

Алгоритмы и структуры данных

Префикс-функция, поиск бора

CS Center, Новосибирск

Периодичность и borders

Пусть $\alpha = \beta \beta \dots \beta \beta'$ и $\beta' \text{ is_prefix } \beta$

- значит $|\beta|$ - период α

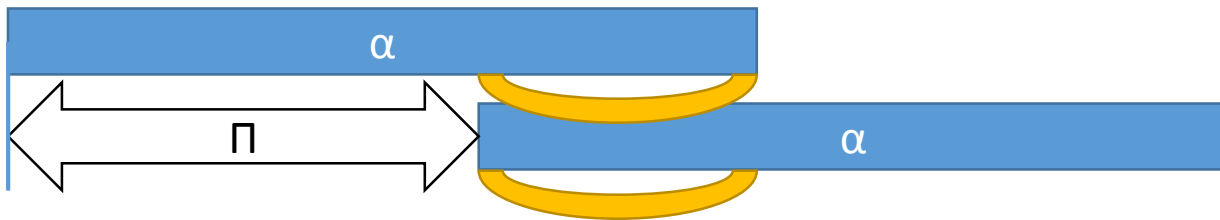
Π – период α , если:

- $\alpha[\Pi :) = \alpha[: |\alpha| - \Pi)$

$\gamma \text{ is_border } \alpha$, если $\gamma \text{ is_prefix } \alpha$ и $\gamma \text{ is_suffix } \alpha$

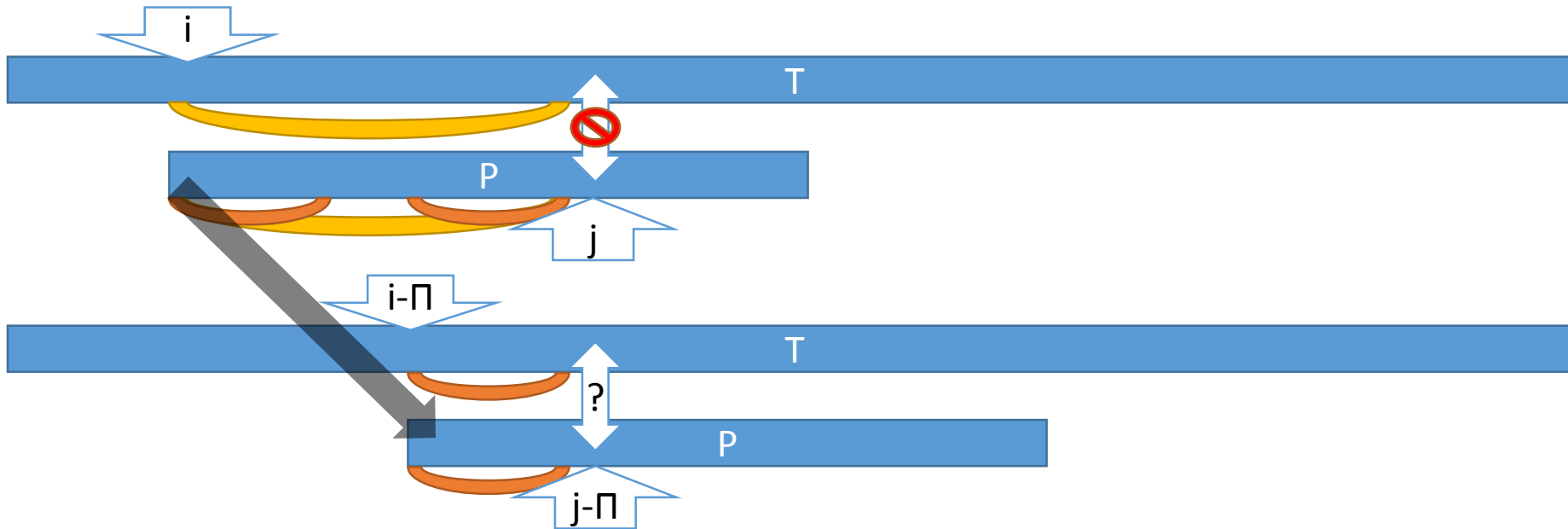
- $\gamma = \alpha[\Pi :) = \alpha[: |\alpha| - \Pi)$

$$|\gamma| = |\alpha| - \Pi$$



Снова поиск шаблона в тексте

$P, T \in \Sigma^*$, $P \text{ is_sub } T$?



Π – минимальный собственный период $P[: j)$

- Можно предподсчитать для все $\Pi[j]$
 - тогда поиск за $O(T)$

Префикс-функция строки

- $P_{\alpha}[j] = \text{длина макс. бордера } \alpha[:j], 1 < j \leq |\alpha|$

Пример:

$\alpha = \text{abacababaca}$

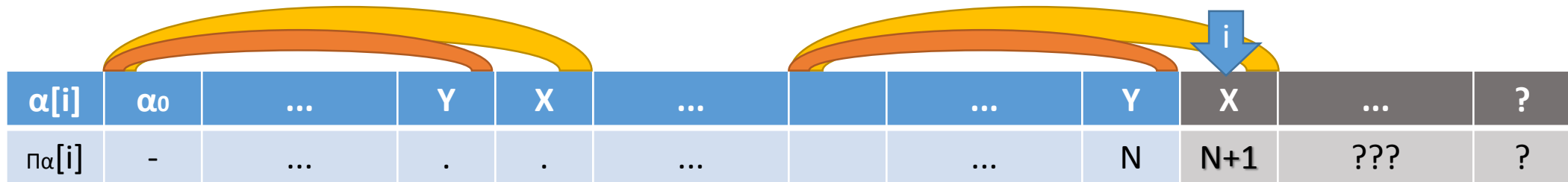
$\alpha[i]$	a	b	a	c	a	b	a	b	a	c	a
$P_{\alpha}[i]$	0	0	1	0	1	2	3	2	3	4	5

Diagram illustrating the prefix function calculation for the string $\alpha = \text{abacababaca}$. The table shows the string characters and their corresponding prefix function values. Colored arrows below the table indicate the matching prefixes:

- Orange arrows: $\alpha[1] = \text{a}$ matches $\alpha[3] = \text{a}$.
- Yellow arrows: $\alpha[1..2] = \text{ab}$ matches $\alpha[6..7] = \text{ab}$.
- Green arrows: $\alpha[1..4] = \text{abac}$ matches $\alpha[9..12] = \text{abaca}$ (Note: the diagram shows a match of length 4, which is the value in the table).

Префикс-функция за линейное время

- $P_{\alpha}[j] = \text{длина макс. бордера } \alpha[:j], 1 < j \leq |\alpha|$



- $\gamma \text{ is_border } \beta, \beta \text{ is_border } \alpha \Rightarrow \gamma \text{ is_border } \alpha$
- $\{ P[j-1], P[P[j-1]-1], P[P[P[j-1]-1]-1], \dots, 0 \}$ – длины всех бордеров $\alpha[:j]$
 - перебираем эти значения, пока не выполнится $\alpha[p] = \alpha[i]$

Префикс-функция за линейное время

$\Pi[0] = 0$

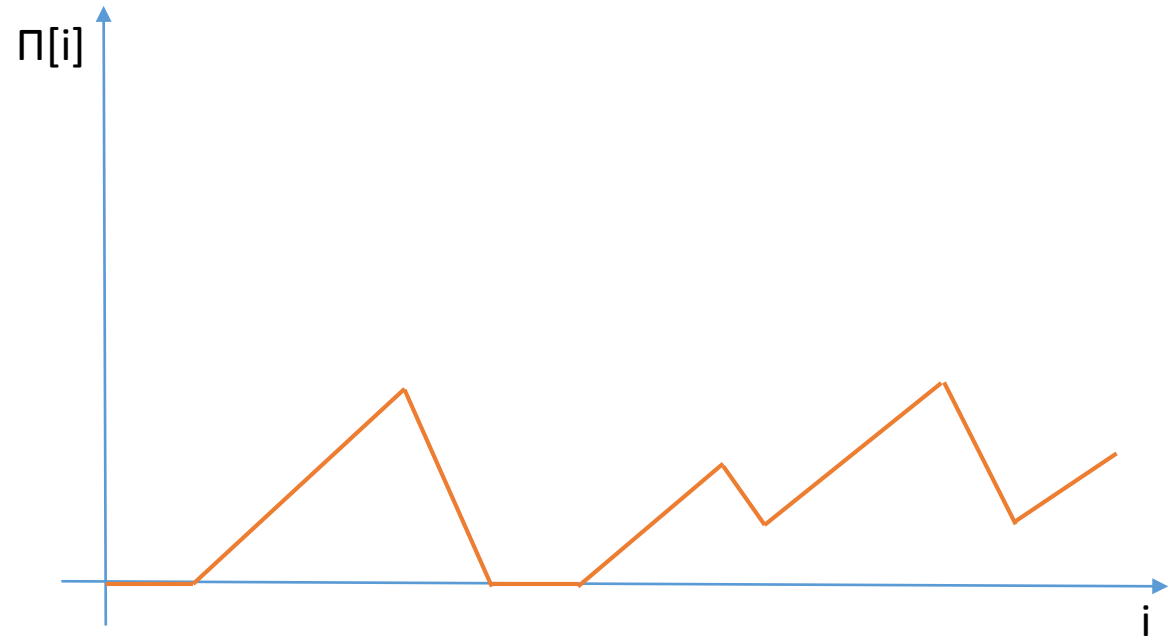
for $i \in [1 : |\alpha|)$

$p = \Pi[i-1]$

 while $p > 0 \ \&\& \ \alpha[p] \neq \alpha[i]$

$p = \Pi[p-1]$

$\Pi[i] = (\alpha[p] == \alpha[i]) ? p+1 : 0$



- $O(N)$ итераций внешнего цикла
- в сумме $O(N)$ итераций внутреннего цикла
 - Π растёт не более N раз и уменьшается не более N раз

Алгоритм Кнута-Морриса-Пратта

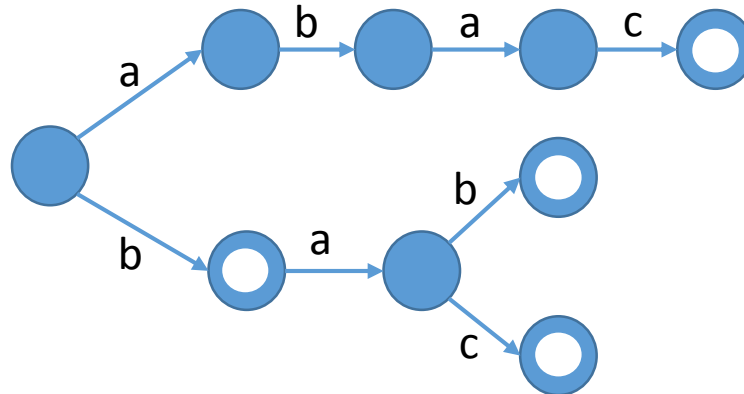
- $\alpha = P \$ T$, $\$$ - sentinel
 - $P_\alpha[j] = |P| \Rightarrow$ вхождение с позиции $j - 2 * |P|$
- Все значения P_α не превосходят $|P|$
 - Можно ограничиться $O(|P|)$ дополнительной памяти
 - Вычисление первых $|P|$ значений P_α – предобработка шаблона

Поиск множества образцов

- Желаемая сложность: $\Sigma |P_i| + |T| + A$ (A – число ответов)
- Бор (trie) – корневое дерево для хранения набора строк
 - рёбра помечены символами
 - помечены конечные вершины

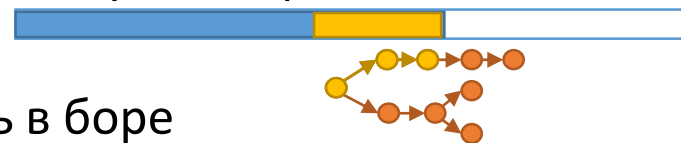
- Пример:

- { abac, b, bab, bac } ->

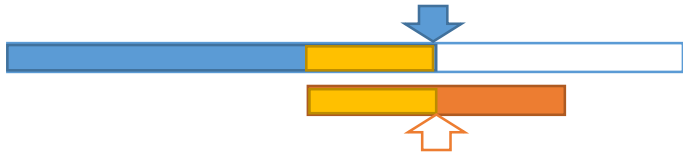


Алгоритм: идея

- Построим бор для заданных шаблонов
- Проходим по тексту и перемещаемся по бору
 - Вспомним: КМП вычисляет длиннейший префикс шаблона равный суффиксу прочитанной части текста
 - Аналогия для множества шаблонов:
 - длиннейший префикс шаблона -> длиннейший путь в боре
 - нужно на каждом шаге вычислять такой путь в боре



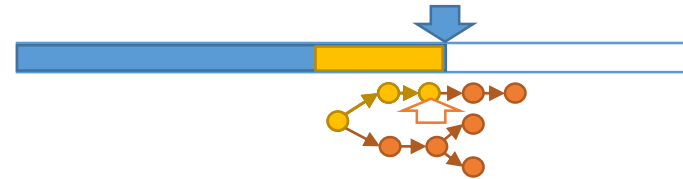
Алгоритм: идея



Один шаблон – КМП

Состояние – позиция в шаблоне

- буква совпадает => продвигаемся по шаблону
- иначе => откатываемся по шаблону, используя префикс-функцию



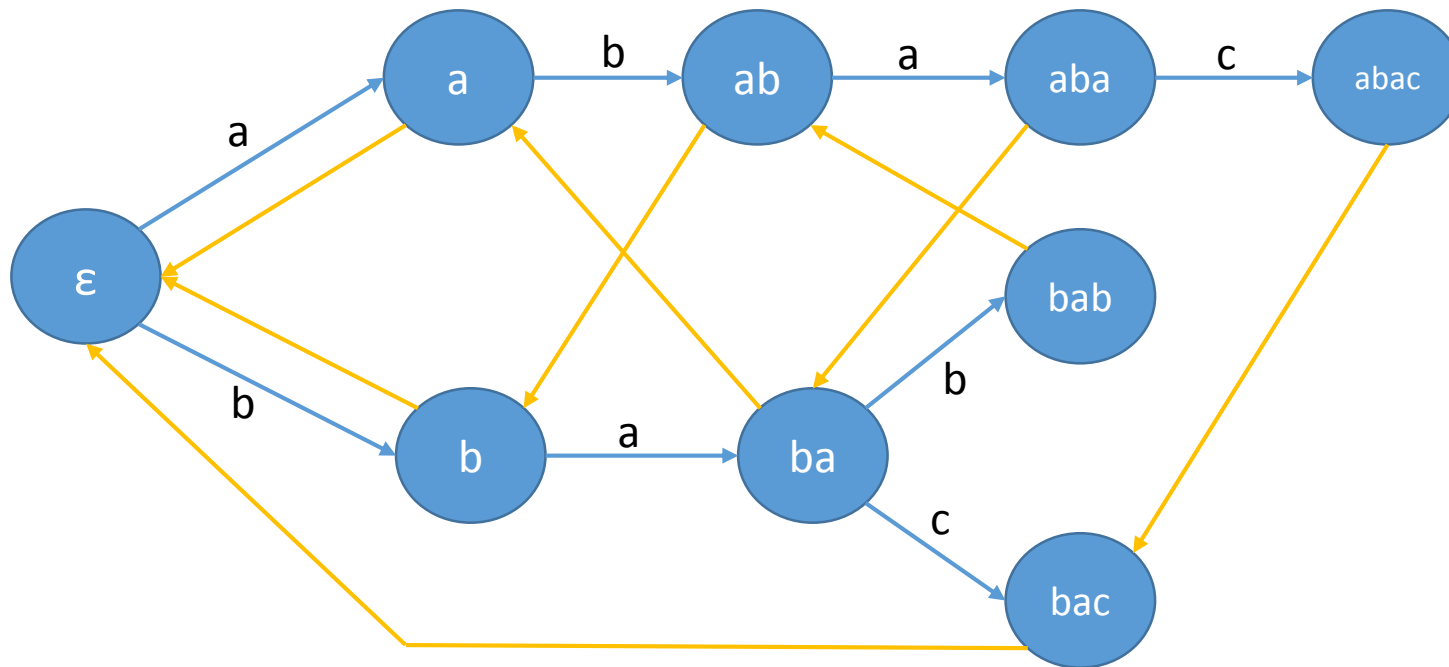
Множество шаблонов

Состояние – вершина бора

- есть ребро, помеченное буквой => продвигаемся по бору
- иначе => откатываемся по бору, используя префикс-функцию на боре

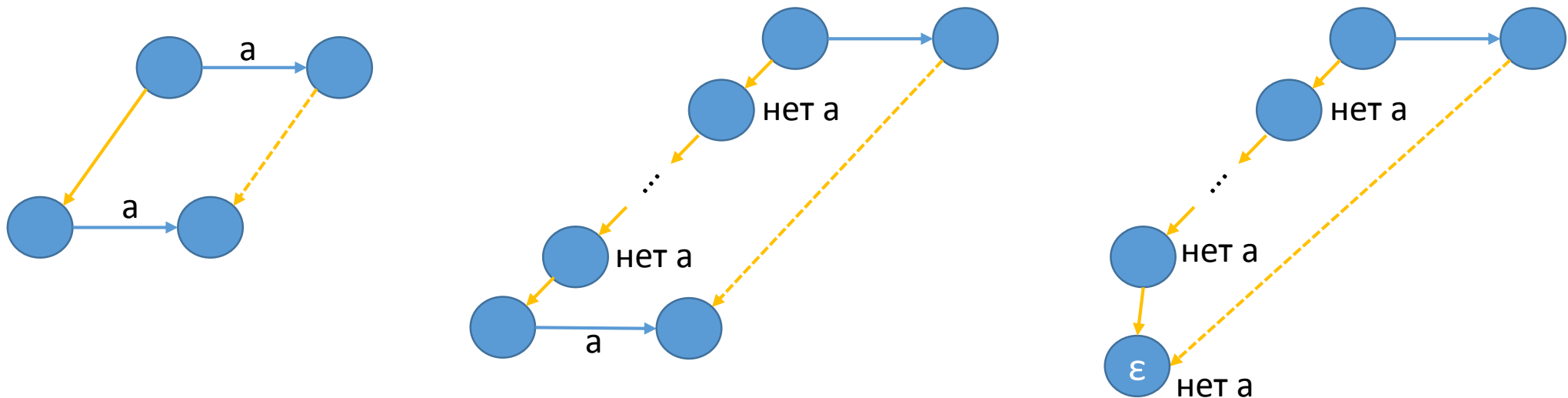
Префикс-функция на боре

- Соответствие: вершины бора \leftrightarrow строки
- $P(\alpha) = \beta$: β есть в боре, β is_suff α , $|\beta| < |\alpha|$, $|\beta|$ максимальна



Префикс-функция на боре: вычисление

- $\{ \alpha, \Pi(\alpha), \Pi(\Pi(\alpha)), \dots, \varepsilon \} = \{ \beta \text{ есть в боре, } \beta \text{ is_suff } \alpha \}$
- Порядок вычисления – по увеличению длины строки
 - порядок обхода дерева в ширину



Префикс-функция на боре: СЛОЖНОСТЬ ВЫЧИСЛЕНИЯ

- $e = (u, v)$
 - $\text{cost}(e)$ – сколько времени вычисляли $P(v)$
- $\sum_e \text{cost}(e) = T$ – сложность вычисления
- $\sum_{v\text{-leaf}} \text{cost}(\text{path}(e, v)) = S \geq T$
- Докажем, что $S \leq \sum |P_i|$:
 - $\text{cost}(\text{path}(e, v)) \leq |v|$ - аналогично префикс-функции для строки
- Значит $T \leq \sum |P_i|$ – линейная сложность вычисления

Алгоритм Ахо-Корасик

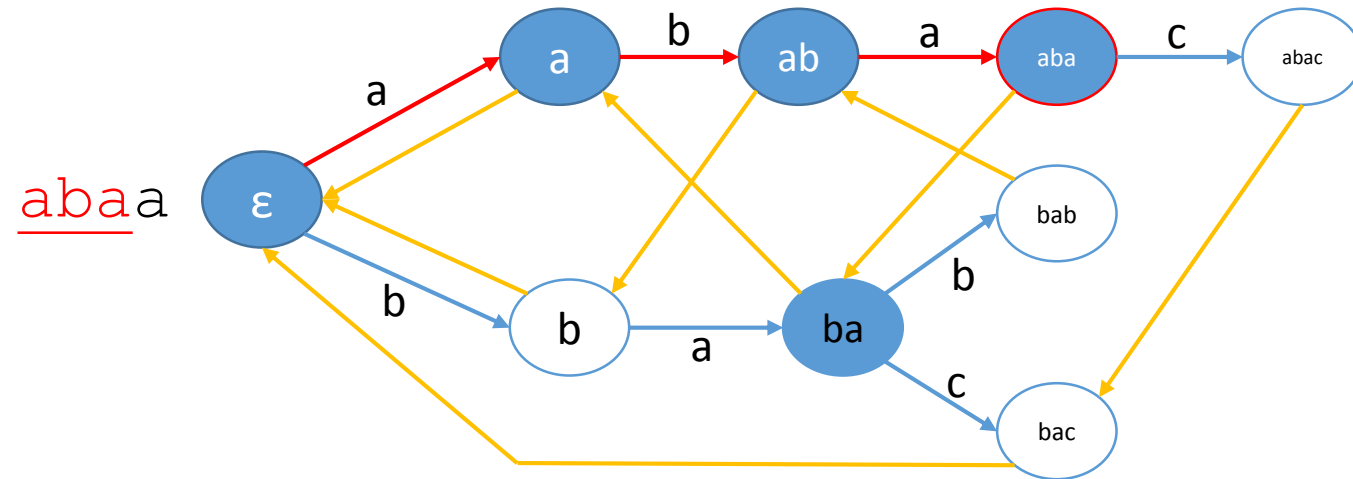
• P:

abac

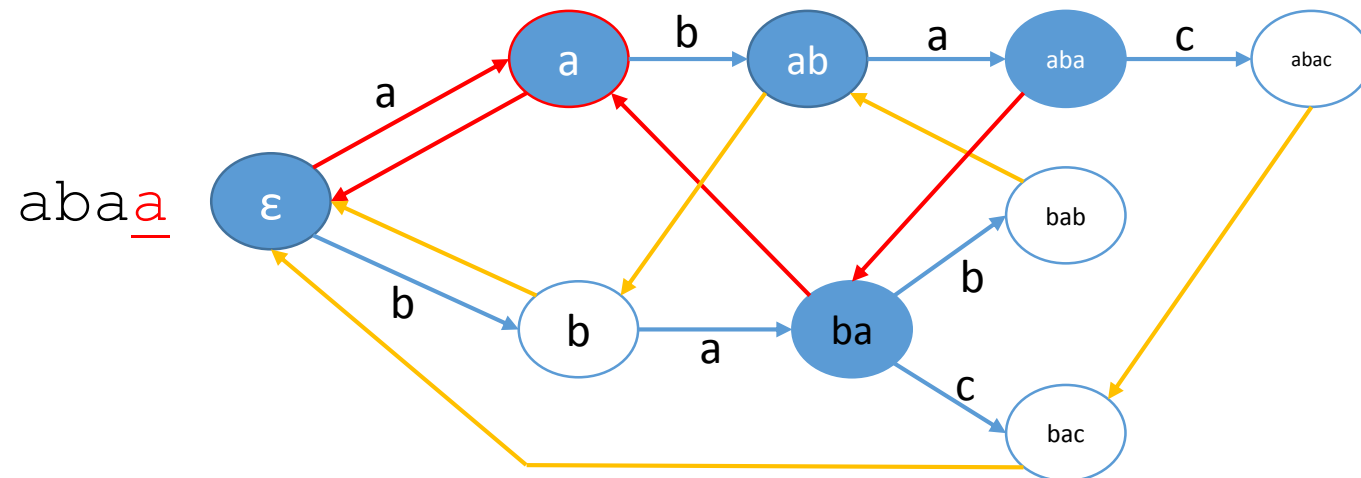
b

bab

bac



• T = abaa



Алгоритм Ахо-Корасик

- Подсчитать общее количество вхождений
 - В каждой вершине бора предподсчитаем число шаблонов, являющихся суффиксами строки
- Найти вхождения
 - В каждой вершине сделаем ссылку на длиннейший шаблон, являющийся собственным суффиксом строки

Бор: накладные расходы

M – сумма длин шаблонов

- Массив ссылок в вершине:
 - Память $O(|\Sigma| * M)$, время перехода $O(1)$
- Дерево поиска:
 - Память $O(M)$, время перехода $O(\log|\Sigma|)$