

Алгоритмы и структуры данных

Суффиксное дерево и суффиксный массив

CS Center, Новосибирск

Суффиксное дерево для задачи LCP

Даны $\alpha, \beta \in \Sigma^*$

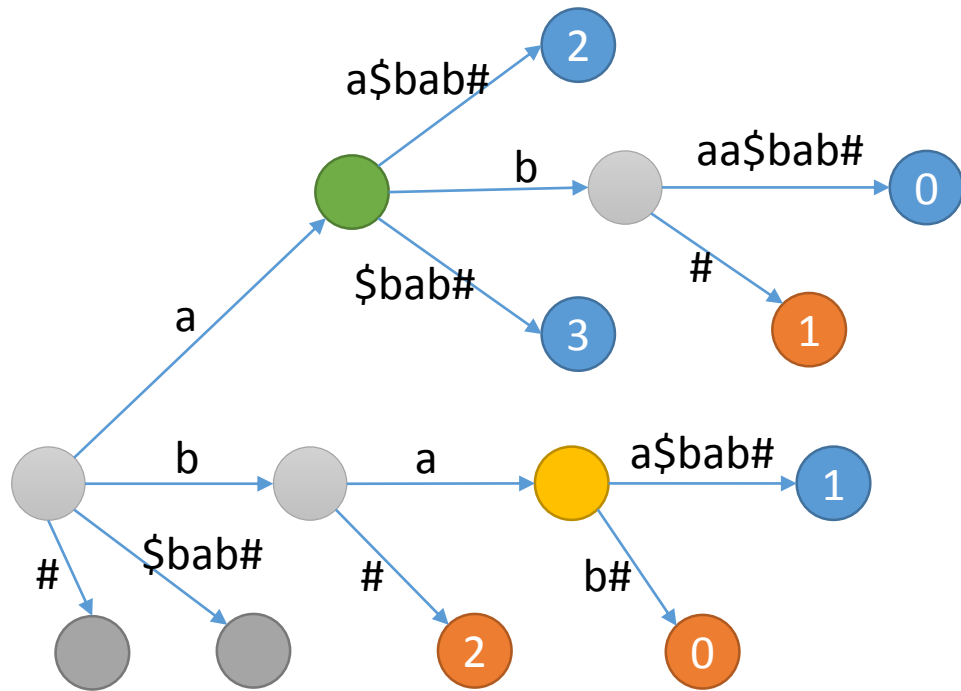
- Запрос $(i, j) \rightarrow |\text{LCP}(\alpha[i:], \beta[j:])| = \max\{k \mid \alpha[i:i+k) = \beta[j:j+k)\}$

Возьмём $S = \alpha\$\beta\#$ ($\$, \#$ нет в α и β)

- $S[i:] = \alpha[i:)\$\beta\#$
- $S[j+|\alpha|+1:] = \beta[j:)\#$
- $\text{LCP}(\alpha[i:], \beta[j:]) = \text{LCP}(\alpha[i:)\$\beta\#, \beta[j:)\#) = \text{LCP}(S[i:], S[j+|\alpha|+1:])$
- Построим суффиксное дерево для S
- Пусть $v[i] := \text{node}(S[i:])$
- $|\text{LCP}(\alpha[i:], \beta[j:])| = \text{CharDepth}(\text{LCA}(v[i], v[j+|\alpha|+1]))$

Суффиксное дерево для задачи LCP: пример

$\alpha = \text{abaa}$, $\beta = \text{bab}$, $S = \text{abaa}\$bab\#$



$LCP(i=1, j=0) = 2$ // baa , bab

$LCP(i=2, j=1) = 1$ // aa , ab

Суффиксное дерево: затраты

- При $|\Sigma| = \text{const}$
 - $O(|T|)$ времени
 - $O(|T|)$ памяти
 - $O(|P|)$ на поиск шаблона
- Зависит от $|\Sigma|$
- Что-то более компактное?
 - Суффиксный массив

Лексикографический порядок

- $S1 < S2$ если
 - $(S1 = \varepsilon \text{ и } S2 \neq \varepsilon)$
 - или
 - $(S1 = a\alpha \text{ и } S2 = b\beta \text{ и } a < b)$
 - или
 - $(S1 = a\alpha \text{ и } S2 = a\beta \text{ и } \alpha < \beta)$
- это линейный порядок на Σ^*

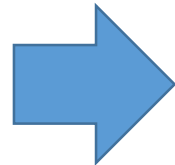
Суффиксный массив

- Выпишем все $\langle T[i:], i \rangle, i = [0 : N)$
- Упорядочим по $T[i:]$
 - Получившаяся перестановка i — суффиксный массив

• Пример: $T=abacaba$

Рассмотрим $\langle T[i:], i \rangle$

- $\langle abacaba, 0 \rangle$
- $\langle bacaba, 1 \rangle$
- $\langle acaba, 2 \rangle$
- $\langle caba, 3 \rangle$
- $\langle aba, 4 \rangle$
- $\langle ba, 5 \rangle$
- $\langle a, 6 \rangle$



Упорядочим лексикографически

- $\langle a, 6 \rangle$
- $\langle aba, 4 \rangle$
- $\langle abacaba, 0 \rangle$
- $\langle acaba, 2 \rangle$
- $\langle ba, 5 \rangle$
- $\langle bacaba, 1 \rangle$
- $\langle caba, 3 \rangle$

Чтобы не запутаться два смысла чисел:

- pos — позиция в строке
- $rank$ — индекс в суффиксном массиве

SA: $rank \rightarrow pos$

ISA: $pos \rightarrow rank$

Суффиксный массив: $[6, 4, 0, 2, 5, 1, 3]$

Суффиксный массив для поиска шаблона

- Подстрока – префикс суффикса
 - Найти суффикс, префикс которого совпадает с шаблоном
 - Идея: бинпоиск // $T[SA[l] :) < P \leq T[SA[r] :)$
- Пример $T=abacaba$, $P=abac$

$\langle a, 6 \rangle$	$\langle a, 6 \rangle$	$\langle a, 6 \rangle$
$\langle aba, 4 \rangle$	$\langle \color{red}{aba}, 4 \rangle$	$\langle aba, 4 \rangle$
$\langle abacaba, 0 \rangle$	$\langle abacaba, 0 \rangle$	$\langle \color{red}{abacaba}, 0 \rangle$
$\langle \color{red}{acaba}, 2 \rangle$	$\langle acaba, 2 \rangle$	$\langle acaba, 2 \rangle$
$\langle ba, 5 \rangle$	$\langle ba, 5 \rangle$	$\langle ba, 5 \rangle$
$\langle bacaba, 1 \rangle$	$\langle bacaba, 1 \rangle$	$\langle bacaba, 1 \rangle$
$\langle caba, 3 \rangle$	$\langle caba, 3 \rangle$	$\langle caba, 3 \rangle$

- Сложность $O(|P| * \log |T|)$
 - Требуется улучшить

Суффиксный массив для поиска шаблона: ОПТИМИЗАЦИЯ

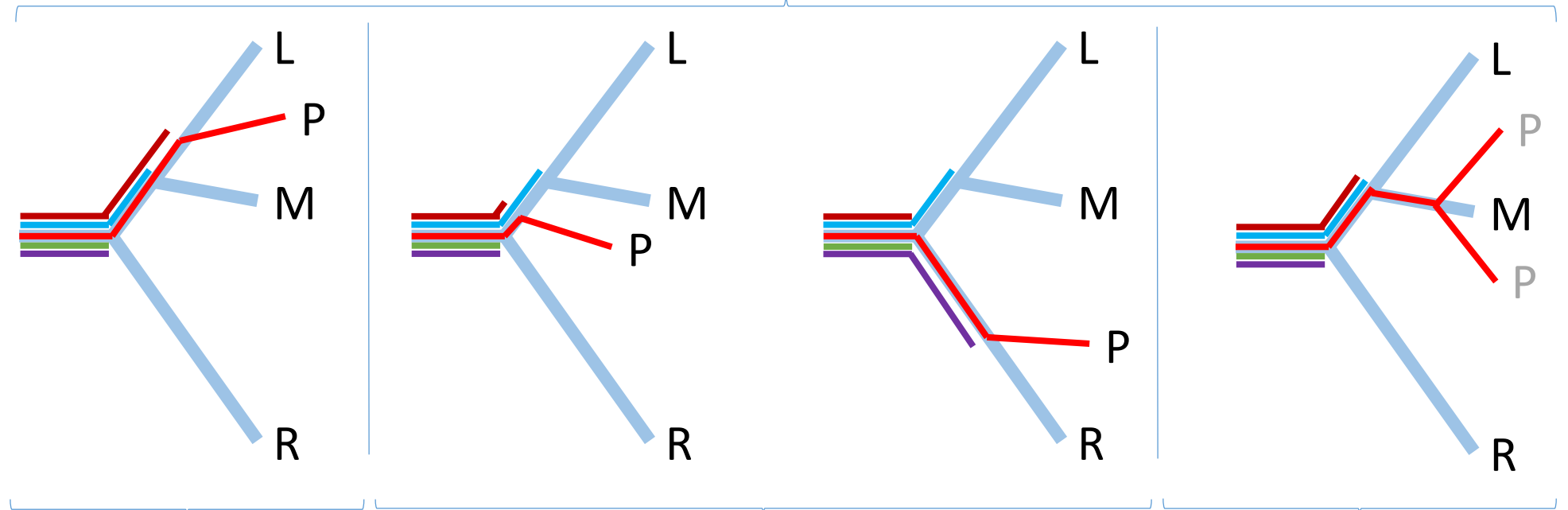
$L = T[SA[l] :)$
 $R = T[SA[r] :)$
 $M = T[SA[m] :)$

- $L_P = |LCP(L, P)|$
- $R_P = |LCP(R, P)|$
- $L_M = |LCP(L, M)|$
- $R_M = |LCP(R, M)|$

Вычислить заранее

Сложность $O(|P| + \log|T|)$

$L_M > R_M$



$L_P > L_M$
 $R := M$
 $R_P := L_M$

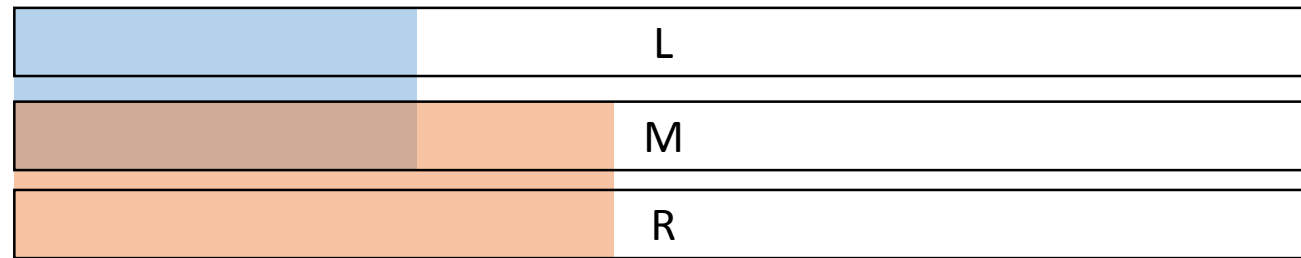
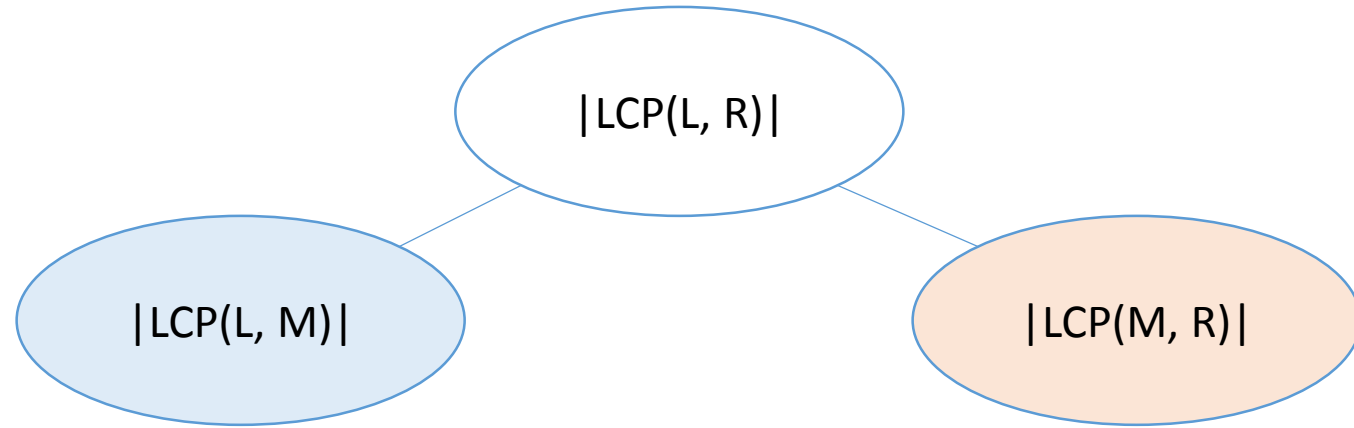
$L_P < L_M$
 $L := M$

$L_P = L_M$
 Сравниваем P и M,
 начиная с позиции L_P
 L или $R := M$
 L_P или $R_P := matched$

LCP-массив

- $L_M = |LCP(L, M)|$
- $R_M = |LCP(R, M)|$

Вычислить заранее



$$|LCP(L, R)| \geq \min(|LCP(L, M)|, |LCP(M, R)|)$$

$$L < M < R$$

$$|LCP(L, R)| = \min(|LCP(L, M)|, |LCP(M, R)|)$$

$$LCP[i] = |LCP(T[SA[i] :), T[SA[i+1] :)) |$$

Суффиксный массив: построение (KMR)

- $L(\beta) \in \{0, \dots, |\beta|\}$ – метка
 - $L(\beta) < L(\gamma) \Leftrightarrow \beta < \gamma$
- $\alpha' := \alpha \text{ \$ } ^{|\alpha|}$
- Требуется отсортировать $\{ \alpha'[i : i+|\alpha|), i \in [0 : |\alpha|) \}$
- Шаг 0: сортируем $\{ \alpha'[i : i+1), i \in [0 : |\alpha|) \}$ и присваиваем им метки L_0
- Шаг $n+1$: сортируем $\{ \alpha'[i : i+2^{n+1}), i \in [0 : |\alpha|) \}$, используя L_n
 - $\alpha'[i : i+2^{n+1}) < \alpha'[j : j+2^{n+1})$, если $(L_n[i], L_n[i+2^n]) < (L_n[j], L_n[j+2^n])$
 - Используем цифровую сортировку и присваиваем метки L_{n+1}
- $\log |\alpha|$ шагов
 - $O(|\alpha|)$ на каждый шаг, $O(|\Sigma|)$ на первый шаг
- Сложность $O(|\alpha| * \log |\alpha| + |\Sigma|)$

Суффиксный массив: построение – пример

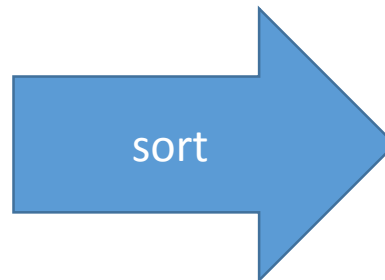
$\alpha = \text{abacaba}\$ \$ \$ \$ \$ \$ \$$

$L_0 = 121312100000000$

$L_1 = 243524100000000$

char	code	#
\$	0	7
a	1	4
b	2	2
c	3	1

i	$\alpha[i:2)$	$L(\alpha[i:i+1))$	$L(\alpha[i+1:i+2))$
0	ab	1	2
1	ba	2	1
2	ac	1	3
3	ca	3	1
4	ab	1	2
5	ba	2	1
6	a\$	1	0



i	$\alpha[i:2)$	$L(\alpha[i:i+1))$	$L(\alpha[i+1:i+2))$	$L(\alpha[i:i+2))$
6	a\$	1	0	1
0	ab	1	2	2
4	ab	1	2	2
2	ac	1	3	3
1	ba	2	1	4
5	ba	2	1	4
3	ca	3	1	5

Суффиксный массив: построение – пример

$\alpha = \text{abacaba}\$ \$ \$ \$ \$ \$ \$$

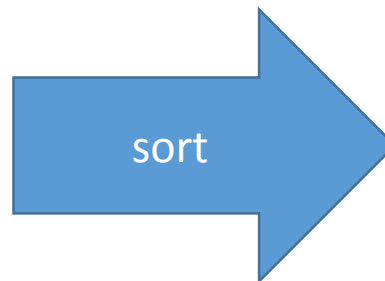
$L_0 = 121312100000000$

$L_1 = 243524100000000$

$L_2 = 364725100000000$

char	code	#
\$	0	7
a	1	4
b	2	2
c	3	1

i	$\alpha[i:4)$	$L(\alpha[i:i+2])$	$L(\alpha[i+2:i+4])$
0	abac	2	3
1	baca	4	5
2	acab	3	2
3	caba	5	4
4	aba\$	2	1
5	ba\$\$	4	0
6	a\$\$\$	1	0

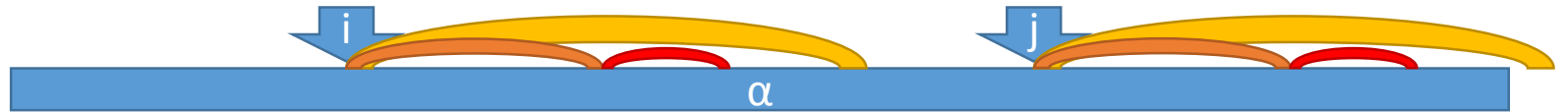


i	$\alpha[i:4)$	$L(\alpha[i:i+2])$	$L(\alpha[i+2:i+4])$	$L(\alpha[i:i+4])$
6	a\$\$\$	1	0	1
4	aba\$	2	1	2
0	abac	2	3	3
2	acab	3	2	4
5	ba\$\$	4	0	5
1	baca	4	5	6
3	caba	5	4	7

SA

Суффиксный массив для нахождения LCP

- Запрос $(i, j) \rightarrow |\text{LCP}(\alpha[i:], \alpha[j:])| = \max\{k \mid \alpha[i:i+k) = \alpha[j:j+k)\}$
- Идея: метки позволяют избежать сравнения подстрок
 - нужно сохранить метки для всех подстрок длины $2^n - O(|\alpha| \log |\alpha|)$ памяти



$n = 2^{\text{ceil}(\log(|\alpha|))}$

ans = 0

while n > 0

 if $L(\alpha[i+\text{ans} : i+\text{ans}+n]) == L(\alpha[j+\text{ans} : j+\text{ans}+n])$

 ans += n

 n /= 2

return ans

Сложность $O(\log |\alpha|)$