

# Суффиксное дерево

Дугер Вайнер, 1973

mississippi

$T = \text{banana}\$$

$T_1 = \text{banana}\$$

$T_2 = \text{nana}\$$

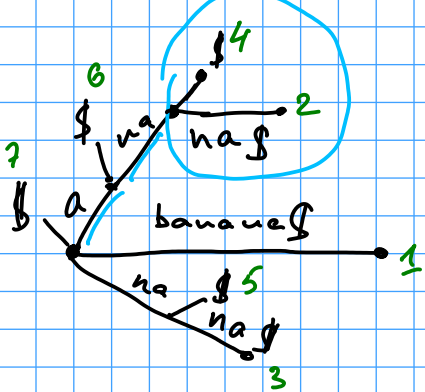
$T_3 = \text{na}\$$

$T_4 = \$$

$T_5 = \text{ana}\$$

$T_6 = \text{aa}\$$

$T_7 = \text{a}\$$



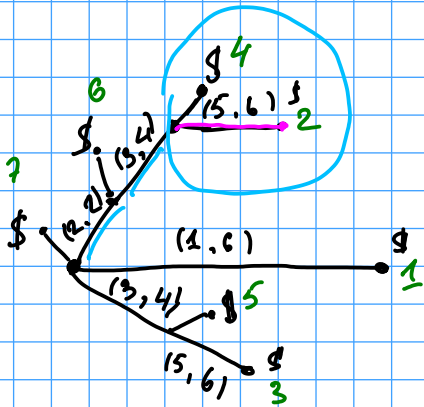
$P = \text{ana}$

Проверка входжений  $P \in T$ ?

$O(|P|)$

Почему -?

## ≡ Схема суффиксного дерева



Помеха  $O(\# \text{ узлов вершин})$

$\# \text{ узлов. вершин} \leq \# \text{ листьев}$

$O(|T|)$

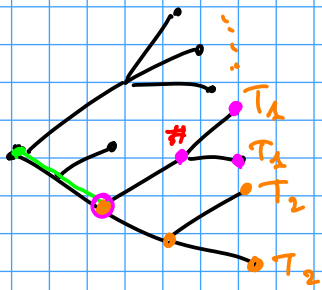
Почему всех входжений  $P \in T$

$O(|P| + \# \text{ входжений})$

Почему наибольшей длины подстроки.

Вход:  $T_1 \sim T_2$

Построим  $ST$  где  $T_1 \neq T_2 \text{ \$}$



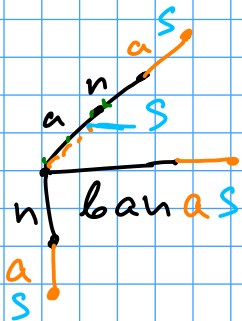
$$O(|T_1| + |T_2|)$$

Книга '20. Гинзбург, 200  $\Omega(n \log n)$

## Алгоритм Уиконена

известное СД - СД где + есть  $\$$

Шаг  $i$ : Построение НСД где  $T[1..i]$   
из дерева где  $T[1..i-1]$



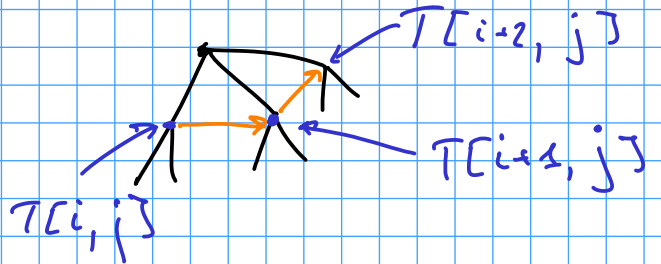
Внутри шага:

пройдем все суффиксы  
 $T_1, T_2, \dots, T_{i-1}$  на одну букву  
и доминировать  $T_i$

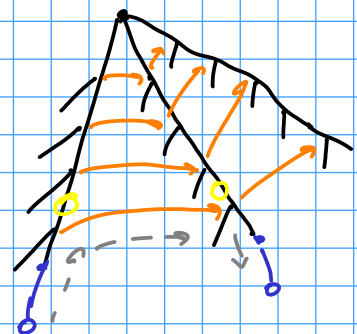
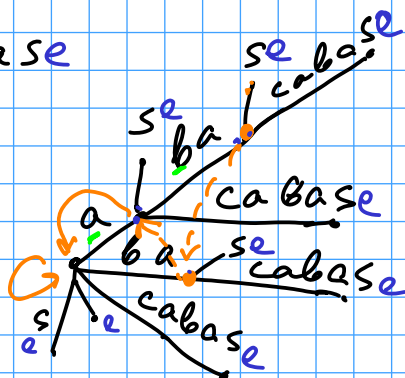
$$\text{Время: } \sum_{i=1}^n \sum_{j=1}^i (i-j) = O(n^3)$$

Типы прохода:

1. проходим места
2. ветвление
3. "пустое прохождение"



abacabase



Переход к след суффиксу:

1. вверх до первой вершины // глубина -1
2. переход по средине // глубина -1
3. вниз до 'места' // глубина ↑↑

$O(n)$  на шаг.

Сложность  $O(n^2)$

**NB:** Важно, что на шаге 3 при спуске вниз по ребрам мы можем найти конец следующего суффикса и количество операций пропорционально не по-вы ребер, а не длине строки. Для этого будем смотреть только на первую букву ребра.

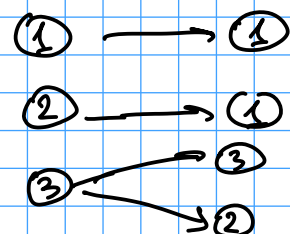
Как поддерживать суффиксные ссылки?

Если произошло ветвление, то запомнить этот факт и парсить ссылку после того, как найдём её конец (при переходе к следующему суффиксу).

Типы прояснения:

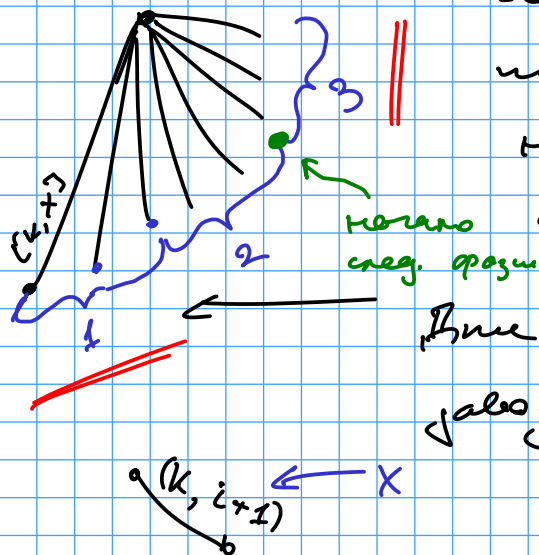
1. прояснение места
2. ветвление
3. "пустое прояснение"

Прояснение:



Для конкретного суффикса

## Внутри графа



Если  $T_j$  предельно  
маленькая иrogenнелая  $\Rightarrow$   
можно перейти к  
след. графу.

Для иrogenнелая мест  
давогда переменной.

$$x = 0$$

$$k = 1$$

for  $i = 1$  to  $n$ :

$$x = x + 1$$

for  $j = k$  to  $i$ :

if ③:

$$k = j$$

break

split

$\Rightarrow$  сложность  $O(n)$

Классификация структурных деревьев:

1. Много ветви с дочерней фрагментацией
2. Зависит от  $|\Sigma|$

# Сурриксный массив

1 бананас  
2 ананас  
3 наанас  
4 ананас  
5 наанас  
6 аанас  
7 санас

ананас  
ананас  
анас  
бананас  
наанас  
наанас  
санас

rank

2  
4  
6  
1  
3  
5  
7

SA  
ISA

4  
1  
5  
2  
6  
8  
7

$P = \text{нап}$

Поиск шаблона бинарным поиском

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

Можно ускорить до  $O(\log |T| + |P|)$

Построение  $\Sigma$  и  $\log u$

бананас  
ананас  
наанас  
ананас  
наанас  
аанас  
санас

ананас  
ананас  
анас  
бананас  
наанас  
наанас  
санас

$\Sigma$   
 $\Sigma^1$   
 $\Sigma \times \Sigma$

1  
1  
2  
3  
4  
4  
5

1  
2  
4  
5  
6

ананас  
ананас  
анас  
бананас  
наанас  
наанас  
санас

$\Sigma^1$   
 $\Sigma \times \Sigma$

1  
2  
3  
4  
5  
6  
2

ананас  
ананас  
анас  
бананас  
наанас  
наанас  
санас

$\leftarrow T_1$   
 $\leftarrow T_3$

$$\Sigma'' \subset \Sigma^1 \times \Sigma^1$$

$$O(u \log u)$$

# Naive binary search

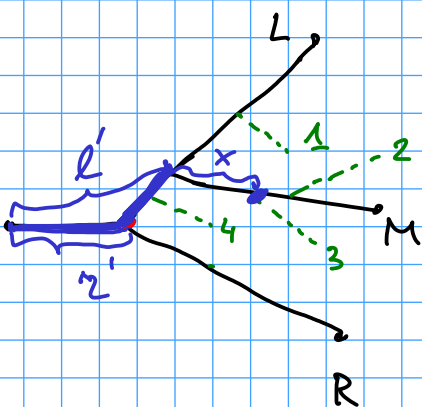
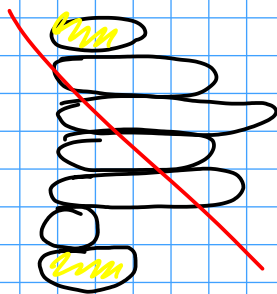
$$O(|P| \cdot \log |T|)$$

Зад. 1

↓	ababc	$s_1$
	abbac	$s_2$
	abcab	$s_3$
	abcd a	$s_4$
	abc ed	$s_5$
	abcd ad	$s_6$

longest common prefix = lcp

$$|lcp(s_1, s_6)| = \min_i |lcp(s_i, s_{i+1})|$$



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

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

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

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

6 уровней  
?

$l \geq z$  : ①, ②, ③, ④

$l' > l \Rightarrow$  ④, go right

$l' < l \Rightarrow$  ①, go left

$\rightarrow l = l' \Rightarrow$  ② vs ③

x символов не имеет  
на этом месте

$$z \leftarrow z' + x$$

переходим

②

$l < z$ :

симметрично.

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

$$l \leftarrow l' + x$$

переходим

③

Как считать  $lcp(L, M)$  ?

1. Прегноцирйт  $lcp(T_i, T_j) \quad O(|T|^2)$
2. Прегноцирйт только где  $tex(i, j)$ ,  
которые встречаются в префиксе  $lcp$ .  
 $O(|T|^2)$   
 $O(|T|)$  номер.
3. Добавьте номерам массив  $\overline{LCP}$ :

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

УТВ 1 + RMQ  $\Rightarrow O(1) \Rightarrow$

по  $\overline{LCP}$  считать

$$|lcp(T_{SA[i]}, T_{SA[j]})| = \min_{i \leq k < j} \overline{LCP}[k]$$

abcabc  
 abcdef  
 abcdeg  
 bcab  
 bcad  
 beap

$O(n^2)$

Алгоритмы:  
 [Arimura,  
 Arikawa,  
 Lee,  
 Kosai,  
 Park]

SA  $\dots \overline{[k] \dots [l]} \dots$

$$|lcp(T_{SA[i]}, T_{SA[i+1]})| = L > 0$$

$\Downarrow$

$$|lcp(T_{k+1}, T_{l+1})| = L - 1$$

$$\overline{LCP}[i] = L$$

$\geq L - 1$

SA  $\dots \overline{[k+1] \dots [l+1]}$

$$\overline{LCP}[ISA[k+1]] \geq L - 1$$

$$L = |lcp(T_1, T_{SA[ISA(1)+1]})|$$

for  $i = 2$  to  $n$ :

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

$$\text{while } T_i[L+1] = T_{SA[ISA(i)+1]}[L+1]$$

$$L = L + 1$$

$O(n)$

