

ИВАНОВ А.Б.

# АЛГОРИТМЫ И СТРУКТУРЫ ДАННЫХ

СТРОКИ

# ПОИСК ПОДСТРОК

Дано

- алфавит  $\Sigma$
- образец  $P[1..m] \in \Sigma^*$
- строка  $T[1..n] \in \Sigma^*$

Требуется найти  $0 \leq s \leq n - m$  такое, что  $T[s + i] = P[i], \forall i = 1..m$  (возможно все такие  $s$ ).

# НАИВНЫЙ АЛГОРИТМ

---

```
function find(string, pattern)
  m = len(pattern)
  n = len(string)

  for s = 0, 2, ...n-m
    i = 1
    while (string[s + i] == pattern[i] && i <= m) i++

    if (i == m + 1) return s
  end
end
```

---

СЛОЖНОСТЬ  $O(nm)$

## АЛГОРИТМ РАБИНА-КАРПА

Будем рассматривать строки как числа в  $|\Sigma|$ -чной системе исчисления. Обозначим  $d = |\Sigma|$ . Тогда строка  $P$  вычисляется за  $O(m)$  как

$$p = P[m] + d(P[m-1] + d(P[m-2] + \dots + dP[1]) \dots).$$

Строку  $T[s+1 \dots s+m]$  можно вычислить так (для всех  $s$  за  $O(n)$ ):

$$t_{s+1} = d(t_s - d^{m-1}T[s+1]) + T[s+1+m]$$

Рассматриваем все эти числа по модулю  $q$ . Тогда

$$t_{s+1} = d(t_s - d^{m-1}T[s+1]) + T[s+1+m] \pmod{q}$$

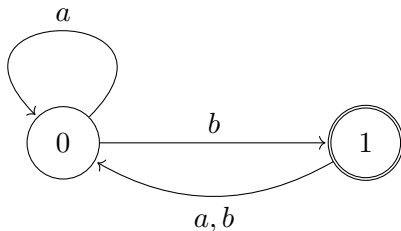
(вычислим  $d^{m-1} \pmod{q}$  заранее)

Сложность  $O(nm)$ , в среднем  $O(n) + O(\frac{mn}{q})$

# КОНЕЧНЫЙ АВТОМАТ

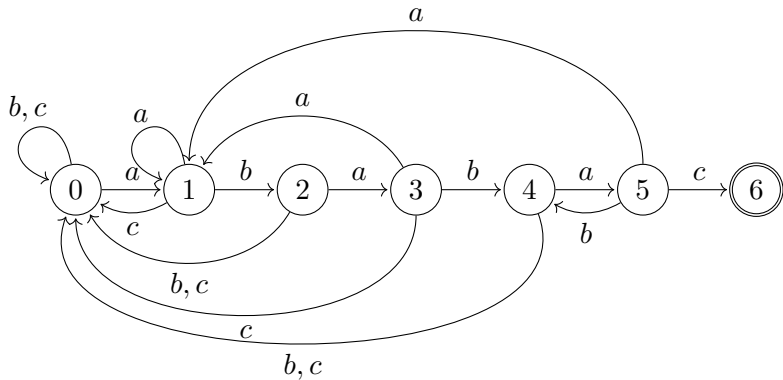
Конечный автомат -  $(Q, q_0, A, \Sigma, \delta)$ :

- $Q$  - множество *состояний*
- $q_0 \in Q$  - *начальное состояние*
- $A \subseteq Q$  - множество *допустимых состояний*
- $\Sigma$  - *алфавит*
- $\delta : Q \times \Sigma \rightarrow Q$  - *функция перехода*



# АВТОМАТ ПОИСКА ПОДСТРОКИ

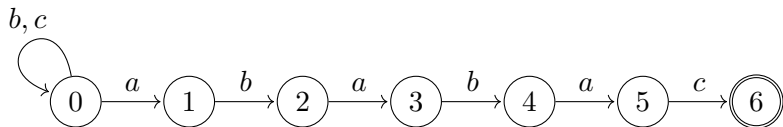
Автомат поиска *ababac*:



Сложность поиска подстроки автоматом -  $O(n)$

# ПОСТРОЕНИЕ АВТОМАТА

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



Сложность построения автомата -  $O(|\Sigma|m)$

# ПОДЫТОГ

Алгоритм:

- Строим автомат за  $O(|\Sigma|m)$
- Поиск подстроки за  $O(n)$

Проблемы (если алфавит  $\Sigma$  большой):

- Автомат занимает память  $O(|\Sigma|m)$
- Общая сложность  $O(|\Sigma|m + n)$

Можно лучше (Алгоритм Кнута-Морриса-Пратта):

- Память  $O(m)$
- Общая сложность  $O(m + n)$

# ПРЕФИКС-ФУНКЦИЯ

Префикс-функция строки  $P$ :

$$\pi(i, P) = \{\text{длина наибольшего префикса } P, \\ \text{являющегося и суффиксом } P[1..i]\}$$

Построение за  $O(m)$ :

---

```
function build_prefix_function(string)
  m = len(string); pi = int[m]
  pi[1] = 0; k = 0
  for i = 0, ..., m
    while (k > 0 && string[k+1] <> string[i]) k = pi[k]
    if (string[k+1] = string[i]) k = k + 1
    pi[i] = k
  end
  return pi
end
```

---

# АЛГОРИТМ КНУТА-МОРРИСА-ПРАТТА (КРМ)

Поиск за  $O(n)$ :

---

```
function match(string, pattern)
  m = len(pattern)
  n = len(string)
  pi = build_prefix_function(pattern)
  q = 0
  for i = 1, ..., n
    while (q > 0 && string[q+1] <> pattern[i]) q = pi[q]
    if (pattern[q + 1] == string[i]) q = q + 1
    if (q == m) return i - m
  end
end
```

---