

# Úsporné datové struktúry

- Chceme pložiť prvky množiny  $C$  (treba jeden ze stromů na  $n$  vrcholech)  
a poskytnout téměř optimální prostor  $\dots$  OPT =  $\lceil \log |C| \rceil$  bitů } entropie při rovnoměrném rozdělení psů

• Přitom chceme rychlé operace

- Variace:
  - implicitní DS - je-li data sama ve vhodném pořadí (setříděné pole, heap) } OPT +  $O(1)$
  - úsporné (succinct) DS - OPT +  $o(OPT)$
  - kompaktní DS -  $O(OPT)$  ... ale pozor, třeba BVS nemusí být kompaktní kvůli pointerům - neproštan slova, ale bity  $\rightarrow$  lze obejít pakováním do natříděného pole

## Problém Reprerendace řetězce nad obecnou abecedou

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

$\hookrightarrow$  limitně se blížíme k  $\log_{10} 10$  b/znak (redundance = délka kódu - entropie)  
 $\dots$  ale ztrácíme lokální dekódovatelnost } jde k 0

Praktické řešení: dělíme na bloky o 9 znacích  $\dots$  délka kódu = 30  
 entropie  $9 \cdot \log_{10} 10 \approx 29.897$   
 $\rightarrow$  redundance  $\approx 0.103$

Ude se vejde do 1 slova  
 $\rightarrow$  potřebujeme  $O(1)$  času i prostoru

$\hookrightarrow$  pro string délky  $n$ :  $30 \cdot \lceil n/9 \rceil \leq \frac{30}{9}n + 30$   
 $\hookrightarrow$  redundance  $\leq \frac{0.103 \cdot n}{9} + 30$   
 $\approx 0.0114n$   
 $\rightarrow$  (a,b)  $\in X \times Y$  kódují po slokách  $\Rightarrow$  délky kódů i entropie se sčítou  $\rightarrow$  redundance také

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

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

Triviální řešení: rozdělíme na bloky o  $b$  bitech, za  $\forall$  blok připsáme 1 bitu značku, do posledního bloku padding  $10^k$   
 $\hookrightarrow$  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

$\hookrightarrow$  převod  $[B+1]^* \rightarrow [B]^*$

• nastavíme  $b \geq 2 \log n + 2 \Rightarrow B \geq 4n^2$

$\hookrightarrow n =$  horní mez na délku stringu

bloky se vejdou do  $O(1)$  slov RAMu

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 domyslíme nulový blok

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

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

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

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

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

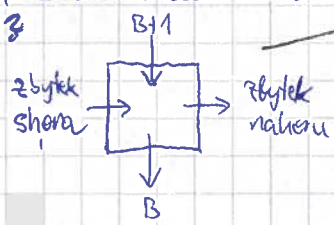
2. průchod:  $(B+3i)(B-3i) = B^2 - 9i^2 \leq B^2$

- redundance  $\leq 3b$  (1 blok zakrouhlení + 1 EOF + 1 extra na lichý # bloků)
- proudové kódování + dekódování
- lokální dekódování i změny v  $O(1)$  na RAMu [potřebují  $\Theta(\log n)$  bitů]  $\hookrightarrow O(1)$  sloz

zakrouhlení na celé bity je přibližně redundancí

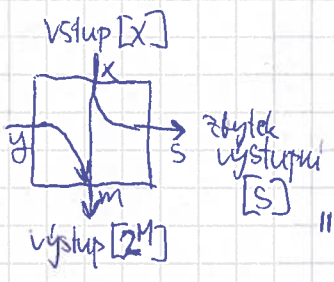
Postupně z toho odvodíme obecnou konverzi abeced...

- proč to vlastně fungovalo?

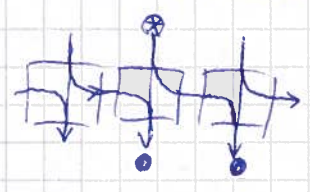


nechám zbytek malé množství redundance a to v příštím kroku přimíchám k většímu, čímž snížím redundanci... zbytek se posílají jen do směrem vztálosti  $\rightarrow$  lokálně dekódovatelné

Obecně: zbytek vstupní [Y]



... při čtení stále lok. dekódovatelné



pro dekódování \* stačí znát obě \*

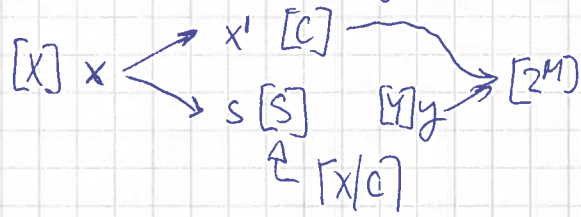
Lemma: Necht  $X, Y \leq 2^w$ . Pak  $\exists M, S$  a zobrazení  $f: [X] \times [Y] \rightarrow [2^M] \times [S]$  t.j.:

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

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

Děs zvolíme  $M$  (předějí)

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



dokud každý má výsledek porovnáme jen 1x, redundance už nemáme sčítat (steleskopují se)

• redundance ~~rozkladu~~ kódování  $x, y \rightarrow m^0$

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

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

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

z toho vyjde opačná nerovnost - správná je v 4 na straně 4

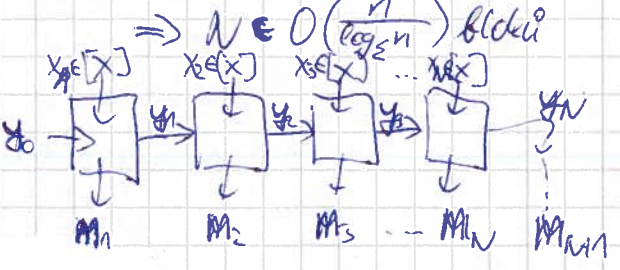
• redundance rozkladu  $x \rightarrow x', s$ :

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

• 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  $X \approx n^2$



$y_i \in [O(\sqrt{x})] = [O(n)]$   
 $m_i \in [O(\sqrt{x} \cdot \sqrt{x})] = [O(n^2)]$

... redundance  $O(1/n)$  v každém řádku

↳ celkově  $O(1)$ , což je milé

... lokální kódování / dekódování

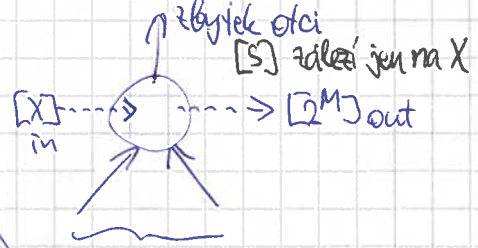
☹☹☹! Každý blok má jiné parametry!

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

↳ takže v sérii nevedlo, protože  $x = 2^{b+1}$ , takže jsme jednotlivá  $y_i$  uměli aproximovat aritmetickou posloupností ... ale dokonce to nebude fu igrat

Záchrana: stromové kódování

Opet zvolíme  $X_B = \Theta(n^2), N = O\left(\frac{n}{\log_{\epsilon} n}\right)$

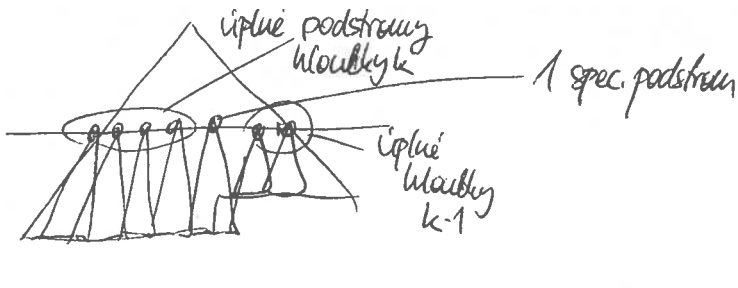


$[y_1] \times [y_2] \cong [y_1 \cdot y_2]$

• dekódování je stále lokální

(z otce stačí out, ten určuje zbytek)

zabraněním na bloky netřeba více dělat, takže pos. blok můžeme ušetřit a vyprázdnit



(4)  
 na každé hladině 3 typy vrcholů  
 pro vě si pamatujeme:  
 • # vrcholů tohoto typu  
 • adresu dat 1. vrcholu  
 • parametry mixéru  
 celkem  $O(\log n)$  konstant (slov)  
 ↓  
 + potřebujeme tabulku  $O(\log n)$  mocin  $|\Sigma|$  na extrakci symbolů = bloků

Dekódování: pozice ve streamu  $\rightarrow$  číslo bloku = pozice ve streamu



vše  $O(1)$

Hladina změna těž  $\sim O(1)$ .

Celková redundance  $O(1)$  [jako u předch. pohusu] +  $O(1)$   
 za kódování bloků

radonabílek posledního bloku (neúplného) (ten kódujeme obvykle)

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

v prostoru  $[n \log |\Sigma|] + O(1)$  bítů

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

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

⊗ Chceme  $\log \frac{1}{1-x} \leq cx$  pro nějaké  $c > 0$  a  $\forall x \in [0, 1/2]$ .

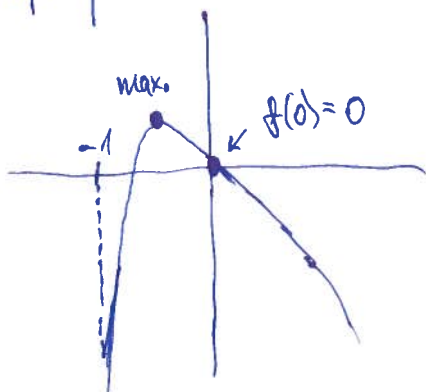
Tedy:  $-\log(1-x) \leq cx$

$\log(1-x) \geq -cx$

$\log(1+x) \geq cx$  pro  $x \in (-1/2, 0]$

Uvažme  $f(x) := \log(1+x) - cx$  ... chceme  $f(x) \geq 0$  pro  $x \in (-1/2, 0]$

Např. pro  $c=2$



Ukážeme, že pro dost velká  $c > 1$  má  $f$  maximum vlevo od  $x=1/2$ .

$$f'(x) = \frac{1}{1+x} - c \quad \dots \quad f'(0) < 0$$

Hledáme  $f'(x) = 0$ :

$$\frac{1}{1+x} - c = 0$$

$$\frac{1}{1+x} = c$$

$$1+x = \frac{1}{c}$$

$$x = \frac{1}{c} - 1$$

a to dostaneme dost velkou  $c$  libovolně blízko  $k-1$

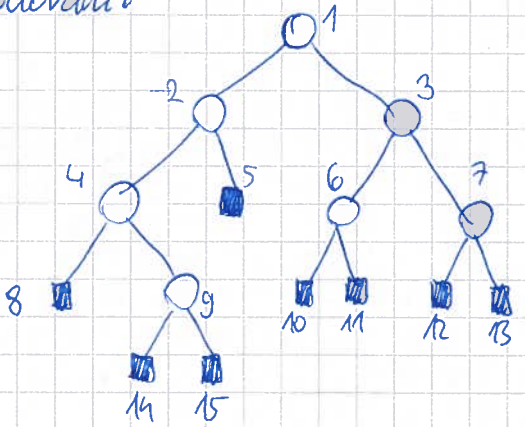
... stačí  $c > 2$

pro max. před  $-1/2$

# Úsporná reprezentace binárních stromů

- chceme uložit strukturu stromu, umět 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 a la halda
  - za  $\forall$  vrchol zapíšeme 1 bit  $\begin{cases} 0 & \text{pro externí} \\ 1 & \text{pro interní} \end{cases}$
- $\downarrow$   
 $2n+1$  bitů

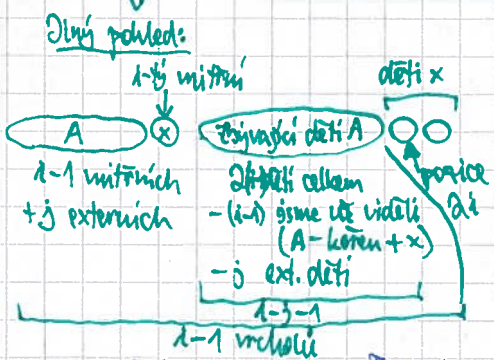
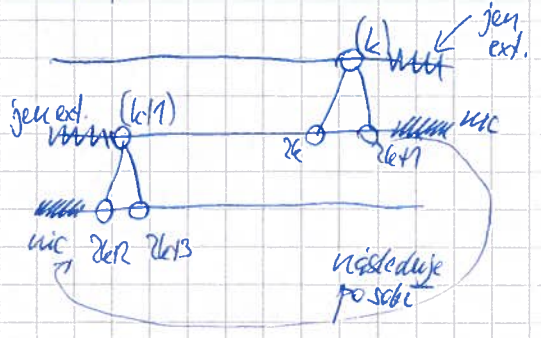
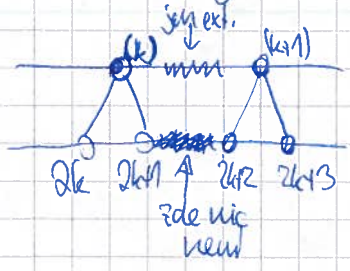
## • adresace:

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

$k \rightarrow k+1$  Indukce podle  $k$ : Pro  $k=1$  triviálně platí.

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

② přes hladinu



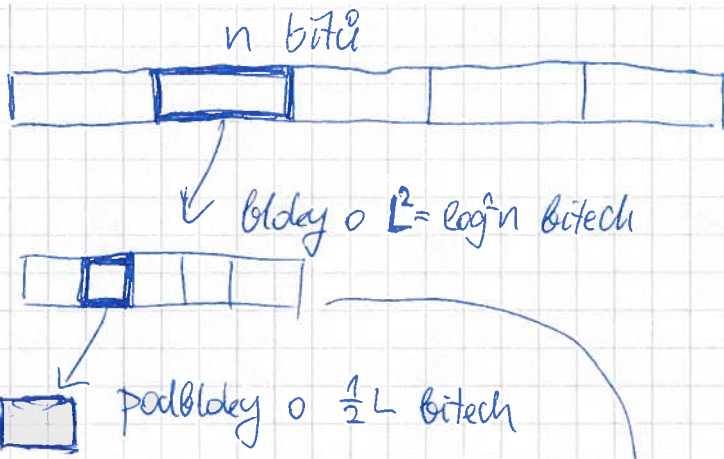
• Redukujeme na  $\text{Rank}_1$  a  $\text{Select}_1$  v posloupnosti  $N = 2n+1$  bitů

$\text{Rank}_1(i) =$   
 počet jedniček  
 na pozicích  $1-i$

$\text{Select}_1(i) =$   
 pozice  $i$ -té  
 jedničky

- Z toho:  $\text{Left}(i) = 2 \cdot \text{Rank}_1(i)$  (vstup i výstup jsou pozice v posl. vrcholů)
- $\text{Right}(i) = \text{Left}(i) + 1$
- $\text{Parent}(i) = \text{Select}_1(\lfloor i/2 \rfloor)$

Rank



pro # bloků předpočítáme #1 v předchozích blocích

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

bloky      bloky na #1

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

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

$\approx LL = \log \log n$

dekompozice...  
předpočítáme odpovědi na všechny dotazy

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

možných obsahů bloků =  $\sqrt{n}$   
možných dotazů  
velikost odpovědi

Select



bloky n bitů s  $L \cdot LL$  jedničkami

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

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

uložíme pozice všech 1

rozdělíme po jedničkách a zaznamenujeme jejich pozice + pointer na substrukturu

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

# částí      pozice + pointer



podbloky  $n'$  bitů s  $LL^2$  jedničkami

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

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

max. # takových bloků      #1 v bloku      pozice

$n' \geq LL^4$

uložíme pozice všech 1 relativně k bloku

$n' < LL^4$

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

# částí celkem      pozice uvnitř bloku + pointer na substrukturu

**2 přidali jsme index**  
Celkem  $\Theta(n)$  bitů navíc  
Čas:  $\Theta(n)$  předvýpočet  
 $\Theta(1)$  dotaz

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

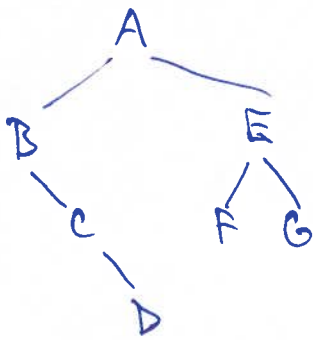
max. # bloků v podbloku      #1      pozice

pro dost velké n je  $n' < \frac{1}{2}L$

předpočítáme pro všechny strany podbloků  
 $2^{\frac{1}{2}L} \cdot LL^2 \cdot LL \in o\left(\frac{n}{LL}\right) \in o(n)$   
trvání dotazů odpovědi

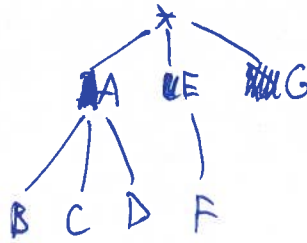
# Jiná reprezentace bin. stromů, kterými umí i počítat velikosti podstromů

bin. strom



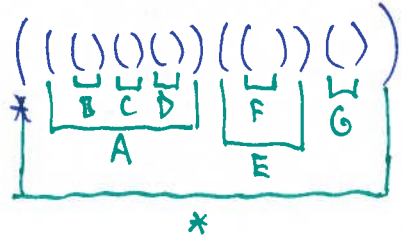
$n$  vrcholů

→ pěstovaný strom (kořen, pořadí dětí)



$n+1$  vrcholů

→ zavorčování



$2n+2$  závorek (bitů)

vrchol	→	vrchol	→	( ... )
levé dítě		první dítě		<del>Už</del> další znak, je-li to (
pravé dítě		další sourozenec		znak po párové ), je-li to (
rodič		předch. sourozenec nebo rodič		předch. znak je ) → párová ( jinak předch. (
velikost podstromu $x$		size(x) + size dalších sourozenců $x$		$\frac{1}{2}$ vzdálenosti do obalující )
je list?		je list bez dalších souroz.?		( )
# listů v podstromu				rank( ) obalující ) - rank( ) zde

Přidáme indexy pro Match ... párová závorka  
a Enclose ... obalující (

index pro Rank  
na posl. bitů generované  
orákulem

Umí se v  $O(n)$  prostoru  
 $O(n)$  čase na Build  
 $O(1)$  čase na dotaz

Opět je to dekompozice, ale výrazně složitější.