

# Суффиксное дерево (Тусар Вайнер, 1973)

$T = \text{bananas}$

$T_1 = \text{bananas}$  §

$T_2 = \text{ananas}$  §

$T_3 = \text{nanas}$  §

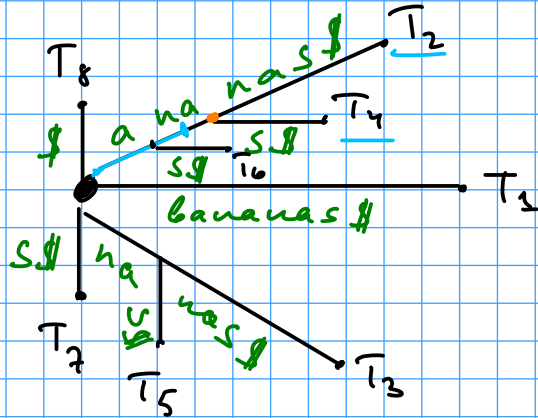
$T_4 = \text{anas}$  §

$T_5 = \text{nas}$  §

$T_6 = \text{as}$  §

$T_7 = \text{s}$  §

$T_8 = \text{}$  §



NB: если добавить вершины только при разветвлении, то потребуются  $O(n)$  памяти

Поиск: ищем образец в доре.

$P = an$  - выводится в доре  
в поддере.

Время  $O(|P| + \# \text{вхождений})$

Факт: Суффиксное дерево можно построить за линейное время.

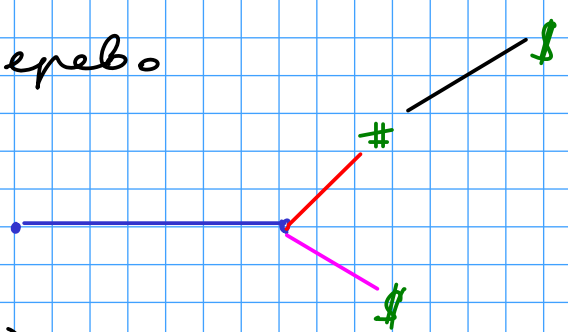
Д. Кнут (1970): // не верно

Задача поиска наибольшей общей подстроки двух строк, то она решается за  $\Omega(n \log n)$

Вход:  $T_1 \sim T_2$  Выход: общая подстрока

$T = T_1 \# T_2 \text{ §}$

# Построение суфр. дерева



⇒ Решение за  $O(n)$

Минусы:

- довольно сложные алгоритмы

- [намерею  $O(|T| \cdot |\Sigma|)$

[поиск  $O(|P| \log |\Sigma|)$  ← алгоритм

## Суффиксный массив

$T = \text{bananas}$

$T_1 = \text{bananas}$

$T_2 = \text{ananas}$

$T_3 = \text{nanas}$

$T_4 = \text{anas}$

$T_5 = \text{nas}$

$T_6 = \text{as}$

$T_7 = \text{s}$

Оптимизируем массив  $T_i$

2461357 SA

SA можно построить за  $O(n)$ .

Например 2/3 С.Д.

Факт

∃ алгоритм построения за  $O(n)$  независимо от  $|\Sigma|$

# Точное время $O(n \log n)$

bananas  
 ananas  
 nanas  
 anas  
 nas  
 as  
 s

сортровка  $O(n)$   
 (возвратом)

Шаг 1

Анализ  $\bar{Z}_1 = \bar{Z}$

abns  $\Rightarrow$   $an \leftrightarrow (1, 3)$   
 $as \leftrightarrow (1, 4)$   
 $i$

Сортируем по  
 второму символу  
 (улучш. сорт.)

Шаг 2

Анализ  $\bar{Z}_2 \subseteq \bar{Z}_1 \times \bar{Z}_1$   
 $|\bar{Z}_2| \leq n$

$T_1$   
 $T_2$

ananas  
 anas  
 as  
 bananas  
 nanas  
 nas  
 s

an  $\leftrightarrow$  1  
 as  $\leftrightarrow$  2  
 bc  $\leftrightarrow$  3  
 na  $\leftrightarrow$  4  
 s  $\leftrightarrow$  5  
 $\leftarrow$  6

ana  $\leftrightarrow$  (1, 1)  
 bana  $\leftrightarrow$  (3, 4)

Шаг 3

Анализ  $\bar{Z}_3 \subseteq \bar{Z}_2 \times \bar{Z}_2$   
 $|\bar{Z}_3| \leq n$

ananas  
 anas  
 as  
 bananas  
 nanas  
 nas  
 s

$T_i$

$\in \bar{Z}_{k+1}$

a b x y z  
 $\uparrow \quad \uparrow$   
 $\bar{Z}_k \quad \bar{Z}_k$

$T_{i+2}^{k-2}$

b x y z

$\log n$  шагов

Время  $O(n \log n)$

# Источники в SA

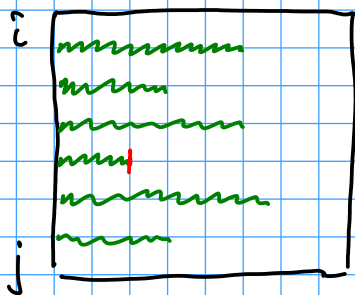
Наибольший источник: для точек  $\Rightarrow$   
 $O(|P| \cdot \log |T|)$

Можно переформулировать  $\Leftarrow$   $O(|P| + \log |T|)$

SA	ISA	
2	4	$SA[ISA[i]] = i$
4	1	$SA[i] = j \Leftrightarrow$
6	5	$ISA[j] = i$
1	2	
3	6	
5	3	
7	7	

$\equiv$   $lcp(S_1, S_2)$  - наибольший общий префикс

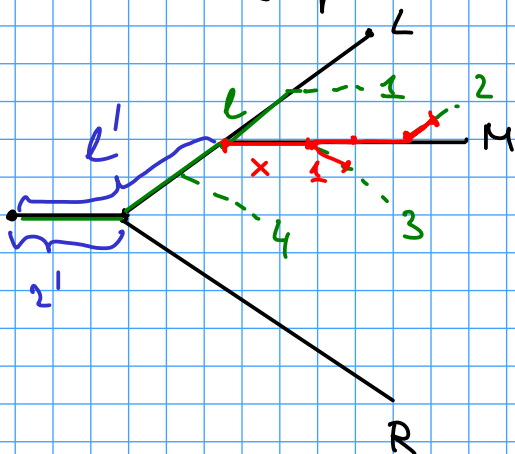
уб:  $|lcp(T_{SA[i]}, T_{SA[j]})| =$   
 $= \min_{i \leq k < j} |lcp(T_{SA[k]}, T_{SA[k+1]})|$



Как устроит точка:

$Find(P, L, R, l, r)$

вспомог. числа



$$M = \lfloor \frac{L+R}{2} \rfloor$$

$$l = |lcp(L, P)|$$

$$r = |lcp(P, R)|$$

с помощью

$$l' = |lcp(L, M)|$$

$$r' = |lcp(M, R)|$$

\*

•  $l \geq r$

•  $l < r$

-  $l' > l$  ④

симметрично

go right

-  $l' < l$  ①

$O(|P| + \log |T|)$

go left

-  $l' = l$  ② vs ③

читаю символы  $P$ , пока

не отклонимся от  $M$

(таким образом получим  $l$  или  $r$   
где следующего шага)

\* как вычислить  $lcp(T_i, T_j)$

1. Посчитать все пары  $lcp(T_i, T_j)$   
 $O(n^2)$

2. Пройтись только по тем  
парам, кот. встречаются в поиске.

$O(n^2)$  времени  
 $O(n)$  памяти

3. Давать вычисления  $lcp$  где  
используются. сгруппировать.

$$\overline{LCP}[i] = |lcp(T_{SACi}, T_{SACi+1})|$$

можно вычислить  $|lcp(T_i, T_j)|$  з/в RMQ  
либо запомнить табл. да  $O(n)$

Алгоритм (Arimura, Arikawa, Lee, Kosai, Park)

$$LCP[1] = l = |lcp(T_{SAC[1]}, T_{SAC[2]})|$$

for  $i=2$  to  $n$ :

$$l = \max(0, l-1)$$

$$\text{while } T_i[l+1] = T_{SAC[i]SAC[i+1]}[l+1]$$

$$l = l+1$$

$$LCP[i] = l$$

# сущности

регулярности  $j = i$

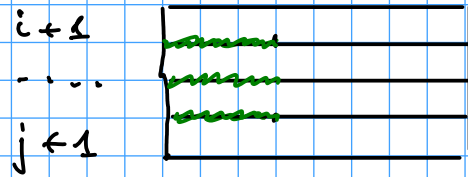
в SA

$$SAC[1] = i$$

$$SAC[2] = j$$

$$|lcp(T_i, T_j)| = l$$

$$|lcp(T_{i+1}, T_{j+1})| = l-1$$



$\overline{LCP}$  за  $O(n)$