

Slovníkový problém

univerzum  $U$ ,  $S \subseteq U$ ,  $|S| = n$

→ chceme reprezentovat  $S$

operace:

- Find(x) / Member(x)
- Insert(x)
- Delete(x)

Naivní řešení: pole velikosti  $U$ , informace o  $x \in U$  je uložena na pozici  $x$ .

lépe: hashození do pole velikosti  $m \geq n$ .

$h: U \rightarrow \{1, \dots, m\}$  hashození fce.

→ prvek  $x$  se uloží na pozici  $h(x)$ .

• typicky  $h$  může být vybráno náhodně z nějaké množiny možných fce:

např.: 1)  $h(x) = ((ax + b) \bmod |U|) \bmod m$

pro náhodně zvolené  $a, b \in U$ .

2)  $h$  je zcela náhodná fce  $U \rightarrow \{1, \dots, m\}$

a čes od čes se změnil. (viz dále)

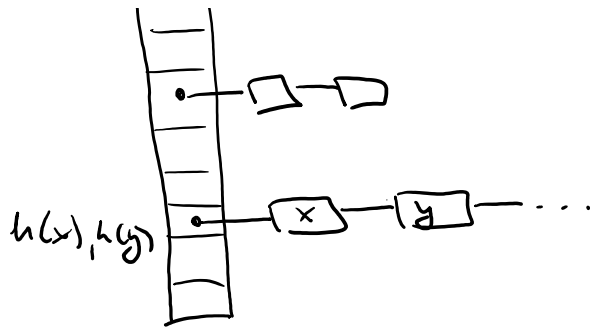
problém: kolize  $x, y \in S$   $x \neq y$   $h(x) = h(y)$

způsoby řešení kolizí:

separování  
řetězce

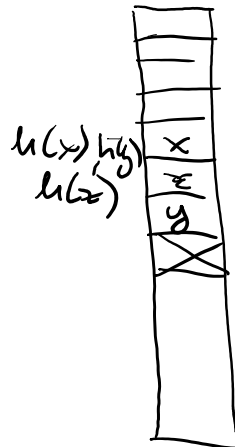


režim



seznam / vyhledávací  
strom

lineární  
přidávání



Insert(x) - najdu nejbližší  
volnou pozici za h(x)  
a tam uložím x.

Find(x) - hledám od h(x)  
do nejbližší volné  
pozice

degitální  
hašování

- podobně jako lineární přidávání, ale  
prvek x ukládám na pozici

$$h_1(x) + i h_2(x) \quad i=0,1,2,\dots$$

kde  $h_1()$  &  $h_2()$  jsou dvě různé nezávislé  
hašovací fu.

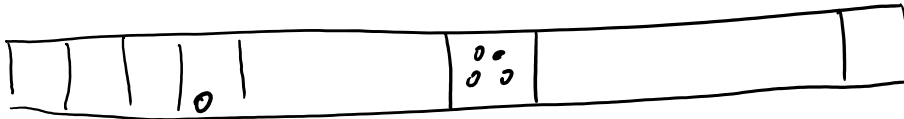
(je potřeba, aby  $h_2(x)$  byla nesoudělná s  $m$ ,  
což je pravda např. pokud  $m$  je prvočíselo  
a  $h_2(x) \in \{1, \dots, m-1\}$ .)

• Fale dalšího přísohu

• DELETE: velmi problematický  
→ označování smazávaných prvků pro Find()

~ jejího využití při Insert(),  
 → pokud je smazávaná prvky příliš mnoho,  
 přecházej vše.

Balls & Bins:  $n$  míček,  $n$  košíků, každý míček hodíme  
 do náhodně zvoleného košíku.



$$\Pr[\text{daný koš je prázdný}] = \left(1 - \frac{1}{n}\right)^n \approx \frac{1}{e}$$

$$\Pr[\text{daný koš obsahuje } k \text{ míčeků}] = \binom{n}{k} \frac{1}{n^k} \left(1 - \frac{1}{n}\right)^{n-k}$$

$$\approx \frac{n^k}{k!} \cdot \frac{1}{n^k} \cdot \frac{1}{e} = \frac{1}{e k!}$$

$k \ll n$

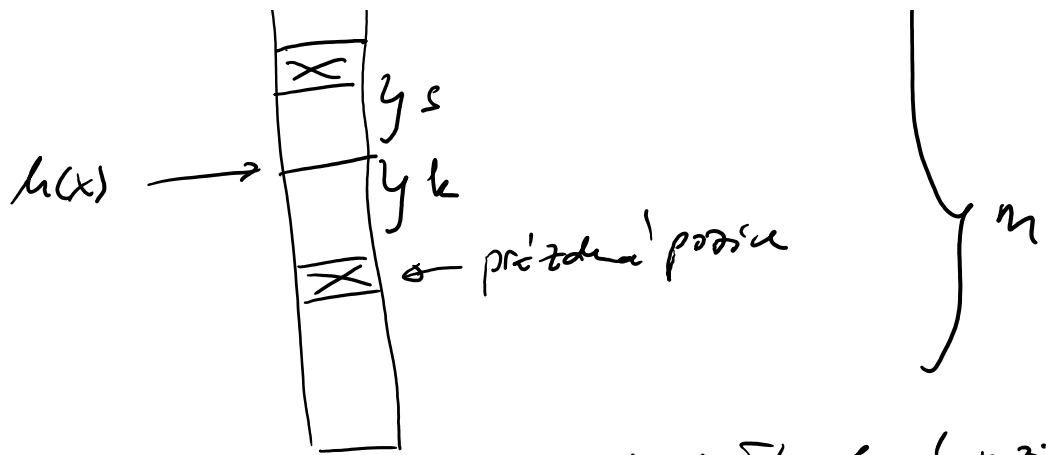
$$\text{pro } k \approx \left(\frac{\lg n}{\lg \lg n}\right) \quad \Pr[\text{existuje koš s } k \text{ míčky}] \leq \frac{1}{n^2}$$

→ s velkou pětí, maximum v libovolném  
 košíku  $\Theta\left(\frac{\lg n}{\lg \lg n}\right)$ .

### lineární přidávání - analýza

klíčik prvek  $x$  a předpokládám, že  $h$  distribuuje  
 prvky zale náhodně.





$P_{k,s}$  ... pravděpodobnost, že nejblíže volené pozice je po  $k$  krocích po  $h(x)$  a před  $h(x)$  je dalších  $s$  prvků

$$P_{k,s} = \binom{n}{k+s} \left(\frac{k+s}{n}\right)^{k+s} = (*)$$

$k+s$  prvků z  $n$  se muselo mapovat do kletky velikosti  $k+s$  z celkového počtu  $n$  prvků.

$$(*) \leq \frac{n^{k+s}}{(k+s)!} \cdot \frac{(k+s)^{k+s}}{n^{k+s}} = \left(\frac{n}{n}\right)^{k+s} \cdot \frac{(k+s)^{k+s}}{(k+s)!} =$$

$$\approx \left(\frac{n}{n}\right)^{k+s} \cdot \frac{1}{\sqrt{k+s}} \cdot e^{(k+s)} \cdot \frac{1}{\sqrt{2\pi}} = (**)$$

Stirlingova aprox:  $a! \approx \sqrt{2\pi a} \left(\frac{a}{e}\right)^a$

Necht'  $m \geq 3n$

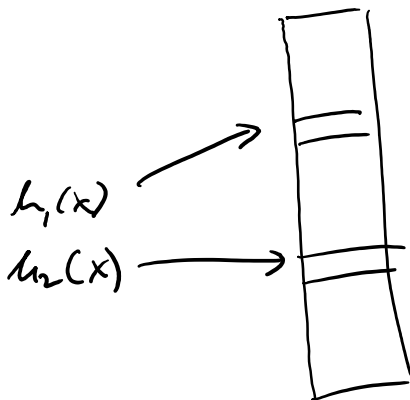
$$(**) \leq \left(\frac{e}{3}\right)^{k+s} \frac{1}{\sqrt{(k+s) 2\pi}}$$

... .. 67

$$\begin{aligned}
 \text{střední doba vyhledávání} &\leq \sum_{k \geq 1} k \cdot \Pr[\text{vyhledávání trvá čas } k] \\
 &\leq \sum_{k \geq 1} k \cdot \sum_{s \geq 0} \left(\frac{e}{3}\right)^{k+s} \frac{1}{\sqrt{(k+s)2\pi}} = O(1) \\
 &= \underbrace{\left(\frac{e}{3}\right)^k \cdot \sum_{s \geq 0} \left(\frac{e}{3}\right)^s}_{O(1)} = O\left(\left(\frac{e}{3}\right)^k\right)
 \end{aligned}$$

• Balls & Bins s volbou:  $n$  míček,  $n$  košíků, pro každý míček zvolím náhodně dva košíky a míček hodím do prázdnějšího.  
 → očekávaní maximální zaplnění košíku  $O(\lg \lg n)$ .

→ Kukací házení [Pagh - Rodler 2004]



$h_1, h_2: U \rightarrow \{1, \dots, m\}$   
 dvě nezávisle vybrané házené fce.

$x$  je buď na pozici  $h_1(x)$   
 nebo na pozici  $h_2(x)$

Insert(x):

if  $T[h_1(x)] = x$  or  $T[h_2(x)] = x$  then return;

```

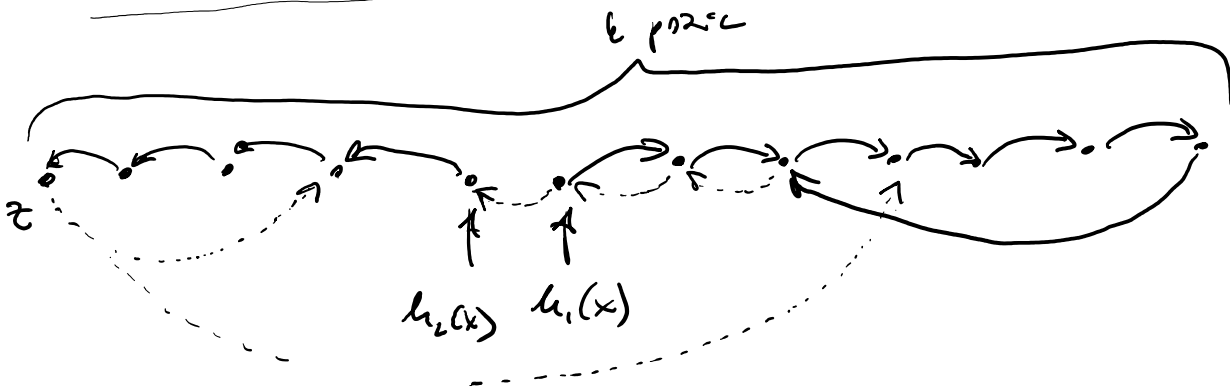
pos = h1(x)
loop n times {
  if T[pos] is empty then { T[pos] = x; return; }
  swap(x, T[pos]);
  if pos = h1(x) then pos = h2(x) else pos = h1(x);
}
rehash(); insert(x);
end

```

- Find  $O(1)$
  - Delete  $O(1)$
  - Insert  $O(1)$
- } v nejhorším případě
- v průměrném případě (přes volbu vhodných fun.)
- Pouze dvě možná pozice pro x.

• Kvalitně kvalitnější funguje dobře pro  $m \geq 2n$ .

• Pohled na Insert(x):



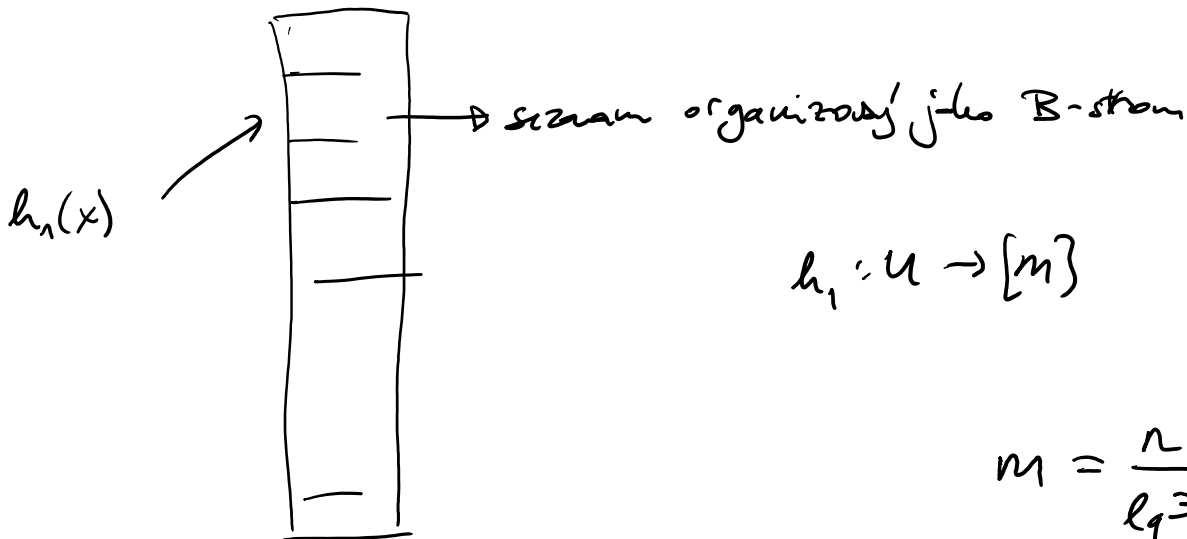
při neúspěšném insertu je pozice  $\geq$  obsazena  
a její obsah se mapuje na již prázdnou

pozici  $\rightarrow$  nekonečný cyklus - prvý je  
 postupný operaci opředm a zpět  
 k pozici pro  $k+1$  místu

- cyklus stáčí zřítizen po  $\approx 6 \lg n$  opakování

Iceberg Hashing [Bender - Conway - Farach-Colton - Kuszman  
 - Tagliarini '21]

$m$  pozic



$$h_1: U \rightarrow [m]$$

$$m = \frac{n}{\lg^3 n}$$

$C$  vhodná konst.

• B-strom pro až  $\lg^3 n + C \lg^2 n$  prvků  
 pokud B-strom přetvoře, prvky se odloží do "sklepa"

•  $\Pr[\text{daná přikrédka přetvoře}] \leq \frac{1}{n^c}$   $C$  libovolná konstanta, závisí  
 číselná rovnost na volbě  $C$

• B-strom organizovaný podle klíče  $h_2(x)$

- B-strom organizovaný podle klíče  $h_2(x)$

$$h_2: U \rightarrow [lg^{10} n]$$

$$B = \sqrt{lg n}$$

klíč  $h_2(x)$  má  $10 lg lg n$  - bitů

celý vchod B-stromu se vejde do jedné paměťové buňky

velikosti  $O(lg n)$  - bitů a lze v něm  
vyhledat za použitím běžných CPU instrukcí  
v čase  $O(1)$ , (viz níže)

hloubka stromu  $\leq 7$ .

→ počet instrukcí  $O(1)$  na vyhledání.

- pokud dojde ke kolizi  $h_2$  v rámci příhradky,  
kolidující prvky se odloží do sklepa

(případně celá příhrádka)

$P_r$  [ v <sup>dané</sup> příhradce nastane kolize ]

$$\leq \binom{lg^3 n + C lg^2 n}{2} \cdot \frac{1}{lg^{10} n} \leq O\left(\frac{1}{lg^4 n}\right)$$

↑  
viz "perfektní  
hesování" dále

s velkou pravděpodobností, počet příhradec

s kolizí  $O\left(\frac{n}{lg^3 n} \cdot \frac{1}{lg^4 n}\right) \rightarrow$  počet prvků



o těchto přehrádkách je  $O(\frac{n}{\lg^4 n})$ .

→ stejný sklop pro  $O(\frac{n}{\lg^4 n})$  prvků  
a pro něj lze použít uopr. kukačkoví  
kašporáci.

panič: přehradka  $\lg^3 n + C \lg^2 n$  prvků po  $\lg |M|$  bitů.  
+ B strom pro  $\lg^3 n + C \lg^2 n$  prvků.  
B strom obsahuje klíče velikosti  
 $O(\lg \lg n)$  - bitů a ukazatele  
velikosti  $O(\lg \lg n)$  - bitů  
[B-strom by měl v rámci přehradky]

$$\text{celkový počet bitů: } \frac{n}{\lg^3 n} \left[ (\lg^3 n + C \lg^2 n) \cdot (\lg |M| + O(\lg \lg n)) \right]$$
$$= \left( 1 + \frac{C}{\lg n} \right) \lg |M| + O(n \lg \lg n).$$

+ sklop  $O(\frac{n}{\lg^4 n} \cdot \lg |M|)$  bitů.

→ d.s. velikosti  $(1 + o(1)) n \lg |M|$   
pro ukládání  $n$  prvků  
s.  $\geq 2 n \lg |M|$  pro kukačkoví kašporáci

→ vsa efektívni

- Find  $O(1)$  operat
- Insert  $O(1)$  s veľkou pravdepodobnosťou (kedykoľvek vejdem do stĺpca)
- $O(1)$  v priemere
- Delete  $O(1)$

Operácie na B-stromoch:

- hľadáme kľúč  $r$ ,  $r$  má  $k = O(\lg n)$  bitov  
 v uzle s  $k$ -bitovými kľúčmi:  $k_1 < k_2 < k_3 \dots < k_B$   
 uloženými jako 

0	$k_1$	0	$k_2$	0	$k_3$	...	0	$k_B$
---	-------	---	-------	---	-------	-----	---	-------

  
 v jednom slove  $O(\lg n)$  bitov  $B = \sqrt{\lg n}$

1) 

0	1	0	1	...	1	0	1
---	---	---	---	-----	---	---	---

 $\cdot$ 

1	5
---	---

 $=$ 

1	r	1	r	...	1	r
---	---	---	---	-----	---	---

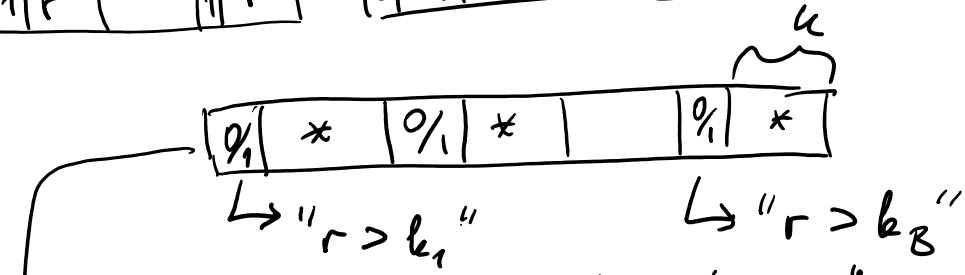
2) 

1	r	1	r	...	1	r
---	---	---	---	-----	---	---

 $-$ 

0	$k_1$	0	$k_2$	...	0	$k_B$
---	-------	---	-------	-----	---	-------

 $=$



3) 

0	*	0	*	...	0	*
---	---	---	---	-----	---	---

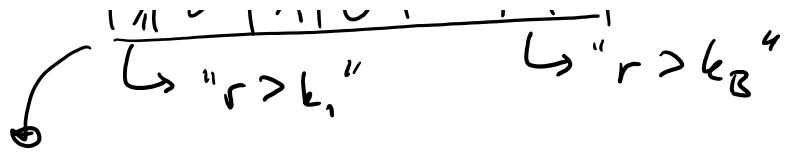
 $\&$ 

1	0	1	0	...	1	0
---	---	---	---	-----	---	---

$=$ 

0	0	0	0	...	0	0
---	---	---	---	-----	---	---

 $\hookrightarrow "r > k_1"$   $\hookrightarrow "r > k_B"$



$$4) \text{ COUNT}(\boxed{\text{||||}}) = \# i, r > k_i$$