

ACH2024

Aula 24 – Hashing em Disco (parte 3) - Hashing Dinâmico Linear

Profa. Ariane Machado Lima



Aulas passadas

Tratamento de colisões

Estratégias:

A) Hashing estático (tamanho da tabela é constante)

1) Encadeamento ou endereçamento fechado – colisões vão para uma lista ligada

1.1) Encadeamento exterior (fora da tabela)

1.2) Encadeamento interior (dentro da tabela)

2) Endereçamento aberto (chaves dentro da tabela, sem ponteiros)

2.1) Tentativa/Sondagem linear

2.2) Tentativa/Sondagem quadrática

2.3) Dispersão dupla / Hash duplo

B) Hashing dinâmico (tabela pode expandir/encolher)

3) Hashing extensível (estrutura de dados adicional)

4) Hashing linear



Hashing Extensível

Pode-se usar uma função de hash mais uniforme, que mapeie as chaves para um intervalo grande (tipicamente inteiro de 32 bits), pois não será criada inicialmente uma tabela desse tamanho...

Diretório: array de 2^i endereços de buckets \longrightarrow

- i : **profundidade global** do diretório
- Cada posição refere-se aos i bits mais significativos de um valor de hash $h(k)$ \rightarrow todos os registros cujas chaves k possuem valores de hash $h(k)$ com os mesmos i primeiros bits são mapeados para a mesma entrada no diretório
- Cada entrada tem um endereço de bucket que contém tais registros
- A vantagem é que diferentes entradas podem apontar para o mesmo bucket ou não
 - Registros com os mesmos i' primeiros bits, $i' < i$ poderiam caber em um mesmo bucket
 - O valor i' depende de cada bucket b (i'_b), e deve ser armazenado com eles: i'_b – **profundidade local**

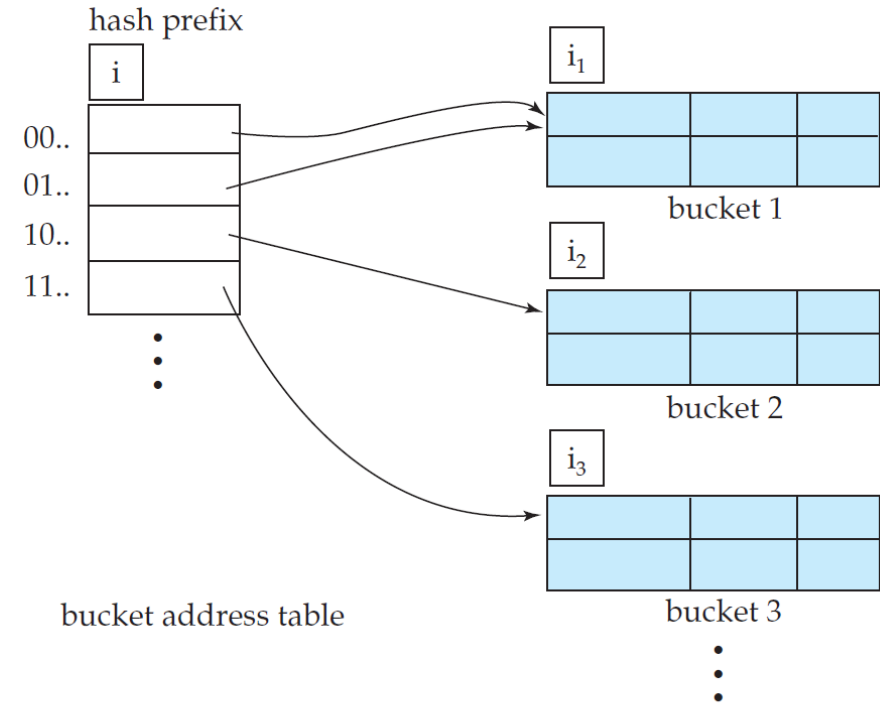


Figure 11.26 General extendable hash structure.

(SILBERSCHATZ, 2011)

Hashing Extensível - Inserção

Inserer(D, k):

Calcula $h(k)$ e pega os i bits mais significativos

Acessa o diretório D **nessa posição** e acessa o bucket aí indicado (b)

se há espaço disponível no bucket b, insere
senão

se $i_b < i$

$e \leftarrow i_b$ bits mais significativos de $h(k)$

$D[e0] \leftarrow$ novo bucket adicional b'

para cada chave k' do bucket apontado por $D[e1]$

move para b' se (i_b+1) -ésimo bit de $h(k') = 0$

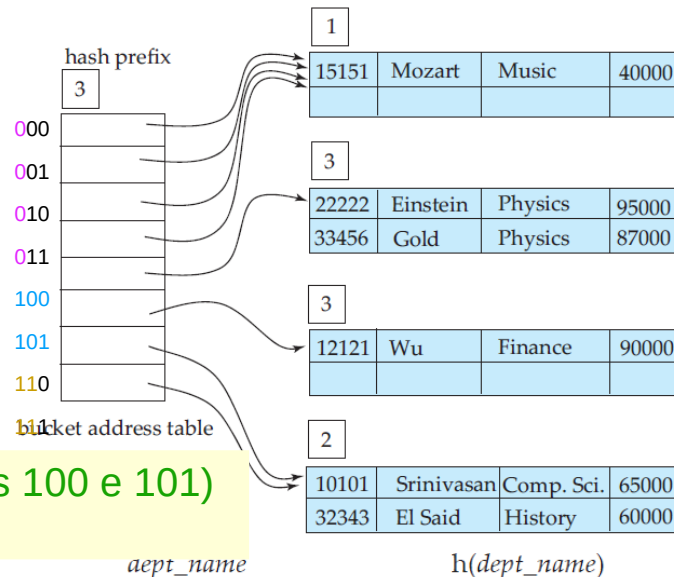
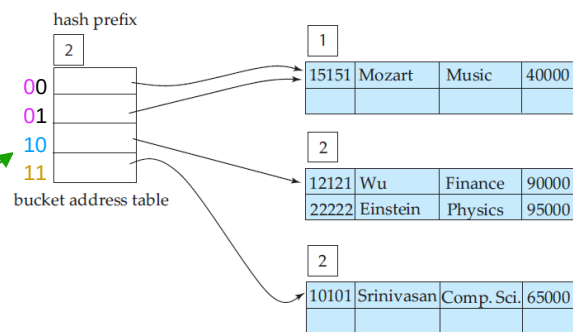
insere(D,k)

senão /* $i_b = i$, ex: 10 */

se b não tiver todas as chaves com os mesmos

$i+1$ primeiro bits

dobra(D) e Insere(D, k)



2 seeks (blocos 100 e 101)
1 novo bloco

dept_name	h(dept_name)
Biology	0010 1101 1111 1011 0010 1100 0011 0000
Comp. Sci.	1111 0001 0010 0100 1001 0011 0110 1101
Elec. Eng.	0100 0011 1010 1100 1100 0110 1101 1111
Finance	1010 0011 1010 0000 1100 0110 1001 1111
History	1100 0111 1110 1101 1011 1111 0011 1010
Music	0011 0101 1010 0110 1100 1001 1110 1011
Physics	1000 0011 1111 1001 1100 0000 0001



Hashing Extensível - Remoção

Remove(D, k):

Calcula $h(k)$ e pega os i bits mais significativos /* ex: 100 */ →

Acessa o diretório D nessa posição e acessa o bucket aí indicado (b)

se está no bucket b principal ou de overflow

remove, traz uma chave do bucket de overflow (se houver)

$e \leftarrow (i_b - 1)$ bits mais significativos de $h(k)$

Se a soma do nr de registros nos blocos apontados por $D[e_0]$ e $D[e_1]$ couber em um único bloco

se $|D[e_0]| + |D[e_1]| \leq r$

acrescenta as chaves do bloco apontado por $D[e_0]$ no bloco apontado por $D[e_1]$

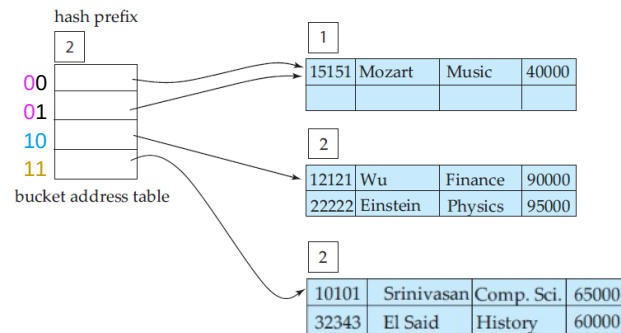
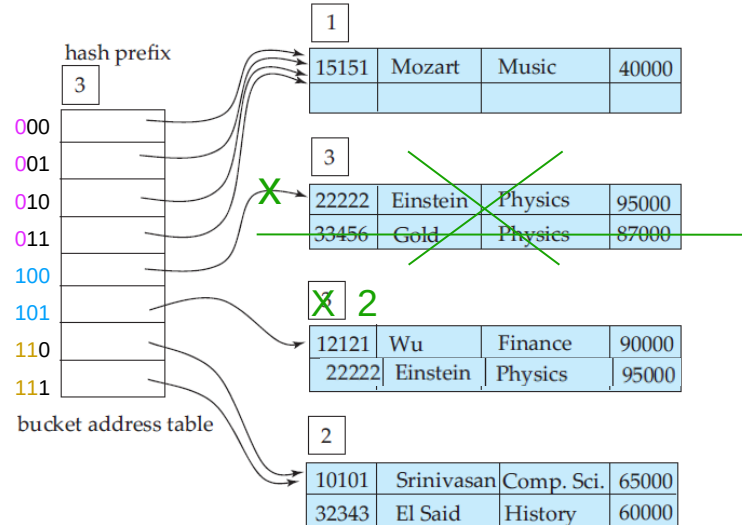
libera bloco apontado por $D[e_0]$

$D[e_0] \leftarrow D[e_1]$

decrementa i' do bloco apontado por $D[e_1]$

se $\max(i') < i$

Divide D pela metade



dept_name	$h(\text{dept_name})$
Biology	0010 1101 1111 1011 0010 1100 0011 0000
Comp. Sci.	1111 0001 0010 0100 1001 0011 0110 1101
Elec. Eng.	0100 0011 1010 1100 1100 0110 1101 1111
Finance	1010 0011 1010 0000 1100 0110 1001 1111
History	1100 0111 1110 1101 1011 1111 0011 1010
Music	0011 0101 1010 0110 1100 1001 1110 1011
Physics	1000 0011 1111 1001 1100 0000 0001

Até 2 seeks
Até 1 bloco a menos

Aula de hoje

Estratégias:

A) Hashing estático (tamanho da tabela é constante)

1) Encadeamento ou endereçamento fechado – colisões vão para uma lista ligada

1.1) Encadeamento exterior (fora da tabela)

1.2) Encadeamento interior (dentro da tabela)

2) Endereçamento aberto (chaves dentro da tabela, sem ponteiros)

2.1) Tentativa/Sondagem linear

2.2) Tentativa/Sondagem quadrática

2.3) Dispersão dupla / Hash duplo

B) Hashing dinâmico (tabela pode expandir/encolher)

3) Hashing extensível (estrutura de dados adicional)

4) Hashing linear



Hashing Linear

Hashing Linear (dinâmico)

- M buckets e função de hash $h_1: x \rightarrow \{0, \dots, M-1\}$
- Colisões que causarem overflow vão para uma lista ligada de registros em buckets de overflow (até agora parece igual ao Hashing Estático..., mas as semelhanças param aqui)
- Uso de um contador (n) de overflows (colisões em buckets lotados)
- À medida que overflows forem ocorrendo (**em quaisquer buckets**), vou partindo em dois os buckets 0, 1, 2, ... **linearmente...**

Hashing Linear - Ideia

Suponha M buckets (0 a $M-1$) e $h_1(k) = k \bmod M$

- $n = 0$
- Primeiro overflow em QUALQUER bucket →
 - Aloca bucket M
 - Divide registros do bucket 0 entre os buckets 0 e M de acordo com $h_2(k) = k \bmod 2M$
 - $n = 1$
- Segundo overflow em QUALQUER bucket →
 - Aloca bucket $M+1$
 - Divide registros do bucket 1 entre os buckets 1 e $M+1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n = 2$
- n -ésimo overflow:
 - Aloca bucket $M+n-1$
 - Divide registros do bucket $n-1$ entre os buckets $n-1$ e $M+n-1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n++$
- (continua no próximo slide...)
- Mas nesse meio tempo... para $n < M$, como faço a busca?
 - Busca(k): $h_1(k) < n$? Se sim, use $h_2(k)$ para saber em qual bucket está, senão use $h_1(k)$

$M = 4; n = 0; h_1(k) = k \bmod 4$

Bucket#	Primary pages	Overflow pages				
0	<table border="1"><tr><td>4</td><td>8</td><td>12</td><td>16</td></tr></table>	4	8	12	16	
4	8	12	16			
1	<table border="1"><tr><td>1</td><td>5</td><td></td><td></td></tr></table>	1	5			
1	5					
2	<table border="1"><tr><td>6</td><td>10</td><td>22</td><td></td></tr></table>	6	10	22		
6	10	22				
3	<table border="1"><tr><td>3</td><td>7</td><td>15</td><td>19</td></tr></table>	3	7	15	19	
3	7	15	19			

DOI: https://doi.org/10.1007/978-0-387-39940-9_742

Hashing Linear - Ideia

Suponha M buckets (0 a $M-1$) e $h_1(k) = k \bmod M$

- $n = 0$
- Primeiro overflow em QUALQUER bucket →
 - Aloca bucket M
 - Divide registros do bucket 0 entre os buckets 0 e M de acordo com $h_2(k) = k \bmod 2M$
 - $n = 1$
- Segundo overflow em QUALQUER bucket →
 - Aloca bucket $M+1$
 - Divide registros do bucket 1 entre os buckets 1 e $M+1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n = 2$
- n -ésimo overflow:
 - Aloca bucket $M+n-1$
 - Divide registros do bucket $n-1$ entre os bucket. $n-1$ e $M+n-1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n++$
- (continua no próximo slide...)
- Mas nesse meio tempo... para $n < M$, como faço a busca?
 - Busca(k): $h_1(k) < n$? Se sim, use $h_2(k)$ para saber em qual bucket está, senão use $h_1(k)$

Ex: $n=0 \rightarrow 1$ inserindo 11

$M = 4; n = 0; h_1(k) = k \bmod 4$

Bucket#	Primary pages	Overflow pages
0	4 8 12 16	
1	1 5	
2	6 10 22	
3	3 7 15 19	

$M = 4; n = 1; h_1(k) = k \bmod 4$
 $h_2(k) = k \bmod 8$

Bucket#	Primary pages	Overflow pages
0	8 16	
1	1 5	
2	6 10 22	
3	3 7 15 19	11
4	4 12	

DOI: https://doi.org/10.1007/978-0-387-39940-9_742

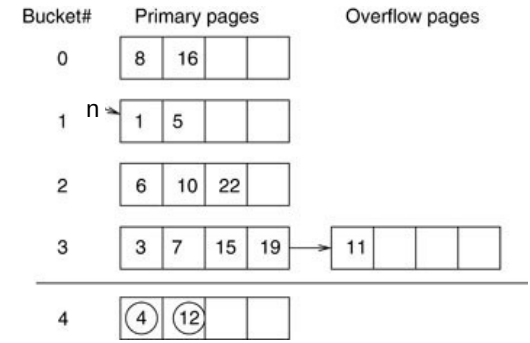
Hashing Linear - Ideia

Suponha M buckets (0 a $M-1$) e $h_1(k) = k \bmod M$

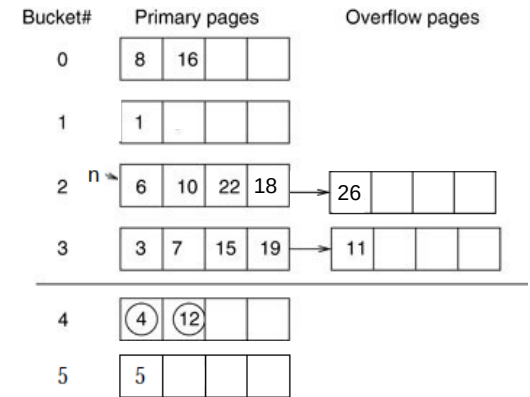
- $n = 0$
- Primeiro overflow em QUALQUER bucket →
 - Aloca bucket M
 - Divide registros do bucket 0 entre os buckets 0 e M de acordo com $h_2(k) = k \bmod 2M$
 - $n = 1$
- Segundo overflow em QUALQUER bucket →
 - Aloca bucket $M+1$
 - Divide registros do bucket 1 entre os buckets 1 e $M+1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n = 2$
- n -ésimo overflow:
 - Aloca bucket $M+n-1$
 - Divide registros do bucket $n-1$ entre os buck. $n-1$ e $M+n-1$ de acordo c/ $h_2(k) = k \bmod 2M$
 - $n++$
- (continua no próximo slide...)
- Mas nesse meio tempo... para $n < M$, como faço a busca?
 - Busca(k): $h_1(k) < n$? Se sim, use $h_2(k)$ para saber em qual bucket está, senão use $h_1(k)$

Ex: $n_1 \rightarrow 2$ inserindo 18, 26

$$M = 4; n = 1; \begin{aligned} h_1(k) &= k \bmod 4 \\ h_2(k) &= k \bmod 8 \end{aligned}$$



$$M = 4; n = 2; \begin{aligned} h_1(k) &= k \bmod 4 \\ h_2(k) &= k \bmod 8 \end{aligned}$$



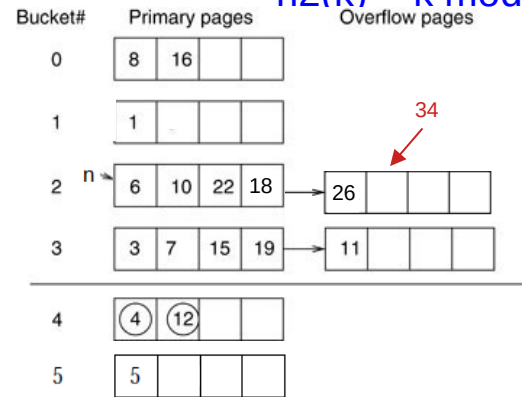
Hashing Linear - Ideia

Suponha M buckets (0 a $M-1$) e $h_1(k) = k \bmod M$

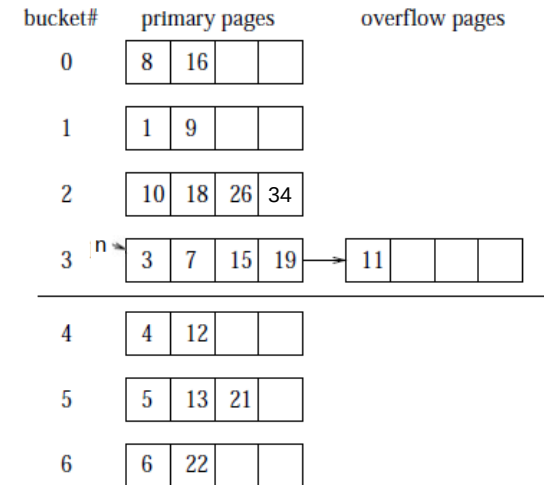
- $n = 0$
- Primeiro overflow em QUALQUER bucket →
 - Aloca bucket M
 - Divide registros do bucket 0 entre os buckets 0 e M de acordo com $h_2(k) = k \bmod 2M$
 - $n = 1$
- Segundo overflow em QUALQUER bucket →
 - Aloca bucket $M+1$
 - Divide registros do bucket 1 entre os buckets 1 e $M+1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n = 2$
- **n -ésimo overflow:**
 - Aloca bucket $M+n-1$
 - Divide registros do bucket $n-1$ entre os buck. $n-1$ e $M+n-1$ de acordo c/ $h_2(k) = k \bmod 2M$
 - $n++$
- (continua no próximo slide...)
- Mas nesse meio tempo... para $n < M$, como faço a busca?

Ex: $n=2 \rightarrow 3$ inserindo 9, 13, 21, 34

$M = 4; n = 2; h_1(k) = k \bmod 4$
 $h_2(k) = k \bmod 8$



$M = 4; n = 3; h_1(k) = k \bmod 4$
 $h_2(k) = k \bmod 8$



Hashing Linear - Ideia

Suponha M buckets (0 a $M-1$) e $h_1(k) = k \bmod M$

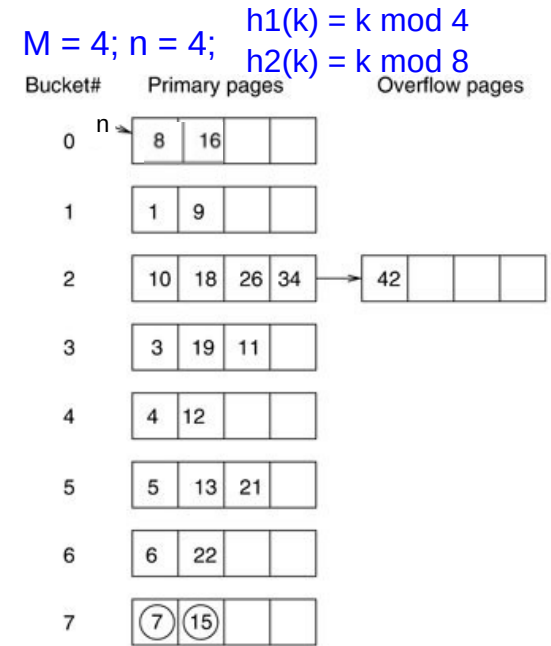
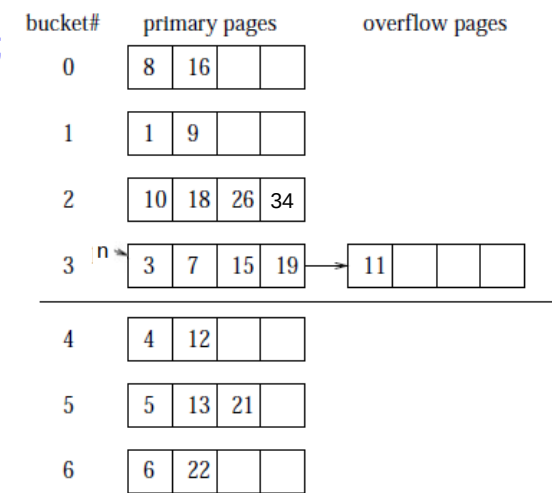
- $n = 0$
- Primeiro overflow em QUALQUER bucket →
 - Aloca bucket M
 - Divide registros do bucket 0 entre os buckets 0 e M de acordo com $h_2(k) = k \bmod 2M$
 - $n = 1$
- Segundo overflow em QUALQUER bucket →
 - Aloca bucket $M+1$
 - Divide registros do bucket 1 entre os buckets 1 e $M+1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n = 2$
- **n -ésimo overflow:**
 - Aloca bucket $M+n-1$
 - Divide registros do bucket $n-1$ entre os buckets $n-1$ e $M+n-1$ de acordo com $h_2(k) = k \bmod 2M$
 - $n++$
- (continua no próximo slide...)
- Mas nesse meio tempo... para $n < M$, como faço a busca?
 - Busca(k): $h_1(k) < n$? Se sim, use $h_2(k)$ para saber em qual bucket está, senão use $h_1(k)$

Ex: $n=3 \rightarrow 4$ inserindo 42

$$M = 4; n = 3;$$

$$h_1(k) = k \bmod 4$$

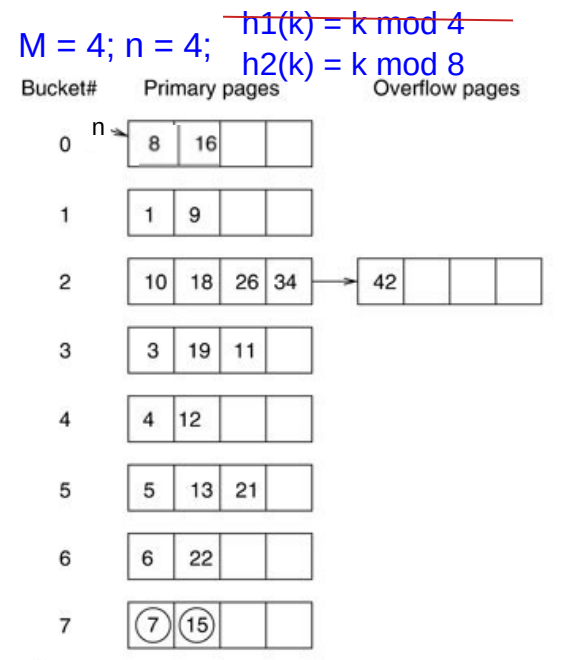
$$h_2(k) = k \bmod 8$$



Hashing Linear

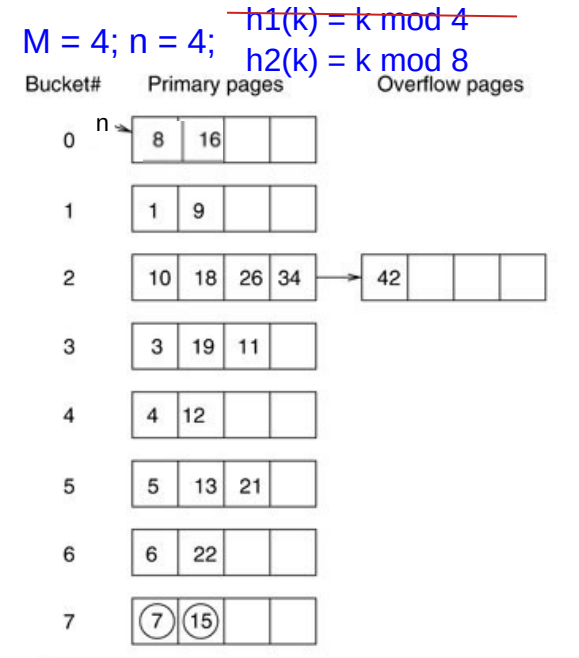
- Processo continua até que $n = M$:
 - todos os M buckets foram divididos
 - Tabela de Hash tem tamanho = $2M$
 - h_1 não é mais necessária
 - $n \leftarrow 0$ e recomeça nova etapa de divisões
 - Funções de hash ativas:
 - $h_2(k) = k \bmod 2M$
 - $h_3(k) = k \bmod 4M$
- ...

- Se $d = nr$ de vezes que a tabela foi inteiramente dobrada:
 - Funções de hash ativas: $h_{d+1}(k)$ e $h_{d+2}(k)$, sendo $h_j(k) = k \bmod 2^{j-1}M$
 - Busca(k):



Hashing Linear

- Processo continua até que $n = M$:
 - todos os M buckets foram divididos
 - Tabela de Hash tem tamanho = $2M$
 - h_1 não é mais necessária
 - $n \leftarrow 0$ e recomeça nova etapa de divisões
 - Funções de hash ativas:
 - $h_2(k) = k \bmod 2M$
 - $h_3(k) = k \bmod 4M$
 - Processo continua até que $n = 2M$
- ...
- Se $d = nr$ de vezes que a tabela foi inteiramente dobrada:



- Funções de hash ativas: $h_{d+1}(k)$ e $h_{d+2}(k)$, sendo $h_j(k) = k \bmod 2^{j-1}M$
- Busca(k): $h_{d+1}(k) < n$? Se sim, use $h_{d+2}(k)$ para saber em qual bucket está, senão use $h_{d+1}(k)$

Hashing Linear

- Estratégia alternativa para controlar dinamismo (ao invés de dividir a cada colisão):
 - Fator de carga $\alpha = N/(b*r)$, N = nr de registros, b = nr de buckets, r = número de registros que cabem em um bucket
 - Definir intervalo aceitável do fator de carga (ex: $0.7 \leq \alpha \leq 0.9$)
 - Se α acima do limite superior, divide buckets
 - Se α abaixo do limite inferior, junta buckets (em um processo inverso, decrementando n)

Hash x árvores B+: diferenças

main
buckets

bucket 0

340	
460	
	record pointer

bucket 1

321	
761	
91	
	record pointer

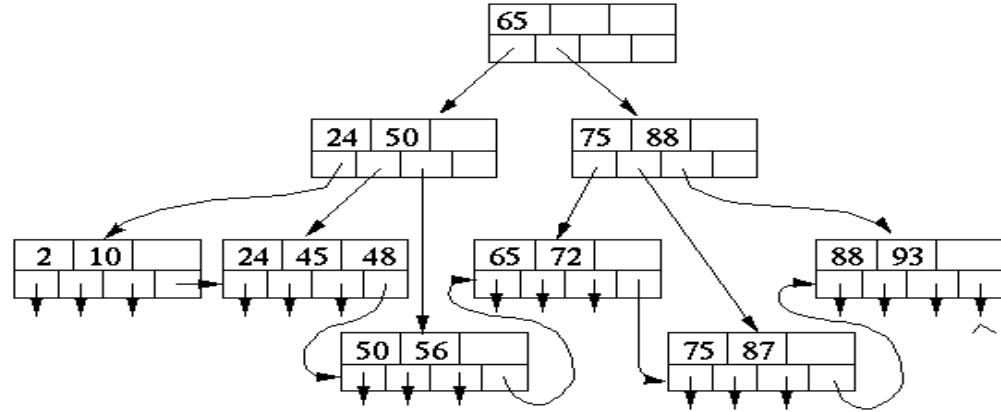
bucket 2

22	
72	
522	
	record pointer

⋮

bucket 9

399	
89	
	record pointer



Hash x árvores B+

- Hash é MUITO rápido independente do tamanho do arquivo (árvore B+ depende da altura)
- Busca ordenada no hash em geral muito ruim
 - Hash mais apropriado para dados tipo dicionário (<chave, valor>, sendo cada entrada “independente” das demais)
- Se precisar fazer um índice para um campo:
 - Busca apenas por valores específicos (utilizando apenas comparação de igualdade)?
 -
 - Precisa ordenar os dados ou buscar conjunto de valores de acordo com uma relação entre os mesmos (<, <=, >, >=)?
 -

Hash x árvores B+

- Hash é MUITO rápido independente do tamanho do arquivo (árvore B+ depende da altura)
- Busca ordenada no hash em geral muito ruim
 - Hash mais apropriado para dados tipo dicionário (<chave, valor>, sendo cada entrada “independente” das demais)
- Se precisar fazer um índice para um campo:
 - Busca apenas por valores específicos (utilizando apenas comparação de igualdade)?
 - Use hash
 - Precisa ordenar os dados ou buscar conjunto de valores de acordo com uma relação entre os mesmos (<, <=, >, >=)?
 - Use árvore B+

Exercícios

Implemente todas essas opções de Hashing (operações de busca, inserção e deleção)

Referências

ELMARIS, R.; NAVATHE, S. B. Fundamentals of Database Systems. 4 ed. Ed. Pearson-Addison Wesley. Cap 13.8. 4 ed. Pearson. 2004

SILBERSCHATZ, A.; KORTH, H. F.; SUDARSHAN, S. Database System Concepts, 6. ed. McGraw Hill, 2011.