

Tabelas Hash

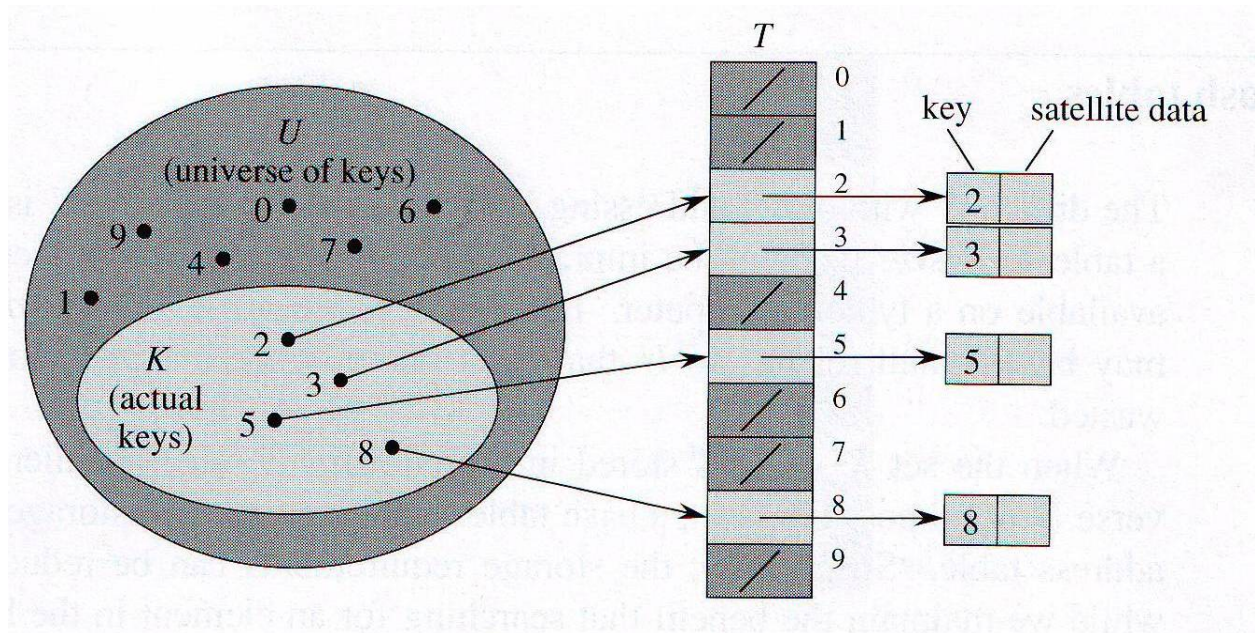
Referência: Cormen et al, Introduction to Algorithms

- Conjuntos que podem ser alterados (aumentados, diminuídos, etc) são chamados **Conjuntos Dinâmicos**.
- Em uma implementação típica de um conjunto dinâmico, cada elemento é representado por um objeto cujos campos podem ser examinados e manipulados, através de um ponteiro para o objeto.
 - Ex: listas ligadas (ordenadas ou não).
- Usualmente, conjuntos dinâmicos assumem que um dos campos do objeto é um campo de **chave** de identificação. Se as chaves são todas diferentes, podemos pensar o conjunto dinâmico como um conjunto de valores de chaves. Os demais campos são considerados **dados satélite**.
- Dicionário: Conjunto dinâmico que suporta as operações de inserção, remoção e busca de elementos no conjunto.
 - Por exemplo, um compilador para uma linguagem de programação mantém uma tabela de símbolos, na qual as chaves dos elementos são strings de caracteres arbitrárias que correspondem aos identificadores na linguagem.

- TAD Dicionário:

- $\text{SEARCH}(S, k)$: Uma consulta que, dado um conjunto S e um valor de chave k , retorna um ponteiro x para um elemento em S tal que $\text{key}[x] = k$, ou NIL se não houver nenhum elemento com esse valor.
- $\text{INSERT}(S, x)$: Insere no conjunto S o elemento apontado por x .
- $\text{DELETE}(S, x)$: Dado um ponteiro x para um elemento do conjunto S , remove x de S .

Tabelas de acesso/endereçamento direto



- Acesso direto: técnica simples que funciona bem quando o universo de chaves U é razoavelmente pequeno.
- Suponha que cada elemento possui uma chave pertencente ao universo $U = \{0, 1, \dots, m - 1\}$, onde m não é muito grande. Assume-se que não haja dois elementos com a mesma chave.

- Para representar o conjunto dinâmico, utiliza-se um arranjo, ou **tabela de endereçamento direto**, denotado por $T[0 \dots m - 1]$, onde cada posição (**slot**), corresponde a uma chave no universo U .
- Operações (Custo de tempo $O(1)$):

Direct-Address-Search(T, k)
 devolva($T[k]$)

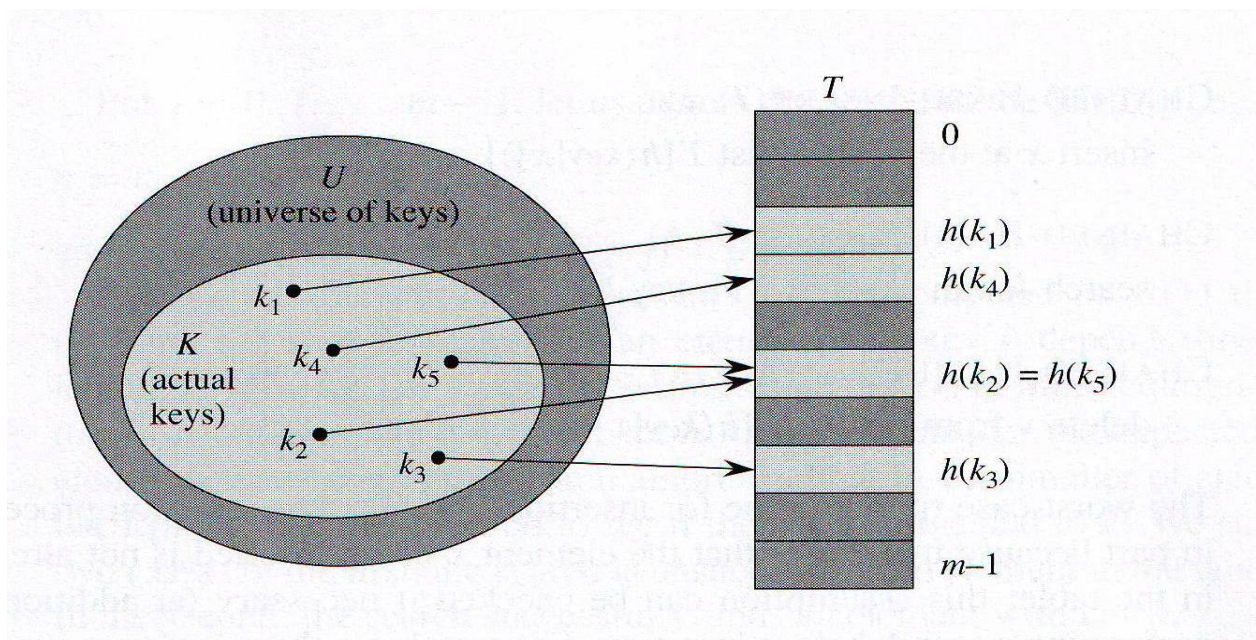
Direct-Address-Insert(T, x)
 $T[\text{key}[x]] := x$

Direct-Address-Delete(T, x)
 $T[\text{key}[x]] := \text{NIL}$

- Para algumas aplicações, pode-se usar o próprio arranjo T para armazenar o dado satélite; a chave não precisa ser guardada, uma vez que corresponde à posição em T . Cuidado para indicar quando a posição está vazia.

Tabelas Hash

- Se o universo U é grande, pode ser impraticável (ou impossível) armazenar uma tabela T de tamanho U .
- Além disso, se o conjunto K de chaves *efetivamente armazenadas* for muito menor do que o universo U , a maior parte do espaço de memória alocada para T será desperdiçada.
- Idéia: Se $|U| \gg |K|$, pode-se reduzir a necessidade de memória para $\Theta(|K|)$, mantendo o tempo (*médio*) das operações $O(1)$.
- Enquanto no acesso direto um elemento com chave k é armazenado na posição k , na tabela Hash esse elemento é armazenado na posição $h(k)$, onde h é uma *função de hash* que calcula a posição a partir do valor da chave, k .



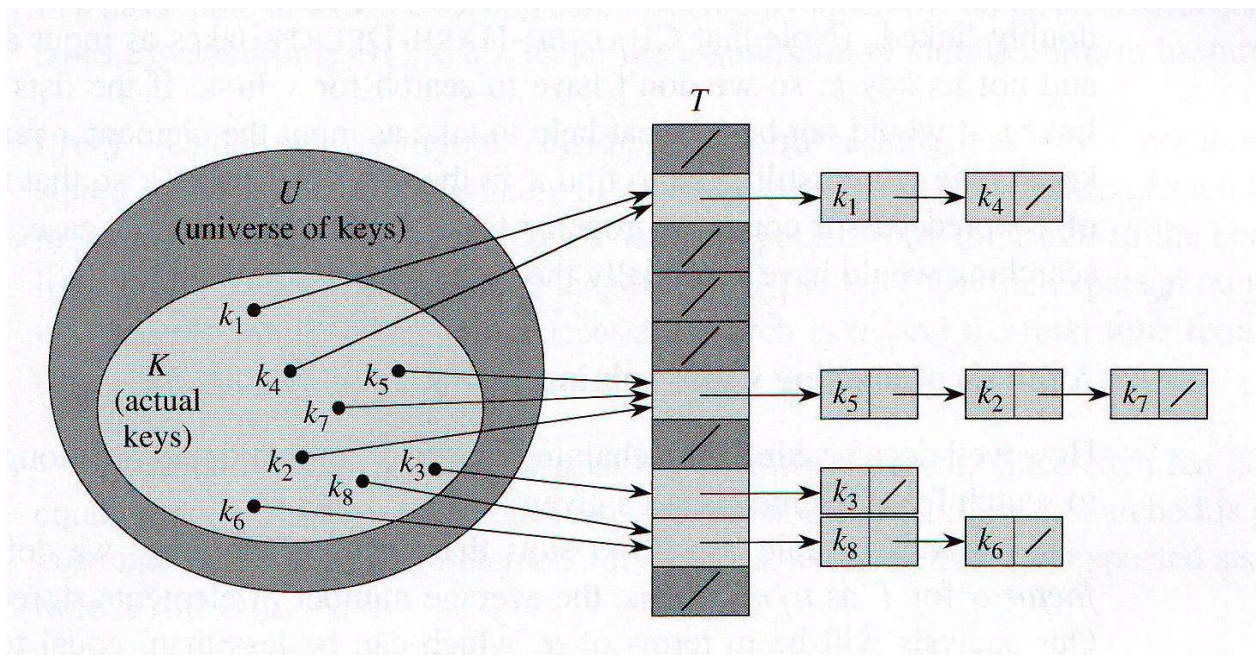
- h mapeia o universo U de chaves nas posições de uma tabela Hash $T[0 \dots m - 1]$:

$$h : U \rightarrow \{0, 1, \dots, m - 1\}.$$

- $h(k)$ não preserva a ordem das chaves.
- **Colisão:** situação em que duas chaves k_1 e k_2 têm o mesmo valor de hashing, ou seja, $h(k_1) = h(k_2)$. Escolha de funções h apropriadas tentam minimizar a probabilidade de ocorrência de colisões.
- Tratamento de colisões:
 - *Resolução por encadeamento* (Hashing aberto)
 - *Endereçamento aberto* (Hashing fechado)

Resolução de colisões por encadeamento (Hashing aberto)

- Na resolução por encadeamento, colocamos todos os elementos com mesmo valor de h em uma lista ligada. Na tabela hash T , cada posição j armazena o ponteiro para o início da lista ligada dos elementos x tais que $h(key[x]) = j$.



- Operações:

`Chained-Hash-Search(T,k)`

procure um elemento com chave k
na lista $T[h(k)]$ e devolva seu ponteiro

`Chained-Hash-Insert(T,x)`

insira x na cabeça da lista $T[h(key[x])]$

`Chained-Hash-Delete(T,x)`

remova x da lista $T[h(key[x])]$

Análise da resolução por encadeamento

- Quanto custa (em média) a procura de um elemento na tabela hash com listas encadeadas?
- Dada uma tabela hash T com m slots armazenando n chaves, definimos o **fator de carga** α para T como $\alpha = n/m$ (número médio de elementos armazenados em cada posição de T).
- Vamos assumir que qualquer elemento tem a mesma probabilidade de ser colocado em qualquer dos m slots, ou seja, $Pr(h(k) = j) = 1/m, j = 0 \dots m - 1$. **Hashing uniforme simples**
- Para $j = 0, 1 \dots m - 1$, denotamos o tamanho da lista $T[j]$ por n_j , de forma que $n = n_0 + n_1 + \dots + n_{m-1}$. Sob a hipótese de hashing uniforme simples, o valor médio de n_j é:

$$E(n_j) = n/m = \alpha.$$

- Assume-se que $h(k)$ é calculado em tempo $O(1)$, e portanto o tempo requerido para procurar um elemento com chave k é linear no comprimento $n_{h(k)}$ da lista $T[h(k)]$. Deve-se então analisar quantos elementos da lista $T[h(k)]$ são acessados, em média, para procurar k .
- **Teorema:** Em uma tabela hash com resolução de colisões por encadeamento, assumindo-se hashing uniforme simples (h.u.s):

a) O tempo esperado para uma busca sem sucesso é $\Theta(1 + \alpha)$.

Dem: Sob a hipótese de h.u.s, qualquer chave k não armazenada na tabela tem probabilidade $1/m$ de ser posicionada em qualquer dos m slots. O tempo esperado para buscar uma chave k sem sucesso é o tempo esperado para percorrer toda a lista $T[h(k)]$, que é $E[n_{h(k)}] = \alpha$. Logo, o número esperado de elementos examinados em uma busca sem sucesso é α , e o tempo total requerido (incluindo o cálculo de $h(k)$) é $\Theta(1 + \alpha)$.

b) O tempo esperado para uma busca bem sucedida é $\Theta(1 + \alpha)$.

Idéia da Demonstração: Seja x o elemento sendo procurado, e $k = \text{key}[x]$. Assumimos que x é igualmente provável de ser qualquer um dos n elementos armazenados na tabela. O número de elementos examinados durante uma busca bem sucedida é 1 mais o número de elementos que aparecem antes de x em sua respectiva lista. Como a posição esperada de x dentro de sua lista é $(n_{h(k)} + 1)/2$ e como $E[n_{h(k)}] = \alpha$, então o número esperado de elementos examinados na busca de x será $1 + (\alpha + 1)/2$, e o tempo total será $\Theta(1 + \alpha/2) = \Theta(1 + \alpha)$.

Endereçamento aberto (Hashing fechado)

- No endereçamento aberto, todos os elementos são armazenados na própria tabela hash. Cada entrada da tabela contém ou um elemento do conjunto dinâmico ou NIL.
- Em endereçamento aberto, a tabela hash pode encher, de forma que não seja mais possível inserir elementos.
- Vantagem: eliminação de ponteiros / eficiência do uso de espaço / implementação em disco
- Para realizar uma inserção, examinamos/testamos sucessivamente a tabela hash até que encontrarmos um slot vazio no qual a chave será inserida. (*rehashing*)
- Para determinar a sequência de posições a serem examinadas, a função de hashing é estendida de maneira a incluir o número do teste como segunda entrada:

$$h : U \times \{0, 1 \dots m - 1\} \rightarrow \{0, 1 \dots m - 1\}$$

- Requer-se que a **seqüência de teste** (*probe sequence*) $\langle h(k, 0), h(k, 1) \dots h(k, m - 1) \rangle$ seja uma permutação de $\langle 0, 1 \dots m - 1 \rangle$, de forma que toda posição da tabela hash seja eventualmente considerada como um slot para uma nova chave à medida em que a tabela enche.

- Operações:

Hash-Insert(T,k)

 i := 0

 repeat

 j := h(k,i)

 if T[j] = NIL or T[j] = DELETED

 T[j] := k

 return j

 else i := i+1

 until i = m

 error "Estouro da tabela hash"

Hash-Search(T,k)

 i := 0

 repeat

 j := h(k,i)

 if T[j] = k

 return j

 i := i+1

 until T[j] = NIL or i = m

 return NIL

```

Hash-Delete(T,k)
  i := 0
  repeat
    j := h(k,i)
    if T[j] = k
      T[j] := DELETED
    i := i+1
  until T[j] = NIL or i = m

```

- Quando o valor especial DELETED é usado, o tempo de busca não é mais dependente do fator de carga α .

Técnicas de rehashing

- Hipótese de hashing uniforme: assume-se que, para cada chave, sua sequência de teste tem a mesma probabilidade de ser qualquer uma das $m!$ permutações de $\langle 0, 1 \dots m - 1 \rangle$.
- **probing linear**: usa uma função de hashing comum, $h' : U \rightarrow \{0, 1 \dots m - 1\}$.

Dada a chave k , o primeiro slot testado é $T[h'(k)]$. Se estiver ocupado, os próximos slots são $T[h'(k) + 1], T[h'(k) + 2], \dots, T[m - 1], T[0], \dots T[h'(k) - 1]$.

Formalmente,

$$h(k, i) = (h'(k) + i) \bmod m, \quad i = 0, 1 \dots m - 1.$$

- Desvantagem do probing linear: *clustering primário* – longas sequências de slots ocupados

Clusters ocorrem porque um slot vazio precedido por i slots ocupados tem probabilidade $(i+1)/m$ de ser ocupado na próxima inserção.

- **probing quadrático:** função da forma

$$h(k, i) = (h'(k) + c_1i + c_2i^2) \bmod m,$$

onde h' é a função de hashing auxiliar, c_1 e c_2 são constantes auxiliares, e $i = 0, 1, \dots, m - 1$.

O rehashing é deslocado por distâncias que dependem quadraticamente de i .

Para que todas as posições da tabela hash sejam exploradas, os valores de c_1 , c_2 e m são críticos.

- Desvantagem do probing quadrático: *clustering secundário* – se duas chaves k_1, k_2 têm a mesma posição inicial de teste, então sua sequência de rehashing será a mesma, uma vez que $h(k_1, 0) = h(k_2, 0)$ implica $h(k_1, i) = h(k_2, i)$.

Assim como no probing linear, a posição inicial determina a sequência inteira, e assim somente m sequências distintas são usadas.

- **hashing duplo:** função na forma

$$h(k, i) = (h_1(k) + ih_2(k)) \bmod m,$$

onde h_1 e h_2 são funções hash auxiliares.

O primeiro slot testado é $T[h_1(k)]$. Cada posição testada posteriormente é um deslocamento de tamanho $h_2(k)$ em relação à posição anterior, módulo m .

Cada par $(h_1(k), h_2(k))$ determina uma sequência \Rightarrow total de $\Theta(m^2)$ sequências.

- **Escolha de m e h_2 com hashing duplo:**

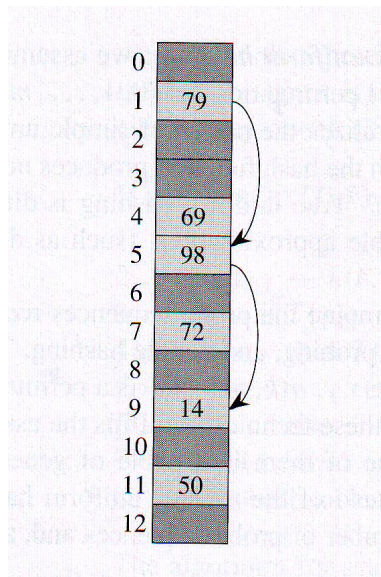
Lema: Se $d \geq 1$ é o máximo divisor comum de m e $h_2(k)$, então uma busca mal sucedida da chave k examinará m/d posições da tabela hash antes de voltar para o slot $h_1(k)$.

Pelo lema acima, idealmente o valor $h_2(k)$ precisa ser primo relativo de m , para que a tabela hash inteira possa ser examinada.

Uma forma possível: toma-se m primo e define-se h_2 que sempre retorna um inteiro positivo menor do que m .

Exemplo: m primo, $h_1(k) = k \bmod m$, $h_2(k) = 1 + (k \bmod m')$ onde m' é ligeiramente menor do que m (p. ex. $m - 1$).

A figura abaixo apresenta um exemplo de inserção por hashing duplo onde $k = 14$, $m = 13$, $h_1(k) = k \bmod 13$, $h_2 = 1 + (k \bmod 11)$.

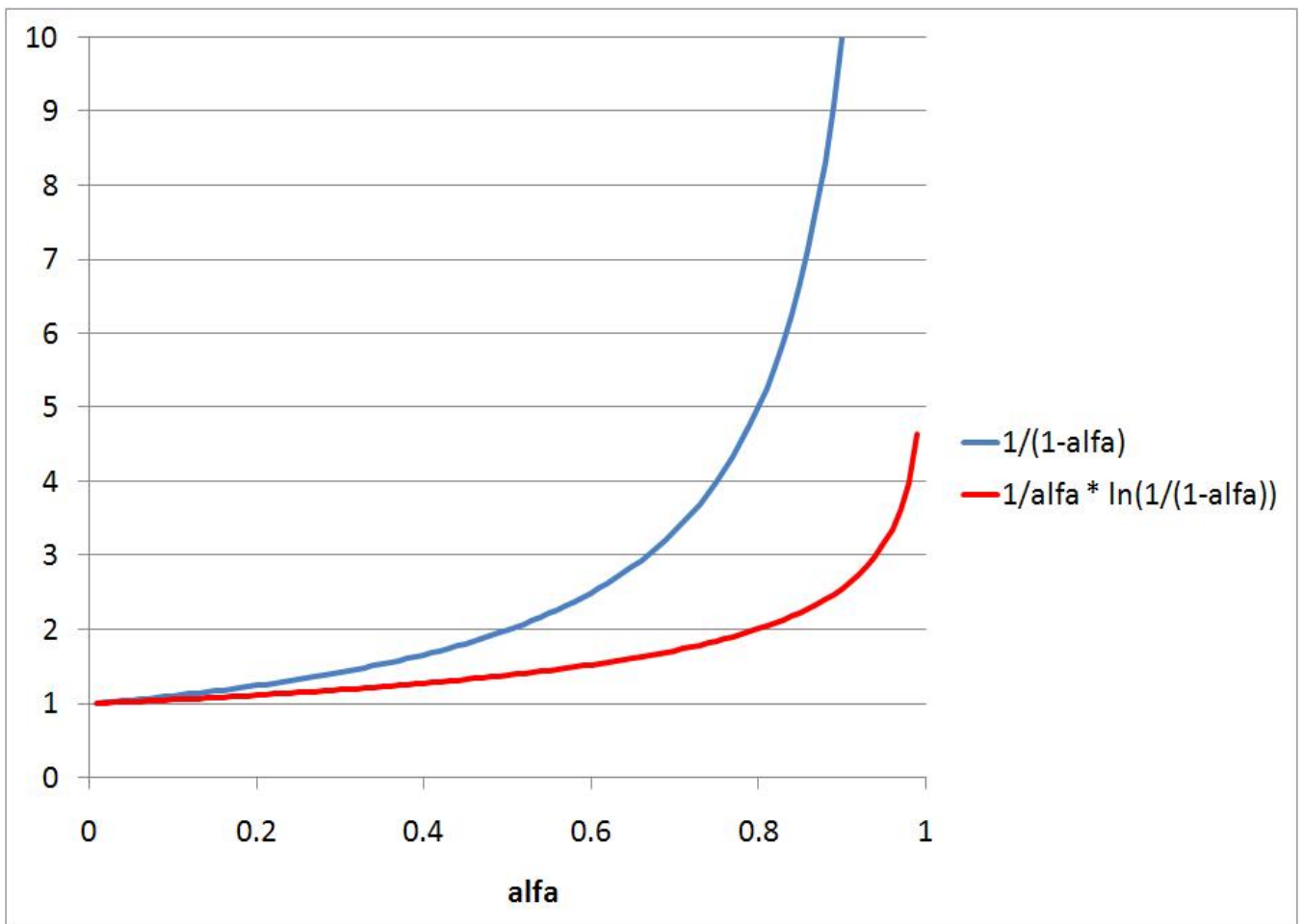


Custo de tempo do hashing de endereçamento aberto

- Análise em termos do fator de carga $\alpha = n/m$ da tabela hash, para valores assintóticos de n, m .

(Note que no hashing fechado tem-se no máximo uma chave por posição, e portanto $n \leq m \Rightarrow \alpha \leq 1$.)

- **Teorema:** Dada uma tabela hash de endereçamento aberto com fator de carga $\alpha = n/m < 1$, o número esperado de slots testados em uma busca sem sucesso é no máximo $1/(1-\alpha)$, assumindo hashing uniforme.
- **Teorema:** Dada uma tabela hash de endereçamento aberto com fator de carga $\alpha = n/m < 1$, o número esperado de slots testados em uma busca bem sucedida é no máximo $1/\alpha \ln[1/(1-\alpha)]$, assumindo hashing uniforme e assumindo que cada chave da tabela tem a mesma probabilidade de ser procurada.



Funções de hash

- Uma boa função de hash deve satisfazer aproximadamente a premissa de hashing uniforme simples: cada chave tem a mesma probabilidade de ser posicionada em qualquer um dos m slots, independentemente da posição indicada para qualquer outra chave.
- *Interpretando chaves como números naturais:*

As funções de hash assumem geralmente que o universo de chaves é o conjunto $\mathbb{N} = \{0, 1, 2, \dots\}$ dos números naturais. Logo, se as chaves não são números naturais, deve-se adotar algum mapeamento do universo de chaves no conjunto dos naturais.

Conversão de cadeias de caracteres para inteiros:

$c_1 c_2 c_3 \dots c_r$ onde c_i é o código do i -ésimo caracter no mapa de caracteres adotado, e cada caracter é codificado com p bits:

$$k = c_1 + c_2 \cdot 2^p + c_3 \cdot 2^{2p} + \dots + c_r \cdot 2^{(r-1)p}.$$

Exemplo: No conjunto de caracteres ASCII, $p = 112$, $t = 116$, e a cadeia pt pode ser convertida para $112 + (116 \cdot 2^7) = 14960$.

- **Método da divisão**

$$h(k) = k \bmod m$$

Restrições sobre m :

- m não pode ser da forma $m = 2^p$ ($p \in \mathbb{N}$), pois nesse caso $h(k)$ corresponderá aos p bits menos significativos de k .
- Se as chaves são cadeias de caracteres convertidas para inteiros conforme método descrito anteriormente, $m = 2^p - 1$ é uma escolha

ruim, pois a permutação dos caracteres de uma cadeia não alteram o valor da função hash:

Suponha duas cadeias c_1c_2 e c_2c_1 , diferindo apenas na permutação de seus caracteres. Logo, $k_1 = c_1 + c_2 \cdot 2^p$ e $k_2 = c_2 + c_1 \cdot 2^p$. Então,

$$\begin{aligned}h(k_1) - h(k_2) &= k_1 \bmod m - k_2 \bmod m \\&= (c_1 + c_2 \cdot 2^p) \bmod (2^p - 1) - (c_2 + c_1 \cdot 2^p) \bmod (2^p - 1) \\&= (c_1 + c_2) \bmod (2^p - 1) - (c_2 + c_1) \bmod (2^p - 1) \\&= 0\end{aligned}$$

- Um número primo não muito próximo de uma potência exata de 2 é usualmente uma boa escolha para m .

Ex: Hashing com encadeamento, $n = 2000$, $m = 701$, $h(k) = k \bmod 701$, $\alpha = 2.85$

• Método da multiplicação

O método opera em dois passos:

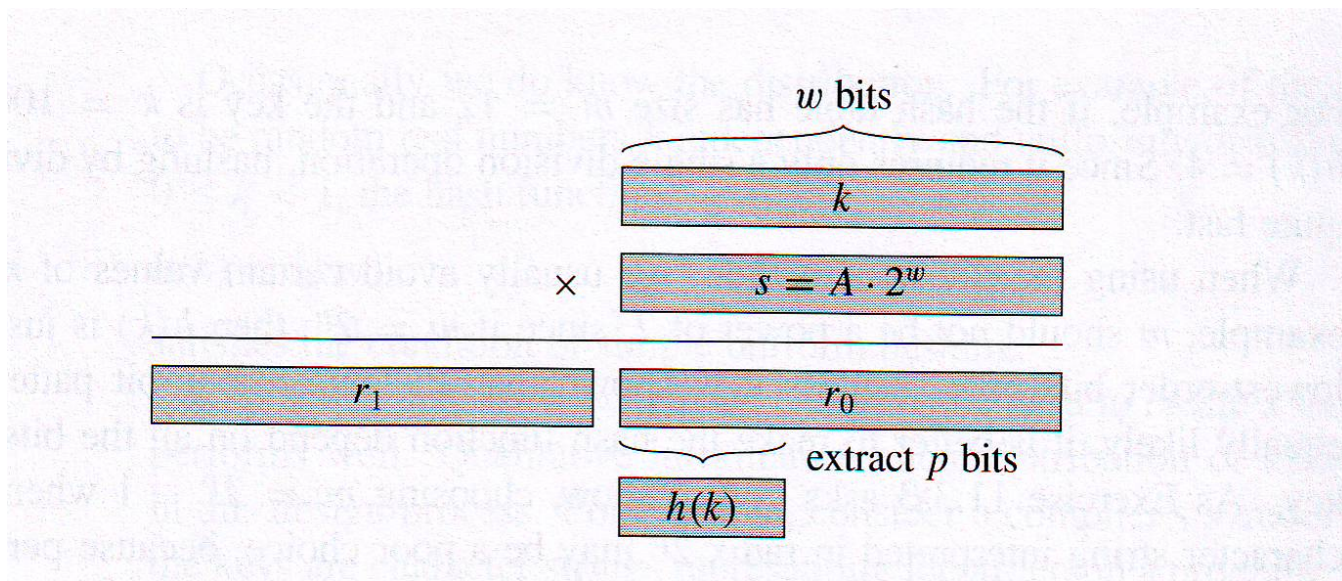
- Primeiro, multiplicamos a chave k por uma constante A no intervalo $0 < A < 1$ e tomamos a parte fracional de kA .
- Então, multiplicamos este valor por m e tomamos o piso do resultado.

$$h(k) = \lfloor m(kA \bmod 1) \rfloor$$

onde $kA \bmod 1 = kA - \lfloor kA \rfloor$.

Vantagem: valor de m não é crítico.

Implementação típica com potências de 2 ($m = 2^p$): Suponha que o tamanho da palavra da máquina seja de w bits e que k cabe em uma única palavra. Restringimos A para ser uma fração na forma $s/2^w$, onde s é um inteiro no intervalo $0 < s < 2^w$.



Primeiro multiplicamos k pelo inteiro de w bits $s = A \cdot 2^w$. O resultado é um valor de $2w$ bits $r_1 2^w + r_0$, onde r_1 é a palavra de ordem mais alta e r_0 é a palavra de ordem mais baixa do produto.

O valor da função de hash desejado consiste dos p dígitos mais significativos de r_0 .

Ex: Suponha que temos $k = 123456$, $p = 14$, $m = 2^{14} = 16384$ e $w = 32$. Suponha que escolhemos A como a fração na forma $s/2^{32}$ que esteja mais próxima de $(\sqrt{5}-1)/2$, e assim $A = 2654435769/2^{32}$. Então $k \cdot s = 327706022297644 = (76300 \cdot 2^{32}) + 17612864$, e assim $r_1 = 76300$ e $r_0 = 17612864$. Os 14 dígitos mais significativos de r_0 resultam no valor $h(k) = 00000001000011_2 = 67$.

• Hashing Universal

Se um adversário malicioso escolher as chaves a serem espalhadas por alguma função de hash fixada, ele poderá escolher n chaves que serão todas designadas para o mesmo slot, resultando em um tempo médio de busca $\Theta(n)$.

Toda função de hash fixa é vulnerável a esse comportamento.

A única forma de diminuir esse risco é escolher a função de hash aleatoriamente de uma forma independente das chaves que serão efetivamente armazenadas. Esta abordagem, denominada *hashing universal*, consegue fornecer performance média comprovadamente boa, não importando quais chaves são escolhidas pelo adversário.

Seja \mathcal{H} uma coleção finita de funções hash que mapeiam um dado universo U de chaves no conjunto $\{0, 1, \dots, m-1\}$. Tal coleção é dita ser *universal* se para cada par de chaves distintas $k, l \in U$, o número de funções hash para as quais $h(k) = h(l)$ é no máximo $|\mathcal{H}|/m$. Em outras palavras, com uma função hash escolhida aleatoriamente de \mathcal{H} , a chance de colisão entre duas chaves distintas k e l não é maior do que a chance $1/m$ de uma colisão se $h(k)$ e $h(l)$ fossem aleatoria e independentemente escolhidas do conjunto $\{0, 1, \dots, m-1\}$.

- **Construção de uma classe universal de funções hashing**

Começamos escolhendo um número primo p grande o suficiente tal que toda chave k esteja no intervalo 0 a $p-1$, inclusive. Denotemos por Z_p o conjunto $\{0, 1, \dots, p-1\}$, e por Z_p^* o conjunto $\{1, \dots, p-1\}$. Uma vez que o universo de chaves é maior do que o número de slots na tabela hash, temos $p > m$.

Definimos a função hash $h_{a,b}$ para todo $a \in Z_p^*$ e todo $b \in Z_p$ usando uma transformação linear seguida por reduções módulo p e então módulo m :

$$h_{a,b} = [(ak + b) \bmod p] \bmod m.$$

P.ex. com $p = 17$ e $m = 6$, temos $h_{3,4}(8) = 5$.

A família de todas as funções hash desse tipo é

$$\mathcal{H}_{p,m} = \{h_{a,b} : a \in Z_p^* \text{ e } b \in Z_p\}.$$

Vantagem: o tamanho m da tabela é arbitrário – não necessariamente primo.

Como há $p - 1$ escolhas para a e p escolhas para b , existem $p(p - 1)$ funções hash em $\mathcal{H}_{p,m}$.