

# HEŠOVÁNÍ (hashing)

- Máme:
- universum  $U$  možných prvků, typicky  $U = [u] = \{0, \dots, u-1\}$
  - množinu přírodních  $\mathcal{P} = [m]$
  - hešovací funkci  $h: U \rightarrow \mathcal{P}$
  - reprezentovanou množinu  $X \subseteq U, |X|=n$
  - v ~~každé~~ přírodně množinu  $P_i = \{x \in X \mid h(x) = i\}$  ... kolize nastane, pokud  $|P_i| > 1$ .  
reprezentovanou jako seznam

Nepráctelské množiny: Pro libovolnou pevnou funkci  $h$  lze najít  $X: h \upharpoonright X = \text{const.}$  ( $\exists i: P_i = X$ )  
... stačí, aby  $U \geq nm$ .  $\Rightarrow$  operace stojí  $\Omega(n)$  w.c.

Radeji zvolíme  $h$  náhodně ... pak můžeme zavést dobrou očekávanou složitost.

Cíl: Nastavíme  $m$ ,  
v každé přírodně  $\mathcal{P}$  const.  
prvků  $\rightarrow$  složitost  $\approx$  const.

Možnosti:  
① Úplně náhodná funkce ... nepraktické (potřebujeme  $\Omega(U \log m)$  bitů paměti)  
Ale chová se dobře:  $\forall x, y \in U, x \neq y, \Pr_n[h(x) = h(y)] = 1/m$ .

② Df: Systém funkcí  $\mathcal{H}$  z  $U$  do  $[m]$  je c-univerzální pro  $c \in \mathbb{R}$ , pokud  
 $\forall x, y \in U, x \neq y: \Pr_{h \in \mathcal{H}}[h(x) = h(y)] \leq c/m$ .

Chceme: Funkce z  $\mathcal{H}$  popsatelem parametry tak, aby stálo  
a) efektivně náhodně vybrat  $h \in \mathcal{H}$   
b) rychle  $h(x)$  vyhodnotit (předp.  $O(1)$ )

Věta: Necht'  $\mathcal{H}$  je c-univ. z  $U$  do  $[m]$ ,  $X \in \binom{U}{n}, y \in U \setminus X$ .  
Potom  $E[\#_{h \in \mathcal{H}} \{x \in X: h(x) = h(y)\}] \leq \frac{cn}{m}$ .  
 $\Rightarrow$  Očekávaná složitost úspěšného hledání je  $O(\frac{n}{m+1})$  pro úspěšné není horší (čas úspěš. hledání = čas Insertu = čas neúsp. hl. v době Insertu)

Důk zavedeme indikatory:  $A_i = \begin{cases} 1 & \text{pokud } h(x_i) = h(y) \\ 0 & \text{jinak} \end{cases}$   
 $E[A] = E[\sum A_i] = \sum E[A_i] \leq n \cdot c/m = c \frac{n}{m}$ .  
ocíslovíme  $X = \{x_1, \dots, x_n\}$   
 $= \Pr[A_i = 1] \leq c/m$  (z c-univerzality)

Df: Necht'  $p \geq m$  je prvočíslo. Pak  $\mathcal{H} := \{h_{a,b} \mid a, b \in [p], a \neq 0\}$ , kde  $h_{a,b}(x) := ((ax+b) \bmod p) \bmod m$ .

Věta: Systém  $\mathcal{H}$  je 1-univerzální.

Důk: Necht'  $x, y \in [U], x \neq y$ .

Nejprve analyzujeme funkce tvaru  $(ax+b) \bmod p$  ... to je výpočet v  $\mathbb{Z}_p$ .

• Pro dané  $(a,b) \in [p]^2$  zavedeme:  
 $r = (ax+b) \bmod p$   
 $s = (ay+b) \bmod p$   
soustava 2 nezávislých lin. rovnic v tělese  $\Rightarrow \exists!$  řešení  $(a,b)$  pro každé  $(r,s)$

- Máme tedy bijekci mezi všemi  $(a,b)$  a  $(r,s)$  (oba  $\in [p]^2$ )
- Požadavek  $a \neq 0$  ~~je~~ je ekvivalentní  $r \neq s$ .

Nyní přidáme modulo  $m$ .

Odhadujeme # spatných dvojic  $(a,b)$ , pro něž  $h_{a,b}(x) = h_{a,b}(y)$ .

Ty odpovídají dvojicím  $(r,s)$  s  $r \equiv s \pmod{m}$ .

rovnoměrně náhodný výběr jednoho dáva r.n. výběr druhého

Pro každé  $r$  spočítáme, kolik je s t.č.  $r \in S$ .

- Pokud  $[p]$  rozdělíme na  $m$ -tice (posl. neúplná), najdeme v 1  $m$ -tici nejvýše 1 takové  $s$ .
- $\#s \leq \lceil p/m \rceil - 1 \leq \frac{p+m-1}{m} - 1 = \frac{p+m-1-1}{m} = \frac{p-1}{m}$ .

Tedy  $\#$  správných dvojic  $\leq p \cdot \frac{p-1}{m}$ , všech dvojic je  $p(p-1) \Rightarrow \text{Pr}[\text{dvojice správná}] \leq 1/m$ .

!  $c$ -universalita (i pro  $c=1$ ) je hodně daleko od skutečné náhodnosti } příklady porady

③ Df: Systém funkcí  $\mathcal{H}$  z  $U$  do  $[m]$  je  $k$ -nezávislý  $\equiv$   
 $\forall a_1, \dots, a_k \in U$  navzájem různá,  $\forall a_1, \dots, a_k \in [m]$   $\text{Pr}_{h \in \mathcal{H}} [\bigwedge_{i=1}^k h(x_i) = a_i] \leq O(\frac{1}{m^k})$ .

Silněji: je  $(k,c)$ -nezávislý  $\equiv \text{Pr}[\dots] \leq c/m^k$ .

- $\mathcal{H}$  je  $k$ -nezávislý  $\Rightarrow \mathcal{H}$  je  $(k-1)$ -nezávislý (pro totéž  $c$ )
- $\mathcal{H}$  je  $(2,c)$ -nezávislý  $\Rightarrow \mathcal{H}$  je  $c$ -universalní
- 1-nezávislost je příliš slabá, splňuje ji třeba systém konstantních funkcí.

Df:  $\mathcal{L} := \{h_{a,b} \mid a,b \in [p]^2\}$ , kde  $h_{a,b}(x) = (ax+b) \bmod p \bmod m$ .

Věta  $\mathcal{L}$  je  $(2,4)$ -nezávislý.

Dk: Opět vyurijeme bijekci mezi  $(a,b)$  a  $(r,s)$ .

Chceme dokázat  $\text{Pr}_{r,s} [r \equiv i \ \& \ s \equiv j] \leq 4/m^2$ . (vše mod  $m$ )

to jsou nezávislé jevy  $\Rightarrow$  stačí  $\text{Pr}_r [r \equiv i] \leq 2/m$

$\hookrightarrow$  čísel  $r \in [p]$  kongruentních s  $i$  je nejvýše  $\lceil p/m \rceil \leq \frac{p+m-1}{m} \leq \frac{2p}{m}$  (kde  $m \leq p$ )  
 takže  $\text{Pr}_r [r \equiv i] \leq \frac{2p}{m} / p = \frac{2}{m}$ .

Věta  $\mathcal{L}$  není 3-nezávislý.

Dk: zvolíme  $x,y,z$ , tak, aby  $x+y=2z$ , ~~například~~  $x,y,z$  různá. Dále  $m=p$ , takže vše v  $\mathbb{Z}_p$ .

Chceme, aby platilo  $ax+b \equiv i$  ① a toho:  $2k \equiv 2az+2b \equiv a \cdot 2z + 2b$   
 (pro dané  $i$  a  $k$ )  $ay+b \equiv j$  ②  $\equiv a(x+y) + 2b \equiv (ax+b) + (ay+b) \equiv i+j$   
 $az+b \equiv k$  ③

Tedy kdykoli ①-③ platí, je  $2k \equiv i+j$ .

Proto  $\text{Pr}_h [h(x)=i \ \& \ h(y)=j \ \& \ h(z)=k] = \underbrace{\text{Pr}_h [h(z)=k \mid h(x)=i \ \& \ h(y)=j]}_1 \cdot \underbrace{\text{Pr}_h [h(x)=i \ \& \ h(y)=j]}_{O(\frac{1}{m^2})}$ , neboť  $\mathcal{L}$  je 2-~~ne~~ nezávislý.  
 pro libovolné  $i,j$  a  $k = (i+j) \cdot 2^{-1}$   
 $= O(1/m^2) \neq O(1/m^3)$ .

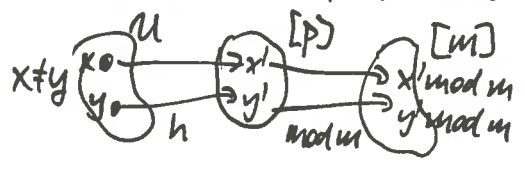
**SKLÁDÁNÍ FUNKCÍ**

díky universalitě  $\mathcal{L}$  / 2-nezávislosti  $\mathcal{L}$  fungoval tak, že jsme nejspíše dokázali 2-ne. v tělese  $\mathbb{Z}_p$  (tedy u funkcí  $[p] \rightarrow [p]$ ) a pak vzhledli, že se modulem  $m$  moc nepokazí.

$\Downarrow$   
zkusíme o tom dokázat něco obecnějšího

Věta: Necht'  $\mathcal{H}$  je  $(2, c)$ -nezavislý systém funkci z  $\mathbb{F}_p^u$  do  $[p]$  a  $m \leq p$ .  
 Potom  $\mathcal{H}^* := \{h \bmod m \mid h \in \mathcal{H}\}$  je  $(2, c)$ -universalní a  $(2, c)$ -nezavislý.

Důk:  $U \xrightarrow{h \in \mathcal{H}} [p] \xrightarrow{\bmod m} [m]$



① universalita:

$$\begin{aligned} & \Pr_h [h(x) \bmod m = h(y) \bmod m] \\ &= \Pr [h(x) = h(y) \vee (h(x) \neq h(y) \mid h(x) \neq h(y))] \\ &\leq \Pr [h(x) = h(y)] + \Pr [h(x) \neq h(y) \mid \dots] \\ &\leq \frac{c}{p^2} \approx 2^{-2\log p} \\ &= \Pr_h \left[ \bigvee_{\substack{i, j \\ i \neq j}} h(x) = i \ \& \ h(y) = j \right] = \sum_{i \neq j} \Pr [h(x) = i \ \& \ h(y) = j] \leq p \cdot \frac{p+m-1}{m} \cdot \frac{c}{p^2} = c \cdot \frac{p+m-1}{pm} \\ &= \frac{c}{m} \cdot \frac{p+m-1}{p} \leq \frac{2c}{p} = 2. \end{aligned}$$

② nezavislost:  $\Pr_h \left[ \bigvee_{\substack{i, j \\ i \neq j}} h(x) = i \ \& \ h(y) = j \right] \leq \frac{c}{p^2} \cdot \left[ \frac{p}{m} \right]^2 \leq \frac{c}{m^2} \cdot \left( \frac{p+m-1}{p} \right)^2 \leq \frac{4c}{m^2}$

Zobecnění:  $\mathcal{H}$   $(k, c)$ -nezavislý  $\Rightarrow \mathcal{H}^k$   $(k, c')$ -nezavislý pro  $c' = c \cdot \left( \frac{p+m-1}{p} \right)^k \leq c \cdot 2^k \leq \left( 1 + \frac{m}{p} \right)^k \leq e^{\frac{km}{p}} \leq \text{const}$  pro  $p = \Omega(km)$

Věta: Necht'  $\mathcal{F}$  je  $c$ -universalní z  $U$  do  $[r]$  a  $\mathcal{G}$  je  $(2, d)$ -nezavislý z  $[r]$  do  $[m]$ .

Potom  $\mathcal{H} := \{f \circ g \mid f \in \mathcal{F}, g \in \mathcal{G}\}$  je  $(2, d')$ -nezavislý z  $U$  do  $[m]$ , pro  $d' = \frac{(c+1)d}{mr}$ .

Důk:  $\Pr_h [h(x) = i \ \& \ h(y) = j] = \Pr_{f, g} \left[ \underbrace{g(f(x)) = i \ \& \ g(f(y)) = j}_{A} \mid f(x) = f(y) \right] + \Pr_{f, g} \left[ \underbrace{A}_{\neq} \mid f(x) \neq f(y) \right]$

$$\leq \Pr[A \mid f(x) = f(y)] \cdot \Pr[f(x) = f(y)] + \Pr[A] \leq \frac{d}{mr^2} \approx \text{nezavislosti } \mathcal{G}$$

$$\leq \frac{cd}{m^2} + \frac{d}{m^2} = \frac{(c+1)d}{m^2} \dots \text{šlo by přepsat na } \frac{(c \cdot \frac{m}{p} + 1)d}{m^2}$$

**POLYNOMY**

Důk:  $\mathcal{P}_k := \{h_a \mid a \in \mathbb{F}_p^k\}$ ,  $h_a(x) := \left( \sum_{i=0}^{k-1} a_i x^i \right) \bmod p$ . fce z  $[p]$  do  $[p]$

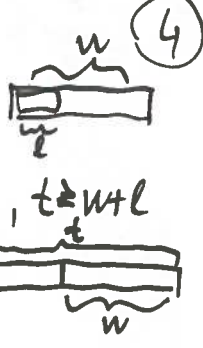
Věta: Systém  $\mathcal{P}_k$  je  $(k, 1)$ -nezavislý.

Důk: Danými  $k$  body prochází právě 1 polynom stupně menšího než  $d$ . (platí v libovolném tělese)  
 Chceme-li hošovat do  $[m]$  pro  $m < p$ , složíme s  $\bmod m$ , pro  $p = \Omega(km)$  je  $(k, \text{const})$ -nezavislý z  $U$  do  $[m]$

MULTIPLY-SHIFT - rychlé hořovací fce  $ax \langle w-l:w \rangle$

Df:  $M := \{h_a \mid a \in [2^w]\}$ , kde  $h_a(x) := \lfloor (ax \bmod 2^w) / 2^{w-l} \rfloor \dots z [2^w]$  do  $[2^l]$

$M' := \{h'_{a,b} \mid a,b \in [2^t]\}$ , kde  $h'_{a,b}(x) := \lfloor (ax^b \bmod 2^t) / 2^{t-l} \rfloor \dots z [2^w]$  do  $[2^l]$ ,  $t \geq w+l$   
 $(ax+b) \langle t-l:t \rangle \bmod 2^l$



Věta:  $M$  je 2-univerzální,  $M'$  je 2-nezávislý.  
 (Dk viz literatura)

TABELACNÍ HOŘOVÁNÍ

Mejme  $U = [2^{dt}]$ . Náhodně zvolíme fce (tabulky)  $T_0, \dots, T_{d-1}: [2^t] \rightarrow [2^t]$ .  
 $h(x) := \bigoplus_{i=0}^{d-1} T_i(x \langle t-i:t \rangle)$

potřebujeme  $dt \cdot 2^t$  bitů paměti

lepší notace:  $x \langle i:j \rangle$  jsou bity i až j-1 čísla x

Věta: Tabelační hořování je 3-nezávislé, ale není 4-nezávislé.  
 (Dk: cvičení)

obecně s čas  $O(d)$  prostor  $O(U^{1/d+\epsilon})$

HOŘOVÁNÍ VEKTORŮ A ŘETĚZCŮ

Df:  $\mathcal{Y}$  z  $\mathbb{Z}_p^k$  do  $\mathbb{Z}_p$ :  $\mathcal{Y} := \{h_{\vec{a}} \mid \vec{a} \in \mathbb{Z}_p^k\}$ ,  $h_{\vec{a}}(\vec{x}) := \vec{a} \cdot \vec{x}$

Věta:  $\mathcal{Y}$  je 1-nezávislý univerzální.

~~obecně s~~

Dk:  $\Pr_{\vec{a}}[h(\vec{x}) = h(\vec{y})] = \Pr_{\vec{a}}[\vec{a} \cdot \vec{x} = \vec{a} \cdot \vec{y}] = \Pr_{\vec{a}}[\vec{a} \cdot (\vec{x} - \vec{y}) = 0] =$

nenulový vektor  $\vec{z}$ , Bůho  $z_k \neq 0$

$$= \Pr_{\vec{a}} \left[ \sum_{i=1}^{k-1} a_i z_i + \underbrace{a_k z_k}_{\neq 0} = 0 \right] = \frac{1}{p}$$

intuice:  $\underbrace{\text{celokv.}}_{\text{hodnota}} + \underbrace{\text{rovnoměrně}}_{\text{uhodně}} = \text{rovnoměrně uhodně}$

pro každou volbu  $a_1 - a_{k-1} \exists! a_k$ , s níž rovnost platí.

Df:  $\mathcal{Y}'$  z  $\mathbb{Z}_p^k$  do  $\mathbb{Z}_p$ :  $\mathcal{Y}' := \{h_{\vec{a},b} \mid \vec{a} \in \mathbb{Z}_p^k, b \in \mathbb{Z}_p\}$ ,  $h_{\vec{a},b}(\vec{x}) := \vec{a} \cdot \vec{x} + b$

Věta:  $\mathcal{Y}'$  je (2,1)-univerzální, 2-nezávislý.

[ (2,1)-univerzálnost plyne z obecné věty ]

Df:  $\mathcal{Q}_k$  z  $\mathbb{Z}_p^k$  do  $\mathbb{Z}_p$ :  $\mathcal{Q}_k := \{h_a \mid a \in \mathbb{Z}_p\}$ ,  $h_a(\vec{x}) := \sum_{i=0}^{k-1} x_i \cdot a^i$

Věta:  $\mathcal{Q}_k$  je  $k$ -univerzální.

Dk:  $\Pr_{\vec{a}}[h(\vec{x}) = h(\vec{y})] = \Pr[\vec{a}(\vec{x} - \vec{y}) = 0]$ , ale to je polynom stupně  $\leq k$ , čili má max.  $k$  kořenů

☺ lze používat i na řetězce proměnlivé velikosti (omezené max. dél), stačí padding znaky, které u jinde neryšují (nejlépe nulami)

↓ složíme s  $(bx + c \bmod p) \bmod m$

$\mathcal{Q}' := \{h_{a,b,c} \mid a,b,c \in [p]\}$ ,  
 $h_{a,b,c}(\vec{x}) := (b \sum_{i=0}^{k-1} x_i a^i + c) \bmod p \bmod m$

Věta:  $\mathcal{Q}'$  je 2-nezávislý pro  $p \geq km$

Dk: z obecné věty o skládání.