

Алгоритмы и структуры данных

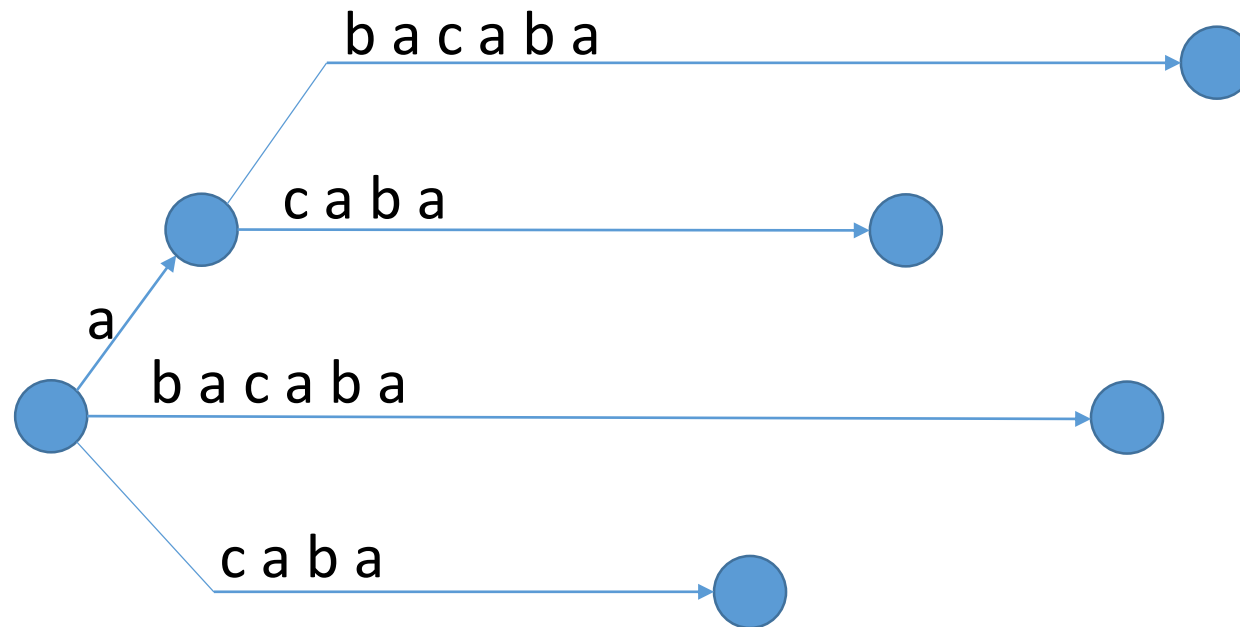
Суффиксное дерево

CS Center, Новосибирск

Online поиск шаблонов в тексте

- Дан текст T
 - Можно его предобработать, $O(|T|)$ на предобработку
- Шаблоны P_i приходят один за другим
 - Хочется ответить на запрос за $O(|P_i|)$
- Как предобработать текст?
 - Подстрока – это префикс суффикса
 - Нужно выяснить: P – префикс какого-либо суффикса T ?
 - Как эффективно закодировать все суффиксы T ?

Идея: сжатый суффиксный бор



- Суффиксное дерево
 - Не более N листьев $\Rightarrow O(N)$ вершин
- Каждое ребро помечено подстрокой $[i : j)$

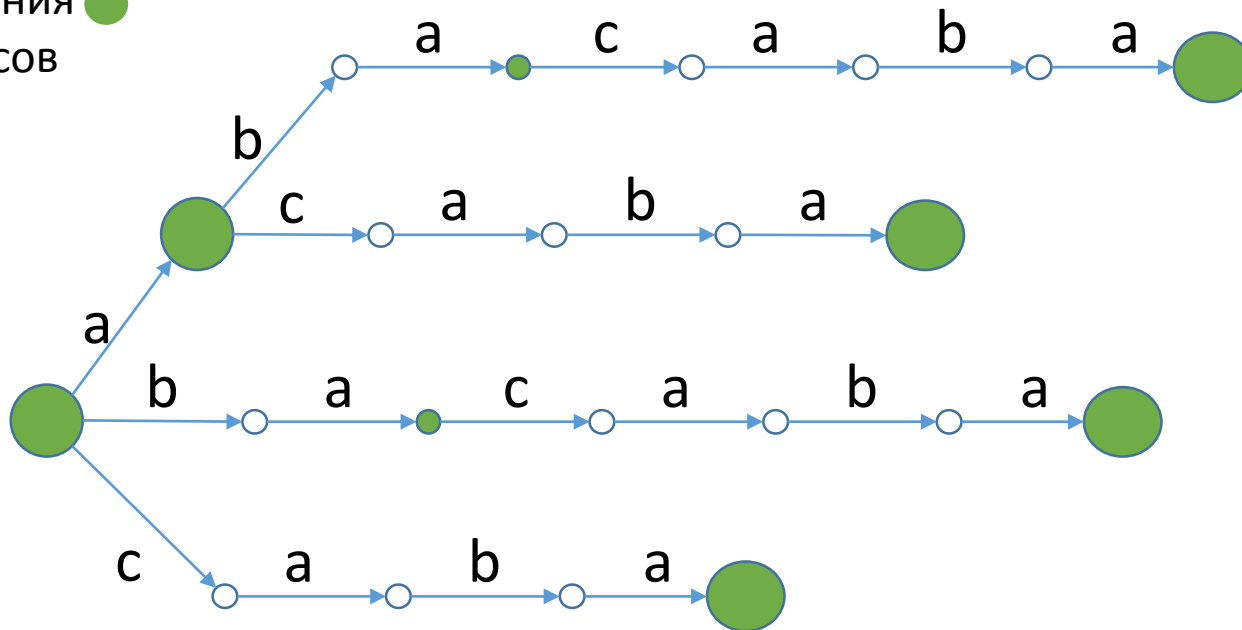
Суффиксное дерево

- Положение в суффиксном дереве:

● Явное – вершина дерева

○ Неявное – позиция внутри строки на ребре

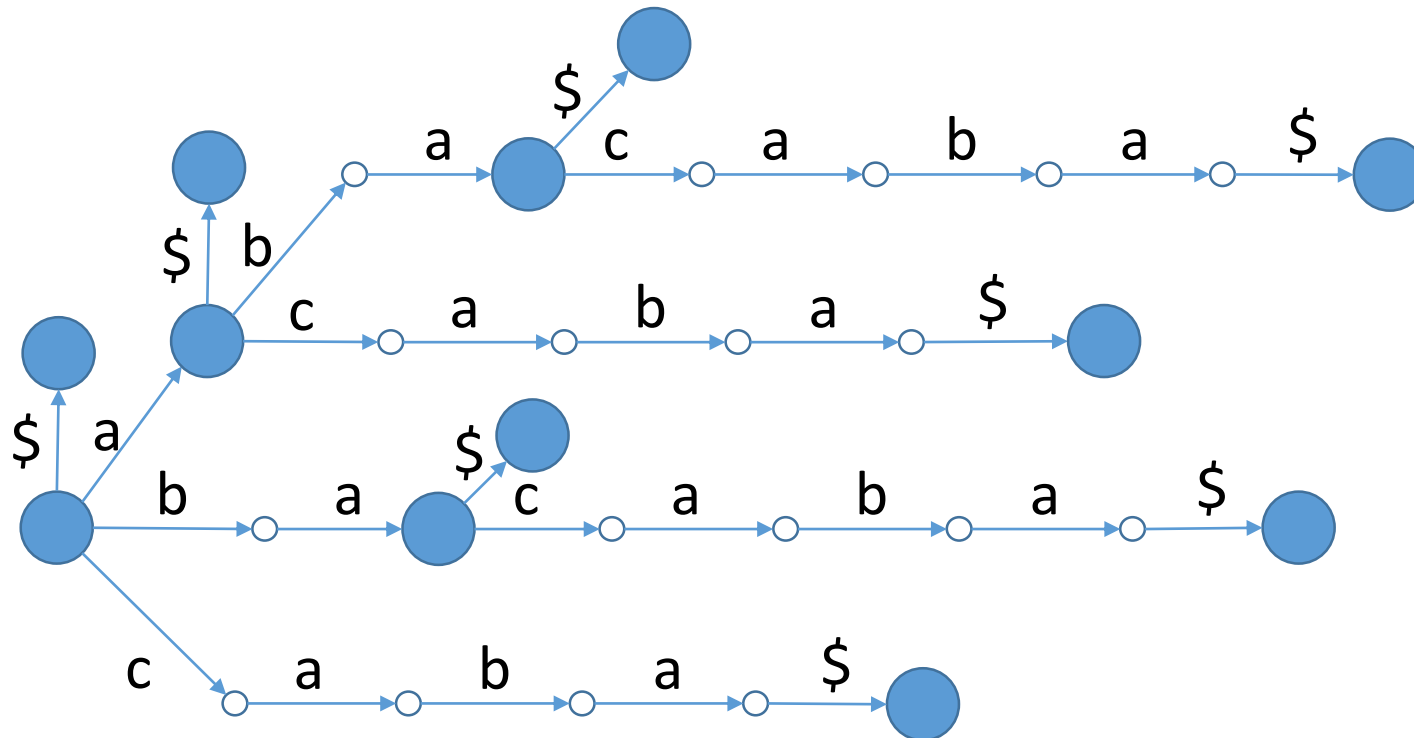
Положения
суффиксов ●



Суффиксное дерево

Сделаем положения суффиксов явными:

- $T' = T + \$$



Алгоритм Укконена

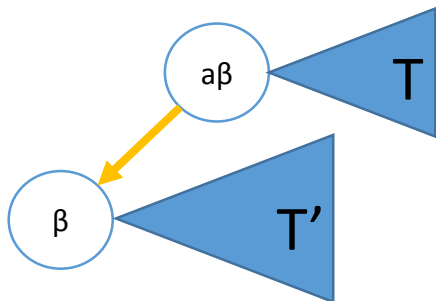
- Начинаем с пустого дерева (для пустой строки)
- На каждом шаге добавляем очередной символ строки
 - На i -ом шаге удлиняются i суффиксов

abac	<u>a</u>bac	<u>ab</u>ac	<u>aba</u>c	<u>abac</u>
ϵ	$\epsilon+a$	$a+b$	$ab+a$	$aba+c$
	ϵ	$\epsilon+b$	$b+a$	$ba+c$
		ϵ	$\epsilon+a$	$a+c$
			ϵ	$\epsilon+c$
				ϵ

- N таких шагов

Алгоритм Укконена: типы положений

- **Типы:**
 - 1 $\xrightarrow{+a}$ Лист
 - Пометка ребра изменяется с $[x : i)$ на $[x : i+1)$
 - Можно удлинить строку на ребре неявно
 - Не требуется изменять дерево
 - 2 $\xrightarrow{\text{нет } a}$ Положение без символа a
 - Добавляет новую вершину
 - $O(N)$ суммарно по всем итерациям
 - 3 \xrightarrow{a} Положение с символом a
 - Не требуется изменять дерево
- Перебираем положения в порядке уменьшения глубины
 - Номера типов положений не убывают: «1, 1, ... 1, 2, 2, ... 2, 3, 3, ... 3»

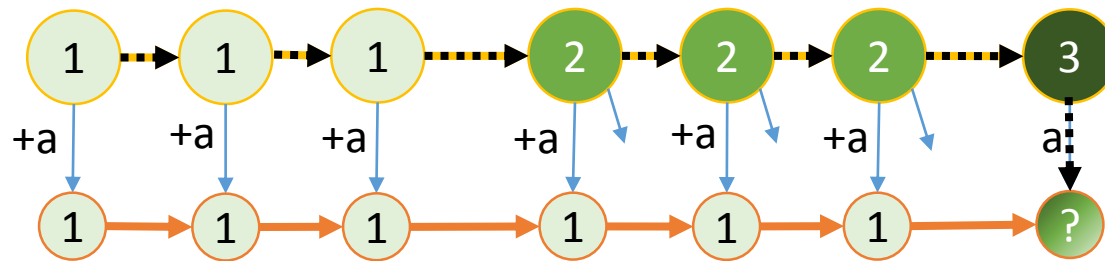


$$\gamma \in T \Rightarrow \gamma \in T'$$

- Пришли в положение типа 3 =>
все следующие положения не требуют изменений
- можно заканчивать итерацию

Алгоритм Укконена: изменение типов положений

Положения на итерации i



Положения на итерации $i+1$

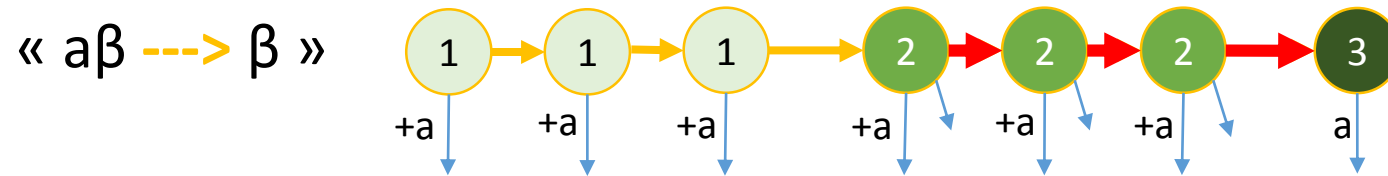
Потенциал := символная глубина

Здесь можно
начать
итерацию $i+1$

- уменьшается на 1 при переходе к следующему положению
- растёт на 1 при переходе на следующую итерацию
 - Растёт не более N раз \Rightarrow уменьшается не более N раз
 - $O(N)$ положений посетим в сумме по всем итерациям

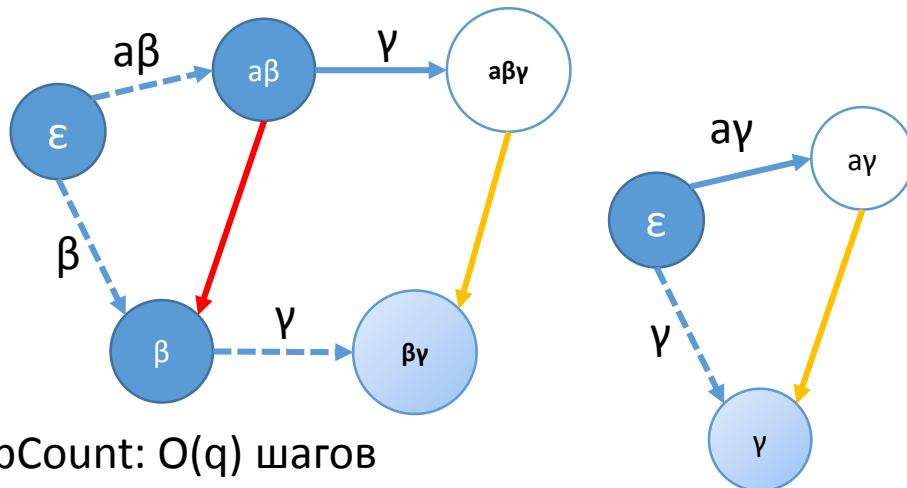
Алгоритм Укконена: переход от положения к положению

- Храним суффиксные ссылки для внутренних вершин дерева



- При создании вершины запомним её суффиксную ссылку

- Как вычислять суффиксные ссылки для произвольных положений?



SkipCount: $O(q)$ шагов

q – число вершин на пути от β до $\beta\gamma$

Потенциал := вершинная глубина

- при SkipCount потенциал растёт не менее, чем на $q-1$
- при переходе на следующую итерацию потенциал не уменьшается
- максимум потенциала = N
- N шагов, на каждом потенциал уменьшается не более, чем на 1
- \Rightarrow сумма q по всем шагам = $O(N)$

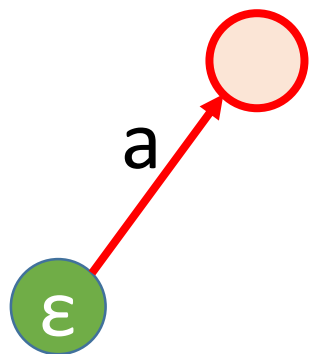
Алгоритм Укконена: пример

$T = \text{abacaba}\$$



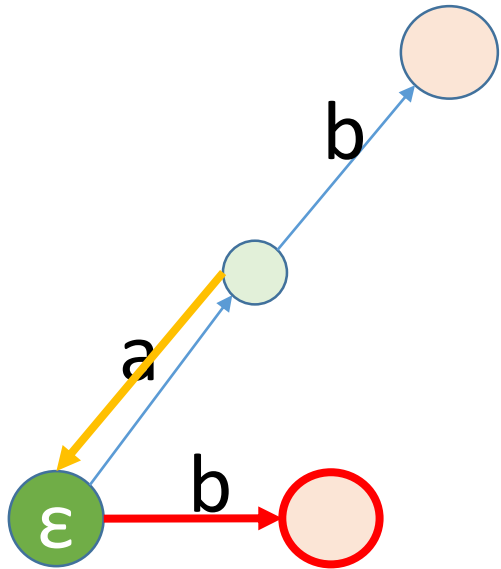
Алгоритм Укконена: пример

T = abacaba\$



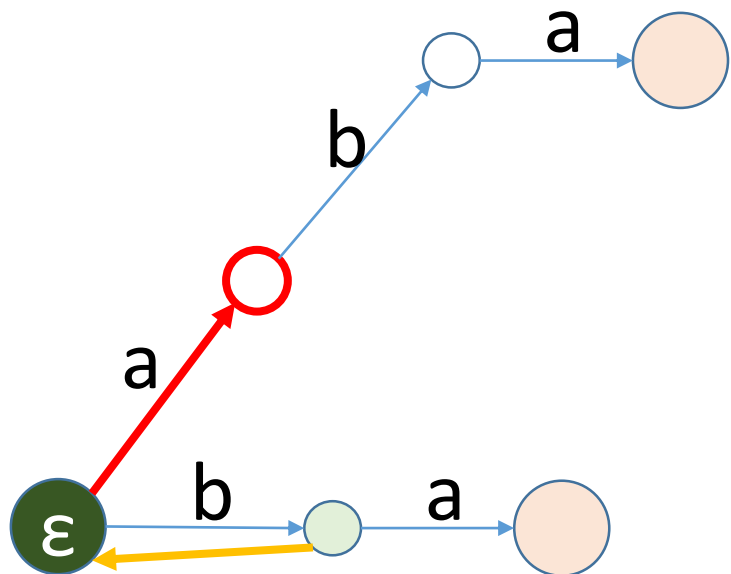
Алгоритм Укконена: пример

$T = \underline{a}bacaba\$$



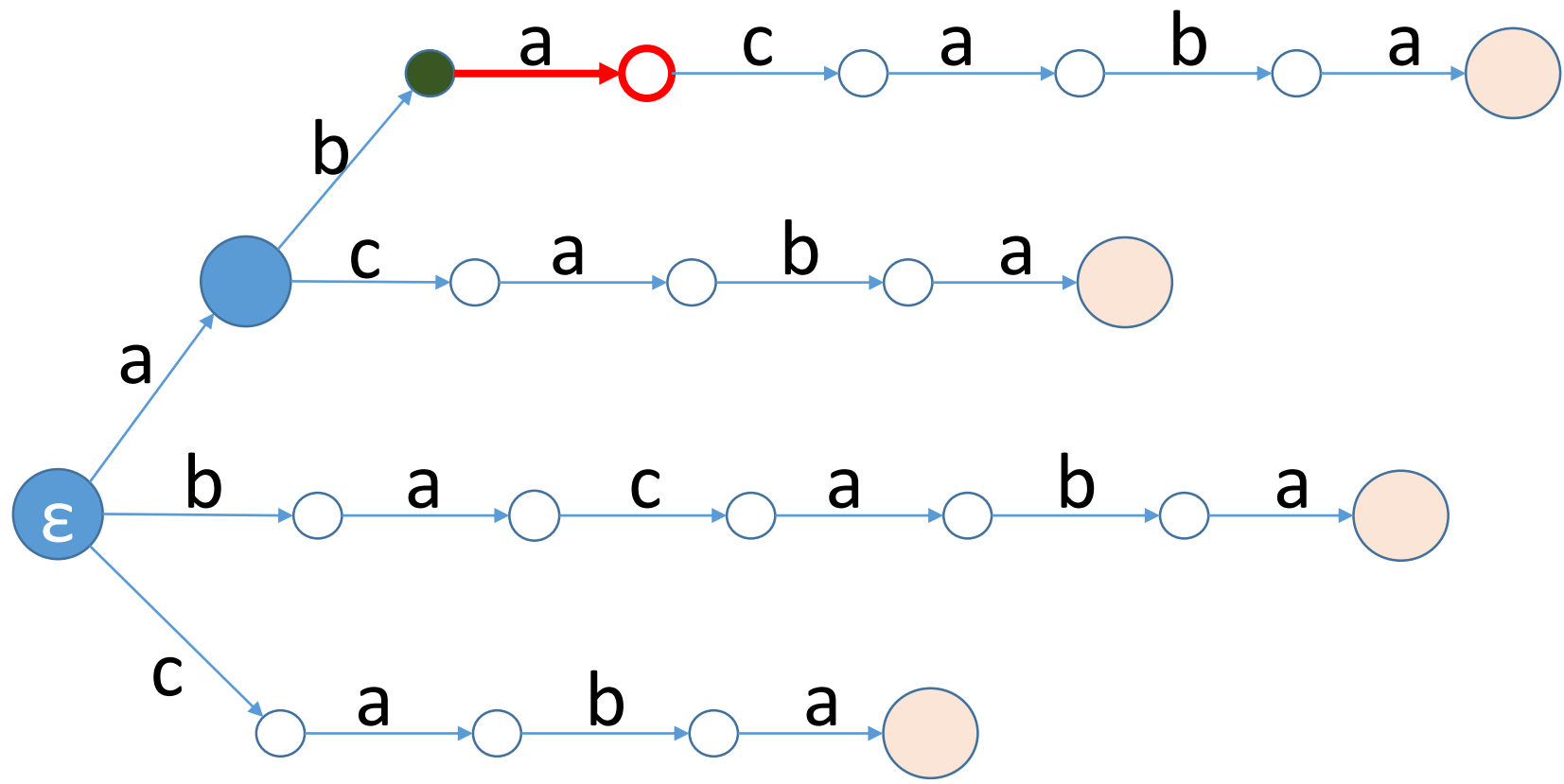
Алгоритм Укконена: пример

$T = \underline{ab}acaba\$$



Алгоритм Укконена: пример

$T = \underline{abacaba} \$$



Алгоритм Укконена: пример

$T = \underline{abacaba} \$$

