

(продолжение)

Строковые алгоритмы и структуры данных



БЕЛОРУССКИЙ
ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ

Бор

Суффиксный бор

Суффиксный массив

Задача поиска заданной строки в множестве строке

Предположим, что есть множество строк

$$S = \{S_0, S_1, \dots, S_{n-1}\}$$

из букв латинского алфавита Σ и необходимо быстро проверить, есть ли среди них некоторая заданная строка

$$A = (a_0, a_2, \dots, a_{l-1}).$$

Если поместить все строки из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$ в **массив** и осуществлять поиск строки $A = (a_0, a_2, \dots, a_{l-1})$ последовательно просматривая элементы массива, то **время поиска строки A** —

$$O(|A| \cdot n),$$

где n — количество строк в S .

При этом требуемая **память** — $O(|S_0| + |S_1| + \dots + |S_{n-1}|)$.

Если для хранения строк из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$ использовать сбалансированное бинарное поисковое дерево, то время **поиска строки A** —

$$O(|A| \cdot \log n),$$

где n — количество строк в S .

Существуют структуры данных, которые позволяют выполнять поиск строки более эффективно, например, за время

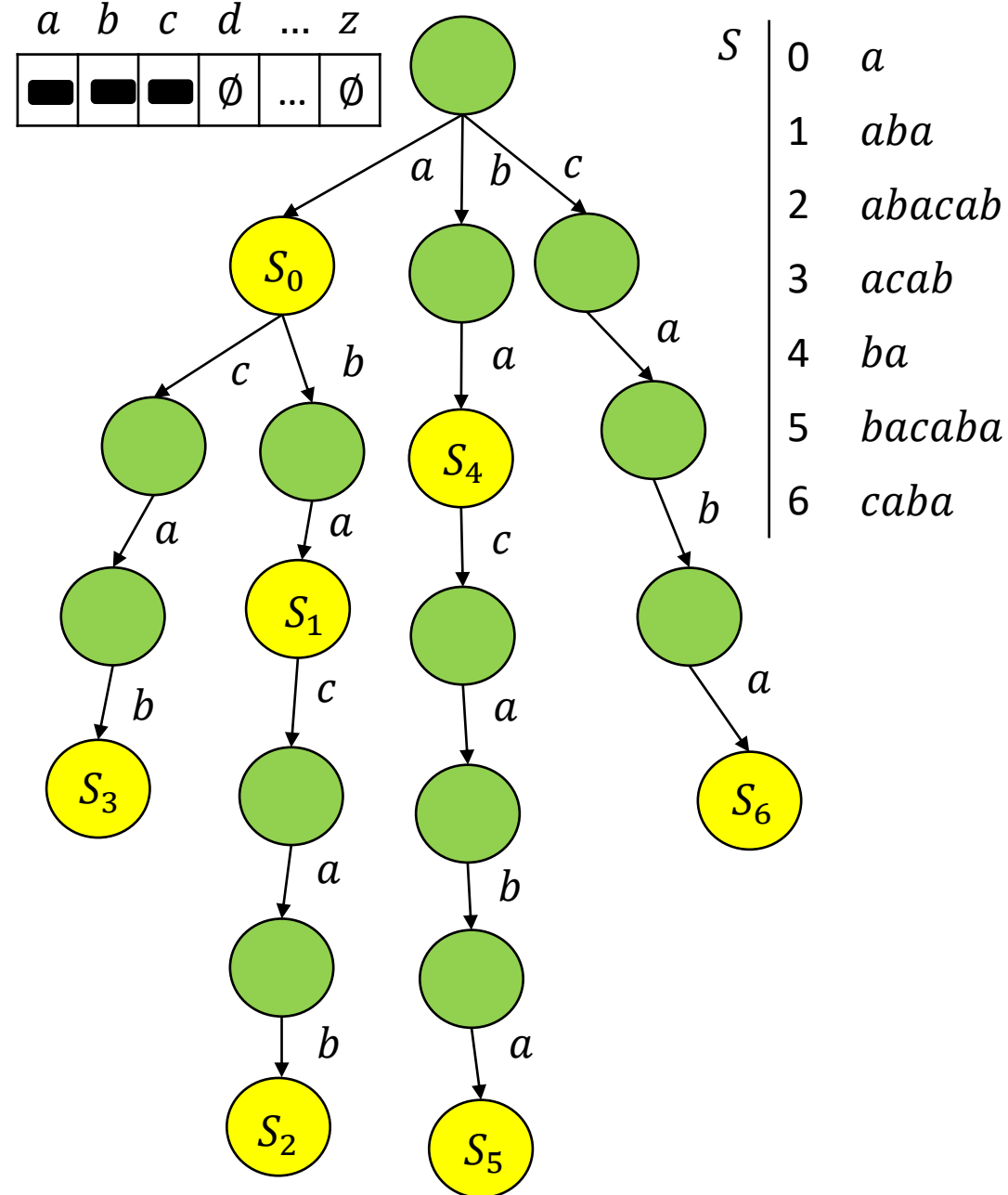
$$O(|A|)$$

(время поиска не зависит от количества строк в S).

Одной из таких структур данных является **бор**.

Бор

(англ. trie, луч, нагруженное дерево)



$A = aba$

Разобьем слова из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$ на группы в соответствии с первыми символами слов.

Пусть S' - группа, у слов которой первый символ совпадает с первым символом строки A . Сокращаем область поиска до множества S' , исключив из рассмотрения те строки, которые туда не попадают.

Заведем массив, в качестве индексов которого будут выступать символы алфавита, а в самом массиве будут храниться сами строки.

$$S \begin{array}{l|l} 0 & a \\ 1 & aba \\ 2 & abacab \\ 3 & acab \\ 4 & ba \\ 5 & bacaba \\ 6 & caba \end{array}$$

$A = (\underline{a}, b, a)$

S' →

	a	b	c	d
	a aba $abacaba$ $acab$	ba $bacaba$	$caba$	

Удаляем из всех слов множества S' и слова A первый символ и повторяем процедуру.

$$S' \begin{array}{l|l} 1 & ba \\ 2 & bacaba \\ 3 & cab \end{array}$$

$A' = (\cancel{a}, \underline{b}, a)$

	a	b	c	d
		ba $bacaba$	cab	

$$S' \begin{array}{l|l} 1 & a \\ 2 & acaba \end{array}$$

$A' = (\cancel{a}, \cancel{b}, \underline{a})$

	a	b	c	d
	a $acaba$			

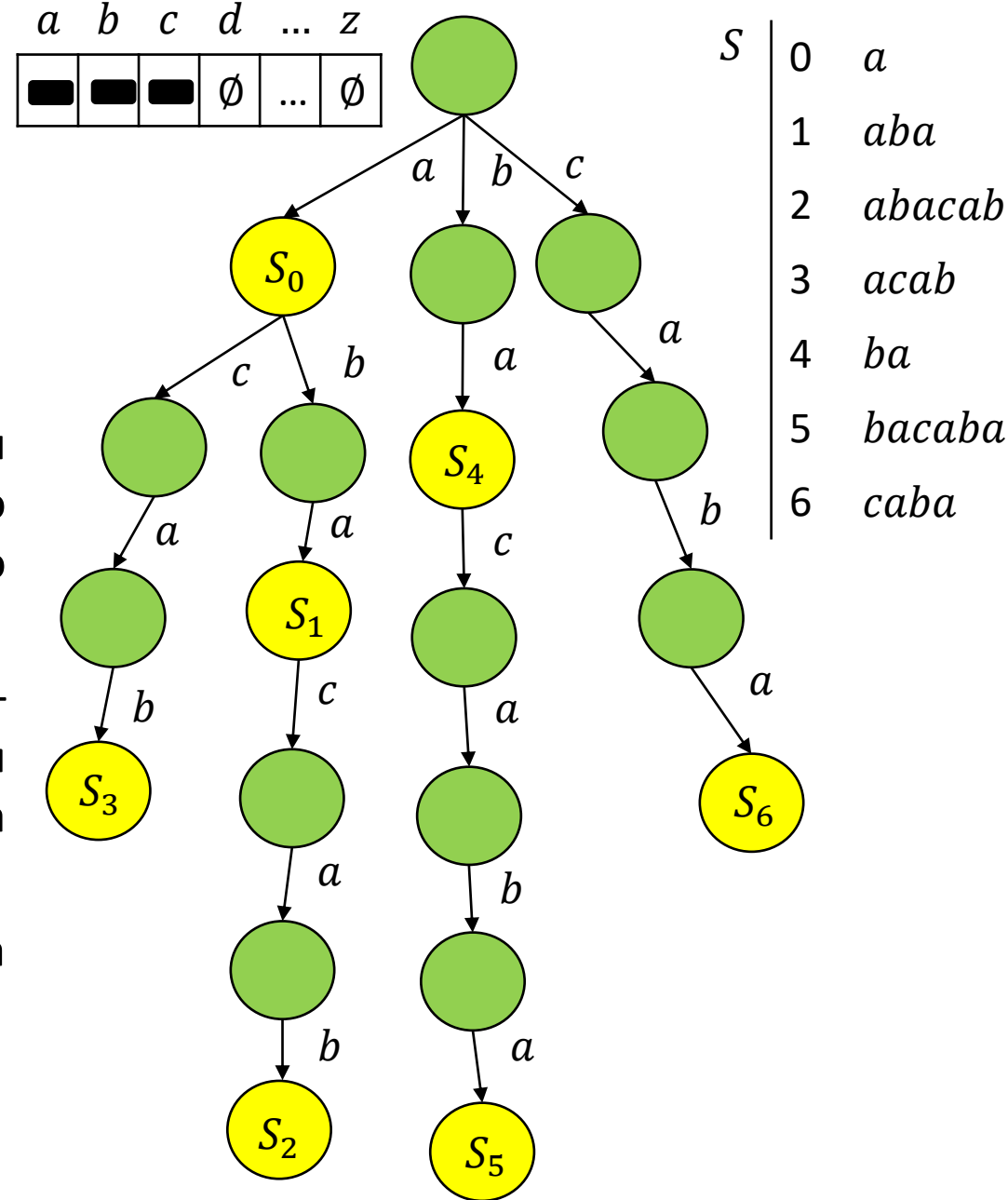
Если искомая строка A' состоит из одного символа, то достаточно проверить наличие в множестве S' строки длины 1.

Приведенное иерархическое разбиение строк на множества можно изобразить в виде древовидной структуры.

Бор

специализированная древовидная структура данных, предназначенная для хранения слов некоторого алфавита Σ .

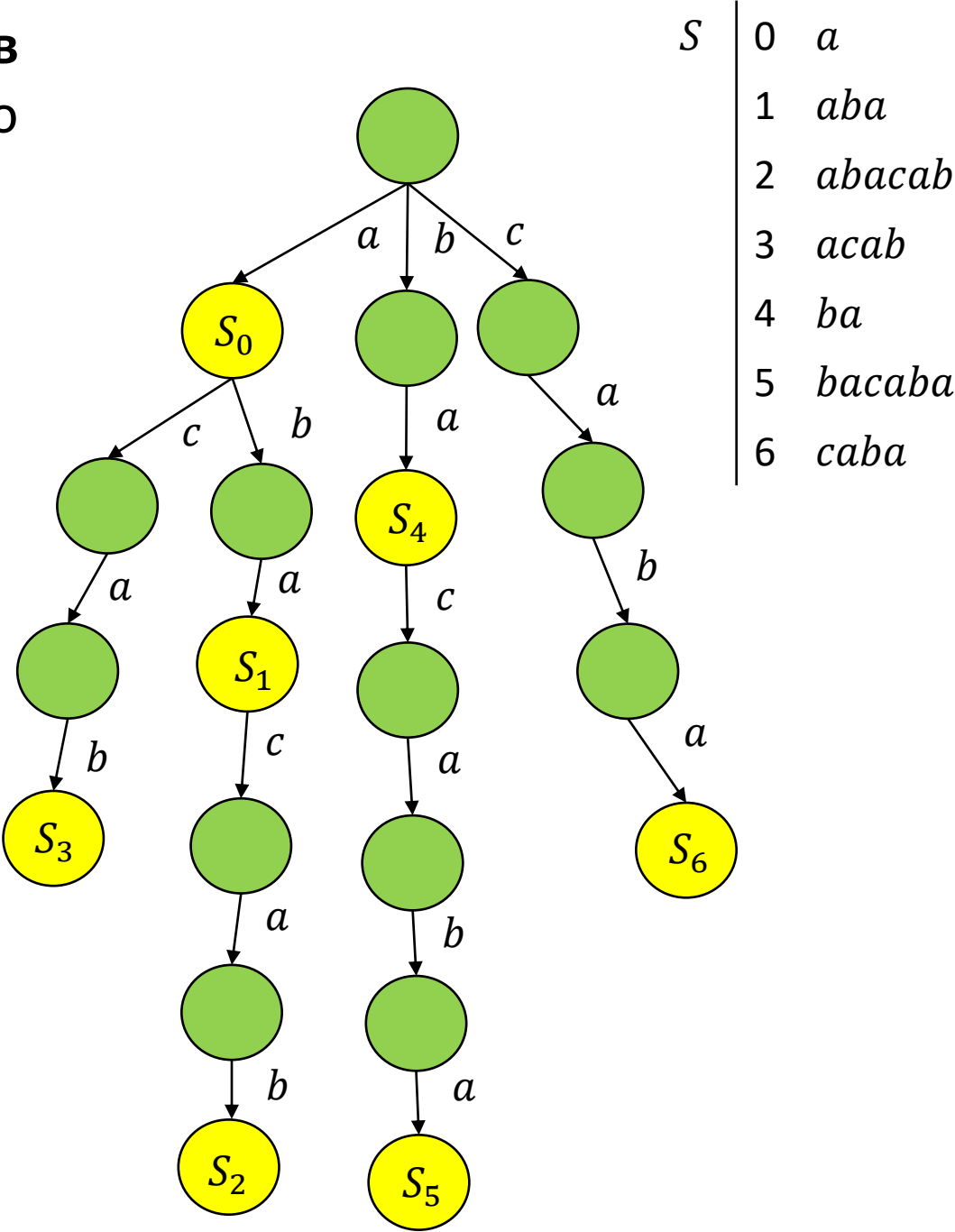
- **вершина дерева** содержит
 - информацию о вершинах, в которые можно перейти по каждому допустимому символу (если переход по дуге с некоторым символом не возможен, то результат перехода будем обозначать \emptyset);
 - пометку: терминальная или нет (вершина v – терминальная, если строка, которая определяется последовательностью символов, встречающихся на дугах в пути от корня к v , есть в множестве S);
 - терминальная вершина может хранить индекс слова в S ;
- **дуге** ставится в соответствие символ строки;



Если **выписать все символы на пути из корня в терминальную вершину**, то получим некоторую строку (слово) из множества S .

Листья дерева соответствуют строкам из S .

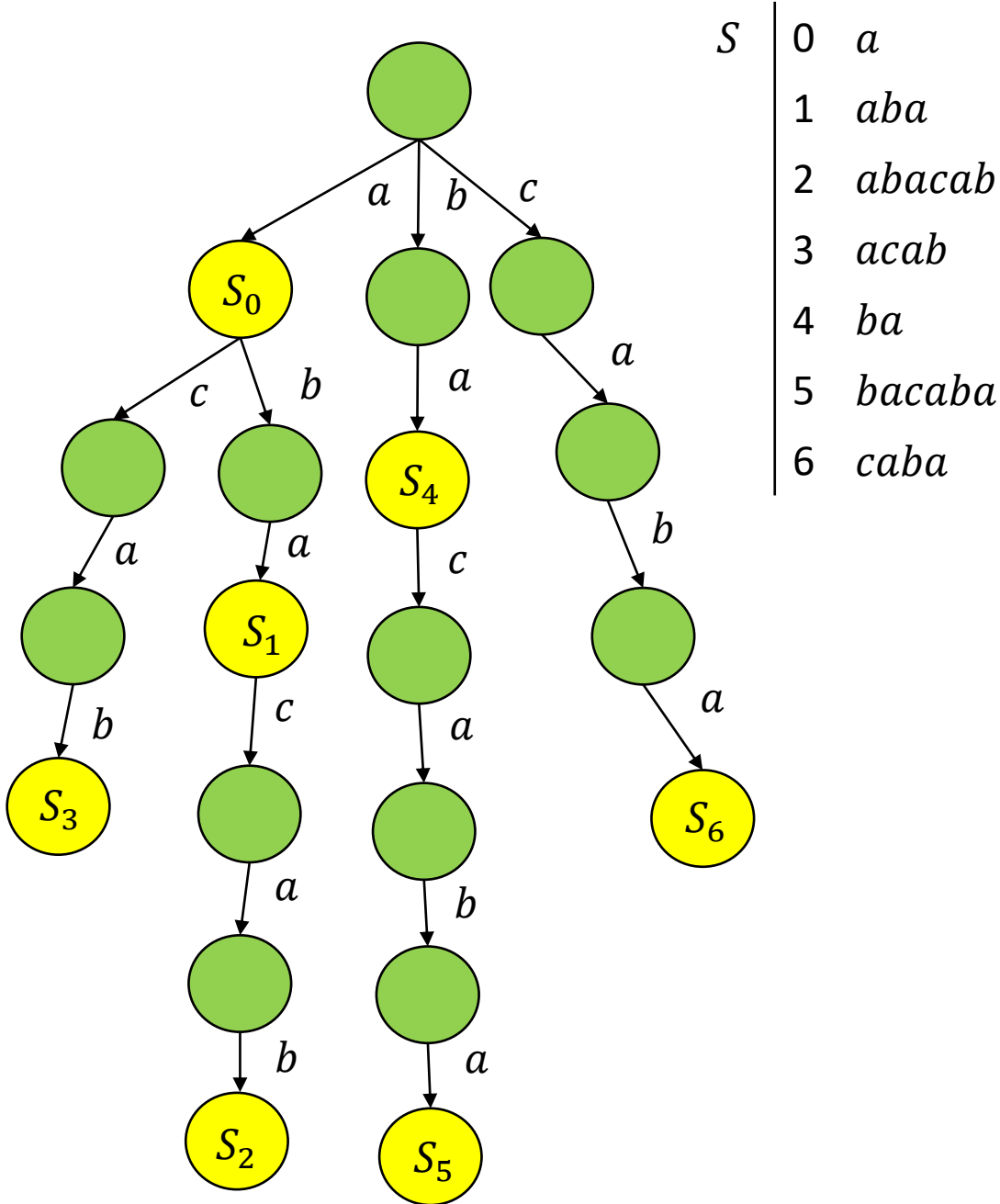
Если некоторая строка из S является префиксом другой строки этого множества, то ей будет соответствовать внутренняя вершина дерева ($S_1 = \mathbf{aba}$, $S_2 = \mathbf{abacab}$).



Для поиска строки A в боре будем последовательно спускаться из корня дерева по дугам, соответствующим символам строки A , пока строка A не закончится и в боре существует дуга, соответствующая текущему состоянию строки A .

Если в результате спуска:

- (1) попадём в терминальную вершину и дойдём до конца строки A , то **искомое слово найдено**;
- (2) если остановимся в нетерминальной вершине либо во время спуска не найдём дуги в дереве, соответствующей текущему символу строки A , то делаем вывод о том, что **строка A в множестве S отсутствует**.

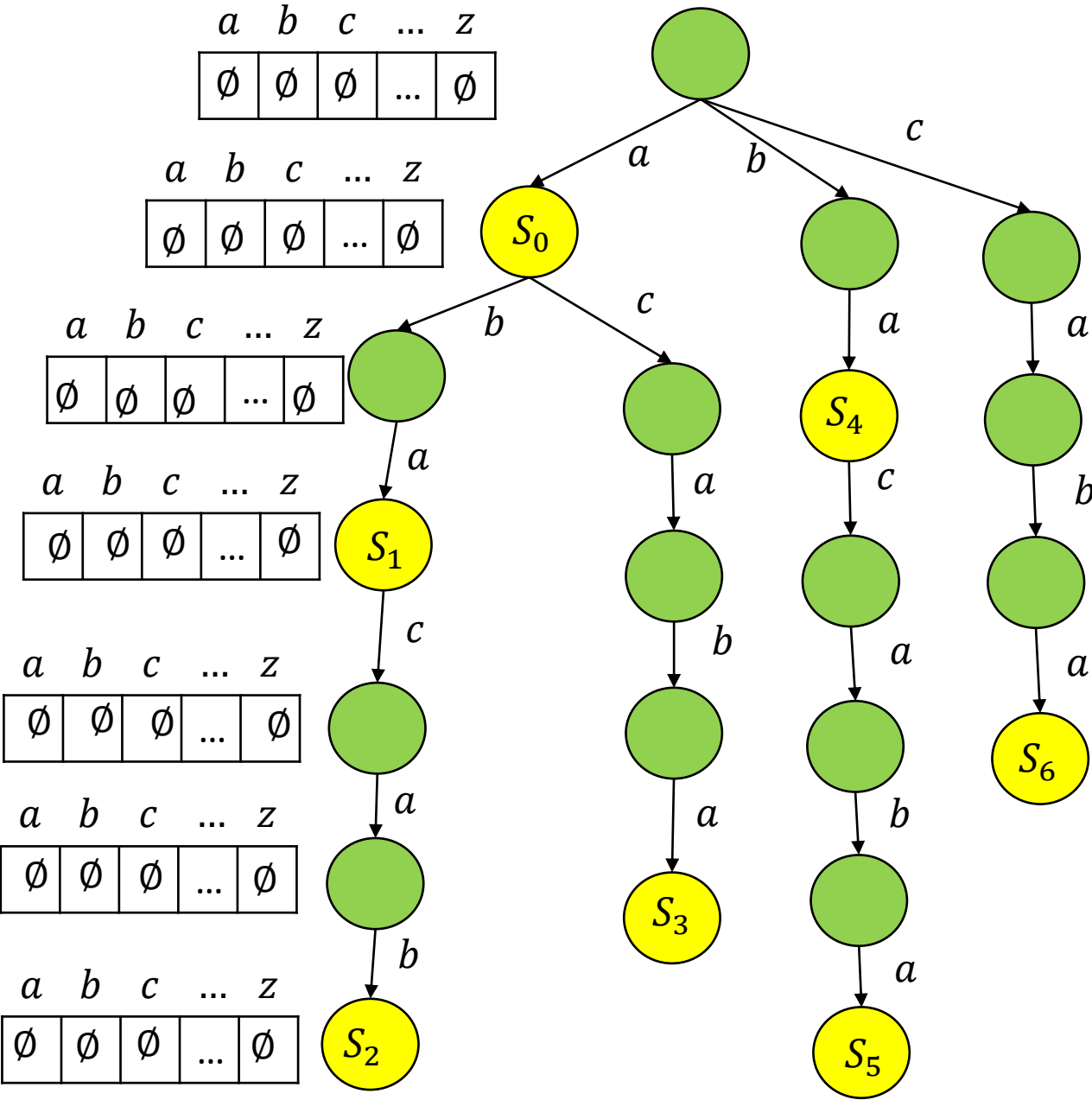


Добавление строки

Как и во время поиска, спускаемся по дереву, начиная с корня. Если во время спуска мы попадаем в ситуацию, что надо перейти по несуществующей дуге, то добавляем к дереву новую вершину и ведущую к ней дугу с соответствующей пометкой.

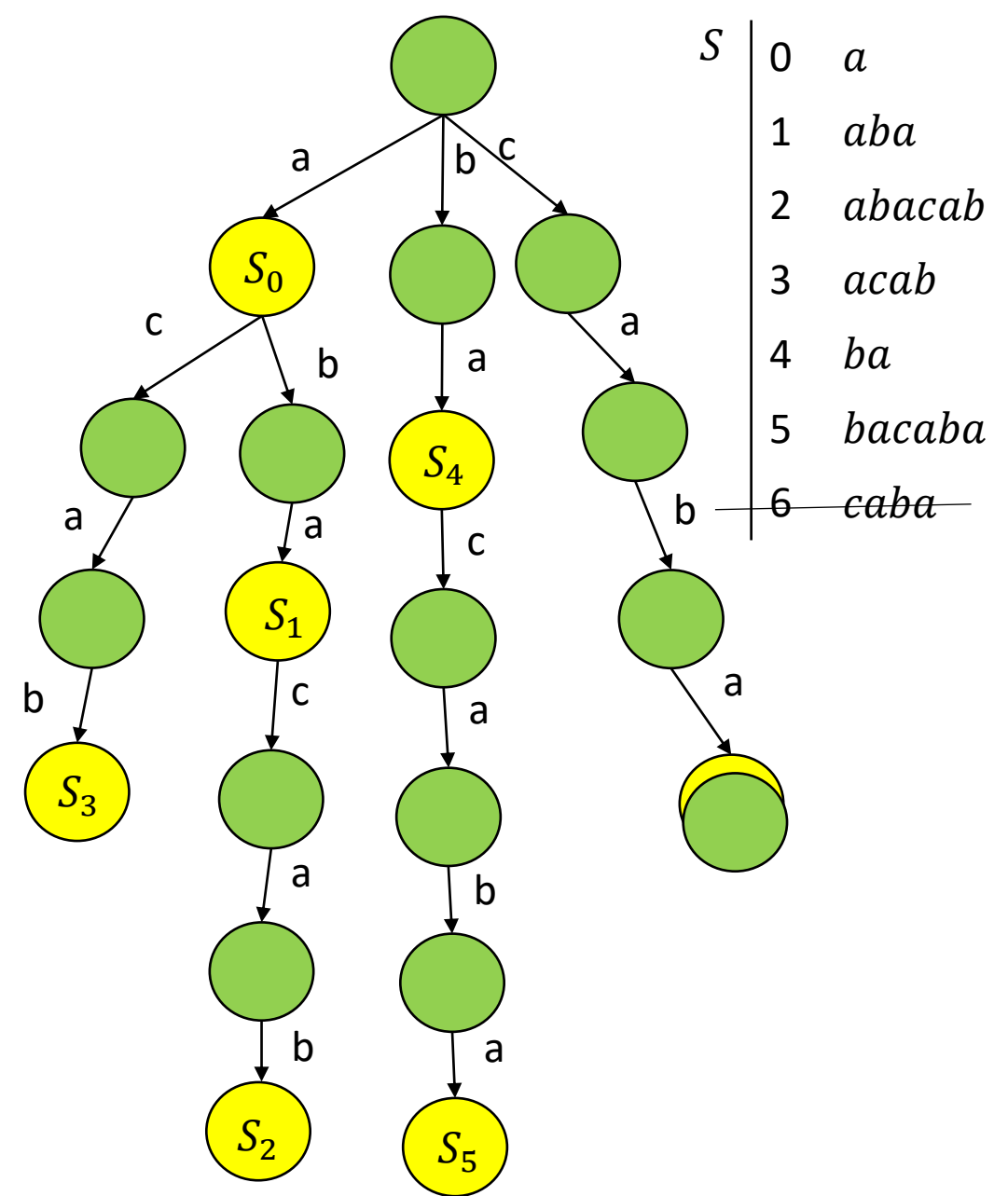
Движение осуществляем, пока не пройдем все символы добавляемой строки.

S	0	a
	1	aba
	2	abacaba
	3	acaba
	4	ba
	5	bacaba
	6	caba

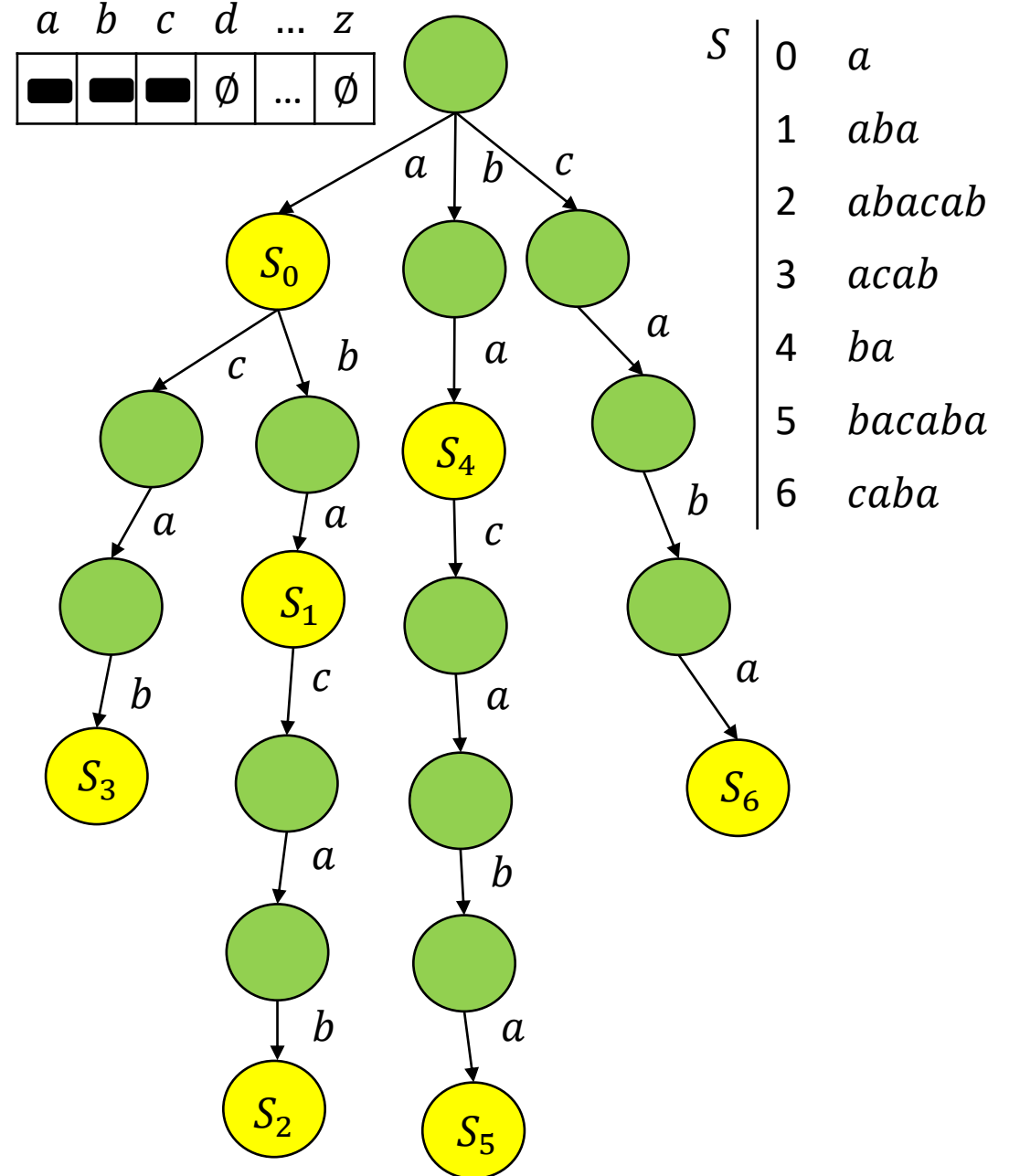


Удалить строку

из бора можно «по ленивому» – снять метку с соответствующей терминальной вершины.



ОЦЕНКИ



Время **построения** структуры данных бор для множества строк из $S = \{S_0, S_1, \dots, S_{n-1}\}$ есть

$$O\left(\sum_{i=0}^{n-1} |S_i|\right),$$

требуемая **память**:

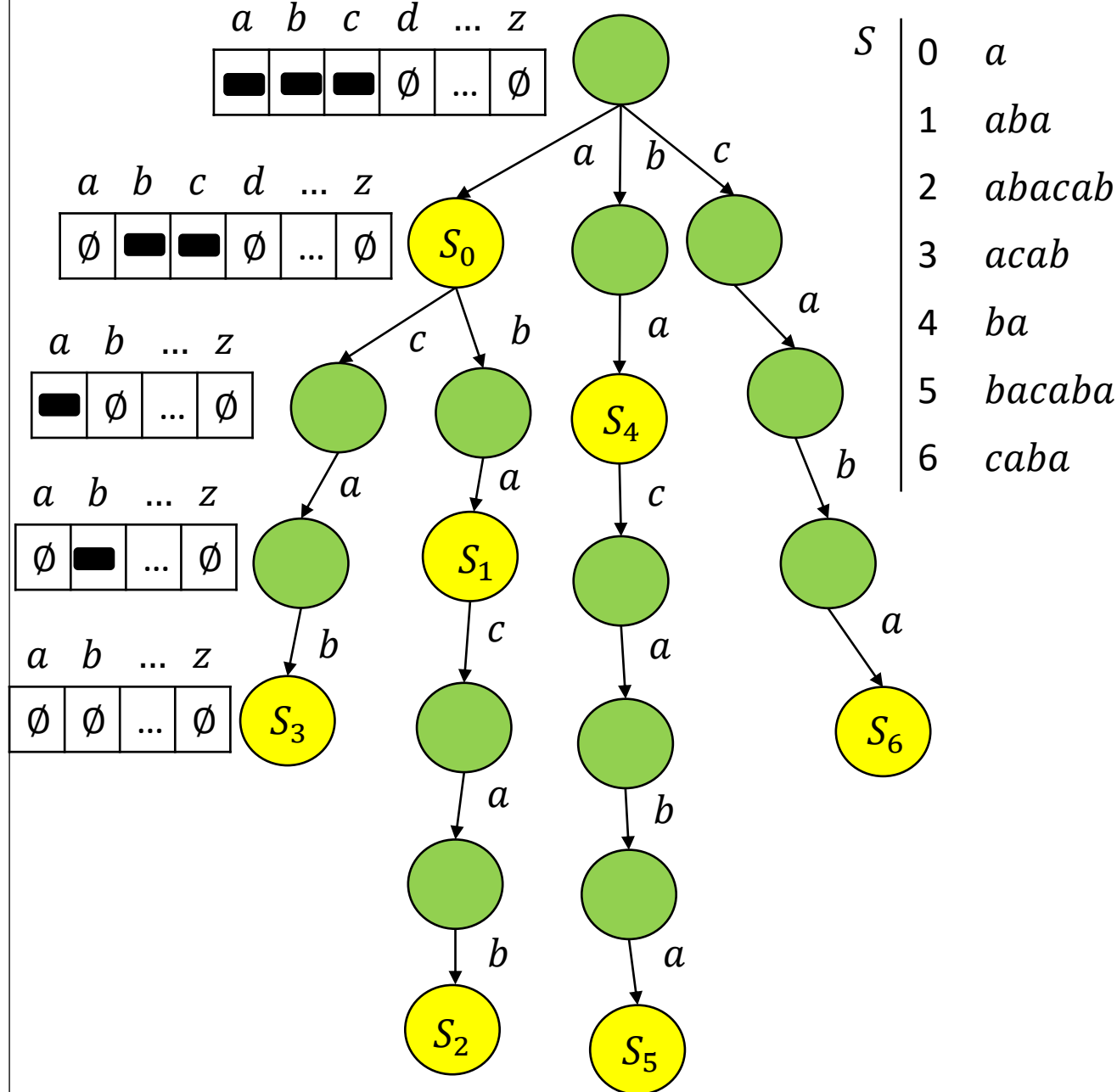
$$O\left(|\Sigma| \cdot \sum_{i=0}^{n-1} |S_i|\right),$$

где n – число строк в множестве S ,
 $|S_i|$ – длина i -ой строки,
 $|\Sigma|$ – размер алфавита.

Время **добавления/поиска** строки

$A = \{a_0, a_1, \dots, a_{l-1}\}$ зависит от длины строки A и не зависит от количества строк в боре:

$$O(l).$$



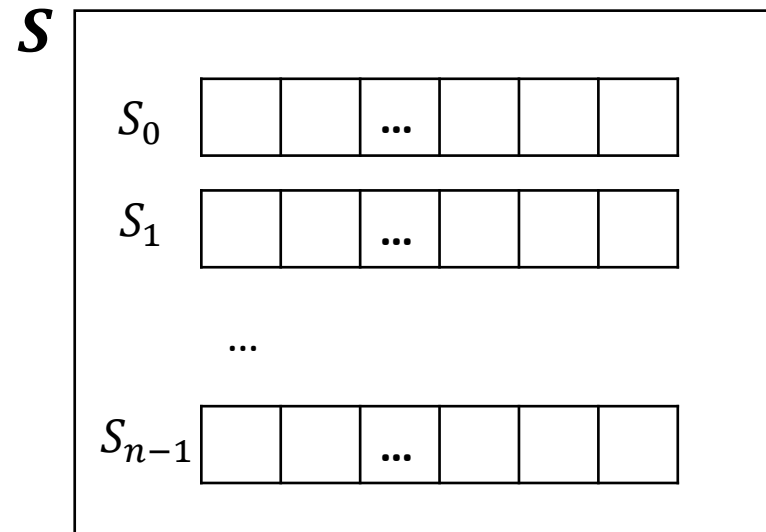
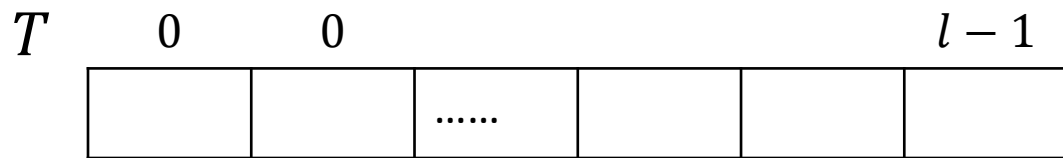
Множество строк $S = \{S_0, S_1, \dots, S_{n-1}\}$	Строка $A = (a_0, a_1, \dots, a_{l-1})$	Определить, есть ли строка A в S?
--	---	--

Структура данных	Время построения	Память	Время поиска
Массив	$O\left(\sum_{i=0}^{n-1} S_i \right)$	$O\left(\sum_{i=0}^{n-1} S_i \right)$	$O(A \cdot n)$
Сбалансированное бинарное поисковое дерево	$O(n \cdot \log n \cdot l_{\max}),$ $l_{\max} = \max_{i=0, \dots, n-1} S_i $	$O\left(\sum_{i=0}^{n-1} S_i \right)$	$O(A \cdot \log n)$
Бор	$O\left(\sum_{i=0}^{n-1} S_i \right)$	$O\left(\Sigma \cdot \sum_{i=0}^{n-1} S_i \right)$	$O(A)$

Задача поиска множества образцов в строке в режиме on-line

- Пусть у нас есть текст $T = (t_0, t_1, \dots, t_{l-1})$, и несколько образцов, которые поступают в режиме *on-line* из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$.

Найти все подстроки текста T , которые совпадали бы с одним из образцов (т.е. необходимо определить, встречается ли образец из S в качестве подстроки в тексте T).



Задача поиска множества образцов в строке в режиме реального времени

Решить задачу можно эффективно,
используя такие структуры данных, как

- ❑ **суффиксный бор**

- ❑ **суффиксный массив**

Задача поиска множества образцов в строке в режиме реального времени

**Суффиксный
бор**

Суффиксный бор - структура данных бор, построенная для всех суффиксов текста $T = (t_0, t_1, \dots, t_{l-1})$.

$T = abacaba$

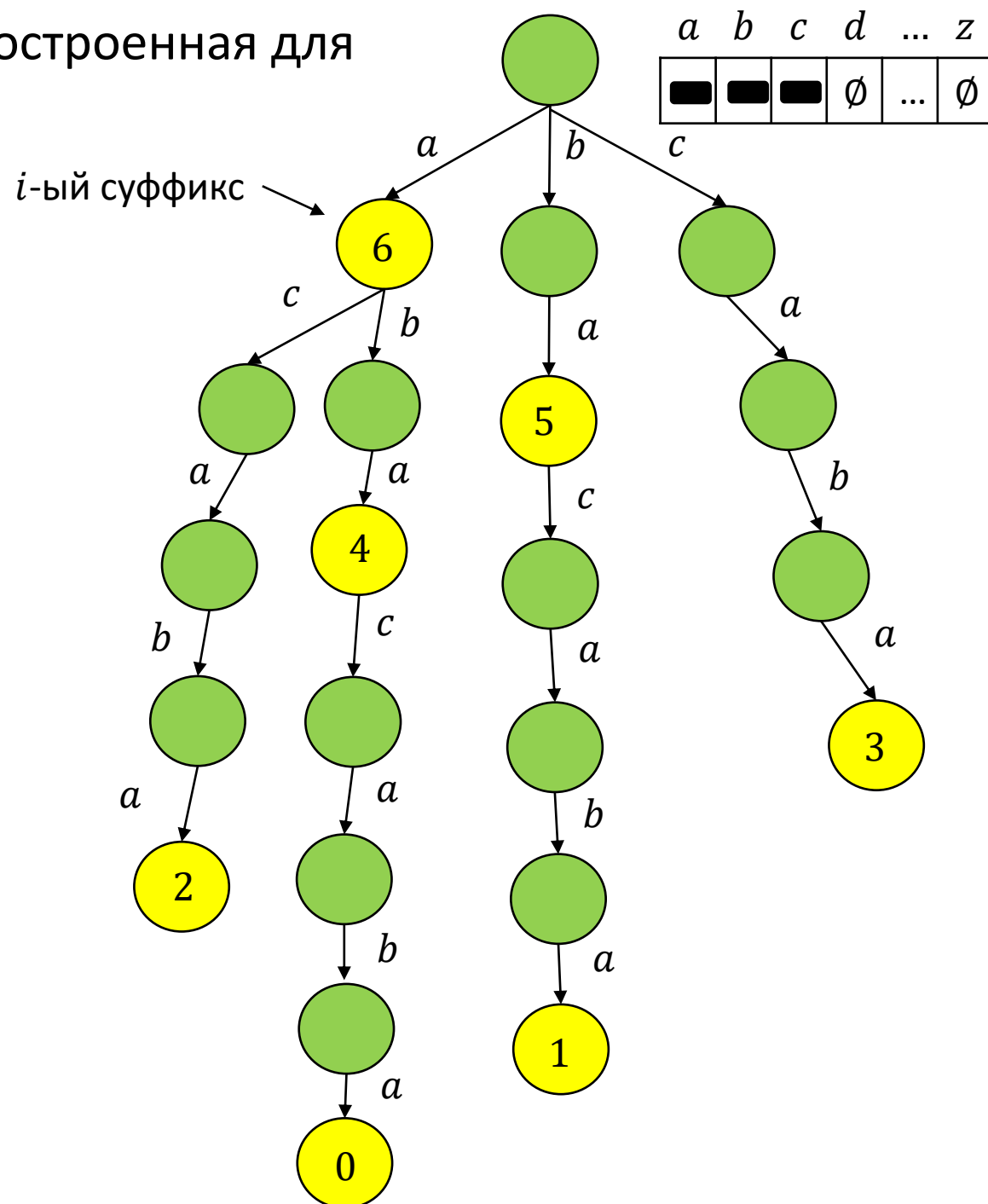
0	1	2	3	4	5	6
a	b	a	c	a	b	a

0-й	<i>abacaba</i>
1-й	<i>bacaba</i>
2-й	<i>acaba</i>
3-й	<i>caba</i>
4-й	<i>aba</i>
5-й	<i>ba</i>
6-й	<i>a</i>

Время построения : $O(l^2)$,

требуемая память: $O(|\Sigma| \cdot l^2)$,

где $l = |T|, |\Sigma|$ - размер алфавита.



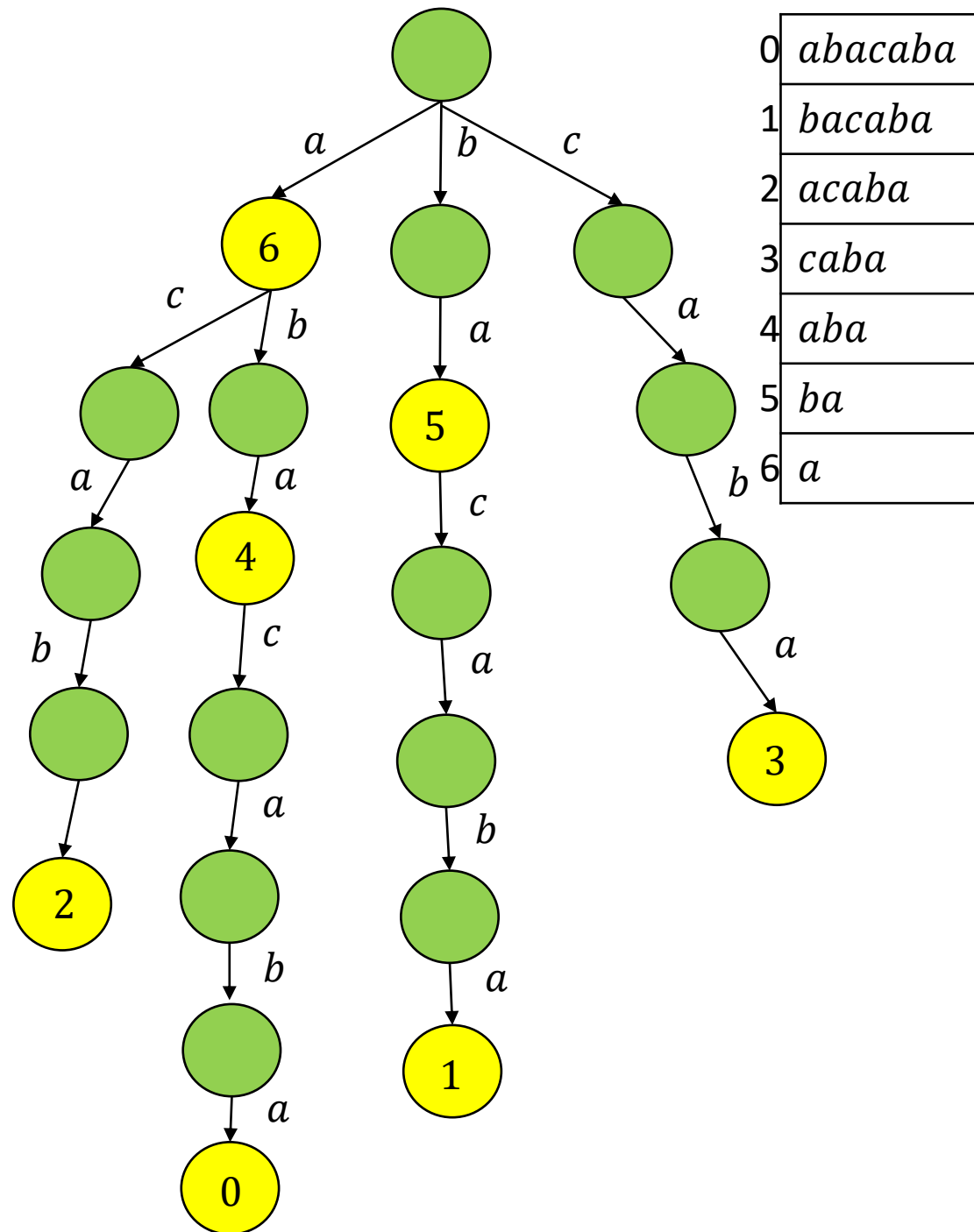
Если в **суффиксном боре** находится некоторый суффикс, то все префиксы этого суффикса также находятся в боре.

Если некоторая строка S_i встречается в качестве подстроки в строке T , то она является префиксом какого-то суффикса строки T . Поэтому суффиксный бор можно использовать для решения задачи поиска множества образцов в строке в режиме on-line.

Время поиска множества образцов в строке:

$$O\left(\sum_{i=0}^{n-1} |S_i|\right) = O(n \cdot l_{max}),$$

где $l_{max} = \max_{i=0, \dots, n-1} |S_i|$.



Задача поиска множества образцов в строке в режиме реального времени

Суффиксный массив

Задан фиксированный текст $T = (t_0, t_1, \dots, t_{l-1})$	On-line поступают образцы из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$	Определить, встречается ли образец из S в качестве подстроки в тексте T	
Структура данных	Время построения	Память	Время поиска
Суффиксный массив	$O(T \cdot \log^2 T)$	$O(T)$	$O\left(\log T \cdot \sum_{i=1}^{n-1} S_i \right)$
Суффиксный бор	$O(T ^2)$	$O(T ^2 \cdot \Sigma)$	$O\left(\sum_{i=0}^{n-1} S_i \right)$

Структура данных **суффиксный массив** была разработана в **1989** году .

Юджин Уимберли «Джин»

Майерс-младший

Eugene Wimberly «Gene» Myers, Jr.



Дата рождения	31 декабря 1953
Место рождения	• Бойсе , США
Страна	• США
Научная сфера	информатика

Уди Манбер

Udi Manber



Место рождения:	Кирьят-Хаим , Израиль
Страна	• США
Научная сфера	информатика
Должность	вице-президент Google по разработкам

Суффиксный массив строки $T = (t_0, t_1, \dots, t_{l-1})$
это последовательность **лексикографически отсортиро-
ванных суффиксов** строки T (очевидно, что достаточно
хранить только индексы суффиксов).

i	0	1	2	3	4	5	6	
T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	$l = 7$

все суффиксы
строки T

0-й	<i>abacaba</i>
1-й	<i>bacaba</i>
2-й	<i>acaba</i>
3-й	<i>caba</i>
4-й	<i>aba</i>
5-й	<i>ba</i>
6-й	<i>a</i>

лексикографическая
сортировка
суффиксов

<i>a</i>
<i>aba</i>
<i>abacaba</i>
<i>acaba</i>
<i>ba</i>
<i>bacaba</i>
<i>caba</i>

6-й	<i>a</i>
4-й	<i>aba</i>
0-й	<i>abacaba</i>
2-й	<i>acaba</i>
5-й	<i>ba</i>
1-й	<i>bacaba</i>
3-й	<i>caba</i>

перестановка индексов суффиксов,
которая задаёт порядок суффиксов
в порядке лексикографической
сортировки

суффиксный
массив
 $P[0..l-1]$

0	6
1	4
2	0
3	2
4	5
5	1
6	3

Пусть есть фиксированный текст $T = (t_0, t_1, \dots, t_{l-1})$:

	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<i>T</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>

На вход в режиме реального времени поступают образцы из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$, и нам необходимо эффективно определять, встречается ли образец $S' \in S$ в качестве подстроки в тексте T (*on-line* версия задачи).

<i>S'</i>	
<i>a</i>	<i>b</i>

Как, используя суффиксный массив, эффективно решить эту задачу?

	0	1	2	3	4	5	6
T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>

	0	1
$S' \in S$	<i>a</i>	<i>b</i>

1) Построили по тексту T суффиксный массив P

0	6	<i>a</i>
1	4	<i>aba</i>
2	0	<i>abacaba</i>
3	2	<i>acaba</i>
4	5	<i>ba</i>
5	1	<i>bacaba</i>
6	3	<i>caba</i>

2) Если строка S' является подстрокой T , то она является префиксом какого-то суффикса, поэтому дихотомией ($LowerBound, UpperBound$) по суффиксному массиву P (осуществляем поиск диапазона суффиксов, у которых префикс совпадает с образцом S').

S'	<i>a</i>	
0	6	<i>a</i>
1	4	<i>aba</i>
2	0	<i>abacaba</i>
3	2	<i>acaba</i>

$L = 0, \quad R = n$
 $L = LowerBound(L, R, s'[0])$
 $R = UpperBound(L, R, s'[0])$
 $L = 0, R = 4$

если в диапазоне у суффиксов, у которых первый символ совпадает, удалить этот символ, то по следующему символу эти префиксы в этом диапазоне также упорядочены лексикографически

S'	<i>a</i>	<i>b</i>
0	6	<i>a</i>
1	4	<i>a</i> <i>ba</i>
2	0	<i>a</i> <i>bacaba</i>
3	2	<i>a</i> <i>caba</i>

$L = LowerBound(L, R, s'[1])$
 $R = UpperBound(L, R, s'[1])$

1	4	<i>a</i> <i>ba</i>
2	0	<i>a</i> <i>bacaba</i>

$L = 1, R = 3$

Время поиска вхождений образца $S' \in S$ в $T = (t_0, t_1, \dots, t_{l-1})$

$$\Theta(|S'| \cdot \log |T|) = \Theta(|S'| \cdot \log l).$$

Как построить суффиксный массив?

Алгоритм построения суффиксного массива для $T = (t_0, t_1, \dots, t_{l-1})$ за время

$$\Theta(|T|^2 + |T| \cdot |\Sigma|)$$

заключается, например, в непосредственном упорядочивании всех суффиксов строки T алгоритмом лексикографической сортировки.

Время работы
алгоритма лексикографической
сортировки

$$\Theta\left(\sum_{i=0}^{l-1} |t^j| + l_{max} \cdot |\Sigma|\right), \text{ где}$$

l — число кортежей,

l_{max} — длина самого длинного кортежа,

$|\Sigma|$ — число различных символов в кортежах.

Как построить суффиксный массив эффективнее?

Если для сортировки суффиксов строки $T = (t_0, t_1, \dots, t_{l-1})$ использовать, например, **сортировку слиянием** (*merge sort*), то время работы алгоритма построения суффиксного массива будет зависеть от того, как будут сравниваться суффиксы.

Непосредственное сравнение двух суффиксов приведет к тому, что время работы алгоритма построения суффиксного массива:

$$O(|T| \cdot (|T| \cdot \log |T|)) = O(|T|^2 \cdot \log |T|).$$

Как научиться сравнивать строки быстрее, например, за константное время?

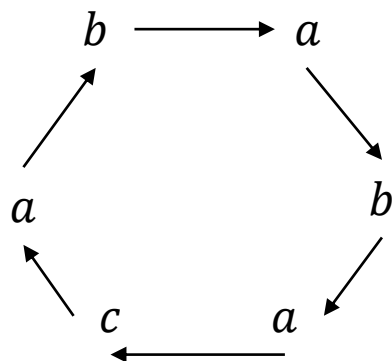
Алгоритм
построения суффиксного массива
строки $T = (t_0, t_1, \dots, t_{l-1})$ за время

$$O(|T| \cdot \log^2 |T|) = \mathbf{O}(l \cdot \log^2 l)$$

Циклический сдвиг строки *T*

это строка, которая получается из исходной строки *T* путем перемещением её первых символов в конец строки.

T = *abacaba*\$



циклический сдвиг длины 1	<i>bacaba</i> <i>a</i>
циклический сдвиг длины 2	<i>acaba</i> <i>ab</i>
циклический сдвиг длины 3	<i>caba</i> <i>aba</i>
циклический сдвиг длины 4	<i>aba</i> <i>abac</i>
циклический сдвиг длины 5	<i>ba</i> <i>abaca</i>
циклический сдвиг длины 6	<i>a</i> <i>abacab</i>
циклический сдвиг длины 6	<i>abacaba</i>

← лексикографически
минимальный
циклический сдвиг

Если к строке T добавить $\$$ (символ, который меньше всех символов строки), выполнить лексикографическую сортировку всех циклических сдвигов $T + \$$, затем удалить из каждого циклического сдвига строки суффикс, который начинается с $\$$, то получим лексикографическую сортировку всех суффиксов строки T .

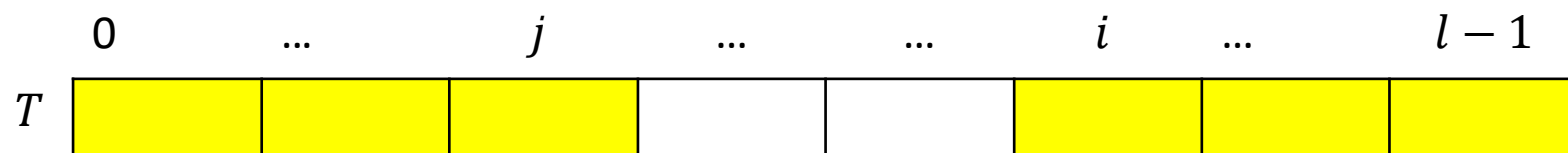
$T = abacaba$

$T + \$ = abacaba + \$$

Циклические сдвиги $T + \$$	Лексикографическая сортировка	Лексикогр. упорядоченные суффиксы строки T
<i>bacaba\$a</i>	<i>\$abacaba</i>	
<i>acaba\$ab</i>	<i>a\$abacab</i>	<i>a</i>
<i>caba\$aba</i>	<i>aba\$abac</i>	<i>aba</i>
<i>aba\$abac</i>	<i>abacaba\$</i>	<i>abacaba</i>
<i>ba\$abaca</i>	<i>acaba\$ab</i>	<i>acaba</i>
<i>a\$abacab</i>	<i>ba\$abaca</i>	<i>ba</i>
<i>\$abacaba</i>	<i>bacaba\$a</i>	<i>bacaba</i>
<i>abacaba\$</i>	<i>caba\$aba</i>	<i>caba</i>

Поэтому научимся эффективно упорядочивать циклические сдвиги строки.

Циклическая подстрока $t[i..j]$, $i > j$ определяется, как $t[i..l-1] + t[0..j]$.



Пусть $l = |T|$, $2^{k-1} < l \leq 2^k$.

Тогда

$$2^{k-1} < l$$

$$k < \log_2 l + 1$$

$$k = \lceil \log_2 l \rceil$$

Выполним $k = \lceil \log_2 l \rceil$ фаз сортировки циклических подстрок.

На каждой фазе выполняется лексикографическая сортировка циклических подстрок одинаковой длины.

На i -ой фазе сортируются циклические подстроки длины 2^i , $0 \leq i \leq k = \lceil \log_2 l \rceil$.

На j — ой фазе поддерживаются следующие массивы:

P^j — перестановка индексов циклических подстрок длины 2^j , которая задаёт порядок циклических подстрок длины 2^j в порядке лексикографической сортировки;

C^j — массив классов эквивалентности (элементы массива - целые числа ≥ 0), где

- ✓ $c^j[i]$ — номер класса эквивалентности для циклической подстроки длины 2^j , начинающейся в позиции i ;
- ✓ если две циклические строки совпадают, то у них один класс эквивалентности;
- ✓ если одна циклическая подстрока лексикографически меньше другой, то номер класса она получит меньший.

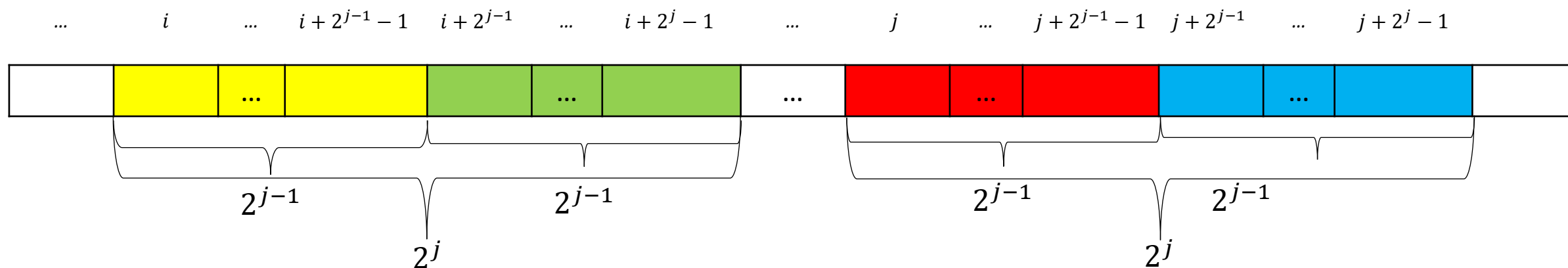
i	0	1	2	3	4	5	6	7
T	a	b	a	c	a	b	a	\$
P^0	7	0	2	4	6	1	5	3
C^0	1	2	1	3	1	2	1	0

на 0-й фазе, когда упорядочиваются подстроки длины 1, сформировать указанные массивы можно устойчивым алгоритмом **сортировки подсчётом** за время $O(l + |\Sigma|)$

На фазах, начиная с 1-й, необходимо лексикографически упорядочивать подстроки длины которых больше 1.

Делать это можно, например, сортировкой слиянием, но надо научиться, используя информацию, полученную на предыдущих фазах, сравнивать две циклические подстроки одинаковой длины за константное время.

Сравнение двух циклических подстрок на j -ой фазе можно **выполнить за константное время**, так как мы знаем для каждой циклической подстроки длины 2^{j-1} к какому классу эквивалентности она отнесена, а любая циклическая подстрока длины 2^j состоит из двух циклических подстрок длины 2^{j-1} .



если $c^{j-1}[i] < c^{j-1}[j]$, то $t[i..i + 2^j - 1] < t[j..j + 2^j - 1]$

если $c^{j-1}[i] > c^{j-1}[j]$, то $t[i..i + 2^j - 1] > t[j..j + 2^j - 1]$

если $c^{j-1}[i] = c^{j-1}[j]$, то

если $c^{j-1}[i + 2^{j-1}] < c^{j-1}[j + 2^{j-1}]$, то $t[i..i + 2^j - 1] < t[j..j + 2^j - 1]$

иначе, если $c^{j-1}[i + 2^{j-1}] > c^{j-1}[j + 2^{j-1}]$, то $t[i..i + 2^j - 1] > t[j..j + 2^j - 1]$

иначе $t[i..i + 2^j - 1] = t[j..j + 2^j - 1]$

Массив C^j формируется на j —ой фазе за время $\Theta(l)$ после формирования массива P^j , проходом по массиву P^j слева направо и сравнением (за константу) двух циклических подстрок длины 2^j на равенство:

- номер класса для $t[p^j[0]..p^j[0] + 2^j - 1]$ полагаем равным 0, т.е. $c[p^j[0]] := 0$;
- для $1 \leq i \leq l$
если $t[p^j[i]..p^j[i] + 2^j - 1] = t[p^j[i-1]..p^j[i-1] + 2^j - 1]$, то номер класса $c^j[p^j[i]]$ для $t[p^j[i]..p^j[i] + 2^k - 1]$ остается таким же, как и для $t[p^j[i-1]..p^j[i-1] + 2^k - 1]$ (т. е. $c^j[p^j[i]] = c^j[p^j[i-1]]$), иначе он увеличивается на 1 (т. е. $c^j[p^j[i]] := c^j[p^j[i-1]] + 1$).

Так как к строке T добавляется \$, то это гарантирует, что каждый суффикс будет иметь свой класс эквивалентности.

2^0

i	0	1	2	3	4	5	6	7	$l = 8$
T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	\$	
C^0	1	2	1	3	1	2	1	0	

2^1

i	0	1	2	3	4	5	6	7
T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	\$
P^1	7	6	0	4	2	1	5	3
C^1	2	4	3	5	2	4	1	0

$t[i..i + 2^1 - 1]$

0	$a_{1\text{кЛ}}b_{2\text{кЛ}}$
1	$b_{2\text{кЛ}}a_{1\text{кЛ}}$
2	$a_{1\text{кЛ}}c_{3\text{кЛ}}$
3	$c_{3\text{кЛ}}a_{1\text{кЛ}}$
4	$a_{1\text{кЛ}}b_{2\text{кЛ}}$
5	$b_{2\text{кЛ}}a_{1\text{кЛ}}$
6	$a_{1\text{кЛ}} \$_{0\text{кЛ}}$
7	$\$_{0\text{кЛ}}a_{1\text{кЛ}}$

P^1

$t[p^1[i]..p^1[i] + 2^j - 1]$

0	7	$\$_{0\text{кЛ}}a_{1\text{кЛ}}$
1	6	$a_{1\text{кЛ}} \$_{0\text{кЛ}}$
2	0	$a_{1\text{кЛ}}b_{2\text{кЛ}}$
3	4	$a_{1\text{кЛ}}b_{2\text{кЛ}}$
4	2	$a_{1\text{кЛ}}c_{3\text{кЛ}}$
5	1	$b_{2\text{кЛ}}a_{1\text{кЛ}}$
6	5	$b_{2\text{кЛ}}a_{1\text{кЛ}}$
7	3	$c_{3\text{кЛ}}a_{1\text{кЛ}}$

$c[p[0]] = 0$

для $1 \leq i \leq l - 1$

если $t[p[i]..p[i] + 2^j - 1] = t[p[i - 1]..p[i - 1] + 2^j - 1]$,

то $c[p[i]] = c[p[i - 1]]$, иначе $c[p[i]] = c[p[i - 1]] + 1$).

2^1	i	0	1	2	3	4	5	6	7
	T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	\$
	C^1	2	4	3	5	2	4	1	0

2^2	i	0	1	2	3	4	5	6	7
	T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	\$
	p^2	7	6	4	0	2	5	1	3
	C^2	3	6	4	7	2	5	1	0

$t[i..i + 2^2 - 1]$

0	$(ab)_{2\text{КЛ}}(ac)_{3\text{КЛ}}$
1	$(ba)_{4\text{КЛ}}(ca)_{5\text{КЛ}}$
2	$(ac)_{3\text{КЛ}}(ab)_{2\text{КЛ}}$
3	$(ca)_{5\text{КЛ}}(ba)_{4\text{КЛ}}$
4	$(ab)_{2\text{КЛ}}(a\$)_{1\text{КЛ}}$
5	$(ba)_{4\text{КЛ}}(\$a)_{0\text{КЛ}}$
6	$(a\$)_{1\text{КЛ}}(ab)_{2\text{КЛ}}$
7	$(\$a)_{0\text{КЛ}}(ba)_{4\text{КЛ}}$

$p^2 \quad t[p^2[i]..p^2[i] + 2^2 - 1]$

0	7	$(\$a)_{0\text{КЛ}}(ba)_{4\text{КЛ}}$
1	6	$(a\$)_{1\text{КЛ}}(ab)_{2\text{КЛ}}$
2	4	$(ab)_{2\text{КЛ}}(a\$)_{1\text{КЛ}}$
3	0	$(ab)_{2\text{КЛ}}(ac)_{3\text{КЛ}}$
4	2	$(ac)_{3\text{КЛ}}(ab)_{2\text{КЛ}}$
5	5	$(ba)_{4\text{КЛ}}(\$a)_{0\text{КЛ}}$
6	1	$(ba)_{4\text{КЛ}}(ca)_{5\text{КЛ}}$
7	3	$(ca)_{5\text{КЛ}}(ba)_{4\text{КЛ}}$

$c[p[0]] = 0$

для $1 \leq i \leq l - 1$

если $s[p[i]..p[i] + 2^j - 1] = s[p[i - 1]..p[i - 1] + 2^j - 1]$,
то $c[p[i]] = c[p[i - 1]]$, иначе $c[p[i]] = c[p[i - 1]] + 1$.

2^2

i

0

1

2

3

4

5

6

7

$l = 8$

T

a

b

a

c

a

b

a

\$

C^2

3

6

4

7

2

5

1

0

2^3

i

0

1

2

3

4

5

6

7

T

a

b

a

c

a

b

a

\$

P^3

7

6

4

0

2

5

1

3

C^3

3

6

4

7

2

5

1

0

$t[i..i + 2^3 - 1]$

0	$(abac)_{3_{\text{КЛ}}}(aba\$)_{2_{\text{КЛ}}}$
1	$(baca)_{4_{\text{КЛ}}}(ba\$a)_{5_{\text{КЛ}}}$
2	$(acab)_{4_{\text{КЛ}}}(a\$ab)_{1_{\text{КЛ}}}$
3	$(caba)_{7_{\text{КЛ}}}(\$aba)_{0_{\text{КЛ}}}$
4	$(aba\$)_{2_{\text{КЛ}}}(abac)_{3_{\text{КЛ}}}$
5	$(ba\$a)_{5_{\text{КЛ}}}(baca)_{6_{\text{КЛ}}}$
6	$(a\$ab)_{1_{\text{КЛ}}}(acab)_{4_{\text{КЛ}}}$
7	$(\$aba)_{0_{\text{КЛ}}}(caba)_{7_{\text{КЛ}}}$

P^3

$t[p^3[i]..p^3[i] + 2^3 - 1]$

\downarrow

\downarrow

0	7	$(\$aba)_{0_{\text{КЛ}}}(caba)_{7_{\text{КЛ}}}$
1	6	$(a\$ab)_{1_{\text{КЛ}}}(acab)_{4_{\text{КЛ}}}$
2	4	$(aba\$)_{2_{\text{КЛ}}}(abac)_{3_{\text{КЛ}}}$
3	0	$(abac)_{3_{\text{КЛ}}}(aba\$)_{2_{\text{КЛ}}}$
4	2	$(acab)_{4_{\text{КЛ}}}(a\$ab)_{1_{\text{КЛ}}}$
5	5	$(ba\$a)_{5_{\text{КЛ}}}(baca)_{6_{\text{КЛ}}}$
6	1	$(baca)_{4_{\text{КЛ}}}(ba\$a)_{5_{\text{КЛ}}}$
7	3	$(caba)_{7_{\text{КЛ}}}(\$aba)_{0_{\text{КЛ}}}$

$c[p[0]] = 0$

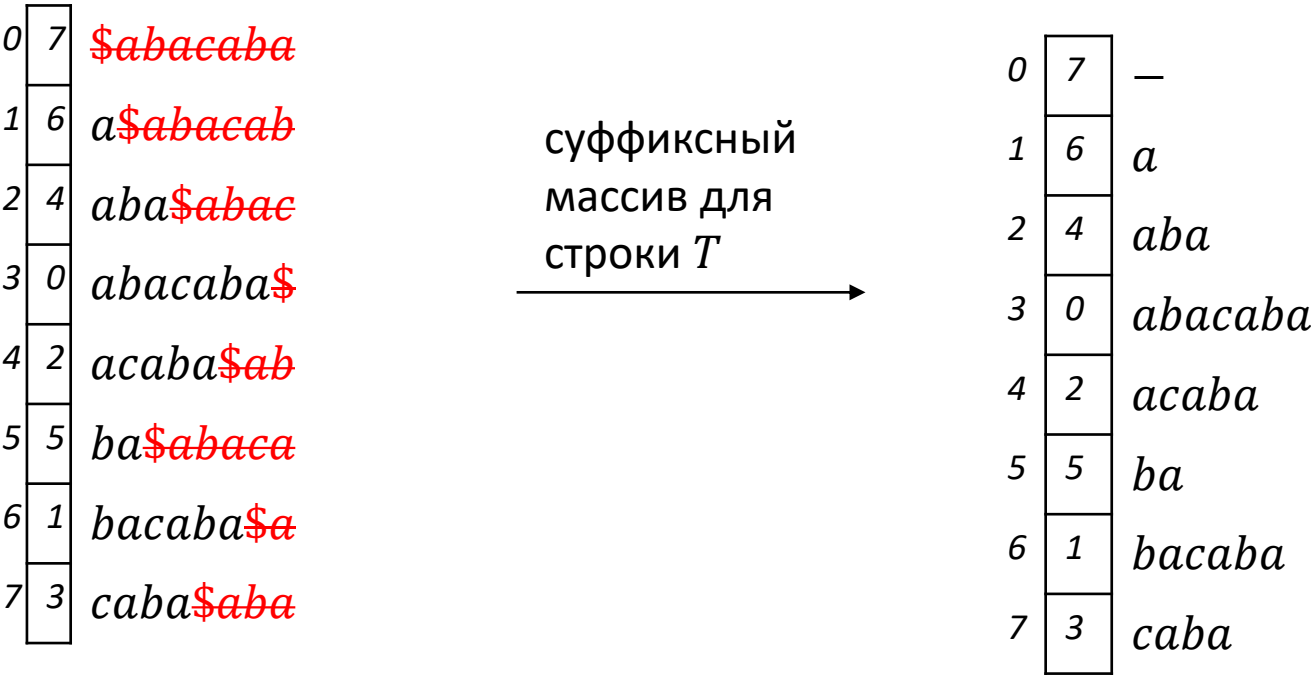
для $1 \leq i \leq l - 1$

если $t[p[i]..p[i] + 2^j - 1] = t[p[i - 1]..p[i - 1] + 2^j - 1]$,

то $c[p[i]] = c[p[i - 1]]$, иначе $c[p[i]] = c[p[i - 1]] + 1$).

После того, как выполнена сортировка циклических сдвигов строки $T + \$$, удалим из каждого циклического сдвига строки её суффикс, который начинается с $\$$, получим суффиксный массив для строки T :

i	0	1	2	3	4	5	6	7	$l = 8$
T	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>	\$	
P^3	7	6	4	0	2	5	1	3	
C^3	3	6	4	7	2	5	1	0	



Если на j -й фазе выполнять лексикографическую сортировку l циклических подстрок длины 2^j ($j < \log_2 |T| + 1$) **сортировкой слиянием**, то так как сравнение двух циклических подстрок мы выполняем за $O(1)$, а массив C формируется за $O(|T|)$, то время выполнения одной фазы

$$O(|T| \cdot \log |T|) + O(|T|) = O(|T| \cdot \log |T|).$$

Учитывая, что число фаз $k = \lceil \log |T| \rceil < \log |T| + 1$, получаем, что **время работы алгоритма лексикографической сортировки циклических сдвигов строки** (построения суффиксного массива строки T)

$$O(|T| \cdot \log^2 |T|).$$

Требуемая память — $O(|T|)$.

Задан фиксированный текст $T = (t_0, t_1, \dots, t_{l-1})$	On-line поступают образцы из множества $S = \{S_0, S_1, \dots, S_{n-1}\}$	Определить, встречается ли образец из S в качестве подстроки в тексте T
---	--	---

Структура данных	Время построения	Память	Время поиска
Суффиксный массив	$O(T \cdot \log^2 T)$	$O(T)$	$O\left(\log T \cdot \sum_{i=1}^{n-1} S_i \right)$
Суффиксный бор	$O(T ^2)$	$O(T ^2 \cdot \Sigma)$	$O\left(\sum_{i=0}^{n-1} S_i \right)$



БЕЛОРУССКИЙ
ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ

Спасибо за внимание!