

Строковые алгоритмы и структуры данных

Основные определения.

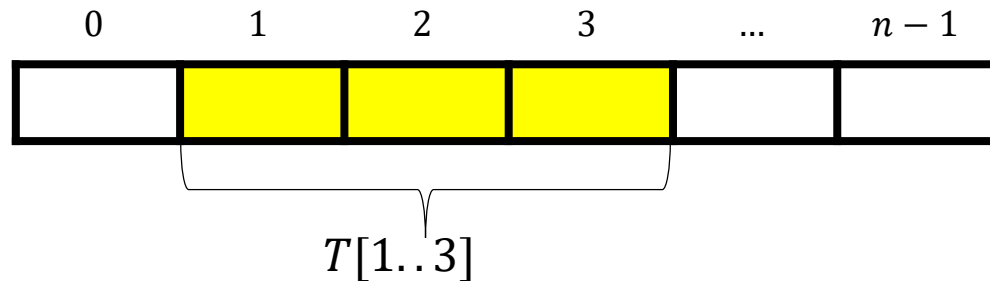
Задано некоторое конечное непустое множество символов Σ , называемое алфавитом ($|\Sigma|$ — мощность алфавита).

Строка - произвольная конечная последовательность символов из алфавита:

$$T = (t_0, t_1, \dots, t_{n-1}), t_i \in \Sigma.$$

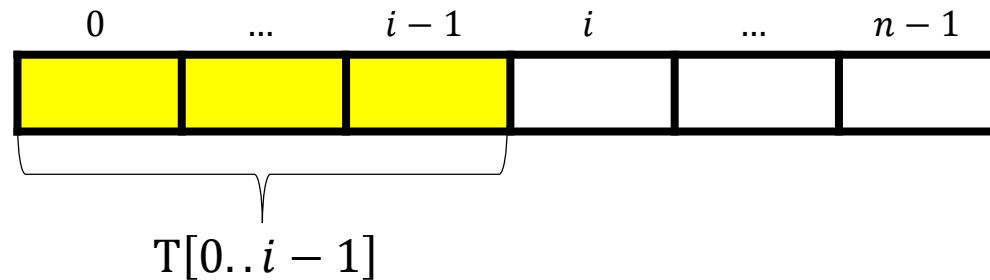
Подстрока – непрерывная последовательность строки

$$T[i..j] = (t_i, t_{i+1}, \dots, t_j), t_i \in \Sigma.$$

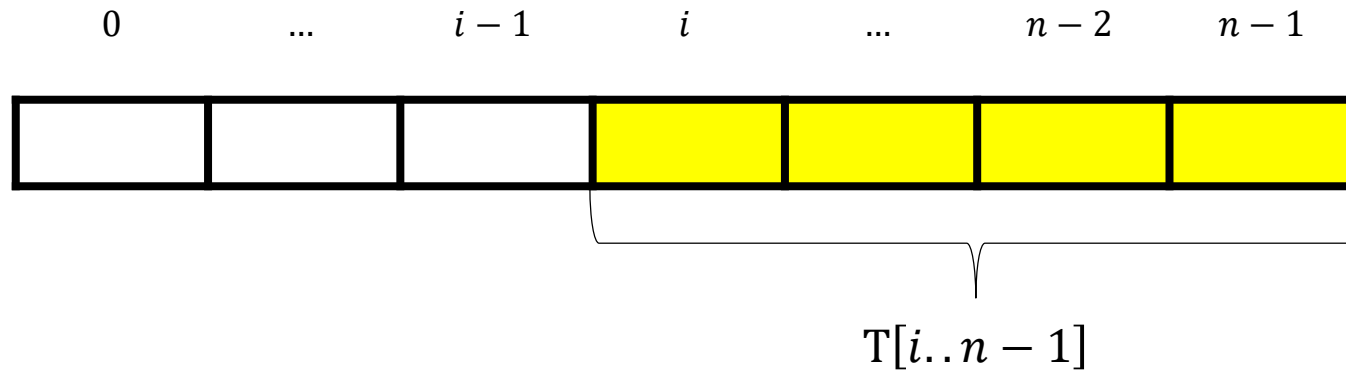


Если не оговорено иное, то при $i > j$ считаем, что подстроки $T[i..j]$ не существует.

Подстрока $T[0..i - 1]$, которая состоит из первых i символов строки, называется **префиксом** ($T[0..i - 1]$ - префикс длины i).



Подстрока $T[i..n - 1]$, которая состоит из последних $n - i$ символов строки, называется i -ым **суффиксом**.



Собственный суффикс/префикс - не совпадающий со всей строкой.

Задача поиска вхождений подстроки в строке

- Пусть есть непустая строка $T = (t_0, t_1, \dots, t_{n-1})$, называемая **текстом**, и непустая строка $S = (s_0, s_1, \dots, s_{m-1})$, называемая **образцом**.

Задача поиска вхождений подстроки в строке состоит в том, чтобы найти все подстроки строки T , которые совпадают с образцом S .

☐ полиномиальное хеширование строк

☐ префиксная функция

☐ Z-функция

Задачу поиска всех подстрок строки $T = (t_0, t_1, \dots, t_{n-1})$, которые совпадают с образцом $S = (s_0, s_1, \dots, s_{m-1})$ можно решить непосредственно – последовательно перебирать в T начальный символ для подстроки, а затем последовательно сравнивать подстроку T с образцом S .

	0	1	2	3		0	1	2	3	4	5	6	7	8	9	10	11	12	13	
S	Ф	П	М	И		T	Б	Ф	П	М	Ф	П	М	И	Г	Ф	П	М	И	У
							Ф	П	М	И	И	И	И	И		Ф	П	М	И	

Трудоемкость поиска такого алгоритма - $O(n \cdot m)$.

Существуют алгоритмы, которые решают эту задачу за время $O(n + m)$, которые мы далее и рассмотрим.

Предположим, что на вход поступает строка:

$$S = (s_0, s_1, \dots, s_n)$$

Можно считать, что элементы строки S – целые числа от 1 до L , где $L = |\Sigma|$ – размер алфавита:

для этого каждый символ строки S заменяем его порядковым номером в алфавите ($s_i \rightarrow a' + 1$).

$$S = (a, b, c, d, b) \rightarrow (1, 2, 3, 4, 2)$$

Прямой полиномиальный хеш (хеширование слева направо)

это число

$$h_{\text{пр.}}(S) = (s_0 + s_1 \cdot k^1 + s_2 \cdot k^2 + \dots + s_n \cdot k^n) \bmod m$$

Обратный полиномиальный хеш (хеширование справа налево)

$$h_{\text{обр.}}(S) = (s_0 \cdot k^n + s_1 \cdot k^{n-1} + s_2 \cdot k^{n-2} + \dots + s_n) \bmod m$$

Рекомендуют выбирать следующие константы:

$k > L$ (натуральное число, чуть больше размера алфавита $L = |\Sigma|$);

m — достаточно большое целое число, например, $m = 2^{64}$;

k и m должны быть взаимно просты, например, если $m=2^{64}$, то k — нечётное число.

Предположим, что алфавит небольшой, $L = |\Sigma| = 9$.

Пусть $k = 10$ и $m = 1\,000\,007$.

$$h_{\text{пр.}}((a, b, c, d, b) \rightarrow (1, 2, 3, 4, 2)) = 24\,321$$

$$h_{\text{обр.}}((a, b, c, d, b) \rightarrow (1, 2, 3, 4, 2)) = 12\,342$$

Две одинаковые строки будут иметь обязательно один хеш, а вероятность коллизий крайне мала.

Вычислить значения функций

$$\mathbf{h}_{\text{пр.}}(S) \text{ и } \mathbf{h}_{\text{обр.}}(S)$$

можно за линейное от длины строки S время $\Theta(n)$.

1) Прямое полиномиальное хеширование:

$$h_{\text{пр.}}(S) = (s_0 + s_1 \cdot k^1 + s_2 \cdot k^2 + \dots + s_n \cdot k^n) \bmod m.$$

Выполним сначала за время $\Theta(n)$ подсчёт нужных степеней.

Обозначим $k[i] = k^i \bmod m$,

тогда:

$$\begin{cases} k[0] = 1 \\ k[i] = (k[i-1] \cdot k) \bmod m, i = 1, \dots, n. \end{cases}$$

Обозначим $h_{\text{пр.}}(0:i)$ - хеш для префикса строки $S[0..i]$:

$$h_{\text{пр.}}(0:i) = \underbrace{(s_0 + s_1 \cdot k^1 + s_2 \cdot k^2 + \dots + s_i \cdot k^i)}_{\text{mod } m} \bmod m.$$

Тогда

$$\begin{cases} h_{\text{пр.}}(0:0) = s_0 \\ h_{\text{пр.}}(0:i) = (h_{\text{пр.}}(0:i-1) + s_i \cdot k[i]) \bmod m, \quad i = 1, \dots, n \end{cases}$$

Так как значение хеша каждого следующего префикса выражается через значение хеша предыдущего префикса, то линейным проходом по всей строке за время $\Theta(n)$ можно вычислить хеши для всех префиксов строки.

2) Обратное полиномиальное хеширование:

$$h_{\text{обп.}}(s) = (s_0 \cdot k^n + s_1 \cdot k^{n-1} + s_2 \cdot k^{n-2} + \dots + s_n) \bmod m$$

$$\left| \begin{array}{l} h_{\text{обп.}} = 0 \\ i = 0, \dots, n \\ h_{\text{обп.}} = (h_{\text{обп.}} \cdot k + s_i) \bmod m, \end{array} \right.$$

$$h_{\text{обп.}} = 0$$

$$i = 0 \quad h_{\text{обп.}} = (0 \cdot k + s_0) \bmod m$$

$$i = 1 \quad h_{\text{обп.}} = (s_0 \cdot k + s_1) \bmod m$$

$$i = 2 \quad h_{\text{обп.}} = (s_0 \cdot k^2 + s_1 \cdot k^1 + s_2) \bmod m$$

...

$$i = n \quad h_{\text{обп.}} = (s_0 \cdot k^n + \dots + s_{n-1} \cdot k^1 + s_n) \bmod m$$

Так как хеш – это значение многочлена, то для многих строковых операций можно быстро пересчитывать хеш результата.

Конкатенация двух строк

$$X = (x_0, x_1, \dots, x_n)$$

$$Y = (y_0, y_1, \dots, y_r)$$

Конкатенация строк X и Y , это строка, которая получается в результате приписывания последовательности символов строки Y после последовательности символов строки X :

$$X \cdot Y = (x_0, x_1, \dots, x_n, y_0, y_1, \dots, y_r)$$

$$h_{\text{пр.}}(X) = (x_0 + x_1 \cdot k^1 + x_2 \cdot k^2 + \dots + x_n \cdot k^n) \bmod m$$

$$h_{\text{пр.}}(Y) = (y_0 + y_1 \cdot k^1 + y_2 \cdot k^2 + \dots + y_r \cdot k^r) \bmod m$$

$$h_{\text{пр.}}(X \cdot Y) = \underbrace{(x_0 + x_1 \cdot k^1 + x_2 \cdot k^2 + \dots + x_n \cdot k^n)}_{h_{\text{пр.}}(X)} + \underbrace{y_0 \cdot k^{n+1} + y_1 \cdot k^{n+2} + \dots + y_r \cdot k^{n+r}}_{h_{\text{пр.}}(Y) \cdot k^{n+1}} \bmod m$$

$$h_{\text{пр.}}(X \cdot Y) = (x_0 + x_1 \cdot k^1 + x_2 \cdot k^2 + \dots + x_n \cdot k^n + y_0 \cdot k^{n+1} + y_1 \cdot k^{n+2} + \dots + y_r \cdot k^{n+r}) \bmod m$$

$$h_{\text{пр.}}(X \cdot Y) = (\underbrace{(x_0 + \dots + x_n \cdot k^n) \bmod m}_{h_{\text{пр.}}(X)} + k^{n+1} \cdot \underbrace{(y_0 + \dots + y_r \cdot k^r) \bmod m}_{h_{\text{пр.}}(Y)}) \bmod m$$

$$h_{\text{пр.}}(X \cdot Y) = (h_{\text{пр.}}(X) + k^{|X|} \cdot h_{\text{пр.}}(Y)) \bmod m$$

Хеш для подстроки $X[l..r] = (x_l, x_{l+1}, \dots, x_r)$

$$h_{\text{пр.}}(l:r) = x_l + x_{l+1} \cdot k^1 + \dots + x_r \cdot k^{r-l}$$

как вычислить быстрее?

Предположим, что $m = 2^{64}$ и значение полинома помещается в 64 бита. Это равносильно тому, что на компьютере можно выполнять все операции над 64-битными числами и не задумываться о том, чтобы брать остаток.

Но в C++ для этого необходимо использовать не `int64_t` (обычно, синоним `long long`), а `uint64_t` (`unsigned long long`), потому что переполнение знаковых целочисленных типов в C++ — опасный пример неопределённого поведения.

$$h_{\text{пр.}}(l:r) = x_l + x_{l+1} \cdot k^1 + \dots + x_r \cdot k^{r-l}$$

$$h_{\text{пр.}}(l:r) \cdot k^l = x_l \cdot k^l + x_{l+1} \cdot k^{l+1} + \dots + x_r \cdot k^r$$

$$h_{\text{пр.}}(X) = \underbrace{x_0 + x_1 \cdot k^1 + \dots + x_{l-1} \cdot k^{l-1}}_{\text{}} + \underbrace{x_l \cdot k^l + x_{l+1} \cdot k^{l+1} + \dots + x_{r-1} \cdot k^{r-1} + x_r \cdot k^r}_{\text{}} + \dots + x_n \cdot k^n$$

$$h_{\text{пр.}}(l:r) \cdot k^l = h_{\text{пр.}}(0:r) - h_{\text{пр.}}(0:l-1)$$

$$h_{\text{пр.}}(l:r) = \frac{h_{\text{пр.}}(0:r) - h_{\text{пр.}}(0:l-1)}{k^l}$$

Из формулы следует, что для вычисления хешей любой подстроки, необходимо знать лишь хеши всех префиксов этой строки.

Теперь, чтобы сравнить на равенство две любые подстроки строки S , можно сравнить соответствующие хеши этих подстрок. Так как предполагается, что значение полинома помещается в 64 бита, то считаем, что сравнение двух хешей будет выполнено за $O(1)$.

Для реализации математической формулы

$$h_{\text{пр.}}(l:r) = \left(\frac{h_{\text{пр.}}(0:r) - h_{\text{пр.}}(0:l-1)}{k^l} \right) \bmod m,$$

необходимо найти обратный элемент для k^l по модулю m .

Для того, чтобы найти обратный элемент для p по модулю m , нужно вычислить значение функции Эйлера $\varphi(m)$ (число взаимно простых с m чисел) и возвести p в степень $\varphi(m) - 1$ ([https://ru.m.wikipedia.org/wiki/Обратное по модулю число](https://ru.m.wikipedia.org/wiki/Обратное_по_модулю_число)).

В случае полиномиального хеша с $m = 2^{64}$ обратным числу p будет число $p^{2^{63}-1}$.

Вычислить величину $p^{2^{63}-1}$ можно бинарным возведением в степень:

$$p^n = \begin{cases} (p^{n/2})^2, n - \text{чётное} \\ p \cdot (p^{n-1/2})^2, n - \text{нечётное} \end{cases}$$

Если предподсчитать степени обратного элемента по выбранному модулю m , то сравнение на равенство двух подстрок можно выполнять за $O(1)$.

Так как деление достаточно сложная операция, то для сравнения подстрок на равенство на практике используют формулу без деления (приведение к одной степени).

Предположим, что $l_2 > l_1$, тогда

$$h_{\text{пр.}}(l_1:r_1) \stackrel{?}{=} h_{\text{пр.}}(l_2:r_2) \\ k^{l_2-l_1} \cdot (h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1-1)) \stackrel{?}{=} (h_{\text{пр.}}(0:r_2) - h_{\text{пр.}}(0:l_2-1))$$

Так как деление достаточно сложная операция, то для сравнения подстрок на равенство на практике используют формулу без деления (приведение к одной степени).

Предположим, что $l_2 > l_1$, тогда для того, чтобы $h_{\text{пр.}}(l_1:r_1) = h_{\text{пр.}}(l_2:r_2)$

надо, чтобы $k^{l_2-l_1} \cdot (h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1 - 1)) = (h_{\text{пр.}}(0:r_2) - h_{\text{пр.}}(0:l_2 - 1))$

Действительно, по определению прямого хеширования

$$(1) \quad k^{l_1} \cdot h_{\text{пр.}}(l_1:r_1) = h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1 - 1)$$

$$(2) \quad k^{l_2} \cdot h_{\text{пр.}}(l_2:r_2) = h_{\text{пр.}}(0:r_2) - h_{\text{пр.}}(0:l_2 - 1)$$

Домножим (1) на $k^{l_2-l_1}$

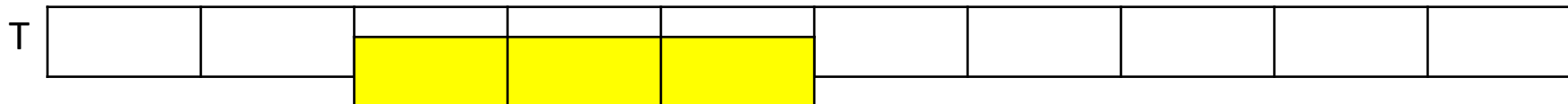
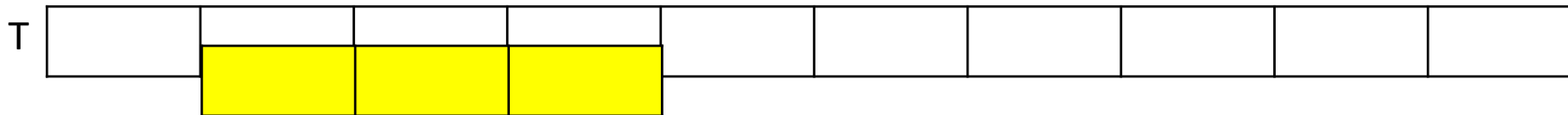
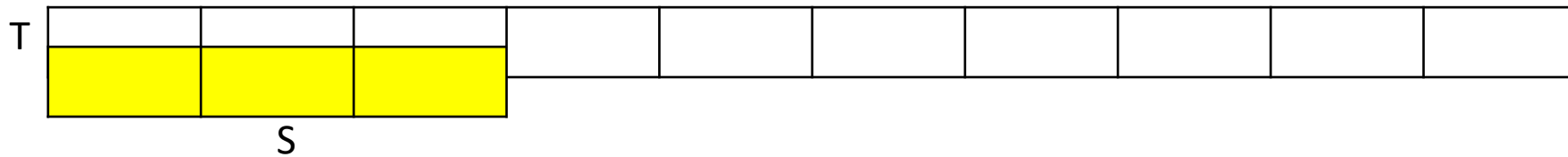
$$k^{l_2-l_1} \cdot k^{l_1} \cdot h_{\text{пр.}}(l_1:r_1) = k^{l_2-l_1} \cdot (h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1 - 1))$$

$$(3) \quad k^{l_2} \cdot h_{\text{пр.}}(l_1:r_1) = k^{l_2-l_1} \cdot (h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1 - 1))$$

Из (2) и (3) получаем, что для того, чтобы $h_{\text{пр.}}(l_1:r_1) = h_{\text{пр.}}(l_2:r_2)$, надо, чтобы $k^{l_2-l_1} \cdot (h_{\text{пр.}}(0:r_1) - h_{\text{пр.}}(0:l_1 - 1)) = (h_{\text{пр.}}(0:r_2) - h_{\text{пр.}}(0:l_2 - 1))$

Применение хеширования строк

- ✓ сравнение двух подстрок на равенство: если выполнить предподсчёт всех префиксов, то сравнение на равенство выполняется за $O(1)$;
- ✓ лексикографически сравнить две строки на больше (меньше): подсчитаем хеши всех префиксов строки S , тогда можно за время $O(\log |S|)$ лексикографически сравнить две строки на больше (меньше): дихотомией ищут такой наибольший индекс i , что префиксы обеих строк длины i равны, после чего сравнивают $i+1$ символ соответствующих строк;
- ✓ за время $O(|T| + |S|)$ можно найти все вхождения образца S в текст T (алгоритм Рабина (Rabin)-Карпа (Karp), 1987 год): сначала для всех префиксов текста T вычислить хеши, вычислить хеш для образца S ; затем, двигаясь слева направо по тексту T , пробуем накладывать окно длины $|S|$.



Префиксная функция

Префиксная функция

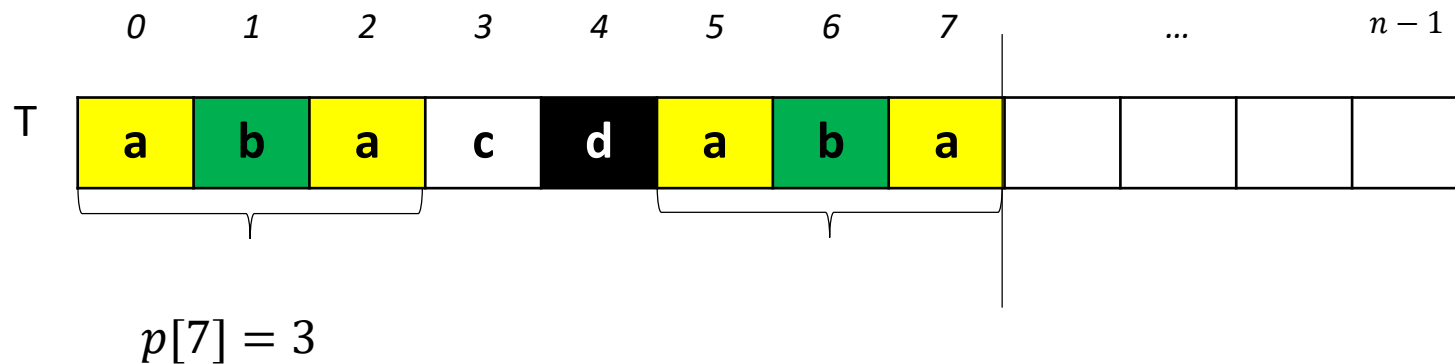
P от строки

$T = (t_0, t_1, \dots, t_{n-1})$

это массив целых чисел $P[0..n-1]$,
где $p[i]$ – длина максимального префикса подстроки $T[0..i]$,
который совпадает с суффиксом этой подстроки, но не равен
всей подстроке

$$\begin{cases} p[i] = 0 \\ p[i] = \max_{0 \leq k \leq i} \{k \mid T[0..k-1] = T[i-k+1..i]\}, i = 1, \dots, n-1 \end{cases}$$

другими словами - наибольшее число $k \leq i$, такое, что первые k символов подстроки $T[0..i]$ совпадают с ее последними k символами;



0 1 2 3 4 5 6

Вычислить префиксную функцию для всех позиций строки T :

a	b	a	c	a	b	a
---	---	---	---	---	---	---

i	0	0	1	0	1	2	0	1	2	3	0	1	2	3	4
T	a	a	b	a	b	a	a	b	a	c	a	b	a	c	a
P	0	0	0	0	0	1	0	0	1	0	0	0	1	0	1

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

Наивный алгоритм вычисления префиксной функции для всех позиций строки T (в основе которого лежит сравнение подстрок) работает за время $\Theta(|T|^3) = \Theta(n^3)$.

Эффективный алгоритм вычисления префиксной функции для всех позиций строки $T = (t_0, t_1, \dots, t_{n-1})$, который работает за линейное от длины строки время

(как основной элемент решения задачи поиска всех вхождений некоторого образца в текст)

был разработан **Д.Э. Кнуттом** и **В.Р. Праттом**,

независимо от них **Д.Х. Мориссом**.

Совместно тремя учёными алгоритм **Кнута, Морисса, Пратта** (**КМП**-алгоритм) был опубликован в **1977** году

Donald Knuth; James H. Morris, Jr, Vaughan Pratt. [Fast pattern matching in strings](#) (англ.) // [SIAM Journal on Computing](#) (англ.)[рус.](#) : journal. — 1977. — Vol. 6, no. 2. — P. 323—350.

Дональд Эрвин **Кнут**

*Donald Ervin
Knuth*



Дата рождения	1938
Место рождения	США ^[1]
Научная сфера	математика , программирование , информатика
Место работы	Стэнфордский университет
Известен как	автор классических трудов « Искусство программирования », « Конкретная математика » и др.

Вон Рональд **Пратт**

*Vaughan Ronald
Pratt*



Дата рождения	1944
Место рождения	Австралия
Научная сфера	математика , программирование , информатика
Место работы	Стэнфордский университет
Известен как	автор сертификата простоты Пратта и синтаксического анализатора Пратта

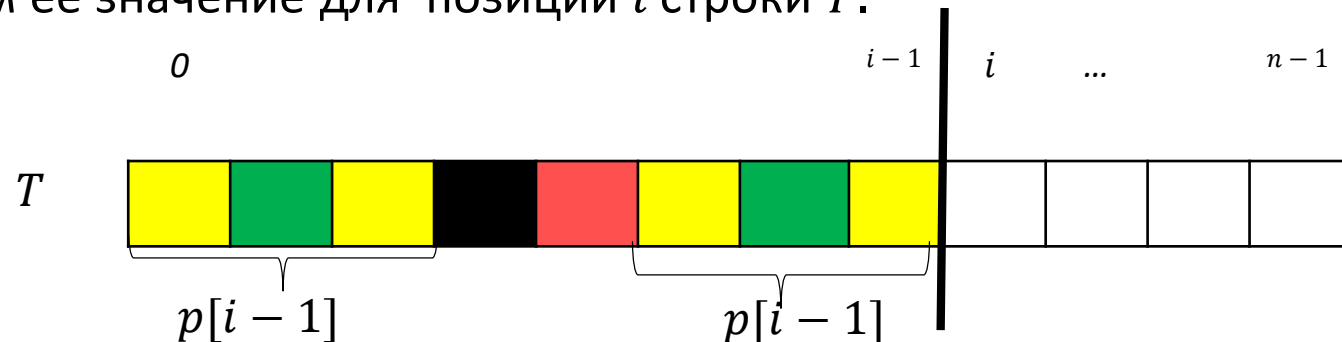
Джеймс Хирам **Моррис**

James Hiram Morris



Дата рождения	1941
Место рождения	США
Место работы	• Университет Карнеги — Меллона

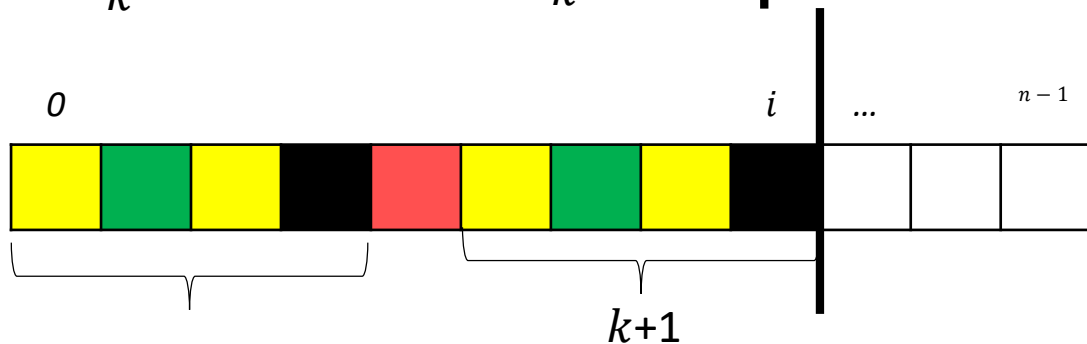
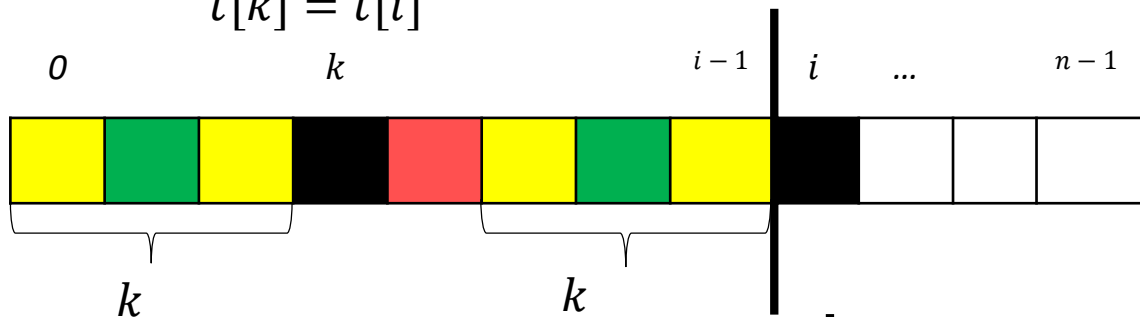
Предположим, что префиксная функция P вычислена для позиций $1, \dots, i - 1$ строки T .
 Вычислим её значение для позиции i строки T .



Случай 1.

$$k = p[i - 1]$$

$$t[k] = t[i]$$



$$k+1$$

$$p[i] = p[i - 1] + 1$$

Покажем, что большее значение для позиции i строки T получить нельзя.

Предположим, что можно получить большее значение префиксной функции в позиции i т.е.

$$p[i] > p[i - 1] + 1.$$

Удалим символ i из подстроки $T[0..i]$, получим, что у подстроки $T[0..i - 1]$ есть суффикс равный префиксу и его длина $p[i] - 1$. По сделанному предположению

$$p[i] - 1 > p[i - 1],$$

т.е. найдено лучшее значение префиксной функции для позиции $i - 1$ строки T – противоречие.

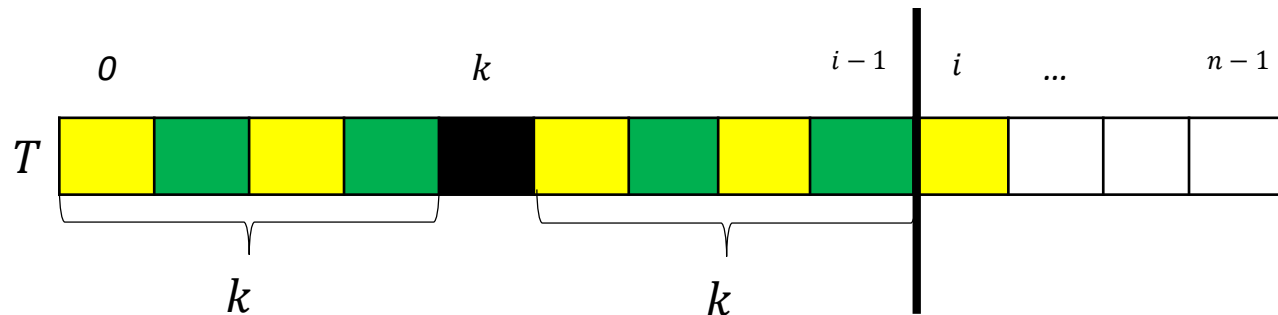
Утверждение 1.

Для каждой позиции i строки T значение префиксной функции может увеличиться максимум на 1 по сравнению с предыдущим значением.

Случай 2.

$$k = p[i - 1]$$

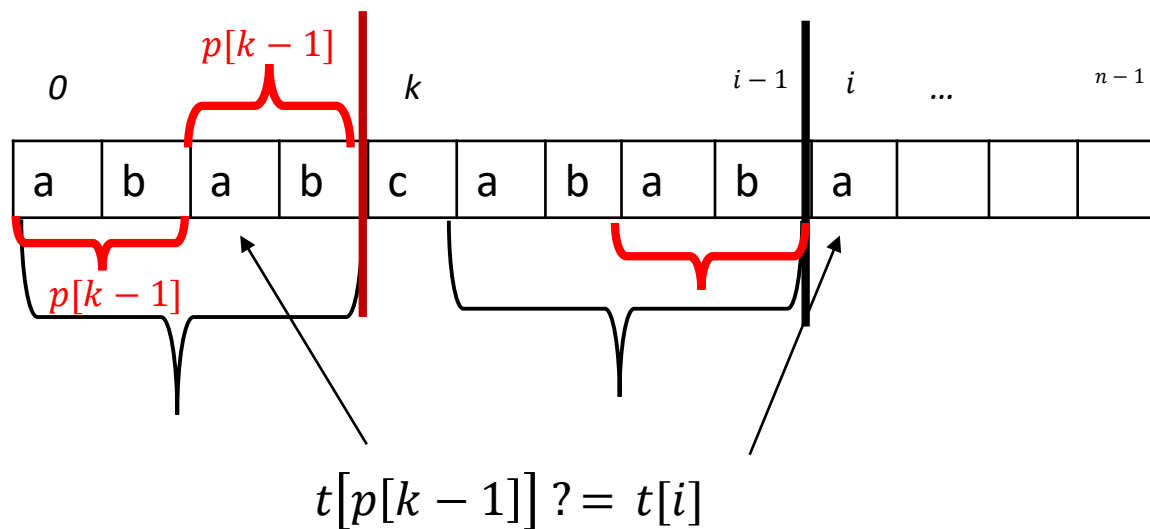
$$t[k] \neq t[i]$$

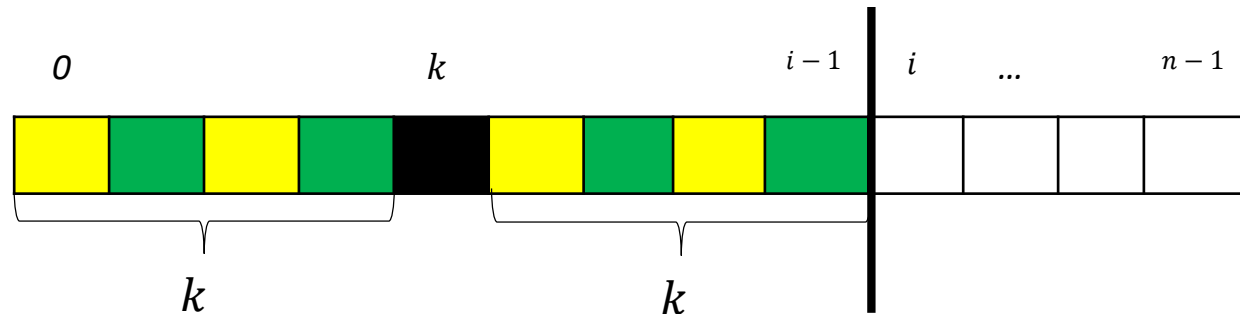


Если $p[i - 1] = 0$, то $p[i] = 0$.

Рассмотрим случай, когда $p[i - 1] > 0$:

значит для подстроки $T[0..i - 1]$ надо искать наибольший по длине префикс, который равен её суффиксу, но более короткий, чем k , а затем попробовать его продолжить





$p[0] := 0$

для всех $i = 1, \dots, n - 1$

$k := p[i - 1]$

пока $(k > 0)$ **и** $(t[k] \neq t[i])$

$k := p[k - 1]$

если $(t[k] = t[i])$, **то** $p[i] := k + 1$, **иначе** $p[i] := k$

Утверждение 2.

При построения префиксной функции P для всех позиций строки $T = (t_0, t_1, \dots, t_{n-1})$ суммарное число уменьшений значения префиксной функции P не превосходит суммарного числа её увеличений и равно $O(n)$.

Обозначим через $t_i \geq 0$ число, на которое уменьшилось значение префиксной функции для позиции i по сравнению с её значением для позиции $i - 1$.

На основании утверждения 1 значение префиксной функции для позиции i по сравнению с её значением для позиции $i - 1$ могло увеличиться максимум на 1. Поэтому справедливы следующие неравенства.

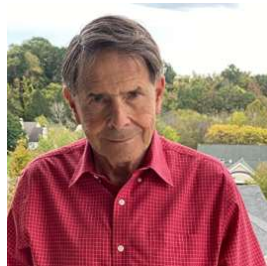
$$\begin{aligned} p[0] &= 0 \\ 0 \leq p[1] &\leq 0 + 1 - t_1 \\ 0 \leq p[2] &\leq (1 - t_1) + 1 - t_2 = 2 - (t_1 + t_2) \\ &\dots \\ 0 \leq p[n - 1] &\leq (n - 1) - (t_1 + t_2 + \dots + t_{n-1}) \end{aligned}$$

$$\sum_{i=1}^{n-1} t_i \leq n - 1$$

Утверждение 3.

Время работы алгоритма построения префиксной функции P для всех позиций строки $T = (t_0, t_1, \dots, t_{n-1})$ есть $O(n)$.

Алгоритм



Кнута-Морриса-Пратта
(сокр. **КМП**)

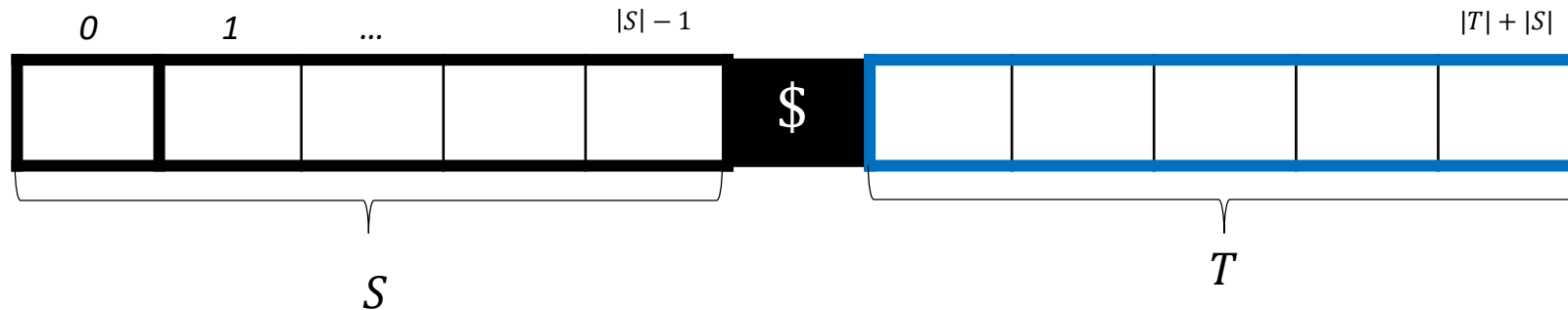
поиска подстроки в строке

Задан текст T и образец S .

Требуется найти все вхождения образца S в текст T .

Эффективный алгоритм **КМП** основан на построении префикс-функции P и работает за время $O(|S| + |T|)$.

- 1) Выполняем конкатенацию образца S и текста T , добавляя между ними символ, которого нет в обеих строках (в нашем примере это \$).

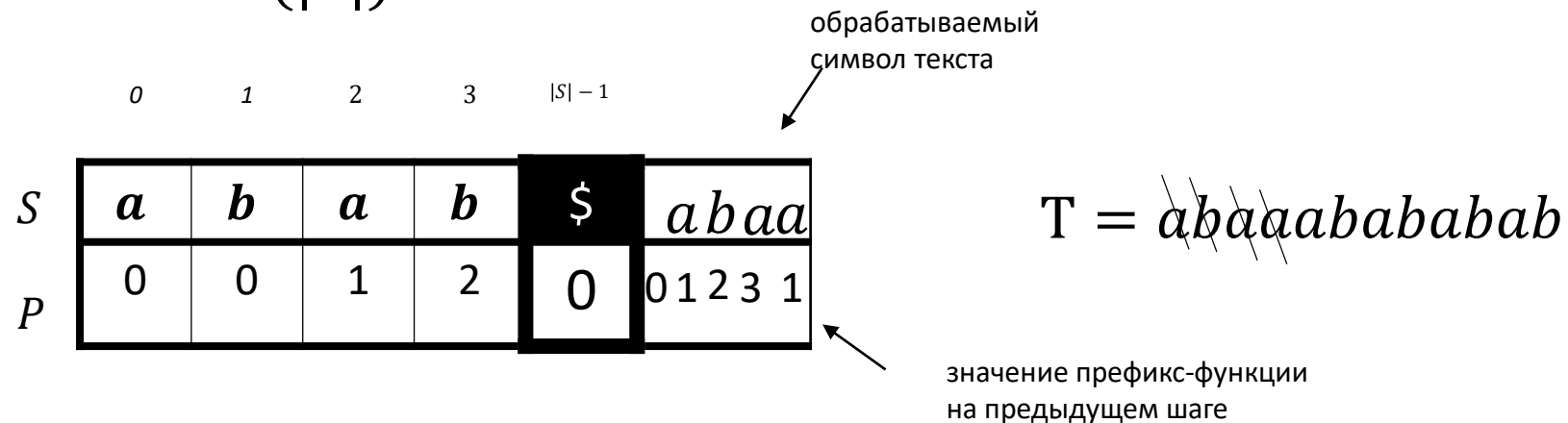


- 2) Вычисляем префиксную функцию P для всех позиций строки $S + \$ + T$.
- 3) Образец встречается в тексте столько раз, сколько раз префиксная функция принимала значение, равное длине образца.

	0	1	2	3		0	1	2	3	4	5	6	7	8	9	10	11	12	13	
S	Ф	П	М	И		T	Б	Ф	П	М	Ф	П	М	И	Г	Ф	П	М	И	У
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
$S + \$ + T$	Ф	П	М	И	\$	Б	Ф	П	М	Ф	П	М	И	Г	Ф	П	М	И	У	
P	0	0	0	0	0	0	1	2	3	1	2	3	4	0	1	2	3	4	0	

Если $p[i] = |S|$, то
образец S встречается в тексте T , начиная с индекса $j = i - 2 \cdot |S|$.

Память - $O(|S|)$.



Так как значение префикс-функции не превосходит $|S|$, то достаточно хранить только строку $S + \$$, значение префиксной функции для строки $S + \$$, последнее вычисленное значение префиксной функции и, считывая по одному очередной символ текста T , и вычислять префикс-функцию для этого символа.

Примеры задач, для которых можно разработать эффективный алгоритм, используя алгоритм построения префиксной функции.

1. Задана строка T , сколько различных подстрок в этой строке. Время работы алгоритма $O(|T|^2)$.
2. Задана строка T . Найти такую строку S наименьшей длины, что T можно представить в виде конкатенации строк S (например, для $T = \text{'аааааааа'}$ строка $S = \text{'а'}$) (т.е. хотим найти самое короткое сжатое представление строки). Время работы алгоритма $O(|T|)$.
3. Циклическим сдвигом строки T на k символов будем называть строку, которая получается из исходной строки переносом первых k символов в конец строки. Необходимо определить, можно ли получить из одной строки (T_1) другую (T_2) при помощи некоторого циклического сдвига. Время работы алгоритма $O(|T_1| +$

Z -функция строки T

(англ. *Z -function*)

z-функция (“зет-функция”) от строки $T = (t_0, t_1, \dots, t_{n-1})$, также, как и префиксная функция, используется в алгоритмах обработки строк

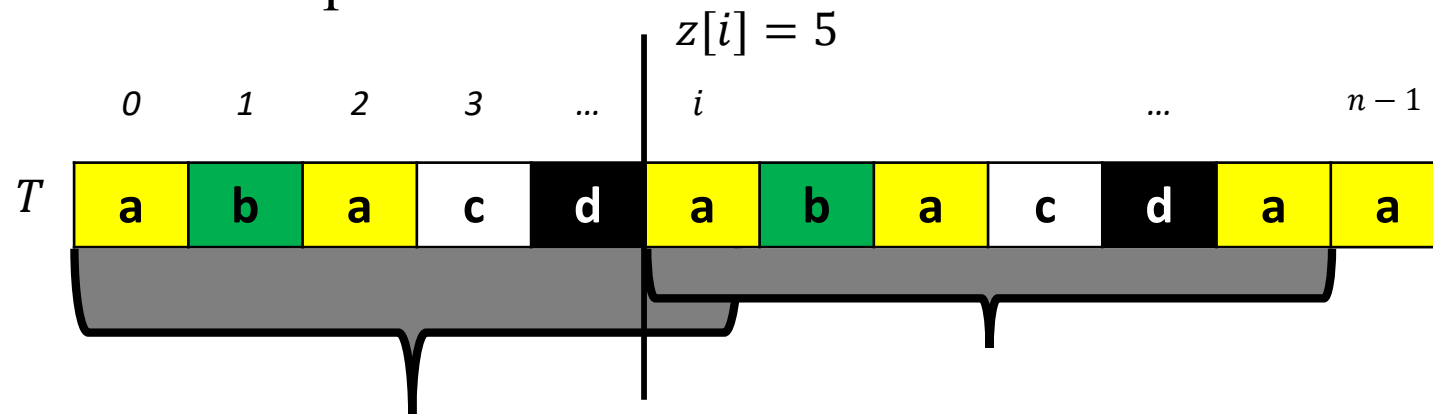
Z-функция

от строки

$$T = (t_0, t_1, \dots, t_{n-1})$$

это массив целых чисел $Z[0..n-1]$,

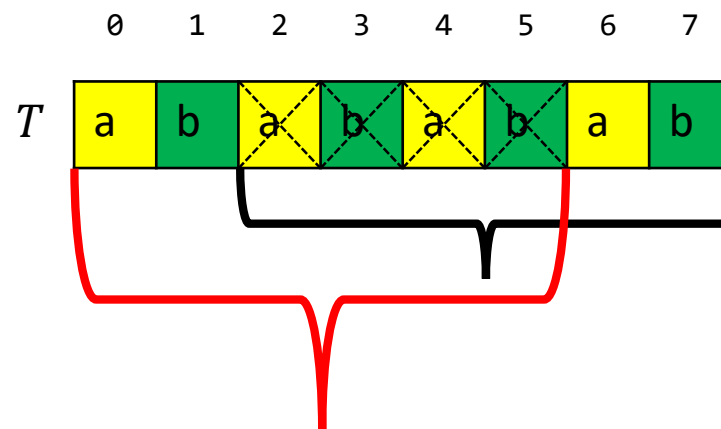
где $z[i]$ – длина наибольшего общего префикса i -ого суффикса строки T и самой строки T



$$\begin{cases} z[0] = 0 \\ z[i] = \max\{k \mid T[i..i+k-1] = T[0..k-1]\}, i = 1, \dots, n-1 \end{cases}$$

другими словами, $z[i]$ – наибольшее число k , такое, что первые k символов подстроки $T[i..n-1]$ совпадают с первыми k символами строки T ;

наибольший общий префикс строки T и подстроки $T[i..n - 1]$ могут пересекаться



$$z[2] = 6$$

i	0	1	2	3	4	5	6
T	a	b	a	c	a	b	a
Z	0	0					

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

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

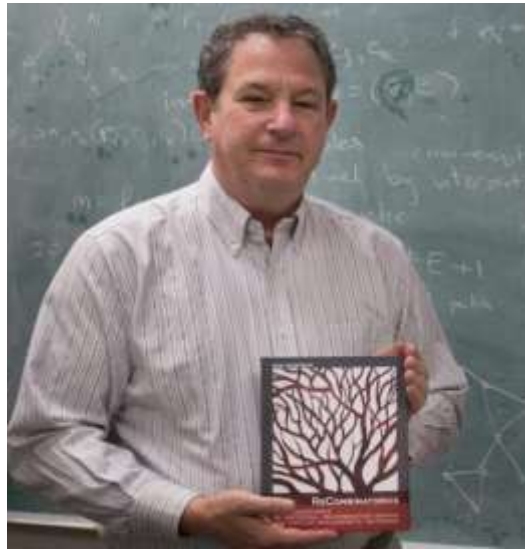
	0	1	2	3	4	5	6
	a	b	a	c	a	b	a
	0	0	1	0	3		

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

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

Наивный алгоритм (в основе которого лежит сравнение подстрок) работает за время $\Theta(|T|^2) = \Theta(n^2)$.

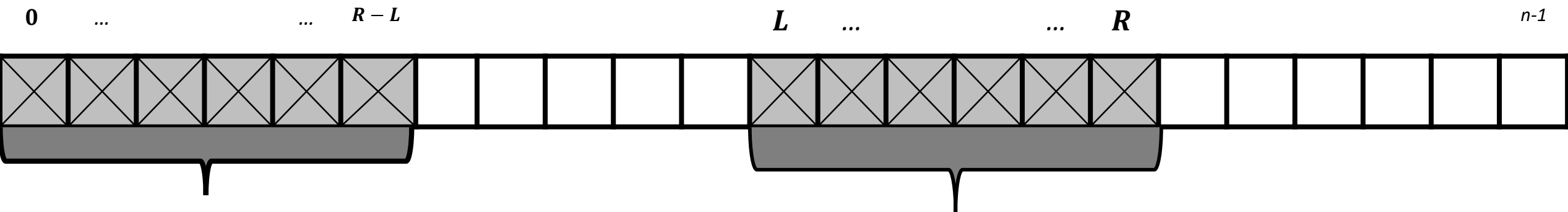
Эффективный алгоритм построения z-функции ($O(|T|) = O(n)$)
был изложен в **1997** году **Дэном Гасфилдом** в его книге
«Строки, деревья и последовательности в алгоритмах. Информатика и вычислительная биология»



Professor Daniel Gusfield
[Computer Science](#)
University of California

Пусть на некоторой итерации алгоритма L и R – индексы левой и правой границы самого «правого» отрезка (блока) совпадения (среди двух блоков с одинаковой правой границей будем брать больший по числу элементов).

Первоначально полагаем L и R равными 0 .



i	0	1	2	3	4	5	6	7
T	a	b	a	b	a	b	c	a
Z	0	0	4	0	2	0	0	1
L, R	0,0	0,0	2,5	2,5	2,5	2,5	2,5	7, 7

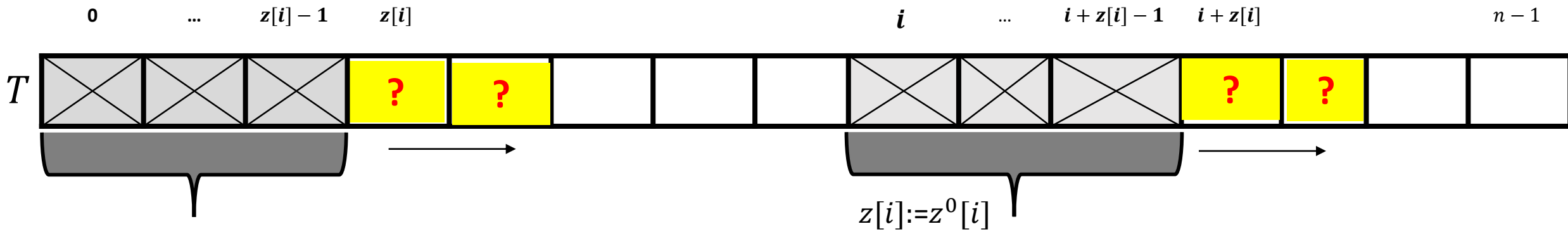
i	0	1	2	3
T	a	a	a	a
Z	0	3	2	1
L, R	0,0	1,3	1,3	1,3

Будем формировать массив Z последовательно от 0 до $n - 1$.

Предположим, что вычислены элементы $z[j], j < i$.

Вычислим $z[i]$.

Сначала попытаемся проинициализировать $z[i]$ некоторым ненулевым значением, используя ранее вычисленные значения: $z[i] := z^0[i]$.



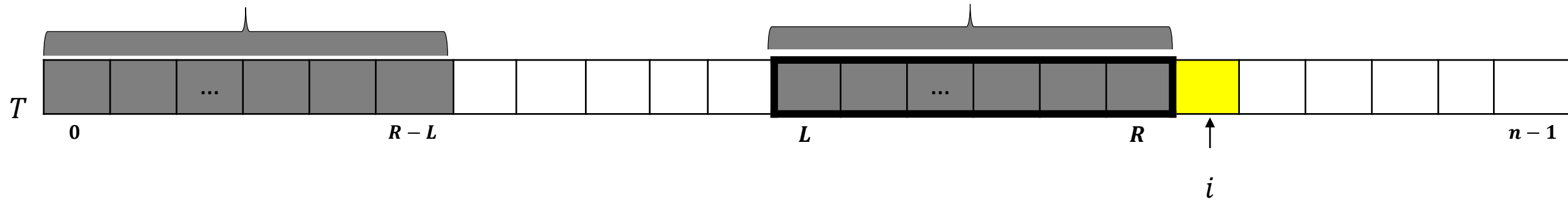
Затем, наивным алгоритмом пробуем увеличить значение $z[i]$, пока идёт совпадение:

$$t[z[i]] \quad ? = \quad t[i + z[i]]$$

$$t[z[i] + 1] \quad ? = \quad t[i + z[i] + 1]$$

...

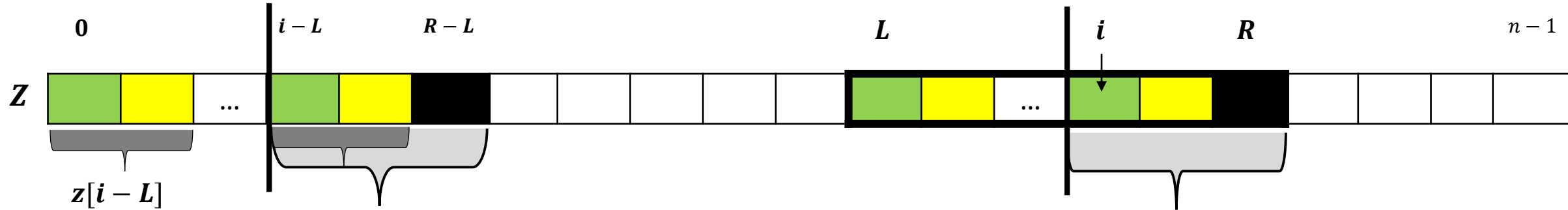
Случай 1: $i > R$ (текущий индекс i вышел за пределы блока).



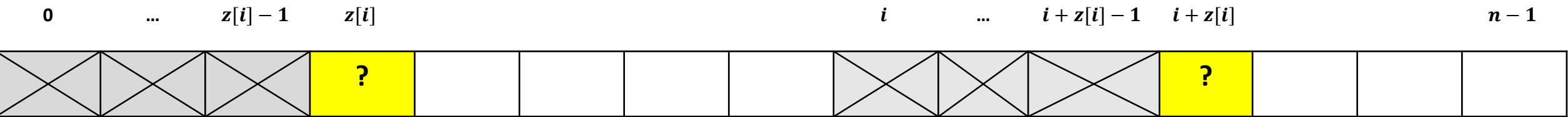
- 1) Проинициализируем $z[i] = z^0[i] = 0$.
- 2) Наивным алгоритмом, пока совпадение, сравниваем $t[0]? = t[i]$, $t[1]? = t[i+1] \dots$, увеличивая $z[i]$ на 1 при каждом совпадении (как только символы не совпадают, останавливаемся).
- 3) Если $z[i] > 0$, то корректируем границы самого правого блока совпадения:
$$L = i, R = i + z[i] - 1.$$

✓ отметим, что при каждом совпадении символов, величина R (правая границы самого правого блока совпадения) будет увеличиваться на 1

Случай 2: $i \leq R$ (текущий индекс i находится в блоке).



- 1) Проинициализируем $z[i] = z^0[i] = \min \{z[i - L], R - i + 1\}$.
- 2) Наивным алгоритмом, **пока совпадение**, сравниваем : $t[z[i]]? = t[i + z[i]], t[z[i] + 1]? = t[i + z[i] + 1] \dots$ увеличивая каждый раз $z[i]$ на 1 при совпадении.



- ✓ если при инициализации $z^0[i] < R - i + 1$, то это начальное значение будет итоговым для $z[i]$;
- ✓ если $z^0[i] = R - i + 1$, то правая граница при каждом последующем успешном сравнении соответствующих символов будет увеличиваться на 1 и число увеличений равно количеству выполненных операций сравнения символов;

для всех $i = 1, \dots, n - 1$

$z[i] := 0$

$L := 0, R := 0$

для всех $i = 1, \dots, n - 1$

если $i \leq R$

$z[i] = \min(R - i + 1, z[i - L])$

пока $(i + z[i] < n)$ и $(t[z[i]] = t[i + z[i]])$

$z[i] := z[i] + 1$

если $(i + z[i] - 1 > R)$

$L := i$

$R := i + z[i] - 1$

При рассмотрении случаев 1 и 2 было показано, что каждый раз, когда происходит **сравнение символов и они оказываются равны** (внутренний цикл пока), значение правой границы **R** увеличивается на 1.

Значение **R** в ходе работы алгоритма не может уменьшиться.

Наибольшее значение, которое может достигнуть **R** , равно длине строки **T** . Следовательно, суммарное число итераций внутреннего цикла **$O(|T|)$** .

Все остальные операции алгоритма выполняются суммарно за время **$O(|T|)$** .

для всех $i = 1, \dots, n - 1$

$z[i] := 0$

$L := 0, R := 0$

для всех $i = 1, \dots, n - 1$

если $i \leq R$

$z[i] = \min(R - i + 1, z[i - L])$

пока $(i + z[i] < n)$ и $(t[z[i]] = t[i + z[i]])$

$z[i] := z[i] + 1$

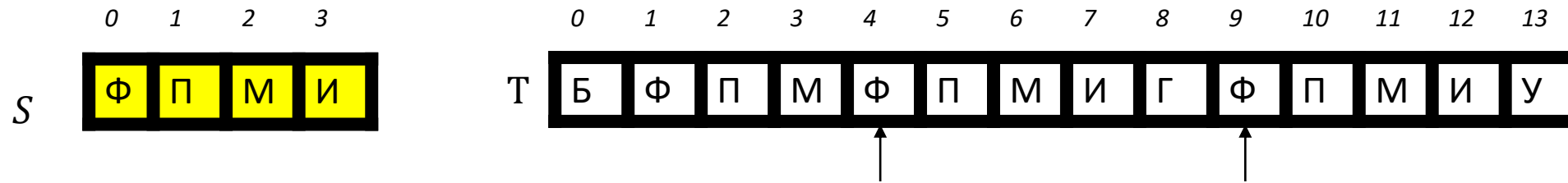
если $(i + z[i] - 1 > R)$

$L := i$

$R := i + z[i] - 1$

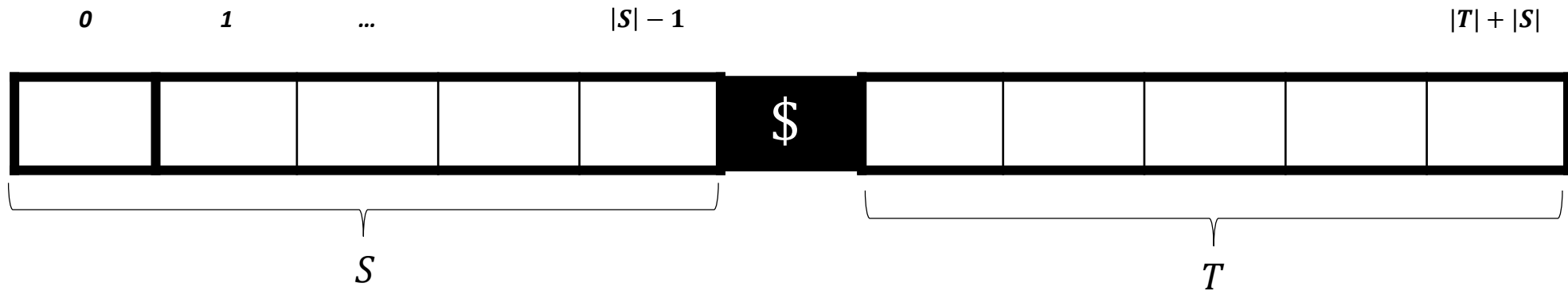
Следовательно время работы алгоритма построения z -функции
 $O(|T|) = O(n)$.

Задан текст T и образец S . Найти все вхождения образца S в текст T .

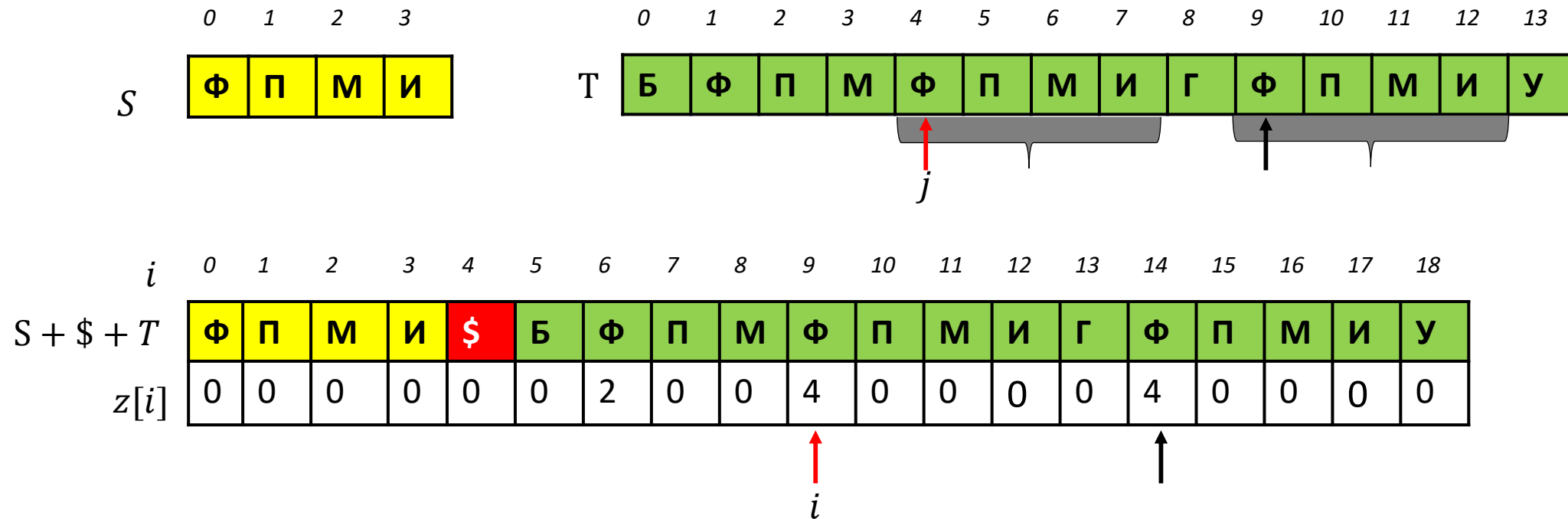


Эффективный алгоритм, основанный на построении z-функции, работает за время $O(|S| + |T|)$.

- 1) Выполняем конкатенацию образца S и текста T , добавляя между ними символ, которого нет в обеих строках (в нашем примере это \$).



- 2) Вычисляем z -функцию для $S + \$ + T$.
- 3) Образец S встречается в тексте T столько раз, сколько раз z -функция принимала значение, равное длине образца.
- 4) Если $z[i] = |S|$, то в тексте T , начиная с позиции $j = i - |S| - 1$, найдено вхождение образца S .



если $z[i] = |S|$, то
образец S встречается в тексте T , начиная с индекса $j = i - |S| - 1$

Задача поиска всех подстрок, которые являются палиндромами

Техника поддержки самого правого блока совпадения, которая использовалась для построения эффективного алгоритма построения z -функции, также используется в задаче **поиска всех подстрок заданной строки T , которые являются палиндромами.**

Эффективный алгоритм решения этой задачи, работающий за время $O(|T|)$, был предложен **Гленном Манакером** (Glenn Manacher) в **1975** г.

Glenn Manacher



Assoc. Professor Emeritus
Mathematics, Statistics, and Computer
Science
Department of Information Engineering,
University of Illinois at Chicago

Простейший же тривиальный алгоритм, который проверяет все возможные подстроки работает за время $O(|T|^2)$.

Для решения задачи нужно сформировать два массива D_1 и D_2 , где

$d_1[i]$ – число палиндромов нечётной длины с центром в i ;

$d_2[i]$ – число палиндромов чётной длины с центром в i
(мнимый центр в i , а подстрока, начиная с индекса i ,
сдвинута на 1 позицию вправо);

	0	1	2	3	4	5	6	7
T	a	b	a	a	a	b	a	c
D_1				4				

abaaaba
baaab
aaa
a

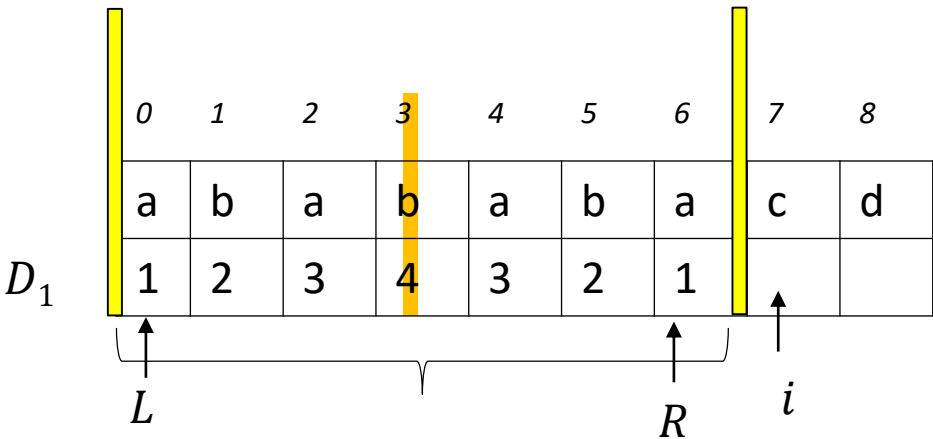
	0	1	2	3	4	5	6	7
T	a	b	a	a	a	b	a	c
D_2				1				

Если для каждой позиции i подсчитать **длину наибольшего палиндрома** чётной и нечётной длины, то легко вычислить значения соответствующих массивов $d_1[i]$ и $d_2[i]$, так как, если строка T является палиндромом, то отбрасывание первого и последнего символов этой строки T также приводит к палиндрому с тем же центром.

Будем поддерживать **левую L** и **правую R** границу самого правого из найденных палиндромов.

Предположим, что вычислены элементы массива $d_1[j], j < i$ и нужно вычислить $d_1[i]$, используя ранее вычисленные значения (первоначально $d_1[j] = 1$ для всех j):

Случай 1 (текущий индекс i вышел за пределы блока) $i > R$



Пока не вышли за пределы строки и есть совпадения, применяем наивный алгоритм, сравнивая

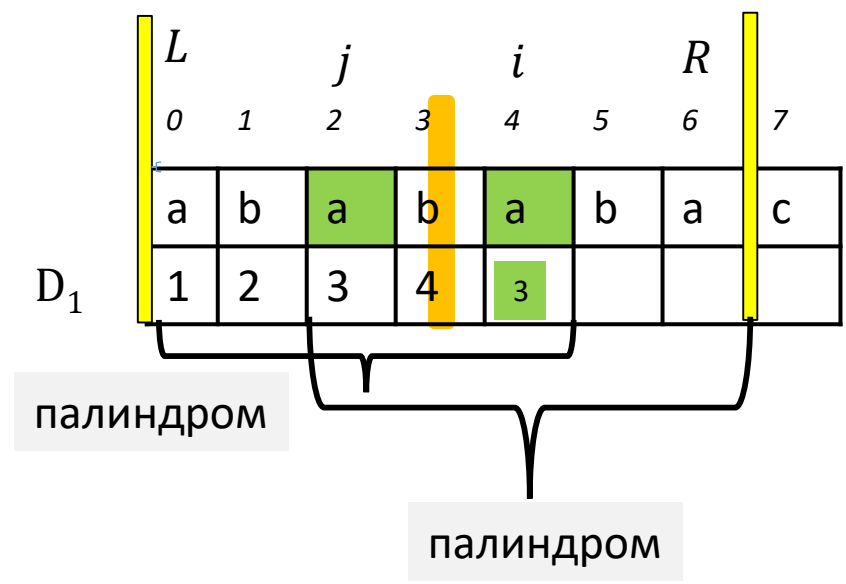
$$t[i + d_1[i]]? = t[i - d_1[i]],$$

увеличиваем значение $d_1[i]$ в случае совпадения.

Корректируем границы самого правого блока.

Случай 2 (текущий индекс i находится в блоке)

$i \leq R$, в этом случае можно проинициализировать $d_1[i]$ через ранее вычисленные значения



при формировании массива d_1
для i симметричный индекс
 $j = (R - i) + L$

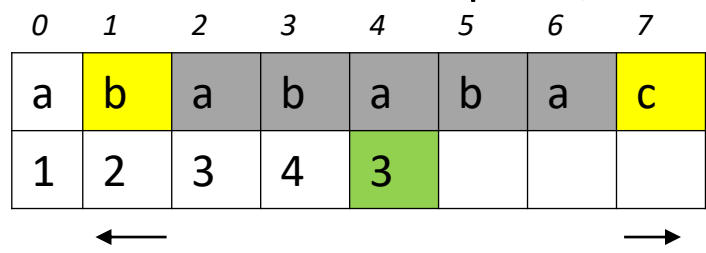
при формировании массива d_2
для i симметричный индекс
 $j = (R - i) + L + 1$

Начальная инициализация $d_1[i] := \min\{d_1[j], R - i + 1\}$, где $j = (R - i) + L$.

Далее, пока не вышли за пределы строки и есть совпадения, применяем наивный алгоритм, сравнивая

$$t[i + d_1[i]] ?= t[i - d_1[i]]$$

и увеличивая $d_1[i]$ в случае совпадения.



При необходимости корректируем границы самого правого блока.

D_1

```

 $L := 0, R := -1$ 
для всех  $i = 0, \dots, n - 1$ 
    если  $i > R$ 
         $k := 1$ 
    иначе  $l := \min(d_1[l + r - i], r - i + 1)$ 
    пока  $(i + k < n)$  и  $(i - k \geq 0)$  и  $(t[i + k] = t[i - k])$ 
         $k := k + 1$ 
     $d_1[i] := k$ 
    если  $(i + k - 1 > R)$ 
         $L := i - k + 1$ 
         $R := i + k - 1$ 

```

 D_2

```

 $L := 0, R := -1$ 
для всех  $i = 0, \dots, n - 1$ 
    если  $i > R$ 
         $k := 1$ 
    иначе  $l := \min(d_2[l + r - i + 1], r - i + 1)$ 
    пока  $(i + k < n)$  и  $(i - k - 1 \geq 0)$  и  $(t[i + k] = t[i - k - 1])$ 
         $k := k + 1$ 
     $d_2[i] := k$ 
    если  $(i + k - 1 > R)$ 
         $L := i - k$ 
         $R := i + k - 1$ 

```

При рассмотрении случаев 1 и 2 было показано, что каждый раз, когда происходит сравнение символов и они оказываются равны, значение правой границы R увеличивается на 1.

Значение R в ходе работы алгоритма не может уменьшиться.

Наибольшее значение, которое может достигнуть R , равно длине строки T .

Следовательно, суммарное число итераций внутреннего цикла **пока** — $O(|T|)$.

Все остальные операции алгоритма выполняться суммарно за время $O(|T|)$.

Следовательно, время работы алгоритма поиска всех подстрок заданной строки T , которые являются палиндромами, есть $O(|T|)$.



Спасибо за внимание!