

Hledání řetězců

Karel Richta a kol.

Katedra počítačů
Fakulta elektrotechnická
České vysoké učení technické v Praze

© Karel Richta a kol., 2018

Datové struktury a algoritmy, B6B36DSA
01/2018, Lekce 14

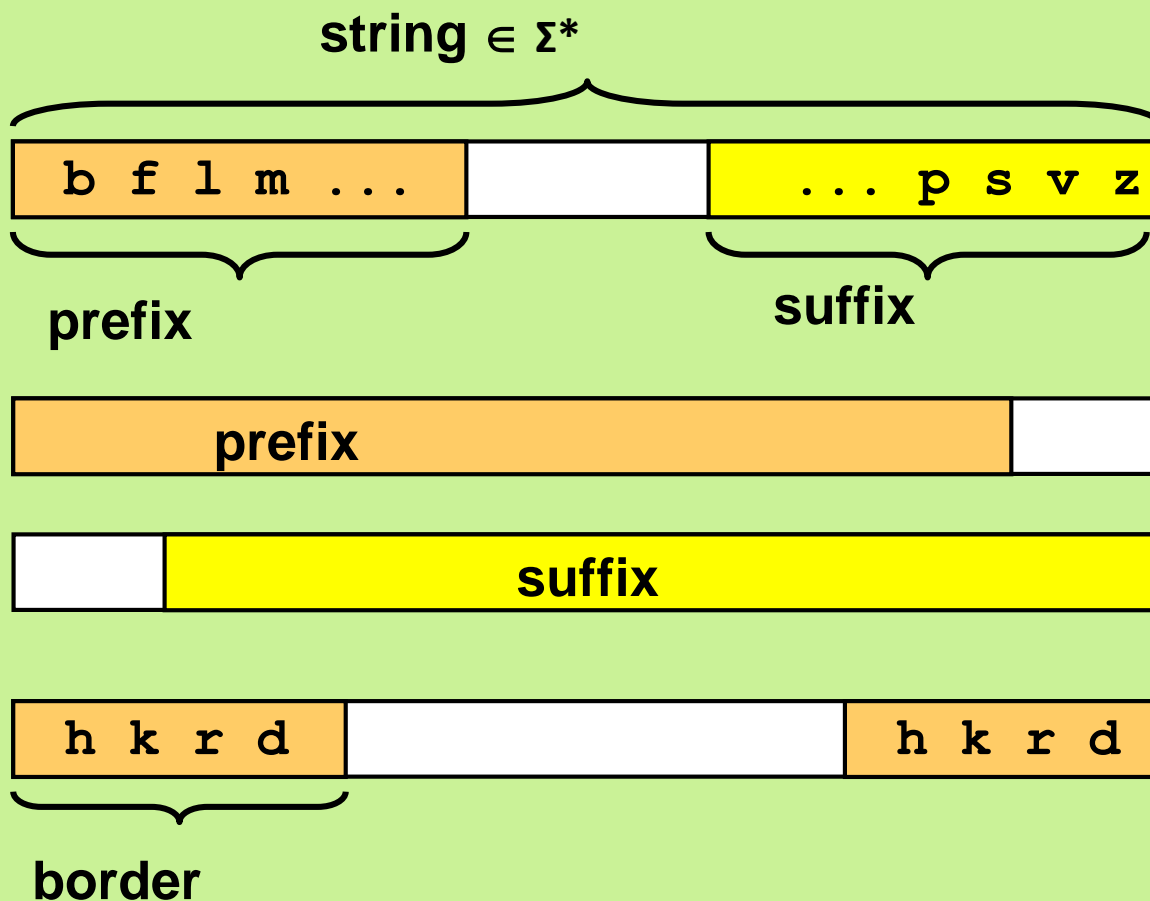
<https://cw.fel.cvut.cz/wiki/courses/b6b36dsa/start>



Evropský sociální fond
Praha & EU: Investujeme do vaší budoucnosti

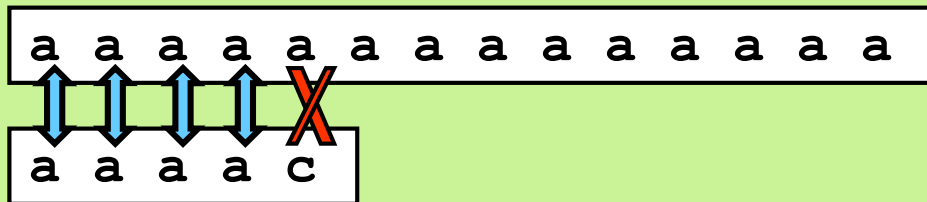
String matching (hledání řetězce v textu)

Σ – alphabet , e.g. $\Sigma = \{a, b, c, d \dots\}$

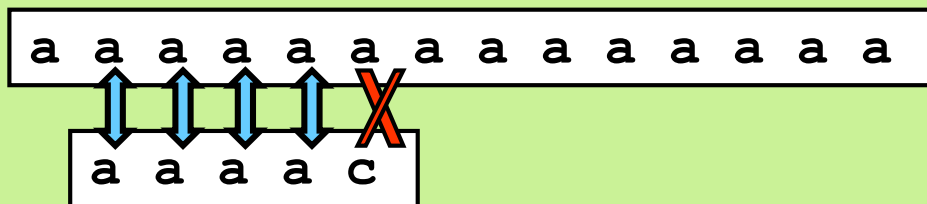


Naive algorithm

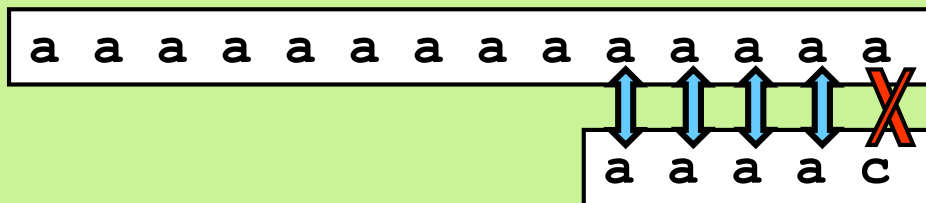
worst case



↕ match
 X mismatch



• • •



Naive algorithm

best case

a a a a a a a a a a a a a a a
X
c a a a a

a a a a a a a a a a a a a a a
X
c a a a a

• • •

a a a a a a a a a a a a a a a
X
c a a a a

↕ match
X mismatch

Karp-Rabin Algorithm

$$t_5: \quad \boxed{5 \ 3 \ 2} \quad \phi(t_5) = 532$$

$$t_6: \quad \boxed{3 \ 2 \ 3} \quad \phi(t_6) = 323$$

$$T: \quad \boxed{1 \ 2 \ 3 \ 4 \ 5 \ 3 \ 2 \ 3 \ 5 \ 6 \ 3 \ 1}$$

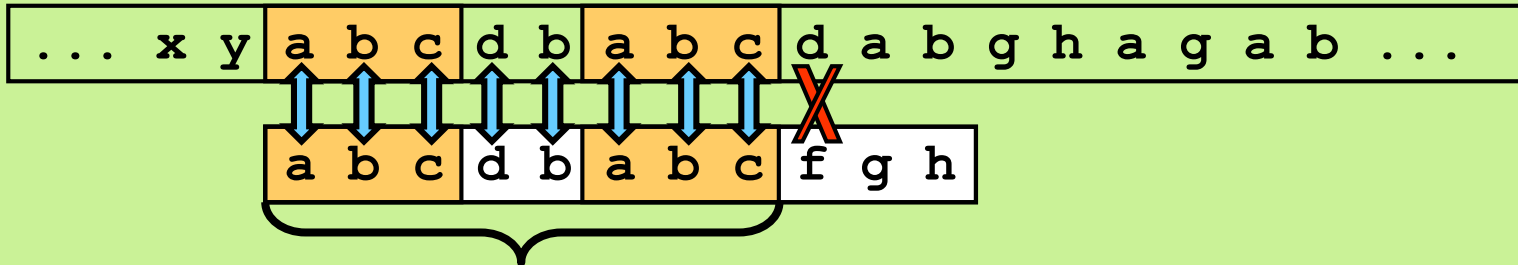
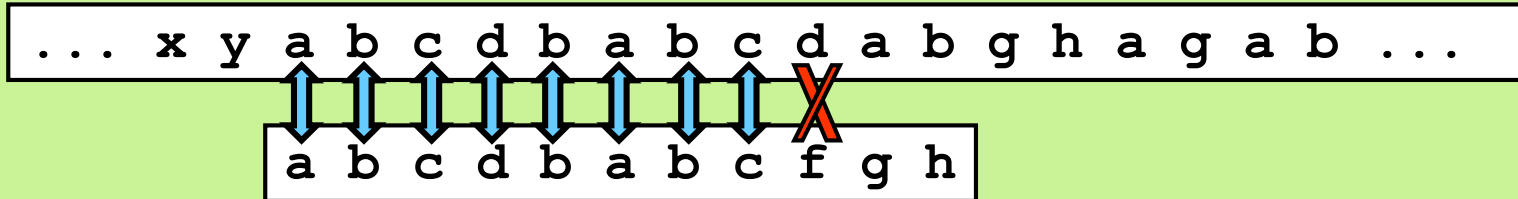
1 2 3 4 5 6 7 8 9 10 11 12

$$\begin{aligned} \phi(t_5) &= t_7 + 10(t_6 + 10(t_5)) = 2 + 10(3 + 10(5)) = 2 + 10(3 + 50) = \\ &= 2 + 10 \cdot 53 = 532 \end{aligned}$$

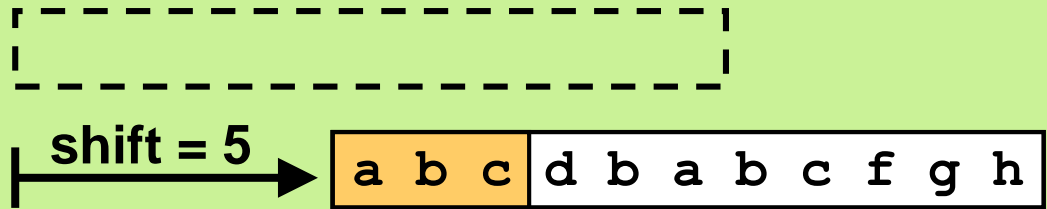
$$\begin{aligned} \phi(t_6) &= 10(\phi(t_5) - 10^2 \cdot t_5) + t_{5+3} = 10(532 - 100 \cdot 5) + 3 = \\ &= 10(532 - 500) + 3 = 10 \cdot 32 + 3 = 323 \end{aligned}$$

$$\begin{aligned} \phi(t_i) &= t_{i+m-1} + 10(t_{i+m-2} + 10(\dots + 10t_i) \dots) \\ \phi(t_{i+1}) &= 10(\phi(t_i) - 10^{m-1}t_i) + t_{i+m} \end{aligned}$$

Knuth–Morris–Pratt algorithm



longest border (=prefix&suffix)



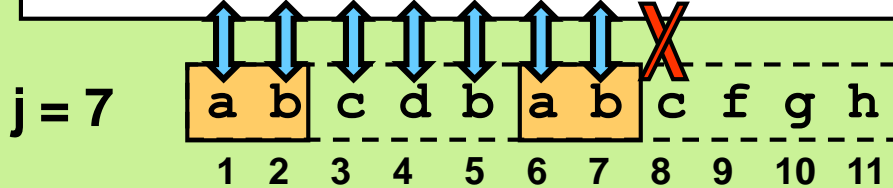
Knuth - Morris - Pratt algorithm

j	1	2	3	4	5	6	7	8	9	10	11
p_j	a	b	c	d	b	a	b	c	f	g	h
$\phi'(p_j)$	0	0	0	0	0	1	2	3	0	0	0

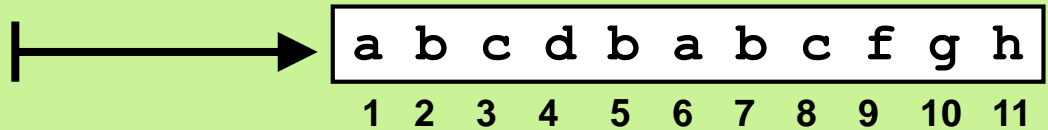
Longest border size
in substr (p_1, \dots, p_j)

shift after mismatch t_i $\times p_{j+1} = j - \phi'(p_j)$

... w a b c d b a b w d a b g h a g a b ...

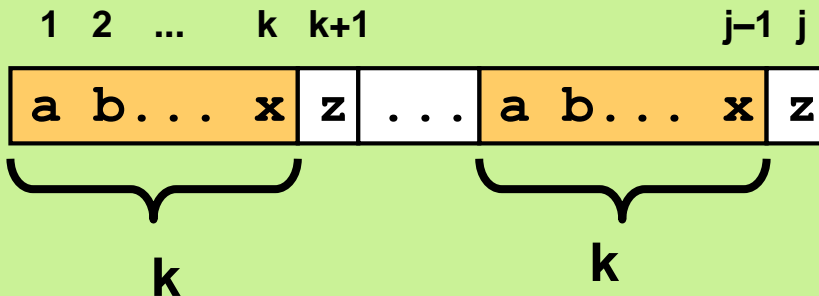


shift = $7 - \phi'(p_7) = 7 - 2 = 5$



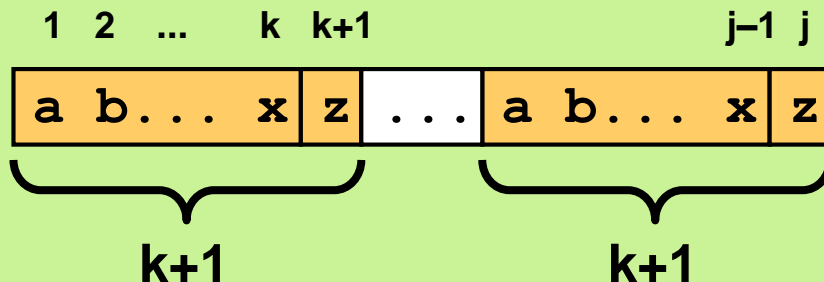
Knuth - Morris - Pratt algorithm

computing $\phi'(p_j)$ from $\phi'(p_{j-1})$



$k = \phi'(j-1) = \text{size of longest border in } (p_1, p_2, \dots, p_{j-1})$

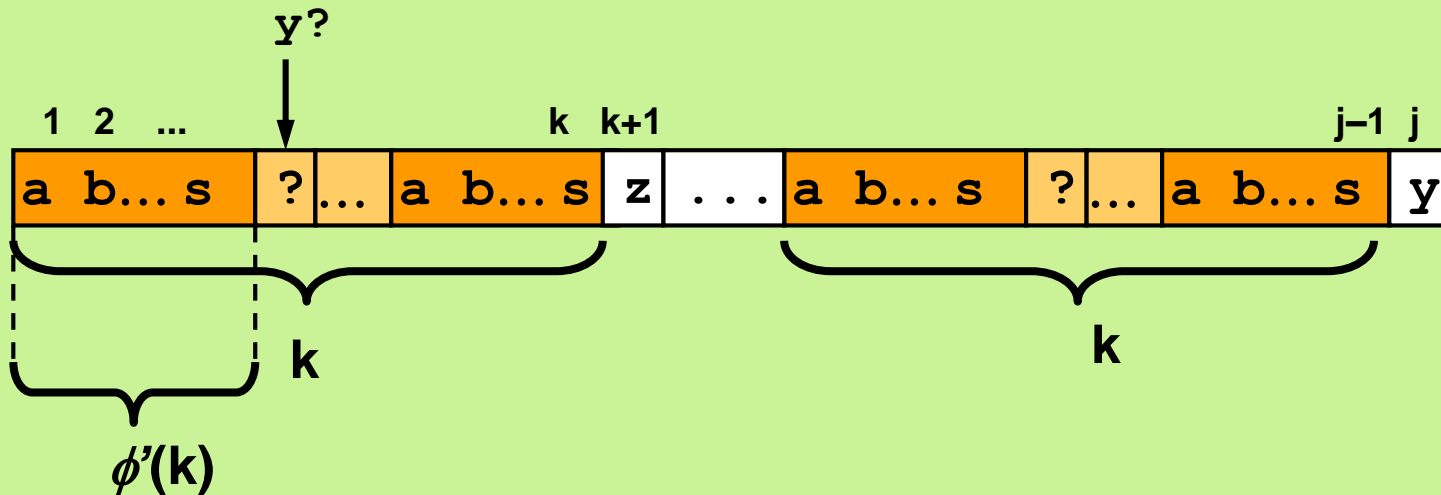
if $p_{k+1} = p_j$ then $\phi'(p_j) = k+1$



Knuth - Morris - Pratt algorithm

computing $\phi'(p_j)$ from $\phi'(p_{j-1})$

if $p_{k+1} \neq p_j$ then longest border of (p_1, \dots, p_j)
might be the longest border of $(p_1, \dots, p_k) \dots$



... and if it is not then continue search for the longest border
in $(p_1, \dots, p_{\phi'(k)})$ recursively

Knuth - Morris - Pratt algorithm

```
Compute shift function  $\phi'$ 
1:  $\phi'(1) \leftarrow 0$ 
2:  $k \leftarrow 0$ 
3: for  $j \leftarrow 2..length(P)$  do
4:   while  $k > 0$  and  $p_{k+1} \neq p_j$  do
5:      $k \leftarrow \phi'[k]$ 
6:   end while
7:   if  $p_{k+1} = p_j$  then  $k \leftarrow k+1$  end if
8:    $\phi'(j) \leftarrow k$ 
9: end for
```

Boyer - Moore algorithm

a b c d b a b c f g h

p
BCS[p]

a	b	c	d	f	g	h	other
5	4	3	7	2	1	11	11

T: ... x y a b c d b a b c d a c g h a g a b ...

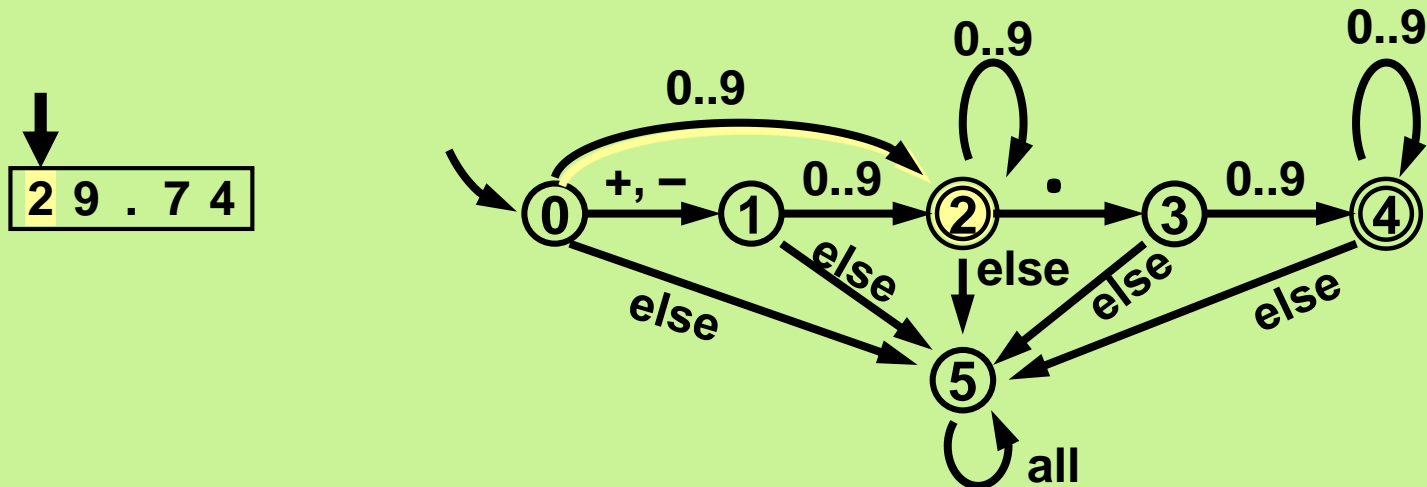
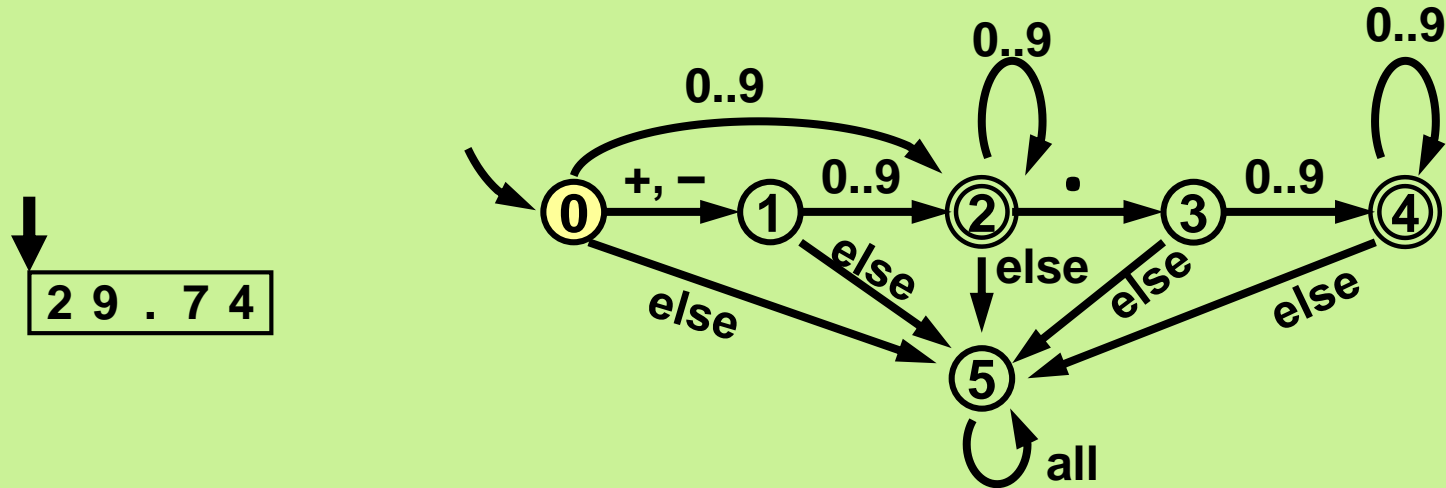
P: a b c d b a b c f g h BCS[b] = 4

4 → a b c d b a b c f g h BCS[c] = 3

3 → a b c d b a b c f g h BCS[a] = 5

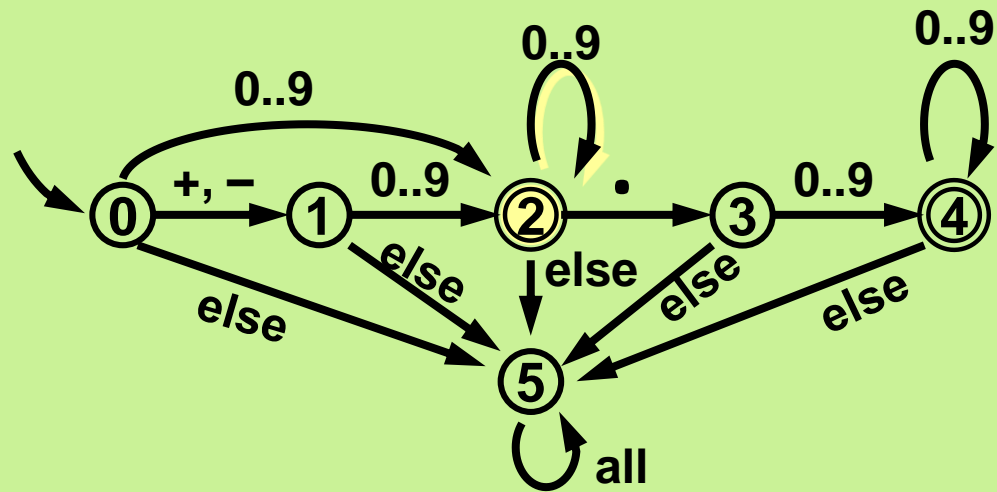
5 → a b c d b a b c f g h

Konečný automat (finite-state automaton)

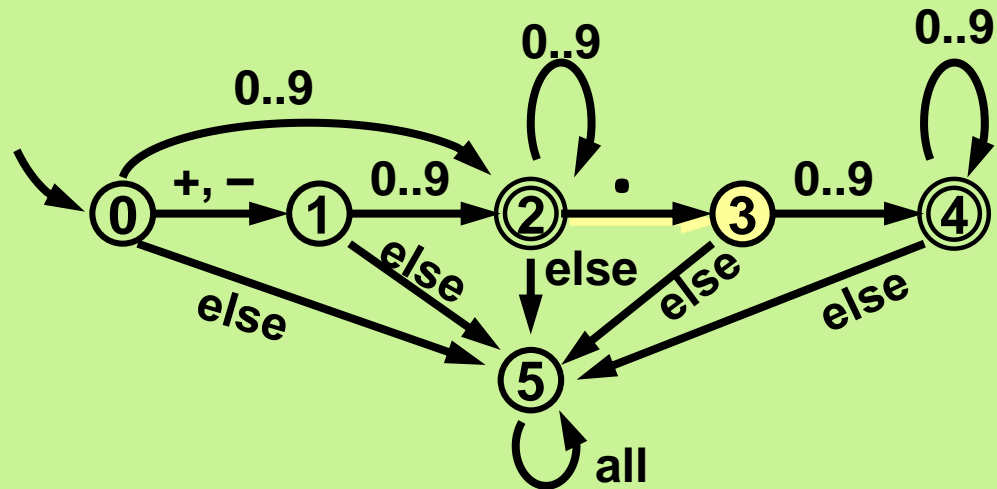


Konečný automat

↓
2 9 . 7 4

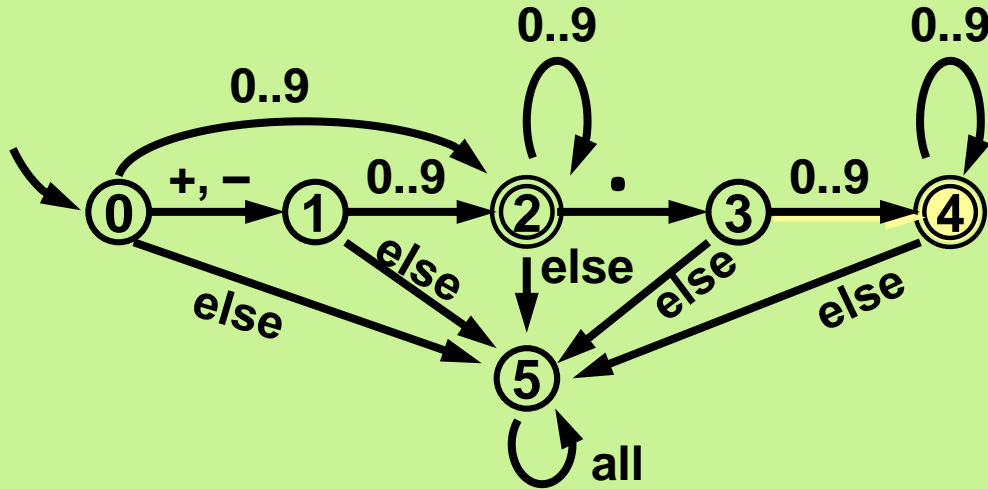


↓
2 9 . 7 4

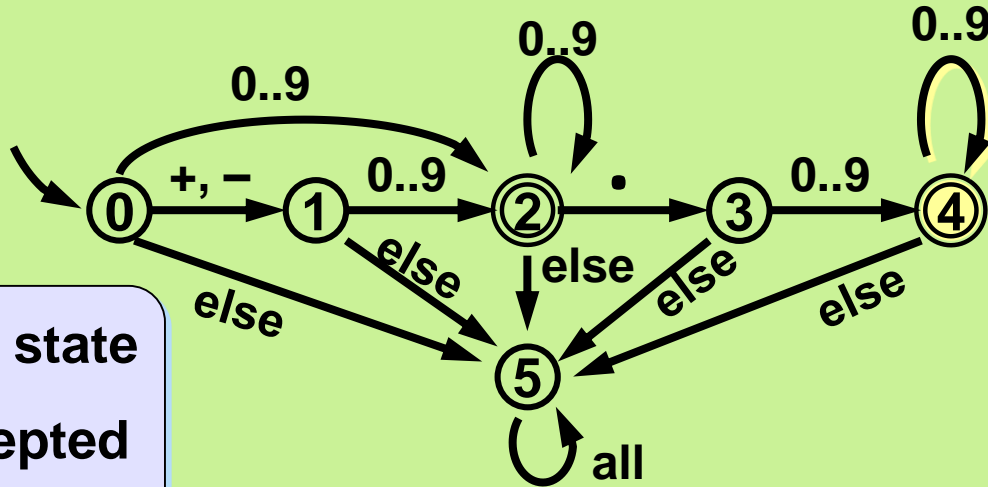


Konečný automat

↓
2 9 . 7 4

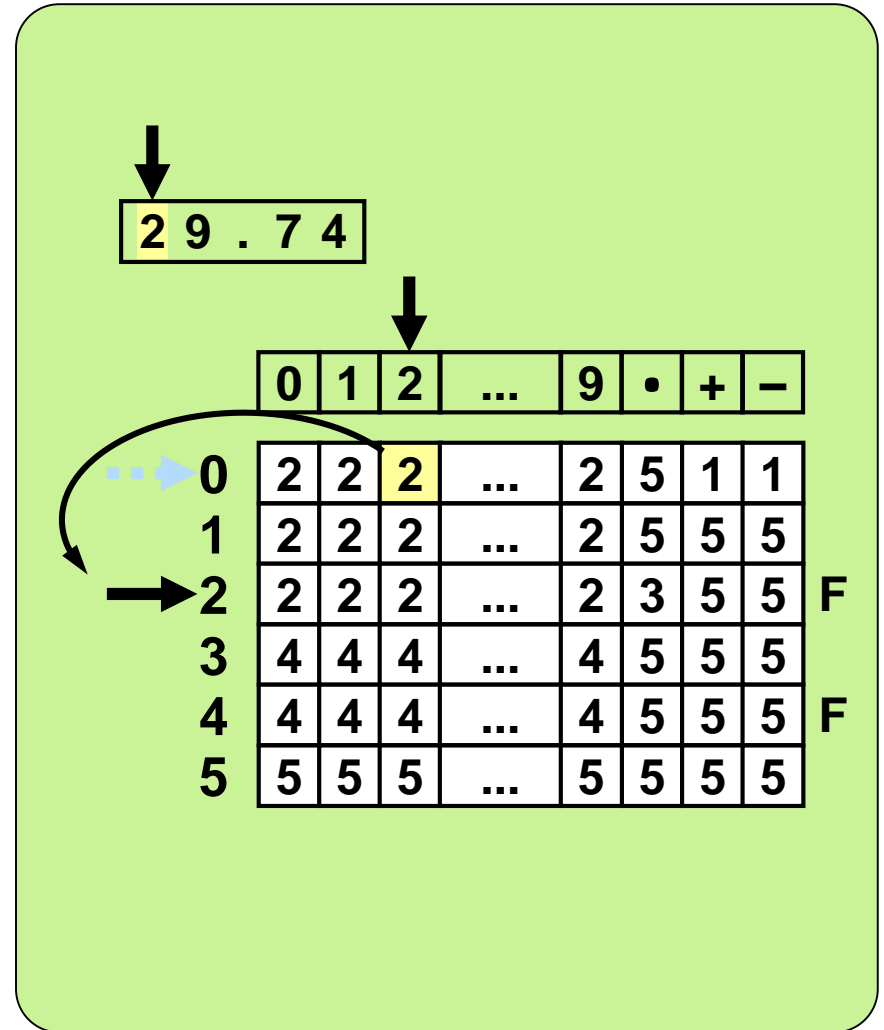
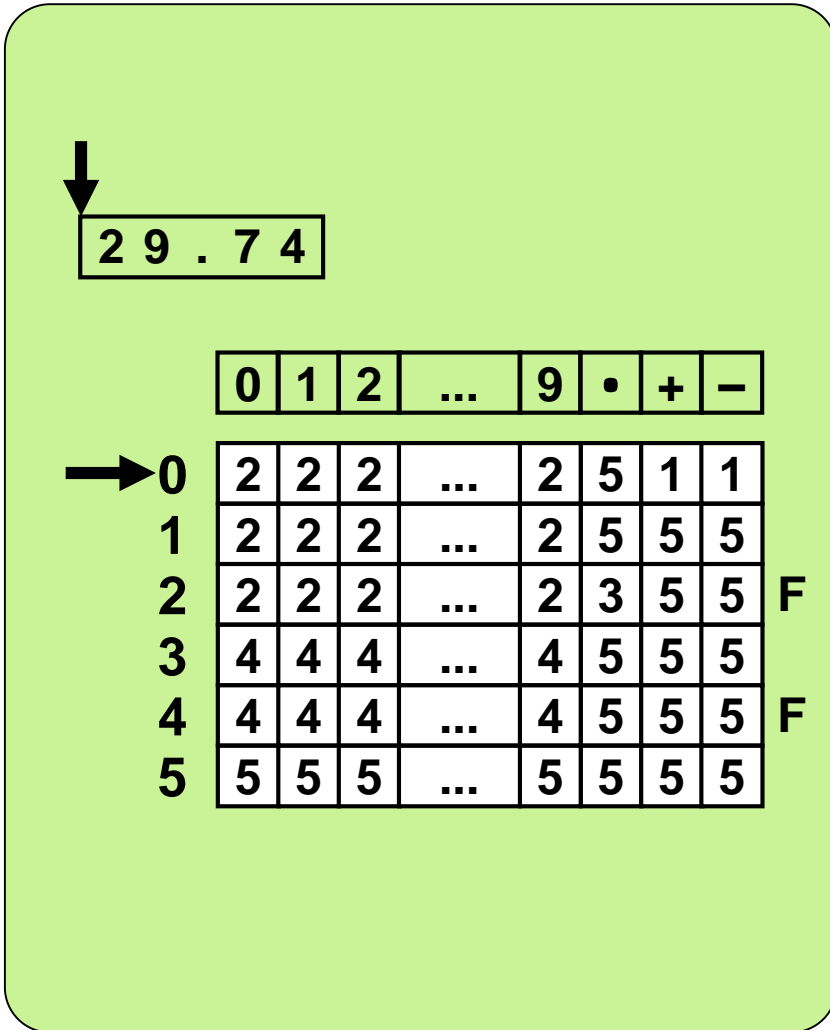


↓
2 9 . 7 4

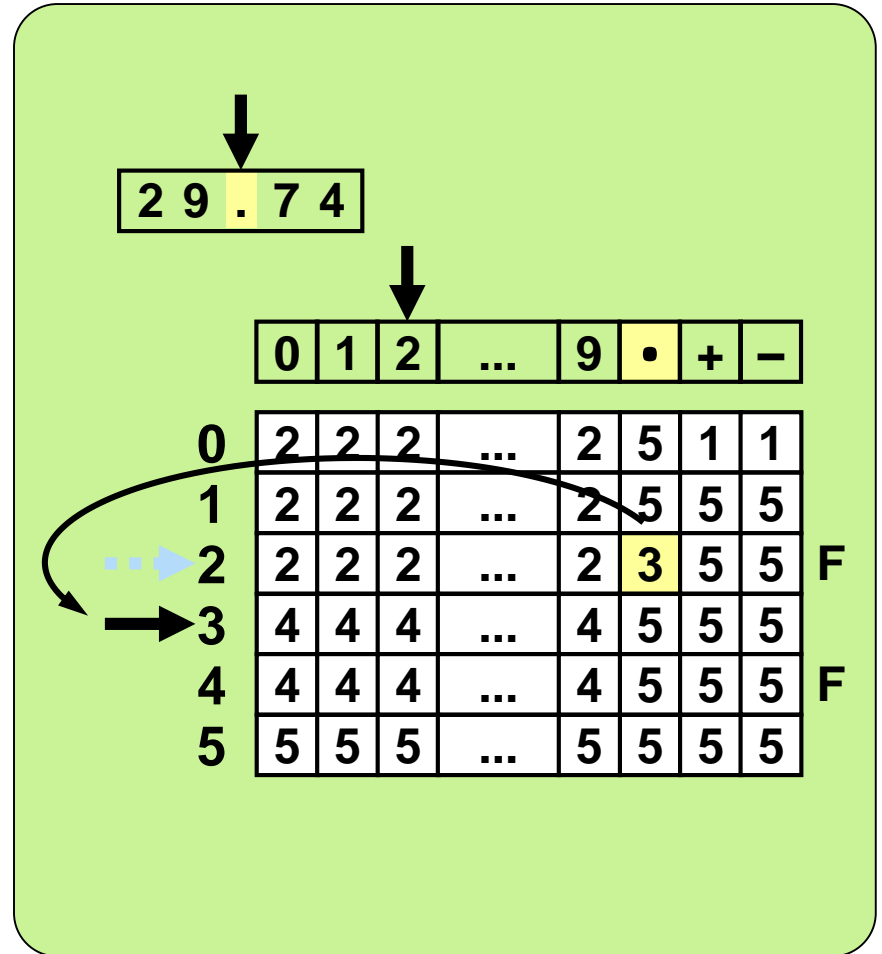
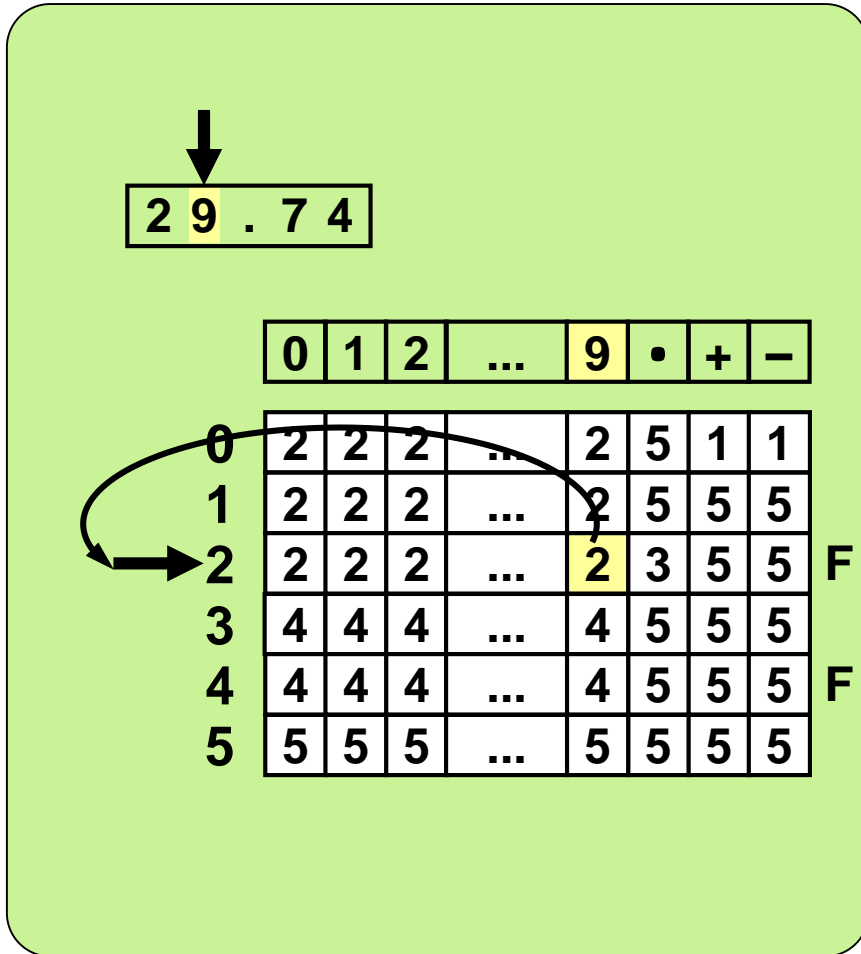


④ is an accepting state
— 29.74 is accepted

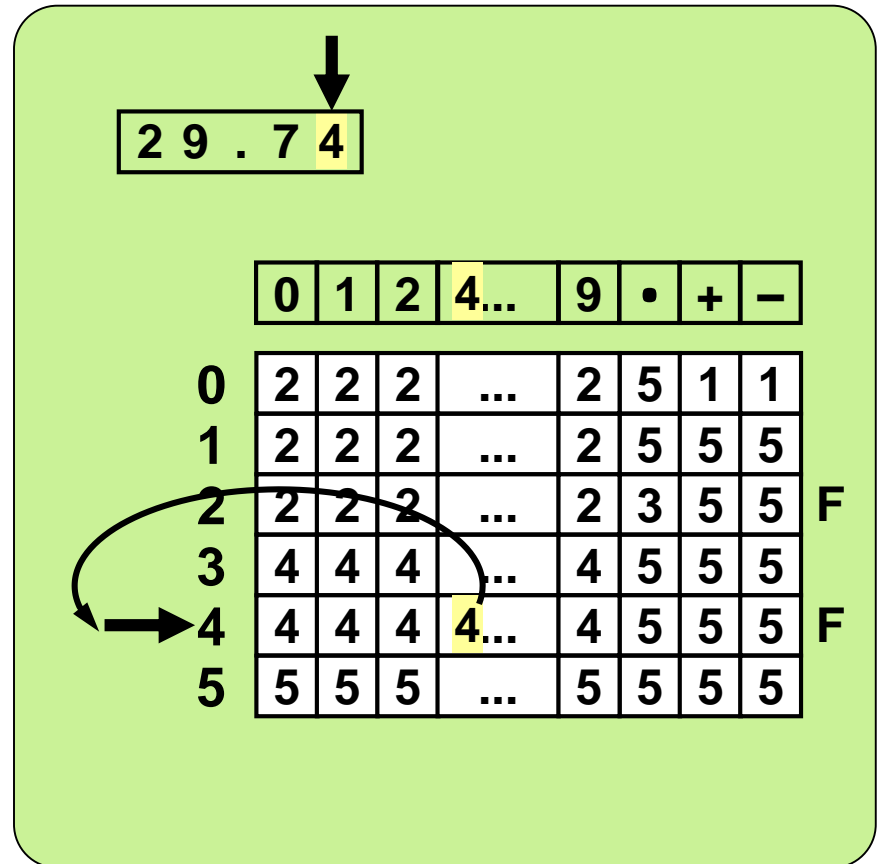
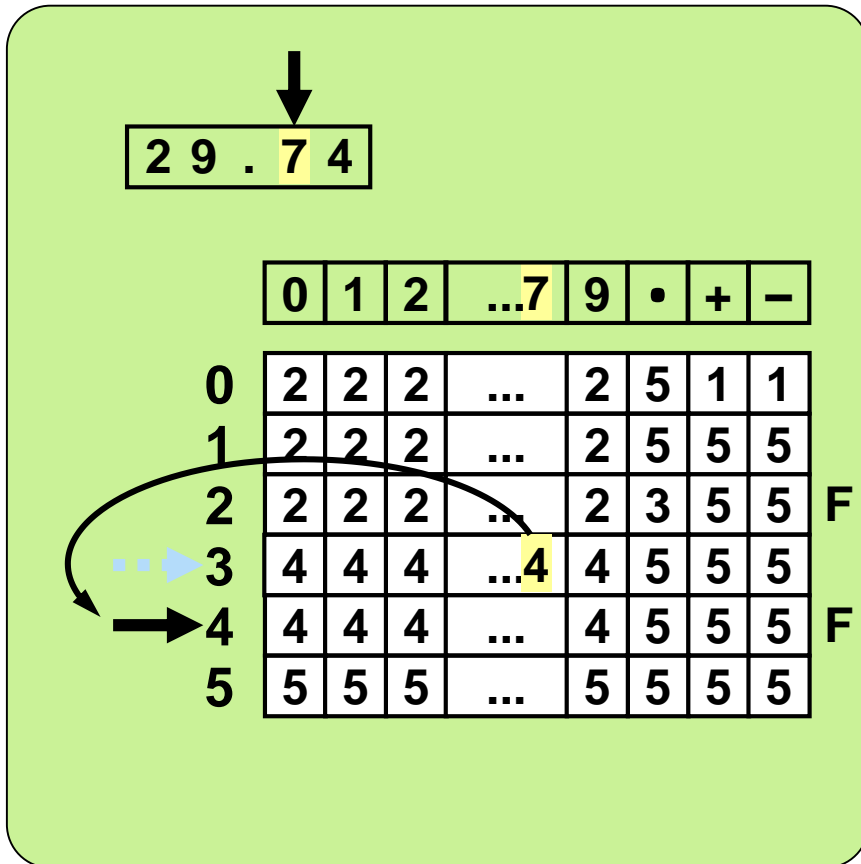
Konečný automat



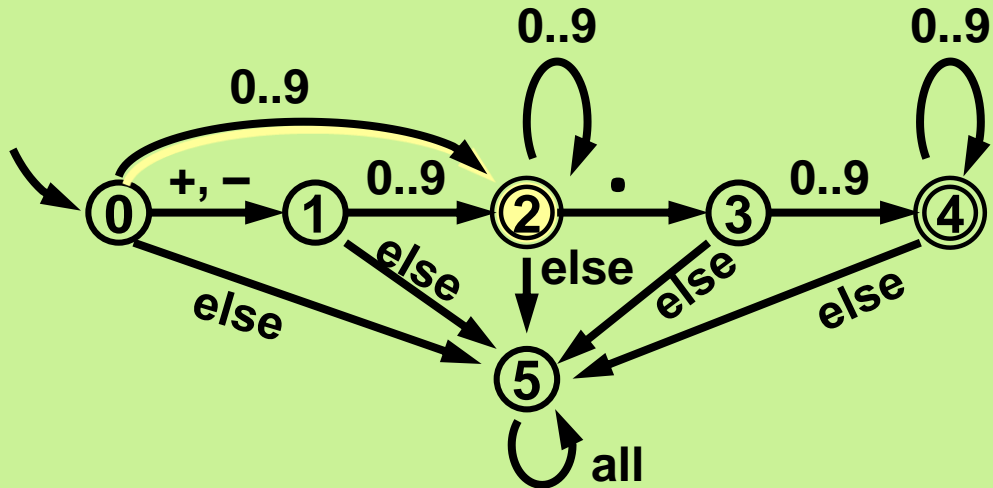
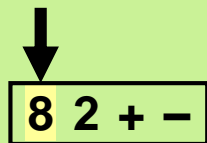
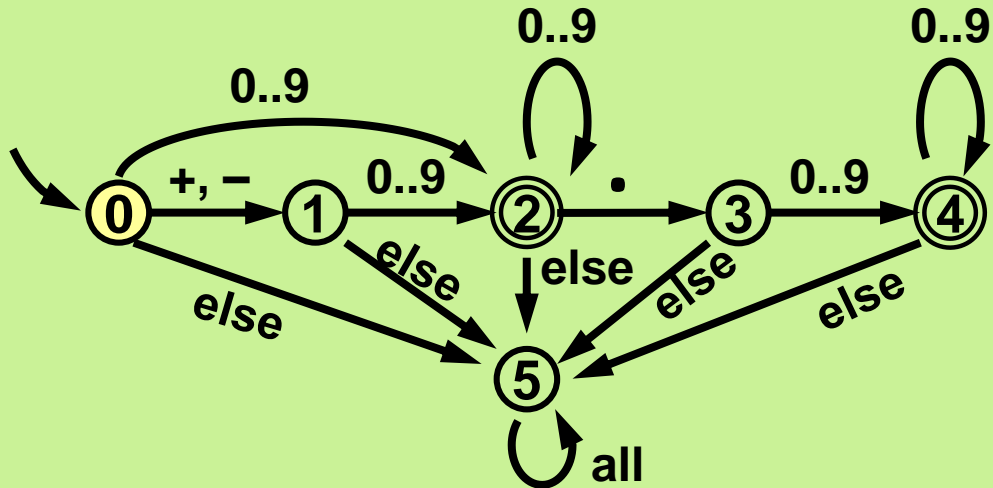
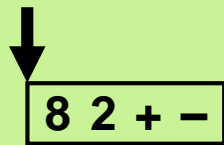
Konečný automat



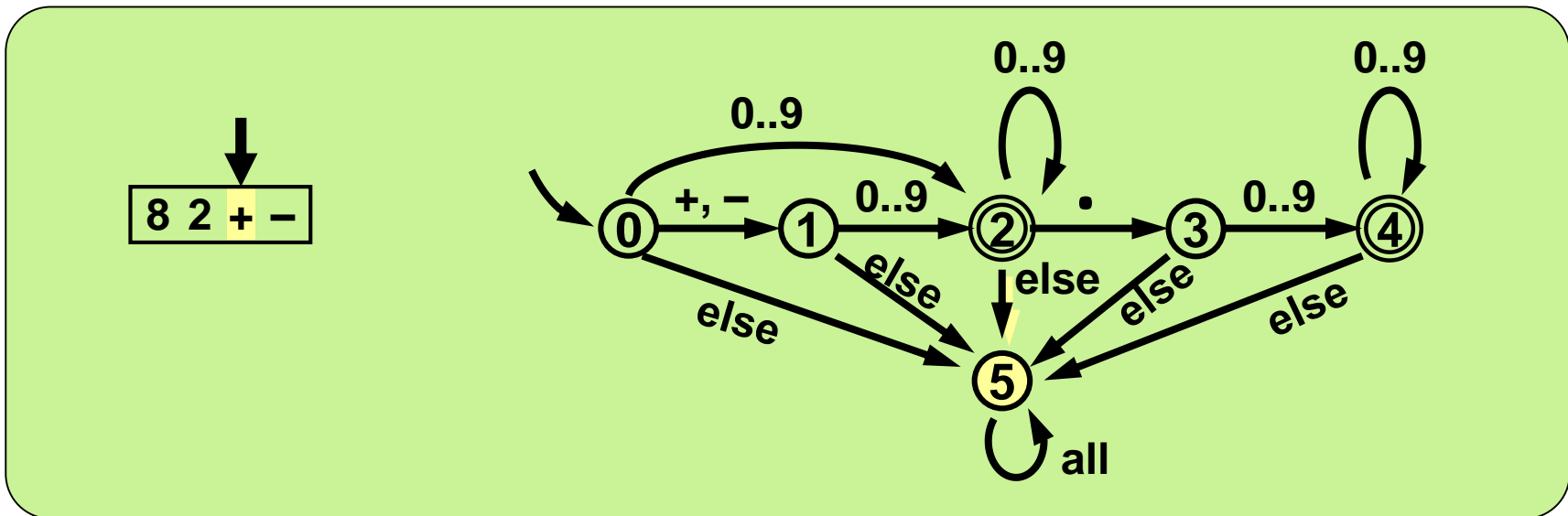
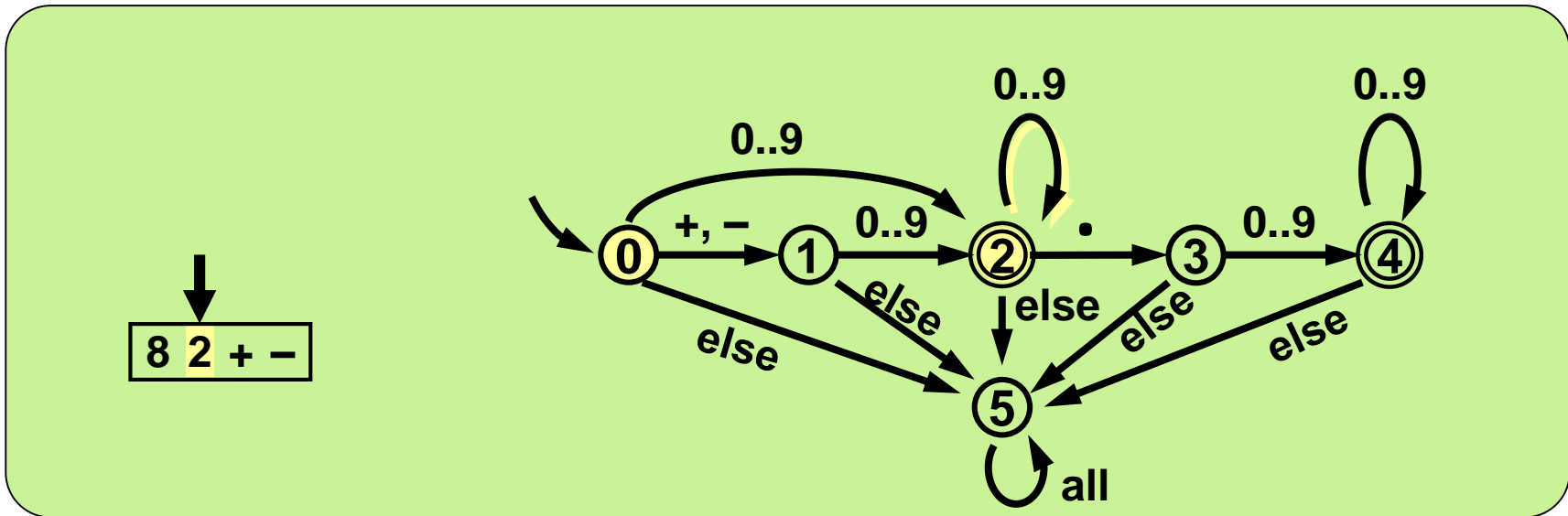
Konečný automat



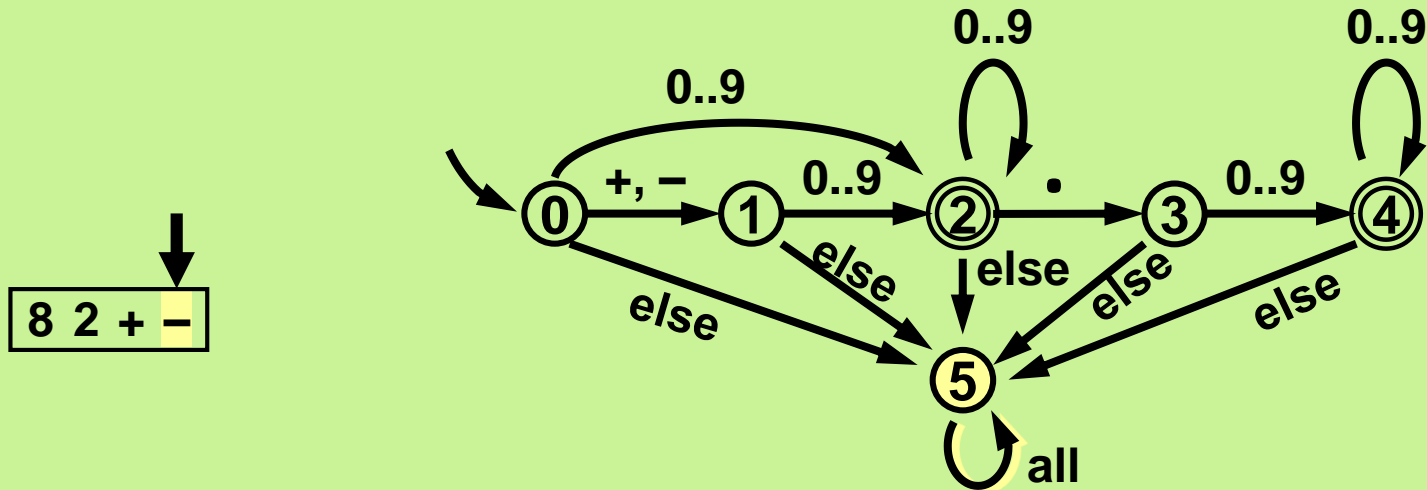
Konečný automat



Konečný automat

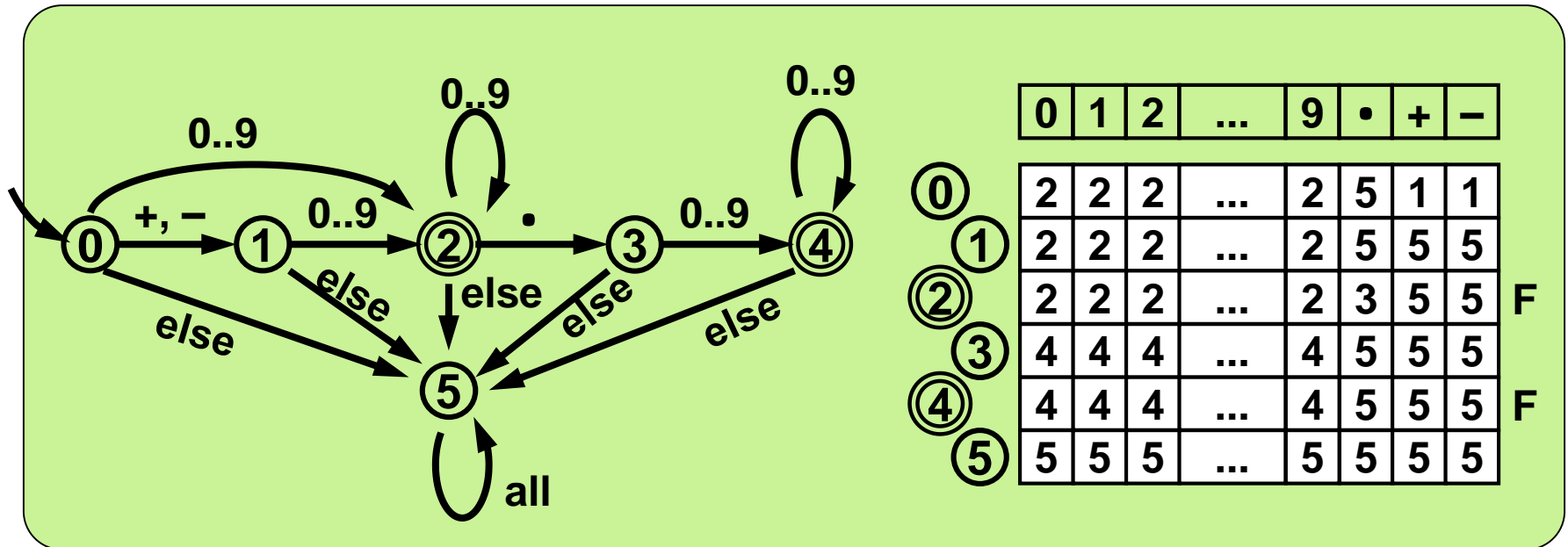


Konečný automat



⑤ is not accepting state — “82+–” is not accepted

Konečný automat



Q: finite set of states ① ② ③ ④ ⑤

Σ : finite alphabet {0,1,2,...,9, •,+, -}

δ : mapping $Q \times \Sigma \rightarrow Q$, all labelled arrows or the whole table

q_0 : initial state



F: set of final (accepting) states ② ④

Konstrukce KA pro „String matching“

		b	a	n	c	d	...
0	0	0	0	0	0	0	...
b	1						...
a	2						...
n	3						...
a	4						...
n	5						...
a	6						...

F

- 1: $M \leftarrow (\{q_0, q_1, \dots, q_m\}, \Sigma, \delta, q_0, \{q_m\})$
- 2: **for** $\forall a \in \Sigma$ **do**
- 3: $\delta(q_0, a) \leftarrow \{q_0\}$ {self-loop of the initial state}
- 4: **end for**

Konstrukce KA pro „String matching“

		b	a	n	c	d	...	$i = 1$
0	1	0	0	0	0	0	...	
b 1	1	0	0	0	0	0	...	$r = \delta(0, b) = 0$
a 2							...	
n 3							...	
a 4							...	
n 5							...	
a 6							...	F

```

5: for  $i \leftarrow 1..m$  do
6:    $r \leftarrow \delta(q_{i-1}, p_i)$ 
7:    $\delta(q_{i-1}, p_i) \leftarrow q_i$  {forward transition}
8:   for  $\forall a \in \Sigma$  do
9:      $\delta(q_i, a) \leftarrow \delta(r, a)$ 
10:  end for
11: end for

```

Konstrukce KA pro „String matching“

		b	a	n	c	d	...	
	0	1	0	0	0	0	...	
b	1	1	2	0	0	0	...	
a	2	1	0	0	0	0	...	
n	3						...	
a	4						...	
n	5						...	
a	6						...	F

$i = 2$

$r = \delta(1, a) = 0$

```

5: for  $i \leftarrow 1..m$  do
6:    $r \leftarrow \delta(q_{i-1}, p_i)$ 
7:    $\delta(q_{i-1}, p_i) \leftarrow q_i$  {forward transition}
8:   for  $\forall a \in \Sigma$  do
9:      $\delta(q_i, a) \leftarrow \delta(r, a)$ 
10:  end for
11: end for

```


Konstrukce KA pro „String matching“

		b	a	n	c	d	...
	0	1	0	0	0	0	...
b	1	1	2	0	0	0	...
a	2	1	0	3	0	0	...
n	3	1	0	0	0	0	...
a	4						...
n	5						...
a	6						...

$i = 3$

$r = \delta(2, n) = 0$

F

```

5: for  $i \leftarrow 1..m$  do
6:    $r \leftarrow \delta(q_{i-1}, p_i)$ 
7:    $\delta(q_{i-1}, p_i) \leftarrow q_i$  {forward transition}
8:   for  $\forall a \in \Sigma$  do
9:      $\delta(q_i, a) \leftarrow \delta(r, a)$ 
10:  end for
11: end for

```

Konstrukce KA pro „String matching“

		b	a	n	c	d	...
	0	1	0	0	0	0	...
b	1	1	2	0	0	0	...
a	2	1	0	3	0	0	...
n	3	1	4	0	0	0	...
a	4	1	0	5	0	0	...
n	5	1	6	0	0	0	...
a	6	1	0	0	0	0	...

$i = 6$

$F \quad r = \delta(5, a) = 0$

```

5: for  $i \leftarrow 1..m$  do
6:    $r \leftarrow \delta(q_{i-1}, p_i)$ 
7:    $\delta(q_{i-1}, p_i) \leftarrow q_i$  {forward transition}
8:   for  $\forall a \in \Sigma$  do
9:      $\delta(q_i, a) \leftarrow \delta(r, a)$ 
10:  end for
11: end for

```

Srovnání metod pro jeden vzorek

Algorithm	Preprocessing time	Matching time ¹
Naive string search algorithm	0 (no preprocessing)	$\Theta((n-m+1) m)$
Rabin–Karp string search algorithm	$\Theta(m)$	average $\Theta(n+m)$, worst $\Theta((n-m+1) m)$
Finite-state automaton based search	$\Theta(m \Sigma)$	$\Theta(n)$
Knuth–Morris–Pratt algorithm	$\Theta(m)$	$\Theta(n)$
Boyer–Moore string search algorithm	$\Theta(m + \Sigma)$	$\Omega(n/m), O(n)$

kde m je délka vzorku a n délka prohledávaného textu

The End