

# Úsporné datové struktury

- Chceme přeložit prvky množiny  $C$  (číslo jeden ze stromů na  $n$  vrcholech)
  - a počítat téměř optimální prostor ...  $OPT = \lceil \log |C| \rceil$  bitů } entropie při rovnoměrné rozdělení psí
  - Přitom chceme rychlé operace
  - Variandy:
    - implicitní DS - jen data sama ve vhodném pořadí (setříděné pole, halda) }  $OPT + O(1)$
    - úsporné (succinct) DS -  $OPT + o(OPT)$
    - kompaktní DS -  $O(OPT)$  ... ale pozor, třeba BVS nemusí být kompaktní kvůli požadování - například slova, ale bitů
- ↑  
rounding  
&c.

## Problém Reprerentace řetězce nad obecnou abecedou

- Příklad:  $[10] \rightarrow 4$  bitů ( $\log 10 \approx 3.322$ )  
 $[10]^2 \dots$  mohu uložit jako  $2 \times 4$ , nebo jako  $[100] \rightarrow 7$  bitů  
 $[10]^k \dots$  uložit jako  $[10^k] \rightarrow k \cdot \log_{10} 10 + 1$  namísto  $k \cdot (\log 10 + 1)$

↳ limitně se blíží  $k \log 10$  b/znak (redundance = délka kódu - entropie)  
 ... ale ztrácím lokální dekódovatelnost ] jde k 0

Praktické řešení: dělím na bloky o 9 znacích ... délka kódu = 30  
 entropie  $9 \cdot \log 10 \approx 29.897$   
 $\rightarrow$  redundance  $\approx 0.103$

Ude se vejde do 1 slova  
 $\rightarrow$  počítám s ním  $O(1)$  času i prostoru

↳ pro string délky  $n$   $30 \cdot \lceil n/9 \rceil \leq \frac{30}{9}n + 30$   
 ↳ redundance  $\leq 0.103 \cdot \frac{n}{9} + 30$

⊙  $(a,b) \in X \times Y$  kódují po složkách  $\Rightarrow$  délky kódů i entropie se sečtou  $\rightarrow$  redundance také

Příklad #2: Posíláme proud bitů, ale dopředu nevíme, jak bude dlouhý.  $\rightarrow$  instantní (prefixový) kód

[potřebujeme umět dekódovat značku pro konec, auz by neustále "jeste to pokračuje" zabralo dost místa]

Triviální řešení: rozdělíme na bloky o  $b$  bitech, za  $\forall$  blok připsáme 1 bitovu značku  
 ↳ redundance  $\frac{n}{b} + b + 1$

## SOLE Encoding [Dodis, Patrascu, Thorup 2010]

(Short-Odd-Long-Even)

- vstup rozdělíme na bloky o  $b$  bitech  $\rightarrow$  prvky abecedy  $[B]$  pro  $B = 2^b$
- poslední blok doplníme (10  $\rightarrow$  0 na konec), ale potřebujeme přidat spec. blok označující EOF  
 ↳ převod  $[B+1]^* \rightarrow [B]^*$
- nastavíme  $b \geq 2 \log n + 2 \Rightarrow B \geq 4n^2$  ] bloky se vejdou do  $O(1)$  slov RAMe

číslo bloku	1	2	3	4	5	6	...	k
vstupní abeceda	B	B	B	B	B	B	...	EOF
+ EOF	B+1	B+1	B+1	B+1	B+1	B+1	...	B+1
1. průchod	B	B+3	B-3	B+6	B-6	B+9	...	0
2. průchod	B	B	B	B	B	B	...	

sem si dovolíme nulový blok

1. průchod:  $(B+1)(B+1) \leq (B-3i)(B+3i+3)$

$$B^2 + 2B + 1 \leq B^2 + 3Bi + 3B - 3Bi - 9i^2 - 9i$$

$$-B + 1 \leq -9i^2 - 9i$$

$$B - 1 \geq 9i^2 + 9i$$

$$B \geq 9i^2 + 9i + 1, \text{ neboť } i \leq \frac{n+1}{2} \text{ a } B \geq 4n^2$$

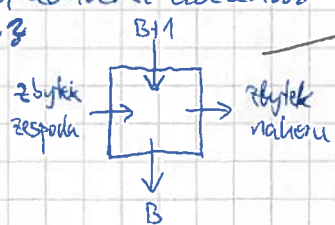
2. průchod:  $(B+3i)(B-3i) = B^2 - 9i^2 \leq B^2$

- redundance  $\leq 3b$  (1 blok zadrůhování + 1 EOF + 1 extra na každý # bloků)
- proudové kódování + dekodování
- lokální dekodování i změny v  $O(1)$  na RAMu [potřebují  $O(\log n)$  bitů]  $\hookrightarrow O(1)$  slov

zadrůhování na celé bits je příliš redundancí

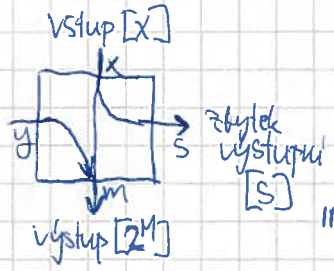
Postupně z toho odvodíme obecnou konverzi abeced...

- proč to vlastně fungovalo?



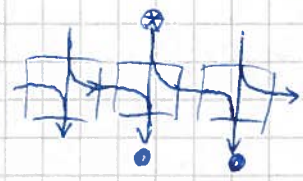
nechám zbytek malé množství informace a to v příštím kroku přimíchám k většímu, čímž snížím redundanci... zbytek se posílají jen do druhé vzdálenosti  $\rightarrow$  lokálně dekodovatelné

Obecně: zbytek vstupní [Y]



"MIXÉR"

... při retrace stále lok. dekodovatelné



pro dekodování stačí znát obě

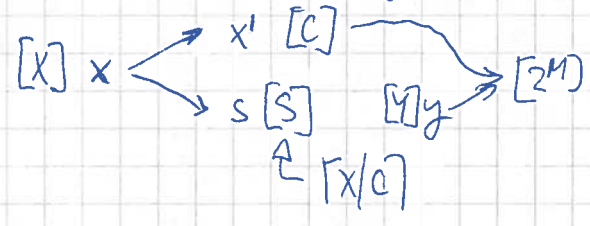
Lemma: Necht  $X, Y \leq 2^W$ . Pak  $\exists M, S$  a zobrazení  $f: [X] \times [Y] \rightarrow [2^M] \times S$  t.č.:

- 1)  $S = O(\sqrt{X})$ ,  $2^M = O(Y \cdot \sqrt{X})$  a  $S, M$  závisí pouze na  $X, Y$
- 2)  $f$  lze vyhodnotit (na RAMu) v čase  $O(1)$
- 3)  $x$  lze dekodovat z  $m, s$  v  $O(1)$
- 4)  $y$  lze dekodovat ze samotného  $m$  v  $O(1)$
- 5) Redundance je  $O(1/\sqrt{X})$  bitů

$$\underbrace{(M + \log S)}_{\text{entropie výstupu}} - \underbrace{(\log X + \log Y)}_{\text{entropie vstupu}}$$

De 8 zvolíme  $M$  (prádeli)

↳  $[2^M]$  musí obsahovat celé  $y$  a nějakou část  $x$  čísla  $x \dots$  jakou?  $Cs = \lfloor 2^M / Y \rfloor$  možnosti



• redundance ~~rozkladu~~ kódování  $x', y \rightarrow m^8$

$$R_1 = M - \log(Y \cdot C) \leq M - \log(2^M - Y) = \log \frac{2^M}{2^M - Y} = \log \frac{1}{1 - \frac{Y}{2^M}} = O\left(\frac{Y}{2^M}\right) = O\left(\frac{1}{C}\right)$$

$\log\left(Y \cdot \frac{2^M}{Y}\right) \geq 2^M - Y$

$e^x \geq 1+x$   
 $x \geq \log(1+x)$   
 $-x \leq \log \frac{1}{1+x}$   
 $x \geq \log \frac{1}{1-x}$

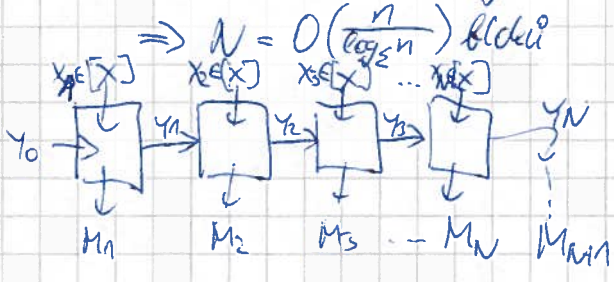
• redundance rozkladu  $x \rightarrow x', s$

$$R_2 = (\log C + \log S) - \log X = \log C + \log \lceil \frac{X}{C} \rceil - \log X = \log \frac{C \cdot \lceil \frac{X}{C} \rceil}{X} \leq \log \frac{X+C}{X} = \log \left(1 + \frac{C}{X}\right) = O\left(\frac{C}{X}\right)$$

• celkem tedy minimalizujeme  $O\left(\frac{1}{C} + \frac{C}{X}\right)$  ...  $C \approx \sqrt{X} \rightarrow S = O(\sqrt{X}), 2^M = O(Y \cdot \sqrt{X})$   
 → redundance  $O(1/\sqrt{X})$ , jak jsme slíbili. [nezávislost na  $Y$ ]

První pokus o reprezentaci stringu

$A \in [E]^n$  ... rozdělíme na bloky o velikosti  $\Theta(\log_{\epsilon} n)$  tak, aby  $X \approx n^2$



... redundance  $O\left(\frac{1}{n}\right)$  v každém bloku

↳ celkové  $O(1)$ , což je milé

... lokální kódování / dekódování

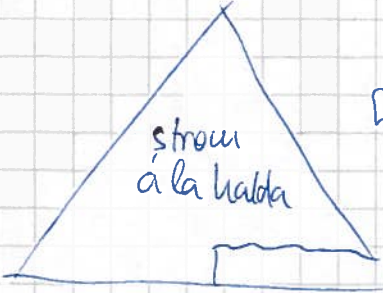
OPPS! Každý blok má jiné parametry!

→ potřebujeme tabulku  $O\left(\frac{n}{\log_{\epsilon} n}\right)$  konstant!

→ tohle v seřadě nevedlo, protože  $X = 2^{\epsilon+1}$ , takže jsou jednotlivá  $y_i$  unesli aproximovat aritmetickou posloupností ... ale obecně to nebude fungovat

Zadání: stromové kódování

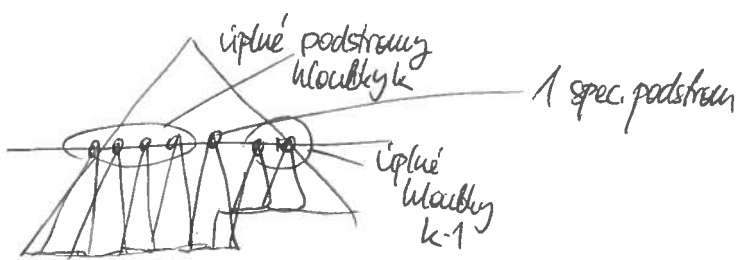
Opět volíme  $X = \Theta(n^2)$ ,  $N = O\left(\frac{n}{\log_{\epsilon} n}\right)$



• dekódování je stále lokální

(z otce stačí out, ten určuje zbytek)

~~Y1~~  $[Y_1] \times [Y_2] \cong [Y_1 \cdot Y_2]$



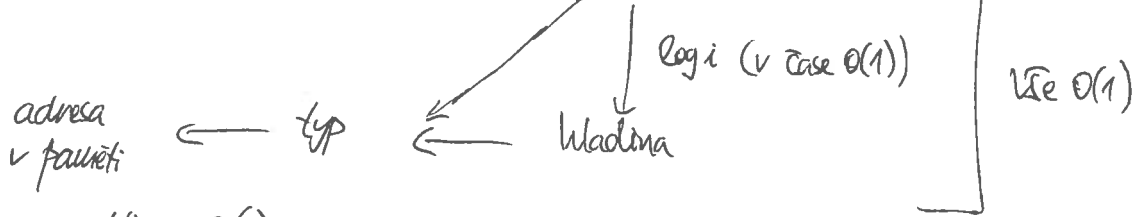
na každé hladině 3 typy vrcholů  
pro ně si pamatujeme:

- # vrcholů tohoto typu
- adresu dat 1. vrcholu
- parametry mixéru

celkem  
 $O(\log n)$   
konstant  
(slov)

4

Decódování s pozice ve streamu  $\rightarrow$  číslo bloku = pozice ve streamu



Hladina změna též v  $O(1)$ .

Cellková redundance  $O(1)$  [jako u předch. polusu] +  $O(1)$

▲ za kódování bloků

▲ zadaná velikost posledního bloku (nepřesného) (ten budujeme odspodu)

Věta: Na word-RAMu lze reprezentovat prvky  $[\Sigma]^n$

v prostoru  $\lceil n \log \Sigma \rceil + O(1)$  bítů

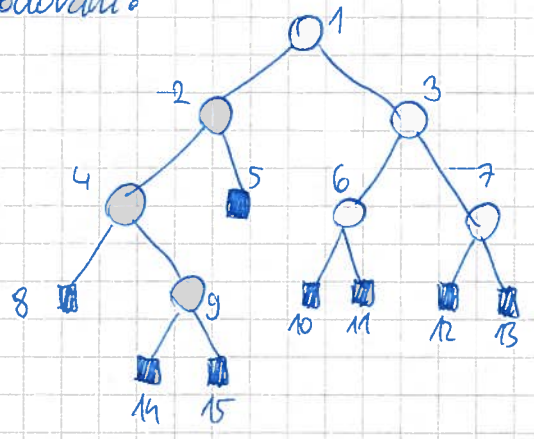
se čtením/zapísem prvku v čase  $O(1)$

s počítáním  $O(\log n)$  konstant závislých na  $n$  a  $\Sigma$ .

# Usporná reprezentace binárních stromů

- chceme uložit strukturu stromu, čímž hledat syny, otce atd.
- bin. stromů s  $n$  vrcholy existuje  $C_n = \frac{1}{n+1} \binom{2n}{n} \approx 4^n / \text{poly}(n)$   
 $\hookrightarrow \approx 2n + o(n)$  bitů

## • kódování:



- přidáme externí vrcholy
- všechny vrcholy očíslováme po hladinách dle kladka
- za  $\forall$  vrchol zapíšeme 1 bit  $\begin{cases} 0 & \text{pro externí} \\ 1 & \text{pro interní} \end{cases}$

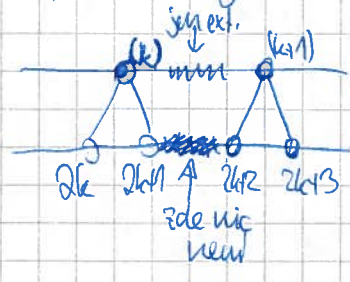
$\downarrow$   
 $2n+1$  bitů

## • adresace:

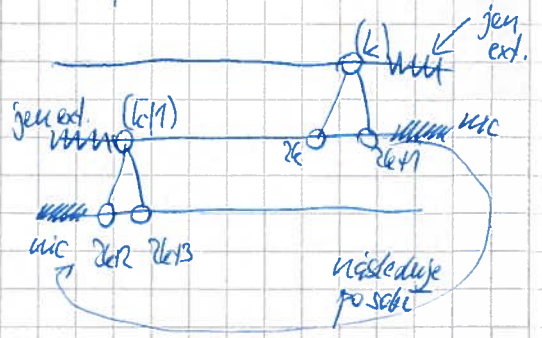
Lemmas: Synové  $k$ -tého vnitřního vrcholu leží na pozicích  $2k$  a  $2k+1$ , opět pod 1  
 počítáno od 1 po hladinách

Dk: Indukcí podle  $k$ : Pro  $k=1$  triviálně platí!

$k \rightarrow k+1$  ... ①  $k, k+1$  na stejné hladině



② přes hladinu



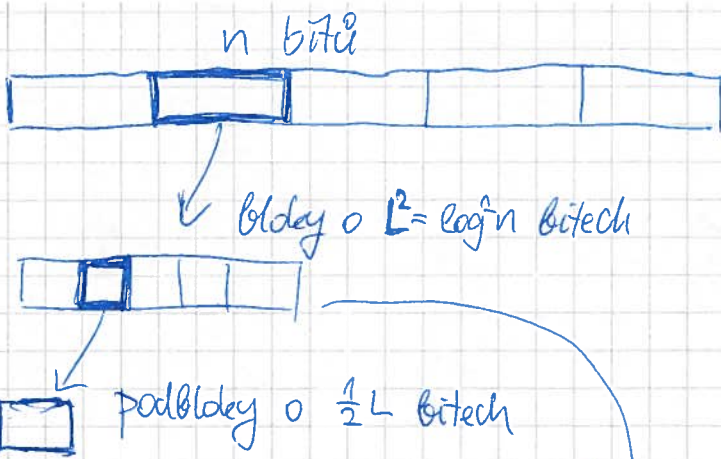
- Redukujeme na  $\text{Rank}_n$  a  $\text{Select}_n$  v posloupnosti  $N = 2n+1$  bitů

$\downarrow$   
 $\text{Rank}_n(i) =$   
 počet jedniček  
 na pozicích  $1-i$

$\downarrow$   
 $\text{Select}_n(i) =$   
 pozice  $i$ -té  
 jedničky

- z toho:  $\text{Left}(i) = 2 \cdot \text{Rank}_n(i)$  (vstup  $i$  vstup jsou pozice v posl. vrcholů)
- $\text{Right}(i) = \text{Left}(i) + 1$
- $\text{Parent}(i) = \text{Select}_n(\lfloor i/2 \rfloor)$

Rank



pro # bloků předpočítáme #1 v předchozích blocích

$$\frac{n}{L^2} \cdot L = \frac{n}{L} \in o(n)$$

↑ bloky      ↑ bloky na #1

dekompozice...  
předpočítáme odpovědi na všechny dotazy

pro # podbloků předpočítáme #1 v předch. podblocích téhož bloku

$$\frac{n}{L} \cdot \log(L^2) \in o(n)$$

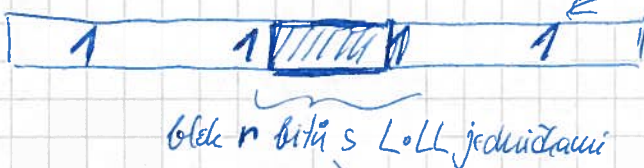
$\approx LL = \log \log n$

$$2^{\frac{1}{2}L} \cdot \frac{1}{2}L \cdot LL \approx \sqrt{n} \cdot L \cdot LL \in o(n)$$

↑ možných obsahů bloků      ↑ možných dotazů      velikost odpovědi

$= \sqrt{n}$

Select



rozdelíme po ~~LL~~ jedničkách a zaznamenujeme jejich pozice + pointery na substruktury

$$\frac{n}{L \cdot LL} \cdot 2L = \frac{2n}{LL} \in o(n)$$

↑ # částí      ↑ pozice + pointer

uložíme pozice všech 1

$r \geq (L \cdot LL)^2$

$r < (L \cdot LL)^2$



ještě 1x totéž, tentokrát korekci  $LL^2$ -ta 1

podblok  $r$  bitů s  $LL^2$  jedničkami

$r \geq LL^4$

$r < LL^4$

uložíme pozice všech 1 relativně ke bloku

$$\frac{n}{LL^2} \cdot o(LL) \in o\left(\frac{n}{LL}\right) \in o(n)$$

# částí celkem      pozice pozice uvnitř bloků + pointer na substruktury

$$\frac{n}{(L \cdot LL)^2} \cdot L \cdot LL \cdot L = \frac{n}{LL} = o(n)$$

↑ max. # takových bloků      ↑ #1 v bloku      ↑ pozice

$$\frac{n}{LL^4} \cdot LL^2 \cdot o(LL) \in o\left(\frac{n}{LL}\right) \in o(n)$$

↑ max. # bloků v podbloku      ↑ #1      ↑ pozice

pro dost velké n je  $r < \frac{1}{2}L$   
předpočítáme pro všechny strany podbloků  
 $2^{\frac{1}{2}L} \cdot \frac{1}{2}L \cdot LL \in o\left(\frac{n}{LL}\right) \in o(n)$   
↑ tvarů      ↑ dotazů      ↑ odpovědí

Celkem  $\Theta(n)$  bitů navíc  
Čas:  $\Theta(n)$  předvýpočet  
 $\Theta(1)$  dotaz