

Succinct data structures.

Rank & Select

▷ rank(i) - #1 up to position i

obě jeho náklady $\approx O(1)$

▷ select(i) - i -th #1 (position)

representace binárního stromu (pouze jeho tvaru)

- každému uzlu přidáme dva potomky
- vyhovíme binární vektor kde přivodíme uzel jako 1, nebo potomci 0
→ indexují prvky po řádkách

$$\text{left_child} = 2 \cdot \text{rank}(i)$$

$$\text{right_child}(i) = 2 \cdot \text{rank}(i) + 1$$

$$\text{parent}(i) = \text{select}\left(\left\lfloor \frac{i}{2} \right\rfloor\right)$$

- reprezentace má $2n+1$ bitů

general trees - každý uzel reprezentován $1^m 0$ kde m je počet potomků

...

Data structure - who knows exactly

- several levels (constant) of indirect access
- logarithmic blocks with superblocks and precalculated table

Wavelet Trees

- konstruācija kark & rekur uz viēnā alfabēda
- visi teksti ir jednom stumū vāstū $O(\lg |\Sigma|)$
- kark stūa dāt - vādu pūvēvām vāstūk jūko vāstū dātū query
- rekur smūkū rēdā mātū

→ compressed suffix arrays

- katrū lēvā pūvām pūvū pūvūvū
- sūkū lēvū SA_k būvū pūvū i ir lēvū $k+1$
- lēvū mājū neighborhood fūktū Φ būvū, kē $SA_k[\Phi_k(i)] - 1$
(k $\Phi(i)$ dātūjū index pūvū jūko $SA = SA[i] + 1$)

Burrows - Wheeler & FM-Index

T = abracadabra P = bdabra

i	1	2	3	4	5	6	7	8	9	10	11
T[i]	a	b	r	a	c	a	d	a	b	r	a
SA	11	8	1	4	6	9	2	5	7	10	3
F	a	a	a	a	a	b	b	c	d	r	r
SA'	10	7	11	3	5	8	1	4	6	9	2
L	r	d	a	r	c	a	a	a	a	b	b

$\sim T[SA[i]]$ = první symboly lexikograficky seřazených řádků
 $\sim SA[i] - 1$ - poslední symboly lex. seřazených řádků
 $\sim T[SA'[i]] = BWT(T)$

C	a	b	c	d	r
	0	5	7	8	9

- kumulativní frekvence symbolů (první výkyt daného symbolu) is F

FM-Index:

 $Occ(X, c, i)$ - # počet výskytů c v prefixu $X[1...i]$

$$Sp_i = C[c] + Occ(L, c, Sp_{i-1}) + 1$$

$$Ep_i = C[c] + Occ(L, c, Ep_{i-1})$$

$$Occ_i = Ep_i - Sp_i + 1$$
 - počet výskytů symbolů (daného řádku)

interval $[Sp_i, Ep_i]$ řádků indexů v SA \rightarrow pozice ve stringu

příklad: $Sp_0 = 1$

$$Ep_0 = 11$$

$$Occ_0 = 11 - 1 + 1 = 11$$
 výskytů předchozího řádku

- pátým řádku od začátku! $\rightarrow c = a$

$$\left. \begin{array}{l} Sp_1 = 0 + 0 + 1 = 1 \\ Ep_1 = 0 + 5 = 5 \end{array} \right\} \text{interval } [1, 5]$$

$$c = r$$

$$\left. \begin{array}{l} Sp_2 = 9 + 0 + 1 = 10 \\ Ep_2 = 9 + 2 = 11 \end{array} \right\} \text{interval } [10, 11]$$