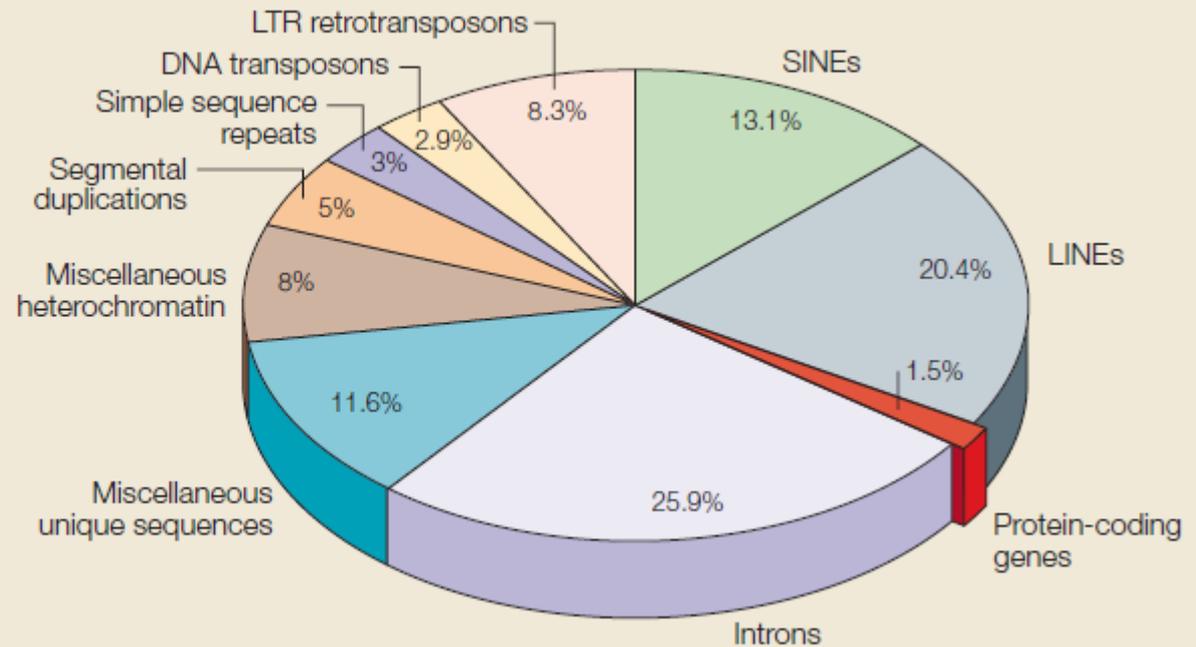
Just because of you don't understand, you can't call us "Junk"!



CORRELAÇÃO ENTRE TAMANHO DE GENOMAS E ELEMENTOS DE TRANSPOSIÇÃO



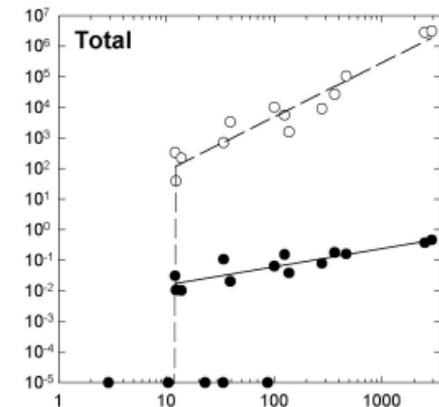
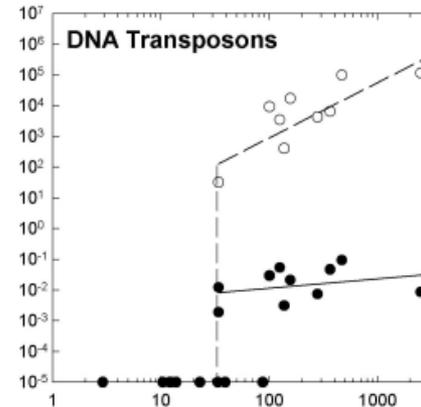
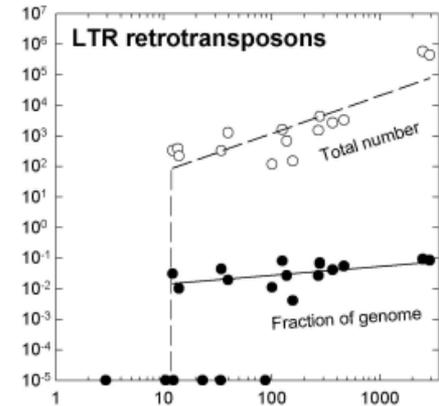
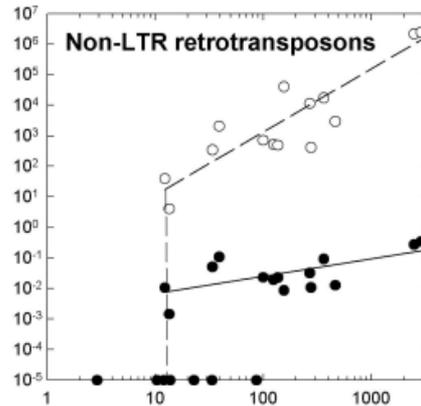
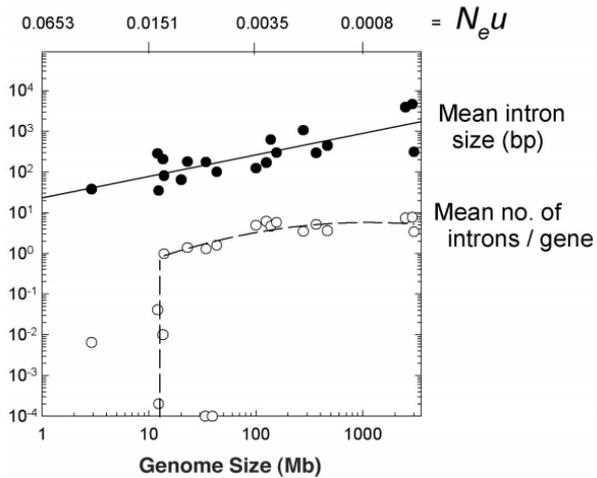
The figure provides a summary of the different components of the human genome. Less than 1.5% of the genome consists of the suspected 20,000–25,000 protein-coding sequences. By contrast, a large majority is made up of non-coding sequences such as introns (almost 26%) and (mostly defunct) transposable elements (nearly 45%). Data are taken from REF. 16.



O surgimento de elementos de transposição tem proporcionado o duplicação do tamanho de genoma em humanos

The Origins of Genome Complexity

Michael Lynch^{1*} and John S. Conery²

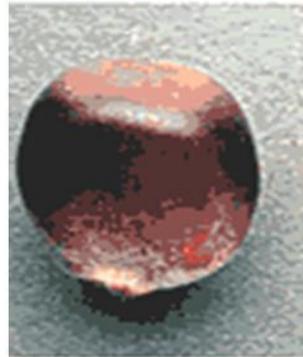


Genome Size (Mb)



Barbara McClintock (1902-1992)

Nobel Prize in 1983



Bz

Forma normal



bz

Mutação de ponto



bz-m

Inserção do
elemento transponível



Transposable elements and the evolution of genome size in eukaryotes

Margaret G. Kidwell

Department of Ecology and Evolutionary Biology, The University of Arizona, Tucson, AZ 85721, USA (Phone: (520) 621-1784; Fax: (520) 621-9190; E-mail: kidwell@azstarnet.com)

Key words: genome size, molecular evolution, transposable element

Novas classes surgem quase diariamente!

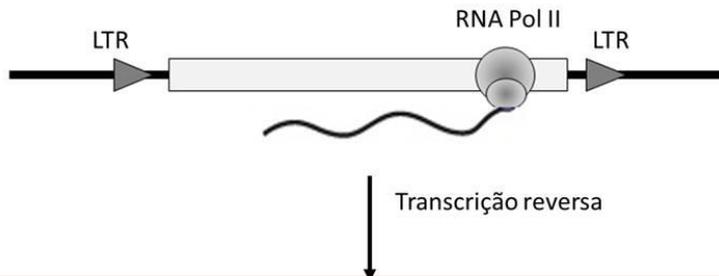
Banco : RepBase -

<http://www.girinst.org/rebase/>

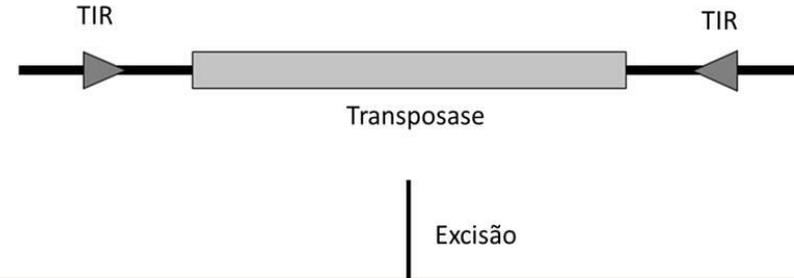
Table 1. Characteristics of some widely distributed types of transposable elements

Class	Subclass	Superfamily	Family examples	Approximate size range (bp)
I. Retroelements RNA-mediated elements	LTR retro-transposons	<i>Ty1-copia</i>	<i>Opie-1</i> in maize	3000–12,000
		<i>BEL</i>	<i>Ce7</i> in <i>C. elegans</i>	3000–20,000
	<i>DIRS-1</i>	<i>DIRS-1</i> in <i>Dictyostelium</i>	~5000	
	<i>Ty3-gypsy</i>	<i>Gypsy</i> in <i>Drosophila</i>	5000–14,000	
Non-LTR retroposons	<i>LINEs</i>	<i>LINE-1</i> in humans;	<i>LINE-1</i> in humans;	1000–7000
		<i>I</i> element in <i>Drosophila</i>	<i>I</i> element in <i>Drosophila</i>	
	<i>SINEs</i>	<i>Alu</i> in humans	100–500	
II. Transposons DNA-mediated elements	Cut and paste transposition DDE signature present	<i>mariner-Tc1</i>	<i>Tc1</i> in <i>C. elegans</i> ;	1000–2000
		<i>Mu</i>	<i>mariner</i> in <i>Drosophila</i>	
	Cut and paste transposition DDE signature absent	<i>MITEs</i>	<i>Mu</i> in maize;	400–20,000
		<i>hAT</i>	<i>MULEs</i> in <i>Arabidopsis</i>	100–500
Rolling circle (RC) transposition	<i>P</i>	<i>Tourist</i> in maize	100–500	
	<i>Helitrons</i>	<i>hobo</i> in <i>Drosophila</i> ;	500–4600	
Unclassified	Rolling circle (RC) transposition	<i>P</i>	<i>Ac</i> in maize;	
		<i>Helitrons</i>	<i>Tam-3</i> in <i>Anthrithum</i>	
Unclassified	Rolling circle (RC) transposition	<i>Foldback</i>	<i>P</i> in <i>Drosophila</i>	5500–17,500
		<i>Mini-me</i>	<i>Helitrons</i> in <i>A. thaliana</i> , <i>O. sativa</i> , and <i>C. elegans</i>	Large, but highly variable
			<i>Galileo</i> in <i>D. buzzatii</i>	500–1200
			<i>Mini-me</i> in many <i>Drosophila</i> species	

**Retrotransposons
Elementos da Classe I**



**Transposons
Elementos da Classe II**



A classe predominante depende do grupo taxonômico:

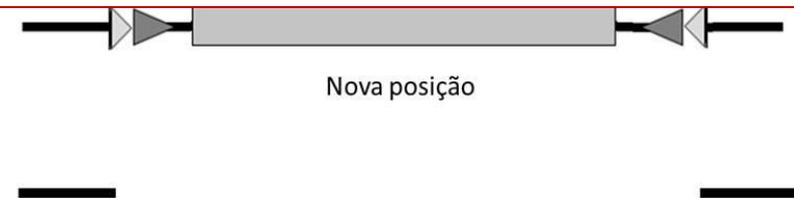
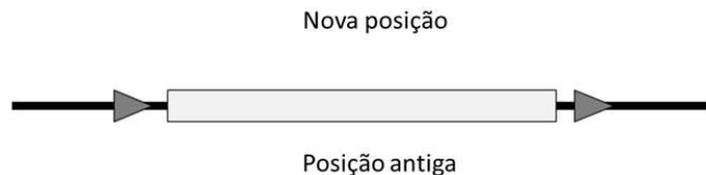
LINE e SINE – em mamíferos.

SINE - presentes em aves mas não ativos

LTR retrotransposons - gramíneas

DNA transposons – nematóides

Em casos muito raros transposon não são encontrados em genomas pequenos



DESAFIO NO GENOMA DE PLANTAS...

“Uma das primeiras e mais proeminentes observações nas propriedades moleculares do material genético de plantas superiores foi a tremenda variação no tamanho do genoma”:

Bennetzen and Kellog, 1997 - The Plant Cell, 9: 1509-1514

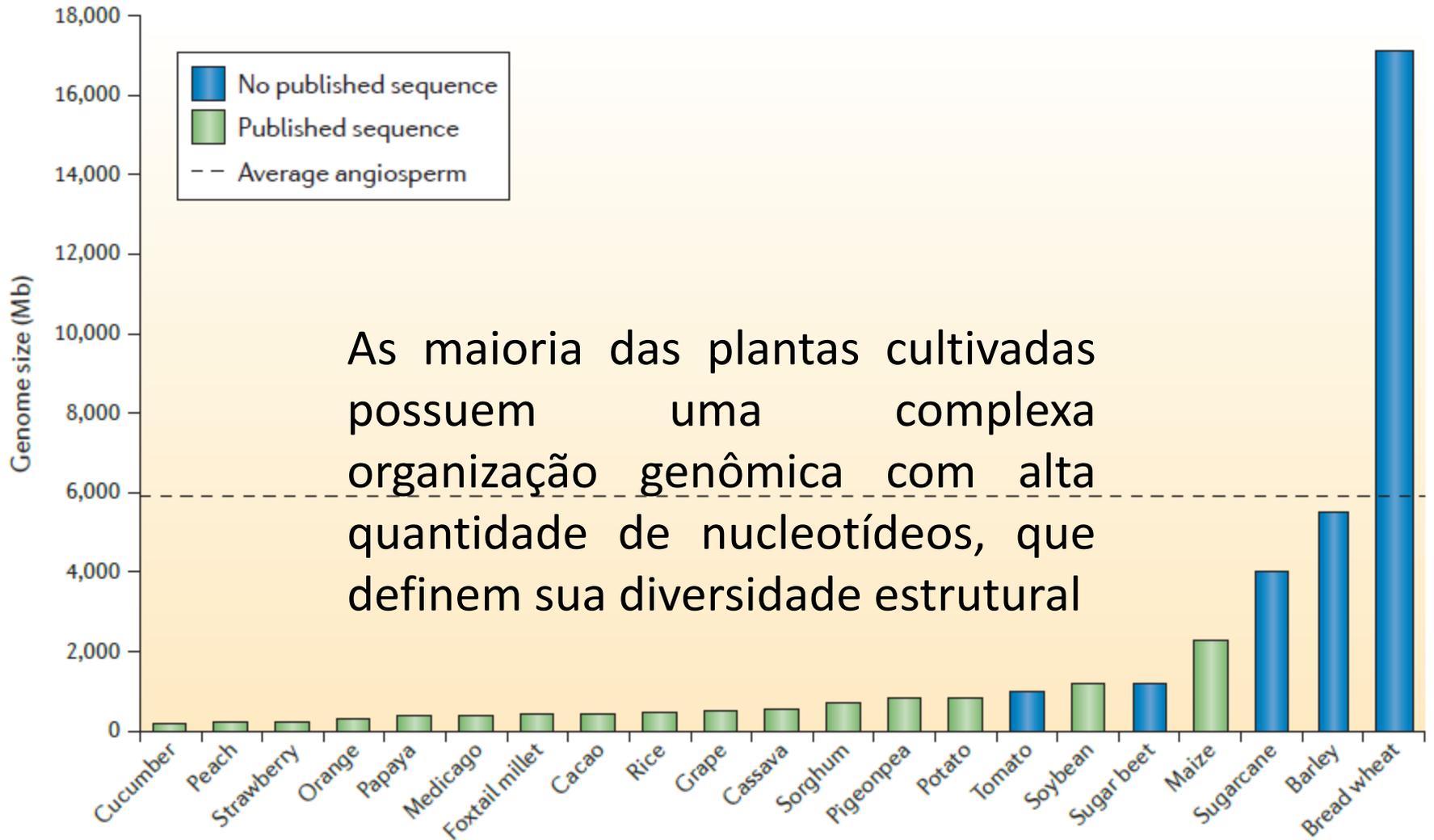
**Muitos modelos estatísticos não
funciona, para plantas..
Muitas famílias multigênicas
Poliplóides
Muitos parálogos**



Fritillaria assyriaca - 110000 Mpb



Arabidopsis thaliana - 110 Mpb



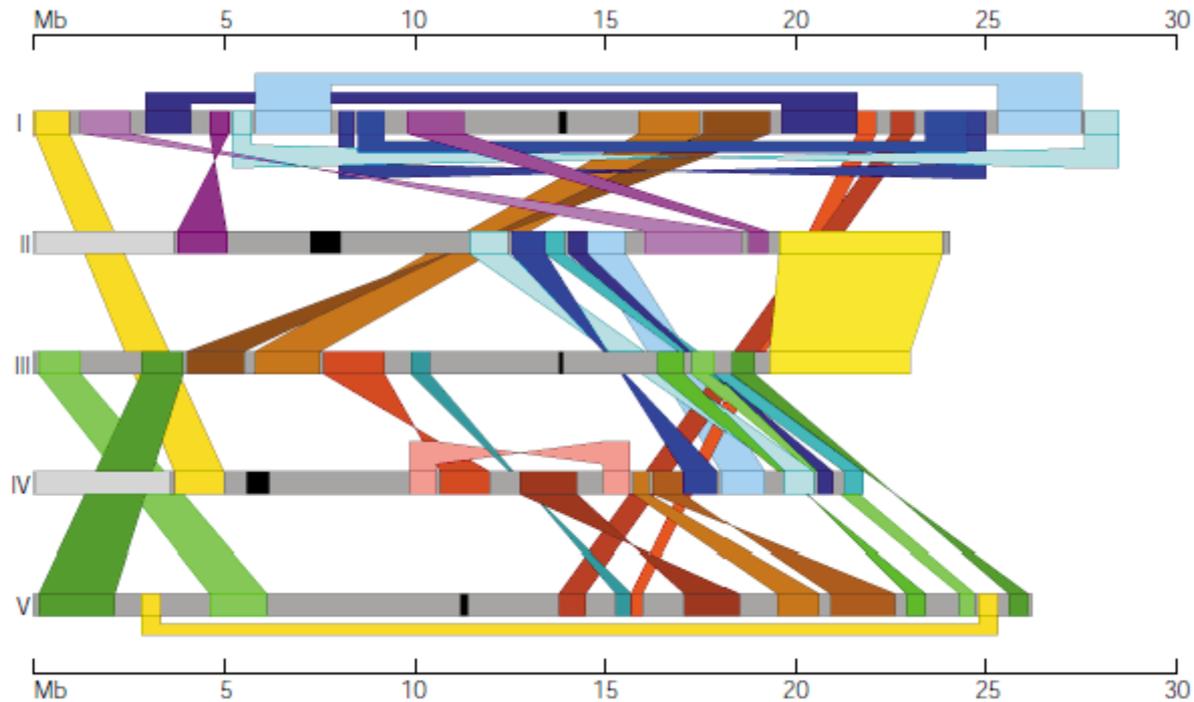
doi: 10.1038/nrg3097

Table 1 | Crop genome characteristics

Crop	Genome size (Mb)	Gene number	Transposable element content (%)	Refs
Cucumber	200	21,000		Phytozome
Peach	230	28,000		Phytozome
Strawberry*	240	35,000	22	16
Orange	320	25,000		Phytozome
Papaya	370	29,000	52	149
Medicago*	375	48,000	30	15
Foxtail millet	410	35,000		Phytozome
Cacao	430	29,000	24	150
Rice	450	41,000	25	151
Grape	490	30,000	41	24
Cassava	530	31,000		Phytozome
Sorghum	730	28,000	63	109
Pigeonpea	833	49,000	52	152
Potato	840	39,000	62	17
Tomato	1,000			Kew
Soybean	1,200	46,000	59	153
Sugar beet	1,200			Kew
Maize	2,300	33,000	85	154
Sugarcane	4,000			Kew
Barley	5,500			Kew
Bread wheat	17,100			Kew
Average angiosperm	5,900			Kew

The table shows the genome size, gene and transposable element content for the world's ten top production crops and all other crops with sequenced genomes. In the 'Refs' column, 'Phytozome' refers to <http://www.phytozome.net> and 'Kew' refers to <http://data.kew.org/cvalues>. *Sequence of a species related to the main crop.

Genoma de *A. thaliana* – muitos parálogos

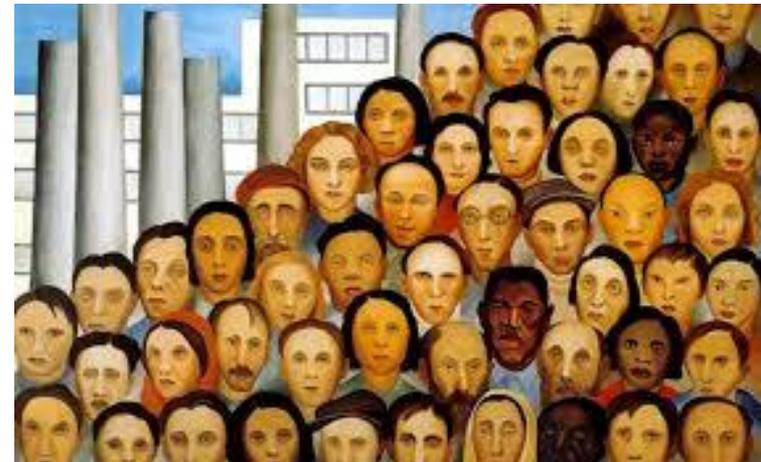


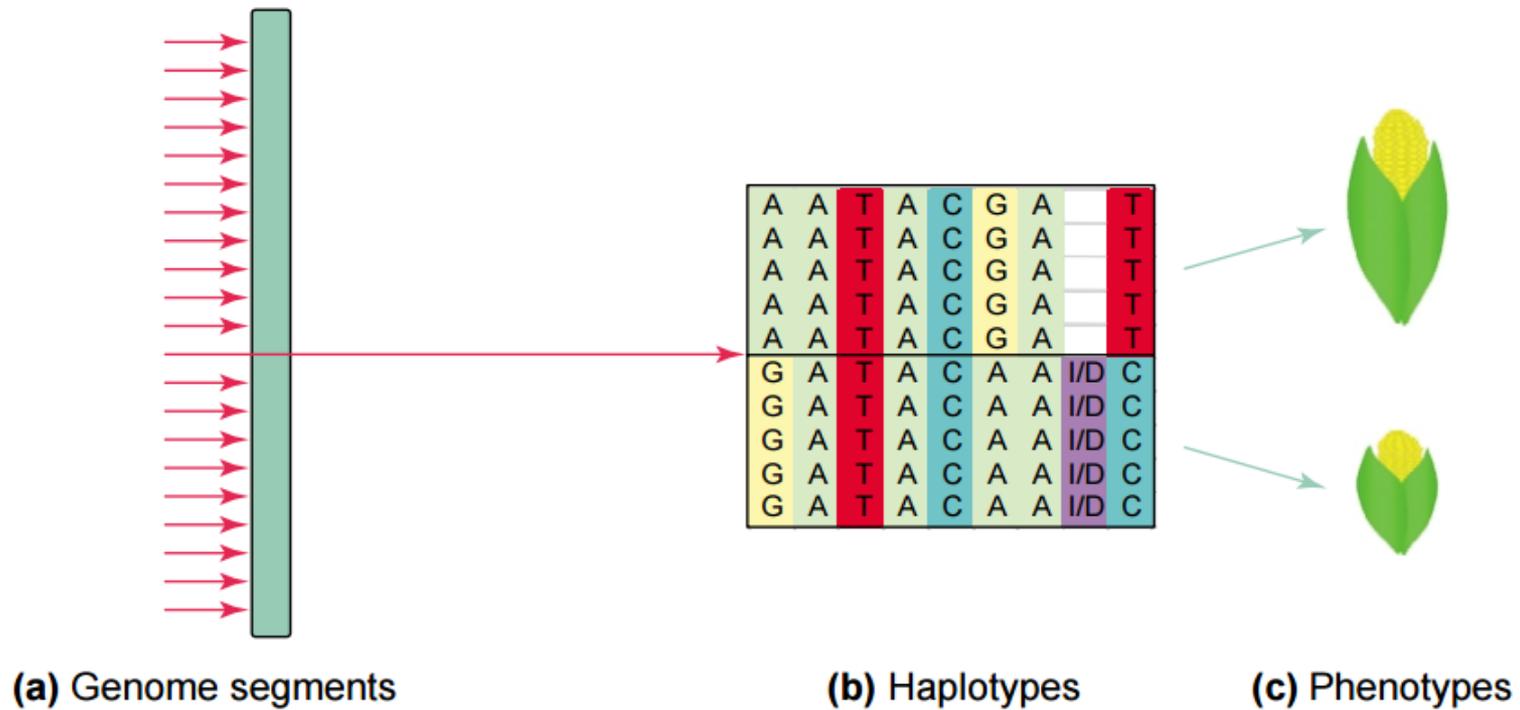
Corn and humans: recombination and linkage disequilibrium in two genomes of similar size

Antoni Rafalski¹ and Michele Morgante²

¹DuPont Crop Genetics, Experimental Station E353, PO Box 80353, Wilmington DE, 19880-0353, USA

²Dipartimento di Produzione Vegetale e Tecnologie Agrarie, Università di Udine, Via delle Scienze 208, Polo Scientifico Rizzi - 33100 Udine, Italy





TRENDS in Genetics

Influência na eficiência de marcadores moleculares devido ao desequilíbrio de ligação

Table 2. Linkage disequilibrium (LD) in different species

Species	LD	Criterion
Human	60 kb	D' half-length, North Europeans
Human	5 kb	D' half-length, Yorubans
Cattle	> 10 cM	D' half length
<i>Arabidopsis thaliana</i>	50–100 kb	r ² , half length
Soybean	> 50 kb	Little LD decay found
Norway spruce	~ 100 bp	r ² half length
	~ 200 bp	r ² = 0.2
Grape	> 500 bp	r ² half length
Maize	~ 400 bp	r ² = 0.2
Maize (inbreds from the USA)	~ 1 kb	r ² = 0.2
Maize	200–1500 bp	r ² = 0.2

^aS. Degli Ivanissevich and M. Morgante, unpublished.

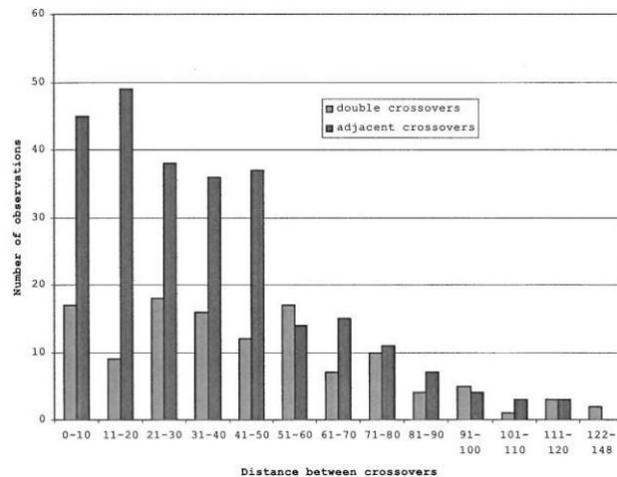
^bF. Cattonaro *et al.*, unpublished.

**A história evolutiva dos genomas
são específicas!!!**

A High-Density Genetic Recombination Map of Sequence-Tagged Sites for *Sorghum*, as a Framework for Comparative Structural and Evolutionary Genomics of Tropical Grains and Grasses

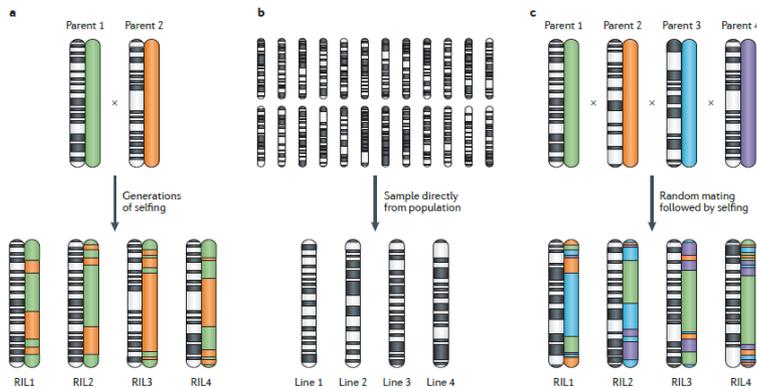
John E. Bowers,* Colette Abbey,† Sharon Anderson,† Charlene Chang,† Xavier Draye,†
 Alison H. Hoppe,† Russell Jessup,† Cornelia Lemke,* Jennifer Lenington,†
 Zhikang Li,† Yann-rong Lin,† Sin-chieh Liu,† Lijun Luo,† Barry S. Marler,*
 Reiguang Ming,† Sharon E. Mitchell,‡ Dou Qiang,† Kim Reischmann,†
 Stefan R. Schulze,* D. Neil Skinner,* Yue-wen Wang,†
 Stephen Kresovich,† Keith F. Schertz†
 and Andrew H. Paterson*^{†,1}

Use de 2050 sondas de RFLP



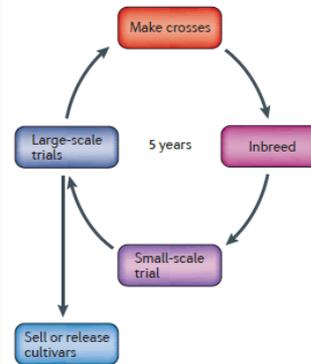
Crop genomics: advances and applications

Peter L. Morrell¹, Edward S. Buckler² and Jeffrey Ross-Ibarra³

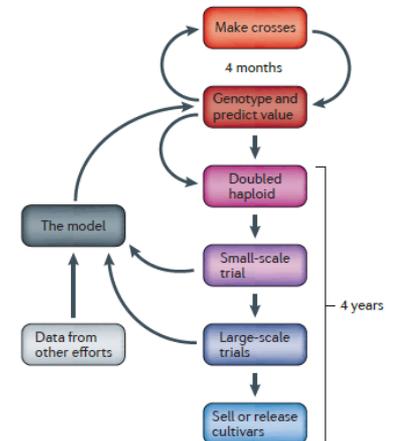


Box 2 | Genomic selection

Standard breeding



Genomic selection



AINDA TEM MUITO O QUE FAZER...

Expectativas de genes em humanos:

antes – 60.000–120.000

draft do genoma – 30.000–35.000

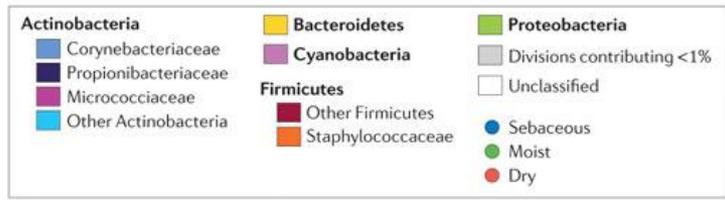
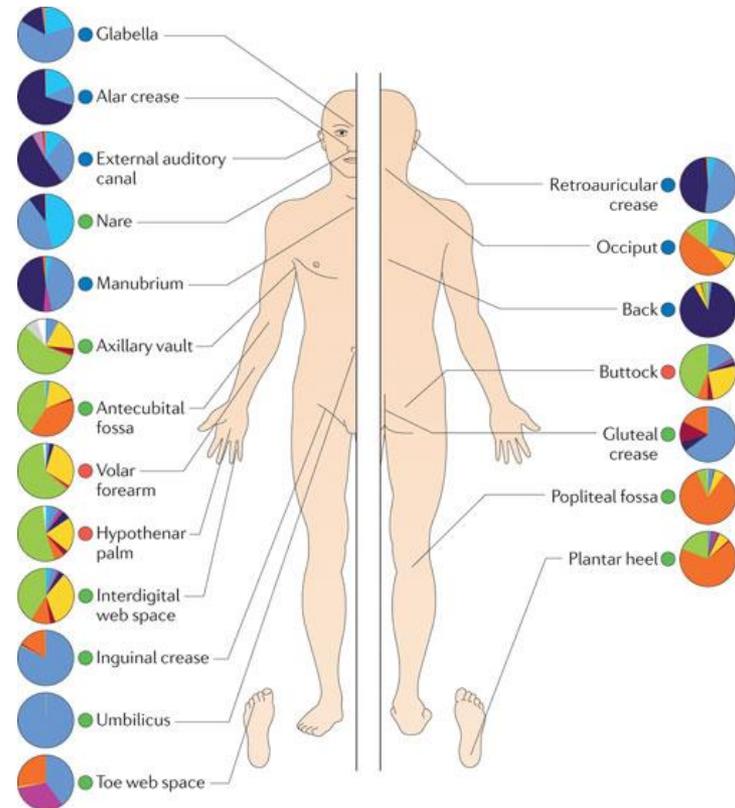
atualmente - 20.000–25.000 na montagem final



QUANTO MAIS SE APRENDE, MAIS PERGUNTAS SURGEM...

- Que tipo de sequências compõem a maioria de regiões não codantes no genoma?
- Como essas sequências são ganhadas e perdidas nos genomas ao longo da evolução?
- Quais efeitos ou talvez funções esse DNA não codante tem no fenótipo?
- Porque alguns genomas são tão enxutos (aves) e outros tão grandes (salamandras)?

O HOMEM NÃO É UMA ILHA...





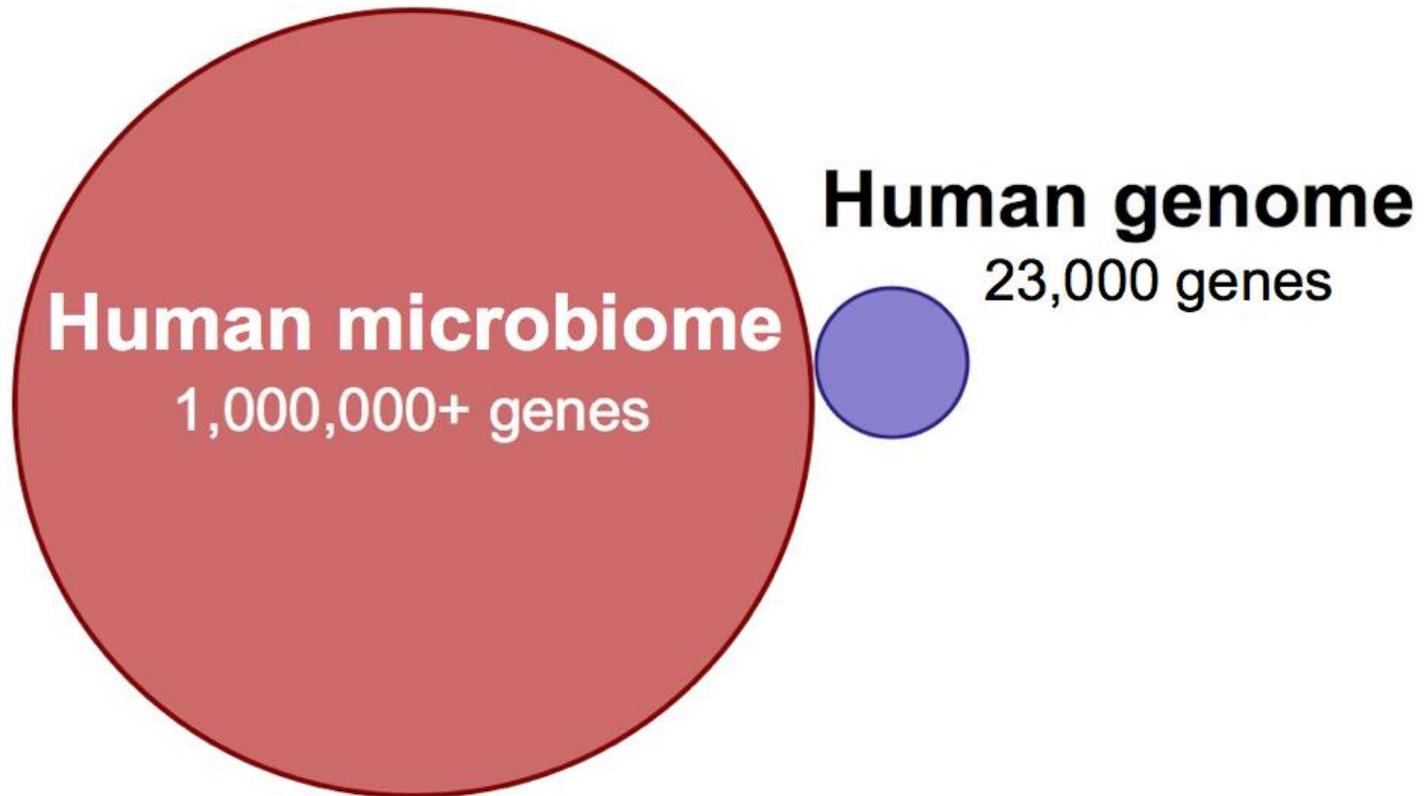
Structure, function and diversity of the healthy human microbiome

The Human Microbiome Project Consortium*

c



Gastrointestinal



**FENÓTIPO = GENÓTIPO +
AMBIENTE?**



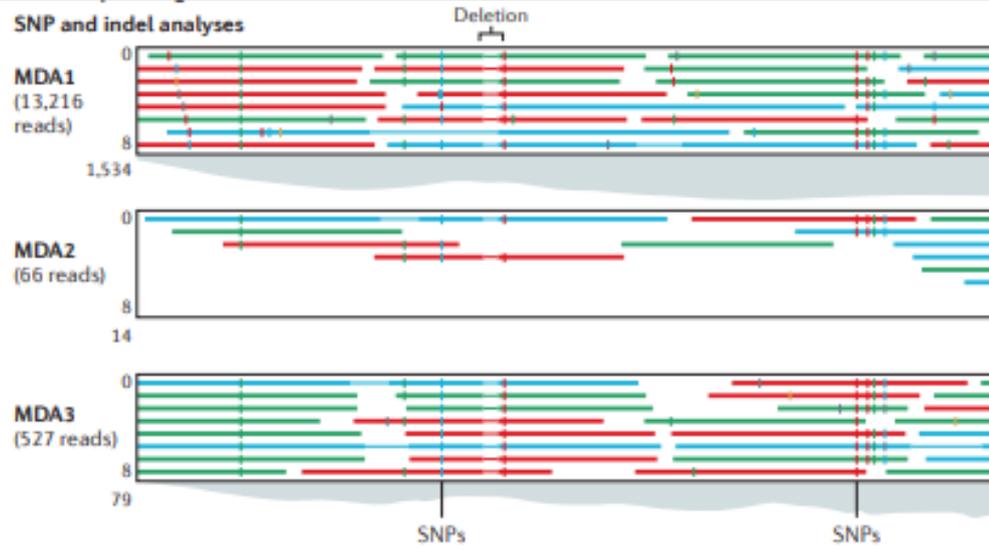
Microbial genome-enabled insights into plant–microorganism interactions

David S. Guttman¹, Alice C. McHardy^{2,3} and Paul Schulze-Lefert^{3,4}

 APPLICATIONS OF NEXT-GENERATION SEQUENCING

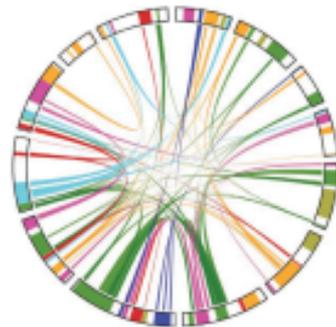
Recent advances in genomic DNA sequencing of microbial species from single cells

Roger S. Lasken and Jeffrey S. McLean

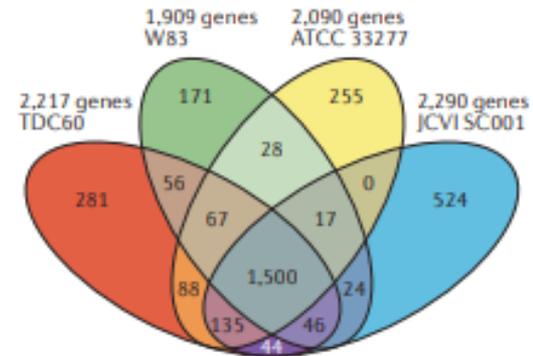


B De novo assembly

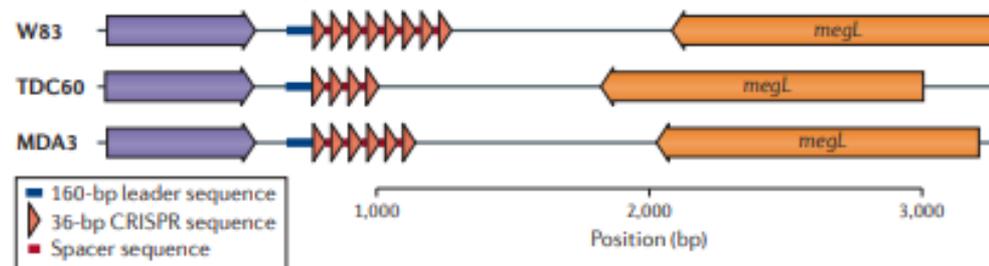
Ba Whole-genome comparisons



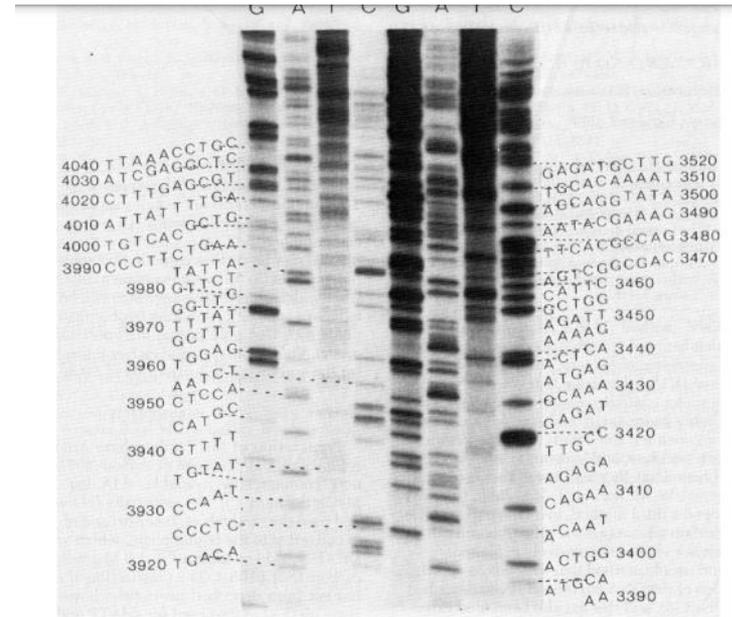
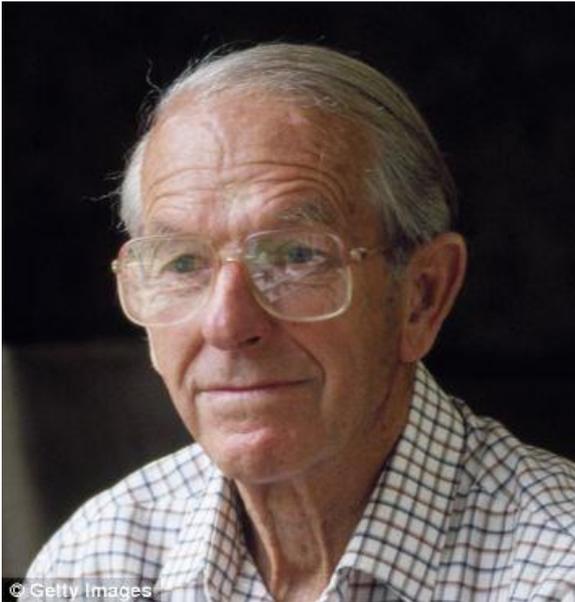
Bb Gene discovery



Bc Structural variation analyses



A CULPA É DO SANGER...



**Nobel de Química
1980**

Proc. Natl. Acad. Sci. USA
Vol. 74, No. 12, pp. 5463–5467, December 1977
Biochemistry

DNA sequencing with chain-terminating inhibitors

(DNA polymerase/nucleotide sequences/bacteriophage ϕ X174)

F. SANGER, S. NICKLEN, AND A. R. COULSON

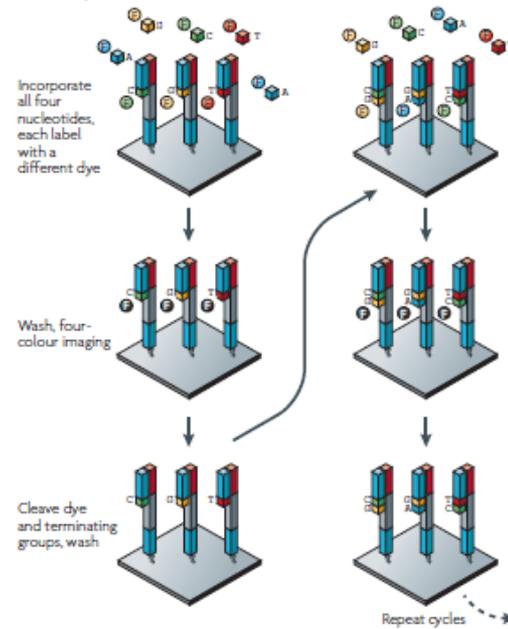
Medical Research Council Laboratory of Molecular Biology, Cambridge CB2 2QH, England

Contributed by F. Sanger, October 3, 1977

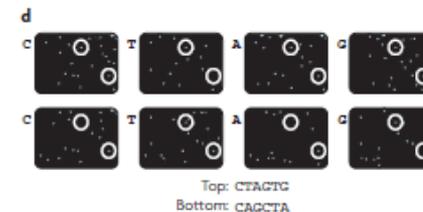
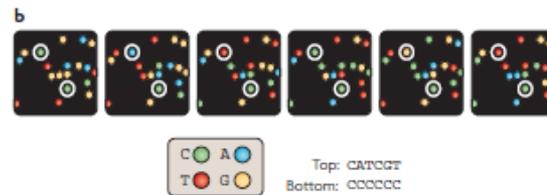
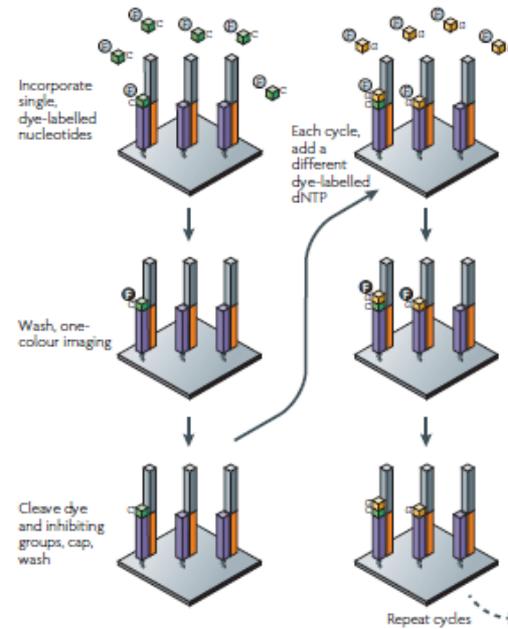
Sequencing technologies — the next generation

Michael L. Metzker*†

a Illumina/Solexa — Reversible terminators



c Helicos BioSciences — Reversible terminators



Repetitive DNA and next-generation sequencing: computational challenges and solutions

Todd J. Treangen¹ and Steven L. Salzberg^{1,2}

Sequencing depth and coverage: key considerations in genomic analyses

David Sims, Ian Sudbery, Nicholas E. Illott, Andreas Heger and Chris P. Ponting

Single-cell genome sequencing: current state of the science

Charles Gawad¹, Winston Koh^{2,3} and Stephen R. Quake^{2,3}

High-throughput functional genomics using CRISPR–Cas9

Ophir Shalem, Neville E. Sanjana and Feng Zhang

APPLICATIONS OF NEXT-GENERATION SEQUENCING

Unravelling the genomic targets of small molecules using high-throughput sequencing

Raphaël Rodriguez^{1,2} and Kyle M. Miller³

Genomes by design

Adrian D. Haimovich^{1,2}, Paul Muir^{1–3} and Farren J. Isaacs^{1,2}

19 de Abril - Aula “Sinalização celular”

Plasmídios e novas tecnologias de clonagem:
Gibson, Gateway e outros