

# Wstęp do Programowania potok funkcyjny

Marcin Kubica

2018/2019

# Outline

## 1 Tablica sufiksowa

# Czym jest tablica sufiksowa

## Definition (Tablica sufiksowa)

- **Dane:** Napis  $x$  złożony z  $n$  liter.  
Zakładamy, że alfabet jest rozmiaru  $O(n)$ .
- Rozpatrujemy wszystkie sufiksy  $y$ .  
Chcemy je uporządkować leksykograficznie.
- **Wynik:**
  - tablica SUF zawierająca początkowe pozycje sufiksów w porządku leksykograficznym.
  - Tablica Rank odwrotna do tablicy SUF.
  - Tablica LCP określająca na ilu znakach są zgodne dwa kolejne sufiksy w SUF.

## Czym jest tablica sufiksowa

## Example

$x = \text{yabbadabbado}$

$i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$x.(i)$	y	a	b	b	a	d	a	b	b	a	d	o	$\emptyset$
Rank. $(i)$	12	1	7	5	3	9	2	8	6	4	10	11	0
SUF. $(i)$	12	1	6	4	9	3	8	2	7	5	10	11	0
LCP. $(i)$	0	5	1	2	0	3	1	4	0	1	0	0	

# Sortowanie pozycyjne

- Sortowanie kubełkowe  $O(n)$ .
- Sortowanie stabilne.
- Sortowanie pozycyjne krotek.
- Sortowanie pozycyjne słów.

# Algorytm Sandersa-Kärkkäinen-Burkhardta

- Dzielimy pozycje w tekście ze względu na reszy mod 3.  
 $B_k = \{i \in \{0, \dots, n\} : i \bmod 3 = k\}$
- Ograniczymy się do obliczenia tablic SUF i Rank dla sufiksów zaczynających się na pozycjach  $B_1 \cup B_2$ .  
 Potem rozszerzymy wyniki o sufiksy zaczynające się na pozycjach  $B_0$ .
- $r_k$  to napis zredukowany do kolejnych trójek znaków zaczynających się w  $B_k$ .  
 $r_k = [|(x.(k), x.(k+1), x.(k+2)); (x.(k+3), x.(k+4), x.(k+5)); \dots|]$   
 Ostatnią trójkę uzupełniamy znakiem  $\emptyset$  (obowiązkowo).

# Algorytm Sandersa-Kärkkäinen-Burkhardta

- Na przykład, dla

$$x = [y' a' b' b' a' d' a' b' b' a' d' o']:$$

$$r_0 = [(y', a', b'); (b', a', d'); (a', b', b'); (a', d', o')]; (\emptyset, \emptyset, \emptyset)$$

$$r_1 = [(a', b', b'); (a', d', a'); (b', b', a'); (d', o', \emptyset)]$$

$$r_2 = [(b', b', a'); (d', a', b'); (b', a', d'); (o', \emptyset, \emptyset)]$$

- Fakt: Kolejność sufiksów o początkach w  $B_1 \cup B_2$  jest taka sama, jak odpowiednich sufiksów  $r_{1,2} = r_1 \hat{ } r_2$ .
- Trójki tworzące  $r_{...}$  sortujemy i zamieniamy na liczby całkowite odpowiadające ich kolejności, tworząc  $r'_{...}$ .

Przykład:

$$r_{1,2} = [abb; ada; bba; do; bba; dab; bad; o]$$

$$r'_{1,2} = [1; 2; 4; 6; 4; 5; 3; 7]$$

## Algorytm Sandersa-Kärkkäinen-Burkhardta

- Rekurencyjnie wyznaczmy  $SUF'_{1,2}$  i  $Rank'_{1,2}$  dla  $r'_{1,2}$ .

$$r'_{1,2} = [1; 2; 4; 6; 4; 5; 3; 7]$$

$$Rank'_{1,2} = [1; 2; 5; 7; 4; 6; 3; 8]$$

$$SUF'_{1,2} = [0; 1; 6; 4; 2; 5; 3; 7]$$

- Przenumerujemy i permutujemy pozycje zgodnie z ich położeniem w  $x$ , otrzymując  $SUF_{1,2}$  i  $Rank_{1,2}$ :

$i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$x.(i)$	$y$	$a$	$b$	$b$	$a$	$d$	$a$	$b$	$b$	$a$	$d$	$o$	$\emptyset$
$Rank_{1,2}$	—	1	4	—	2	6	—	5	3	—	7	8	—
$SUF_{1,2}$	1	4	8	2	7	5	10	11					



# Algorytm Sandersa-Kärkkäinen-Burkhardta

- Znamy względną kolejność leksykograficzną sufiksów zaczynających się na pozycjach z  $B_1 \cup B_2$ .
- Sufiksy  $x$  zaczynające się na pozycjach  $i$  i  $j$  są w takim samym porządku leksykograficznym, jak:  
 $(x.(i), \text{Rank}(i + 1))$  i  $(x.(j), \text{Rank}(j + 1))$ , oraz  
 $(x.(i), x.(i + 1), \text{Rank}(i + 2))$  i  $(x.(j), x.(j + 1), \text{Rank}(j + 2))$ .
- Sufiksy zaczynające się na pozycjach z  $B_0$  sortujemy porównując  $(x.(i), \text{Rank}_{1,2}(i + 1))$ .  
 Przykład:  $r'_0 = [ ('y', 1); ('b', 2); ('a', 5); ('a', 7) ]$   
 $\text{Rank}'_0 = [ |4; 3; 1; 2; 0 | ]$   
 $\text{SUF}'_0 = [ |4; 2; 3; 1; 0 | ]$

## Algorytm Sandersa-Kärkkäinen-Burkhardta

- $\text{Rank}'_0 = [4; 3; 1; 2; 0]$   
 $\text{SUF}'_0 = [4; 2; 3; 1; 0]$
- Przenumerujemy pozycje zgodnie z ich położeniem w  $x$ , otrzymując  $\text{SUF}_0$  i  $\text{Rank}_0$ :

$i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$x.(i)$	$y$	$a$	$b$	$b$	$a$	$d$	$a$	$b$	$b$	$a$	$d$	$o$	$\emptyset$
$\text{Rank}_0$	4	—	—	3	—	—	1	—	—	2	—	—	0
$\text{SUF}_0$	12	6	9	3	0								

# Algorytm Sandersa-Kärkkäinen-Burkhardta

## Scalanie:

Mamy dwa posortowane leksykograficznie ciągi sfiksów do scalenia.

Ciąg zaczynający się na pozycji  $i \in B_0$  porównujemy z:

- $j \in B_1$  — porównując:  
 $(x.(i), \text{Rank}(i + 1))$  i  $(x.(j), \text{Rank}(j + 1))$ ,
- $j \in B_2$  — porównując:  
 $(x.(i), x.(i + 1), \text{Rank}(i + 2))$  i  $(x.(j), x.(j + 1), \text{Rank}(j + 2))$ .

## Algorytm Sandersa-Kärkkäinen-Burkhardta

## Example

$i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$x.(i)$	$y$	$a$	$b$	$b$	$a$	$d$	$a$	$b$	$b$	$a$	$d$	$o$	$\emptyset$
$\text{Rank}_{1,2}$	—	1	4	—	2	6	—	5	3	—	7	8	—
$\text{SUF}_{1,2}$	1	4	8	2	7	5	10	11					
$\text{Rank}_0$	4	—	—	3	—	—	1	—	—	2	—	—	0
$\text{SUF}_0$	12	6	9	3	0								
$\text{Rank}$	12	1	7	5	3	9	2	8	6	4	10	11	0
$\text{SUF}$	12	1	6	4	9	3	8	2	7	5	10	11	0

# Algorytm Sandersa-Kärkkäinen-Burkhardta

- Złożoność czasowa:

$$T(n) = O(n) + T\left(\frac{2}{3}n\right) = O(n)$$

## Wyznaczenie tablicy LCP

- Oznaczenie  $LCP'.(i) = LCP.(Rank.(i))$ .
- Fakt:  $LCP'.(i + 1) \geq LCP'.(i) - 1$ .

$i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$x.(i)$	y	a	b	b	a	d	a	b	b	a	d	o	$\emptyset$
$Rank.(i)$	12	1	7	5	3	9	2	8	6	4	10	11	0
$LCP'.(i)$		5	4	3	2	1	1	0	1	0	0	0	0
$SUF.(i)$	12	1	6	4	9	3	8	2	7	5	10	11	0
$LCP.(i)$	0	5	1	2	0	3	1	4	0	1	0	0	

# Wyszukiwanie wzorców

- Przez bisekcję wyszukujemy wystąpienia pierwszego znaku — spójny fragment SUF.
- Drugi znak — zawężamy fragment; itd.
- Złożoność czasowa  $O(m \log m)$ .

# Kompresja tekstu

- Idea kompresji LZ77.
- Tablica LPF.
- Jak obliczyć LPF na podstawie SUF i LCP.
- Fakt: LPF jest permutacją LCP.