

Úsporné datové struktury

- Chceme počet prvků množiny C (treba jeden ze stromů na n vrcholech)
- a počet též optimální prostor ... $OPT = \lceil \log |C| \rceil$ bitů } entropie při rovnoměrném rozdělení psí
- Přitom chceme rychlé operace
- Variandy:
 - implicitní DS - jen data sama ve vhodném pořadí (seřazené pole, halda) } $OPT + O(1)$
 - úsporné (succinct) DS - $OPT + o(OPT)$
 - kompaktní DS - $O(OPT)$... ale pozor, třeba BVS nemusí být kompaktní kvůli pointerům - u počítačů slova, ale bity

rounding
lc.

Problém Repräsentace řetězce nad obecnou abecedou

Příklad: $[10] \rightarrow 4$ bity ($\log_2 10 \approx 3.322$)
 $[10]^2 \dots$ mohu uložit jako 2×4 , nebo jako $[100] \rightarrow 7$ bitů
 $[10]^k \dots$ uložíme jako $[10^k] \rightarrow k \cdot \log_2 10 + 1$ namísto $k \cdot (\log_2 10 + 1)$

\hookrightarrow limitně se blížíme k $\log_2 10$ 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 $9 \cdot \log_2 10 \approx 29.897$
 \rightarrow redundance ≈ 0.103

Ukážu vejde do 1 slova
 \rightarrow počítám s úniky $O(1)$ času i prostoru

\hookrightarrow pro string délky n $\approx 30 \cdot \lceil n/9 \rceil \leq \frac{30}{9}n + 30$

\hookrightarrow redundance $\leq 0.103 \cdot \frac{n}{9} + 30$

\hookrightarrow (a,b) $\in X \times Y$ kódují po slokách \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 "ještě to pokračuje" zabralo dost místa]

Triviální řešení: rozdělíme na bloky o b bitech, za τ bloků 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 $b \geq 2 \log n + 2 \Rightarrow B \geq 4n^2$

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. přechod	B	B+3	B-3	B+6	B-6	B+9	...	0
2. přechod	B	B	B	B	B	B	...	

sem si domyslíme nulový blok

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

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

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

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

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

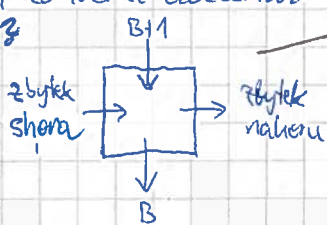
2. přechod: $(B+3i)(B-3i) = B^2 - 9i^2 \leq B^2$

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

zadržování na celé bits je příliš redundancí

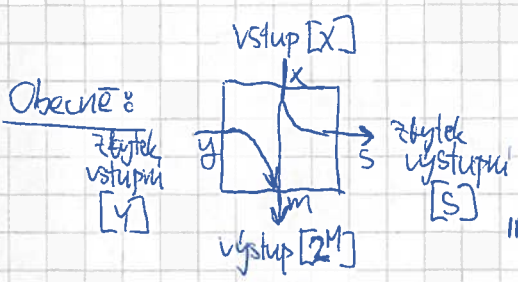
Postupně z toho odvodíme obecnou Lemberzi abecedu

- proč to vlastně fungovalo?



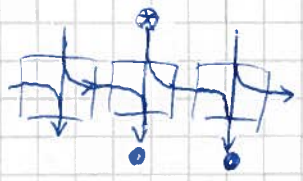
nedám zbytek malý množství informace a to v příštím bloku přimíchám k většině, čímž snížím redundanci

zbytek se posílají jen do směru určité vlnělosti \rightarrow lokálně dekódovatelné



"MIXÉR"

... při překodu stále lok. dekódovatelné



pro dekódování stačí znát obě

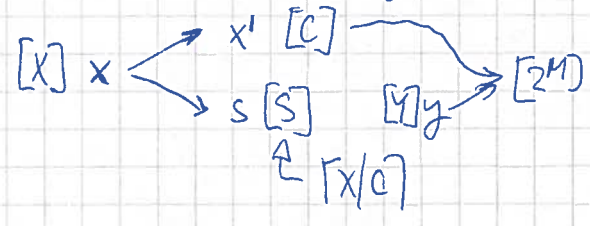
Lemma: Necht $X, Y \leq 2^W$. Pak $\exists M, S$ a zobrazení $f: [X] \times [Y] \rightarrow [2^M] \times [S]$ t.j.:

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

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

Důs 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 možný výsledek považujeme jen $1x$, redundance mírně se sčítá (steleskopují se)

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

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

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

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

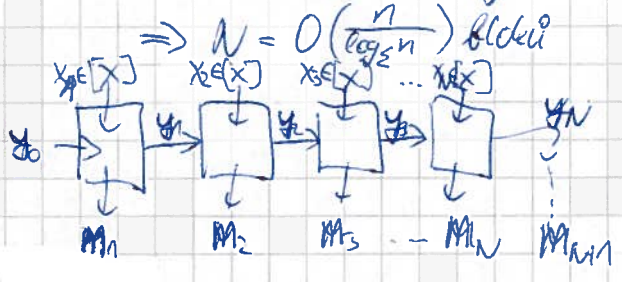
• redundance rozkladu $x \rightarrow x', s$

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

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

První pokus o reprezentaci stringu

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

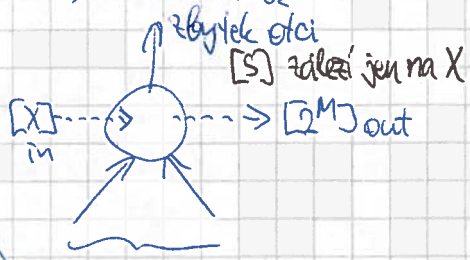
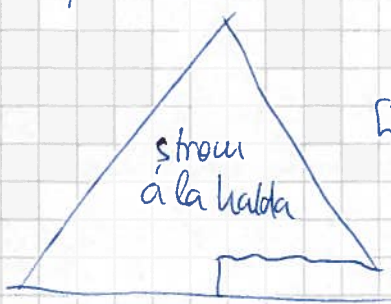


$y_i \in [O(\sqrt{X})] = [O(n)]$
 $m_i \in [O(\sqrt{X} \cdot \sqrt{X})] = [O(n^2)]$

••• redundance $O\left(\frac{1}{n}\right)$ v každém linku
 ↳ celkově $O(1)$, což je milé
 ••• lokální kódování / dekodování
 •••? Každý blok má jiné parametry?
 ↳ potřebujeme tabulku $O\left(\frac{n}{\log n}\right)$ konstant!
 ↳ tchle v serié nevadilo, protože $X = 2^b + 1$, takže jsme jednotlivá y_i uměli aproximovat aritmetickou posloupností ... ale obecně to nebude fungovat

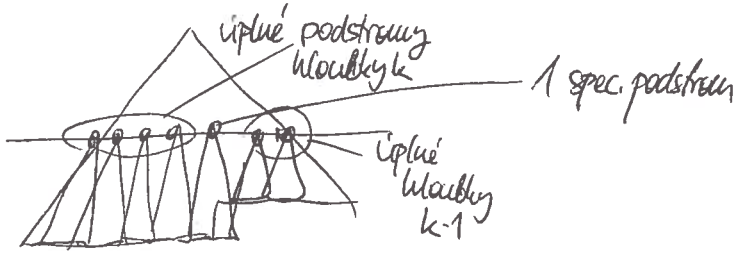
Záchrana: stromové kódování

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



zbytek ze signál
 $[Y_1] \times [Y_2] \cong [Y_1 \cdot Y_2]$

• dekodování je stále lokální (z otcí stačí out, ten určuje zbytek)



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)

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



$\forall i \in O(1)$
 + potřebujeme tabulku $O(\log n)$ mocnin $|\Sigma|$ na extrakci symbolů z bloku

Hladiní změna též $v O(1)$.

Celková redundance $O(1)$ [jako u předch. podusu] + $O(1)$
 A za kódování bloků B zadaná velikost posledního bloku (nepřímého) (ten kódujeme obvykle)

Věta: Na word-RAMu lze reprezentovat prvky $[\Sigma]^n$
 v prostoru $\lceil n \log \Sigma \rceil + O(1)$ bítů
 se čtením/zápisem prvku v čase $O(1)$
 s počítáním $O(\log n)$ konstant závislých na n a Σ .