

Design and Analysis of Algorithms



Lecture 11: String Matching

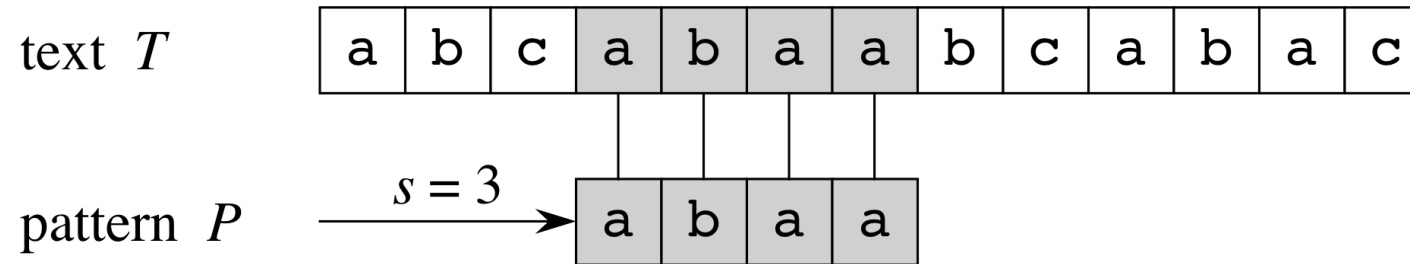
Ahmed Hamdy

Agenda

- Naïve algorithm
- Rabin-Karp algorithm
- Finite Automata (FA) algorithm
- Knuth-Morris-Pratt (KMP) algorithm

String matching

- Simply put, find *all occurrences* of string called *pattern P* (of length *m*) inside another one called *text T* (of length *n*).



- Can be viewed as find the *shift s* ($0 \leq s \leq n - m$) by which *P* appears in *T*.
- Σ denotes the alphabet; the unique characters in *P* (*a* and *b* in the above example).

String matching algorithms

Performance of algorithms

- Σ denotes the alphabet

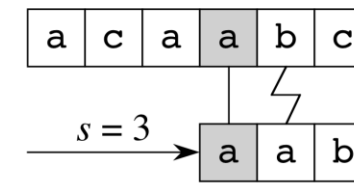
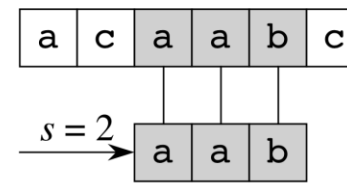
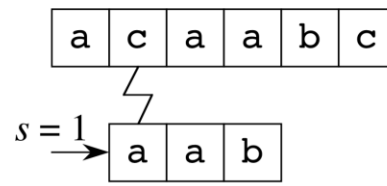
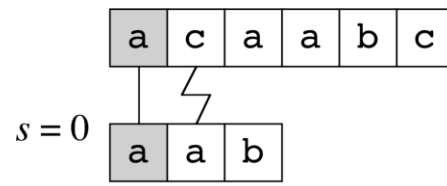
Algorithm	Preprocessing time	Matching time
Naive	0	$O((n - m + 1)m)$
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$
Finite automaton	$O(m \Sigma)$	$\Theta(n)$
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(n)$

- When is each algorithm suitable??

Naive string-matching

NAIVE-STRING-MATCHER(T, P)

```
1   $n = T.length$ 
2   $m = P.length$ 
3  for  $s = 0$  to  $n - m$ 
4      if  $P[1..m] == T[s + 1..s + m]$ 
5          print "Pattern occurs with shift"  $s$ 
```

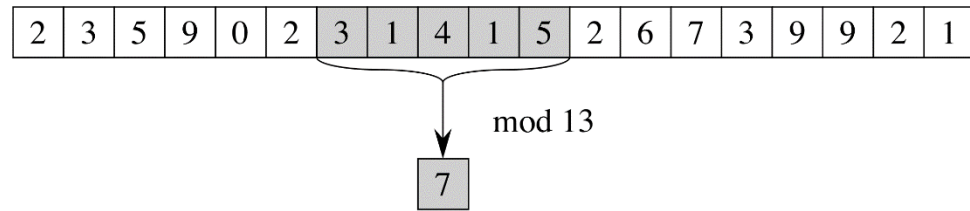


- Worst case running time $O((n - m + 1)m)$ which is $O(n^2)$ if $m = \lfloor n/2 \rfloor$.
- Room for optimization where the algorithm does not make use of information from previous iteration.

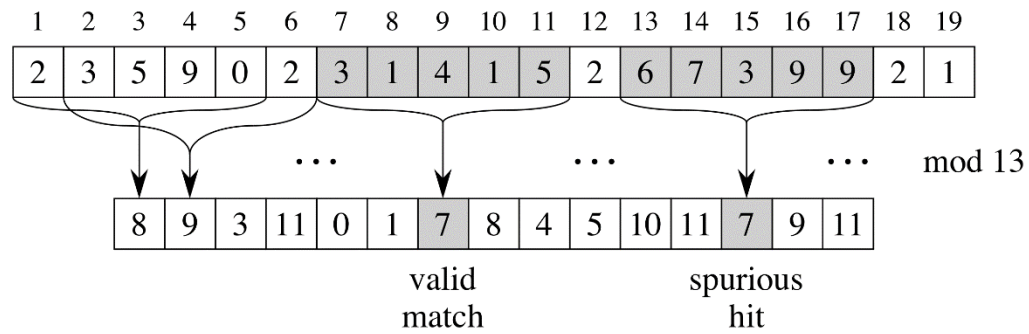
Rabin-Karp algorithm

- Compute the hash for the *pattern* P .
- Compute the hash for the *text* T at each shift $s = 0..(n - m + 1)$.
- Compare $hash(P)$ with each $hash(T[s + 1 : s + m])$
 - If equal, compare character by character to verify the collision of hashes.
- How to optimize the computation of hashes for each s ?

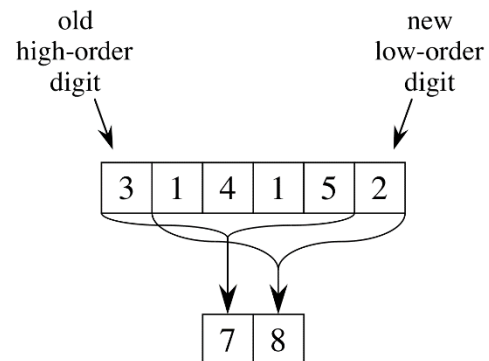
Rabin-Karp algorithm



(a)



(b)



(c)

$$\begin{aligned}
 14152 &\equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13} \\
 &\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13} \\
 &\equiv 8 \pmod{13}
 \end{aligned}$$

$$\begin{aligned}
 &\text{hash}(s + 1) \\
 &= (\text{hash}(s) - T[s + 1] \cdot \text{const1}) \cdot \text{const2} + T[s + m + 1] \pmod{\text{const3}}
 \end{aligned}$$

Rabin-Karp algorithm

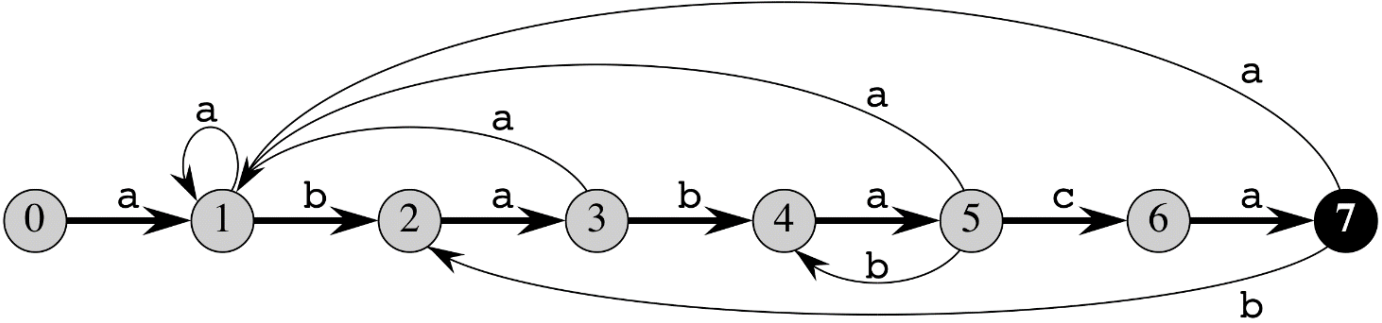
RABIN-KARP-MATCHER(T, P, d, q)

```
1   $n = T.length$ 
2   $m = P.length$ 
3   $h = d^{m-1} \bmod q$ 
4   $p = 0$ 
5   $t_0 = 0$ 
6  for  $i = 1$  to  $m$            // preprocessing
7       $p = (dp + P[i]) \bmod q$  //  $(d \cdot p + P[i]) \bmod q$ 
8       $t_0 = (dt_0 + T[i]) \bmod q$  //  $(d \cdot t_0 + T[i]) \bmod q$ 
9  for  $s = 0$  to  $n - m$        // matching
10     if  $p == t_s$ 
11         if  $P[1..m] == T[s + 1..s + m]$ 
12             print "Pattern occurs with shift"  $s$ 
13     if  $s < n - m$ 
14          $t_{s+1} = (d(t_s - T[s + 1]h) + T[s + m + 1]) \bmod q$ 
```

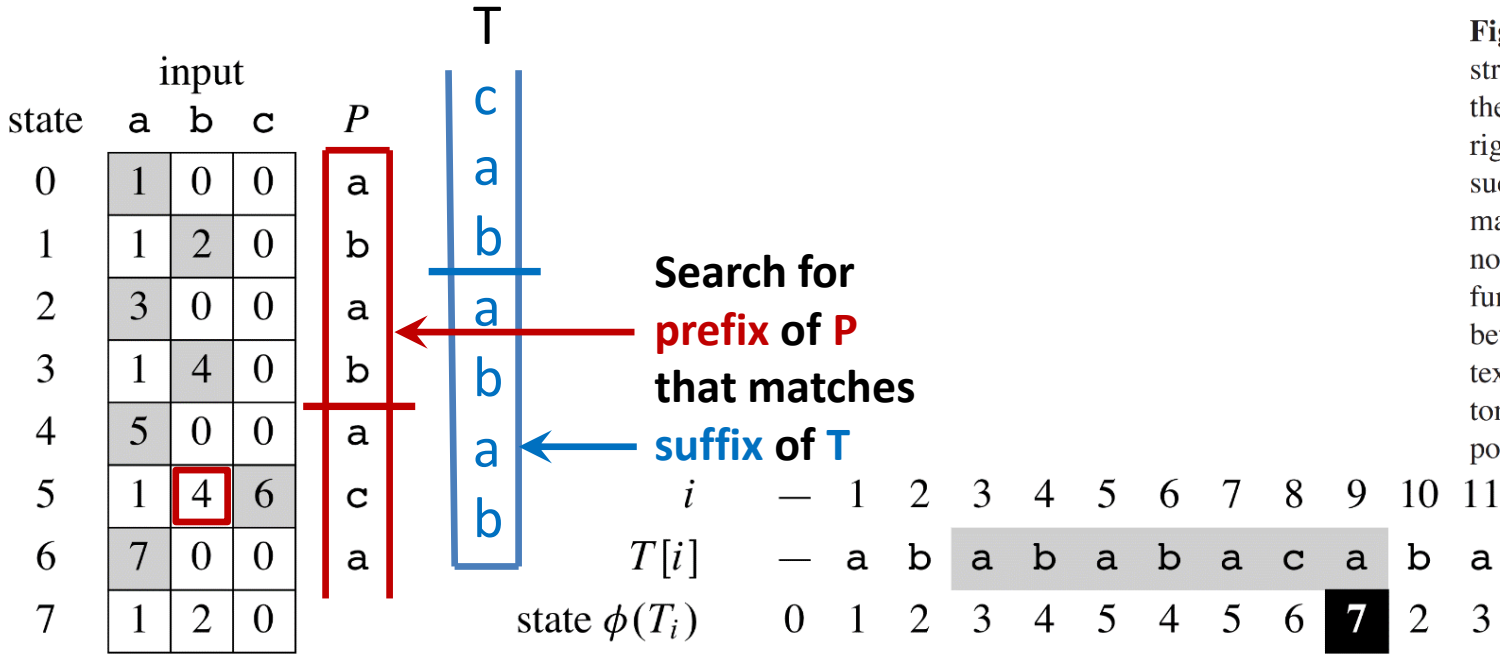
- Though worst case is not better than the naïve algorithm, the average case is much better, typically $O((n - m + 1) + cm) = O(n + m)$

Finite automata

- Complexity: $\Theta(n)$



(a)



(b)

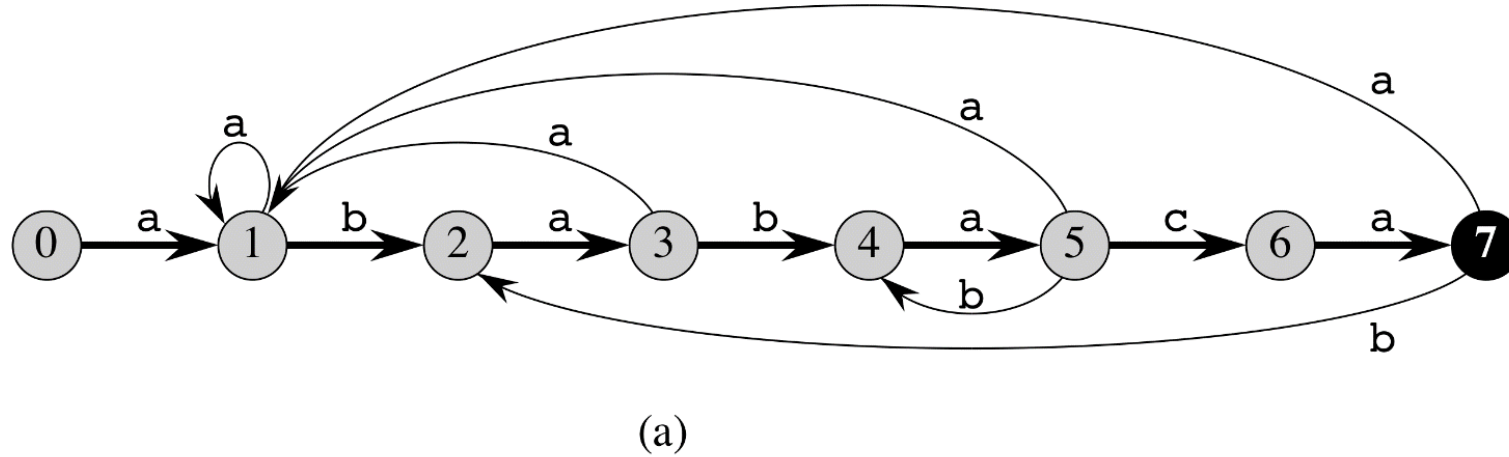
(c)

FINITE-AUTOMATON-MATCHER(T, δ, m)

```
1   $n = T.length$ 
2   $q = 0$ 
3  for  $i = 1$  to  $n$ 
4       $q = \delta(q, T[i])$ 
5      if  $q == m$ 
6          print "Pattern occurs with shift"  $i - m$ 
```

Figure 32.7 (a) A state-transition diagram for the string-matching automaton that accepts all strings ending in the string $ababaca$. State 0 is the start state, and state 7 (shown blackened) is the only accepting state. A directed edge from state i to state j labeled a represents $\delta(i, a) = j$. The right-going edges forming the “spine” of the automaton, shown heavy in the figure, correspond to successful matches between pattern and input characters. The left-going edges correspond to failing matches. Some edges corresponding to failing matches are omitted; by convention, if a state i has no outgoing edge labeled a for some $a \in \Sigma$, then $\delta(i, a) = 0$. (b) The corresponding transition function δ , and the pattern string $P = ababaca$. The entries corresponding to successful matches between pattern and input characters are shown shaded. (c) The operation of the automaton on the text $T = abababacaba$. Under each text character $T[i]$ appears the state $\phi(T_i)$ that the automaton is in after processing the prefix T_i . The automaton finds one occurrence of the pattern, ending in position 9.

Finite Automata



state	input			<i>P</i>	<i>T</i>
	a	b	c		
0	1	0	0	a	c
1	1	2	0	b	a
2	3	0	0	a	b
3	1	4	0	b	a
4	5	0	0	a	b
5	1	4	6	c	a
6	7	0	0	a	b
7	1	2	0		

(b)

When $q = 5$, initially $k = (q + 2) - 1 = 6$.

Search for
prefix of *P*
that matches

suffix of *T*

<i>i</i>	—	1	2	3	4	5	6	7	8	9	10	11
<i>T</i> [<i>i</i>]	—	a	b	a	b	a	b	a	c	a	b	a
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

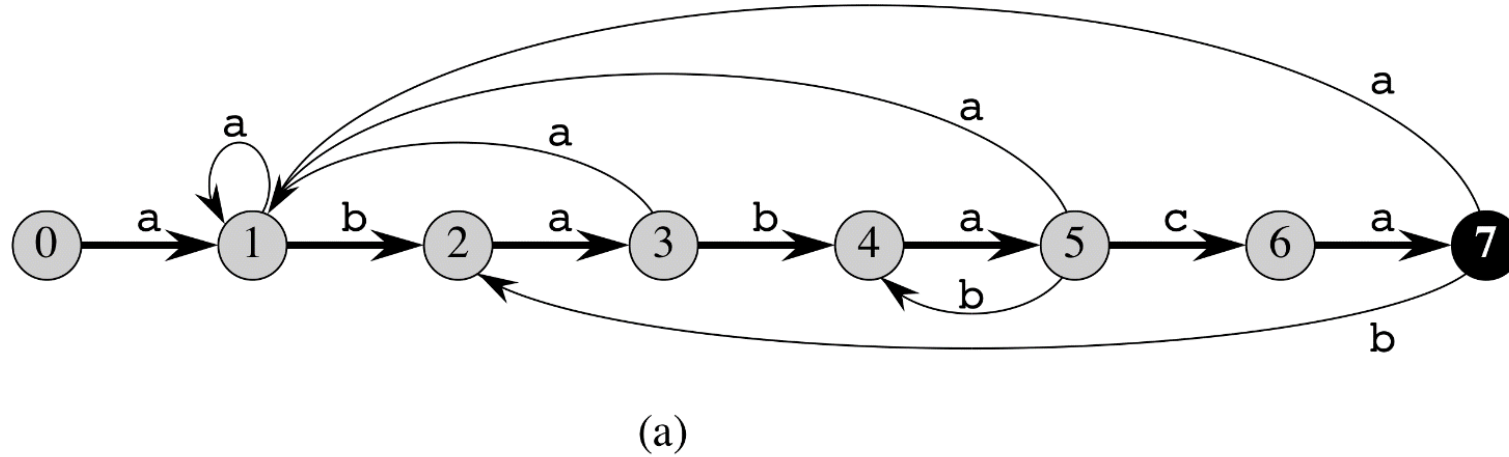
(c)

COMPUTE-TRANSITION-FUNCTION(*P*, Σ)

```

1  m = P.length
2  for q = 0 to m
3      for each character a ∈  $\Sigma$ 
4          k = min(m + 1, q + 2)
5          repeat
6              k = k - 1
7          until  $P_k \sqsubseteq P_q a$ 
8           $\delta(q, a) = k$ 
9  return  $\delta$ 
    
```

Finite Automata



state	input			<i>P</i>	<i>T</i>
	a	b	c		
0	1	0	0	a	c
1	1	2	0	b	a
2	3	0	0	a	b
3	1	4	0	b	a
4	5	0	0	a	b
5	1	4	6	c	a
6	7	0	0	a	b
7	1	2	0		

When $q = 5$, initially $k = k - 1 = 5$.

Search for
prefix of *P*
that matches

suffix of *T*

<i>i</i>	—	1	2	3	4	5	6	7	8	9	10	11
<i>T</i> [<i>i</i>]	—	a	b	a	b	a	b	a	c	a	b	a
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

COMPUTE-TRANSITION-FUNCTION(*P*, Σ)

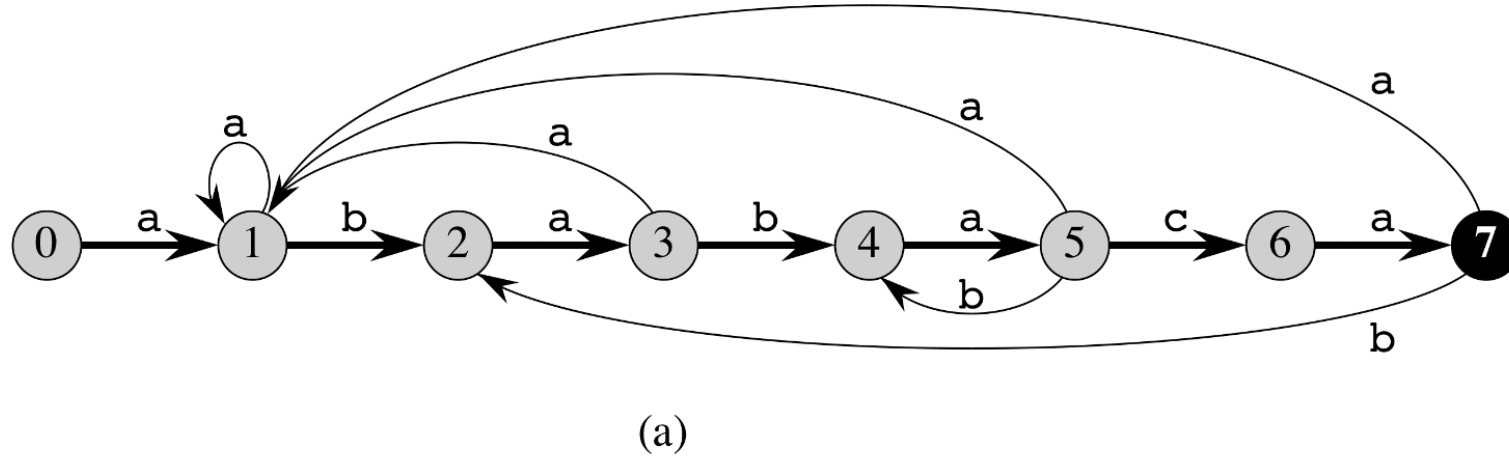
```

1   $m = P.length$ 
2  for  $q = 0$  to  $m$ 
3      for each character  $a \in \Sigma$ 
4           $k = \min(m + 1, q + 2)$ 
5          repeat
6               $k = k - 1$ 
7          until  $P_k \sqsubseteq P_q a$ 
8               $\delta(q, a) = k$ 
9  return  $\delta$ 
    
```

(b)

(c)

Finite Automata



state	input			<i>P</i>
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

(b)

When $q = 5$, initially $k = k - 1 = 4$.

Search for
prefix of *P*
that matches

suffix of *T*

<i>i</i>	—	1	2	3	4	5	6	7	8	9	10	11
<i>T</i> [<i>i</i>]	—	a	b	a	b	a	b	a	c	a	b	a
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

(c)

COMPUTE-TRANSITION-FUNCTION(*P*, Σ)

```

1   $m = P.length$ 
2  for  $q = 0$  to  $m$ 
3      for each character  $a \in \Sigma$ 
4           $k = \min(m + 1, q + 2)$ 
5          repeat
6               $k = k - 1$ 
7          until  $P_k \sqsubseteq P_q a$ 
8               $\delta(q, a) = k$ 
9  return  $\delta$ 
    
```

From FA to KMP

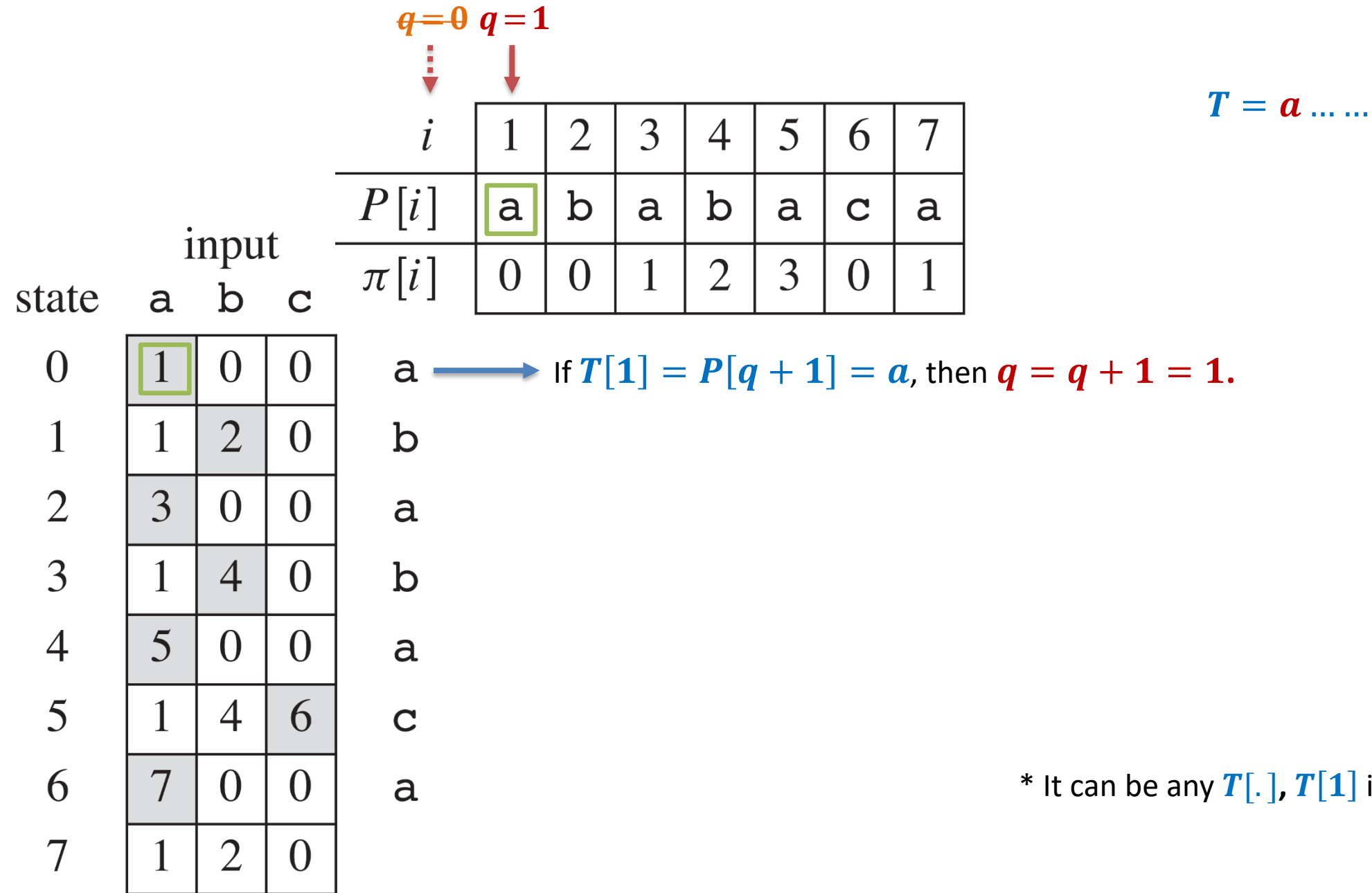
- The table produced by FA can be squeezed to be one column!!
- Don't compute where to go back for each character ('a', 'b', and 'c').
- Instead, go back to what is common and check there.
- Now preprocessing complexity is $\Theta(m)$ instead of $O(m|\Sigma|)$.

From FA to KMP

				$q = 0$	
				↓	
				i	
					1 2 3 4 5 6 7
				$P[i]$	a b a b a c a
				$\pi[i]$	0 0 1 2 3 0 1
state	input				
	a	b	c		
0	1	0	0	a	
1	1	2	0	b	
2	3	0	0	a	
3	1	4	0	b	
4	5	0	0	a	
5	1	4	6	c	
6	7	0	0	a	
7	1	2	0		



- A variable q holds the current state
- Instead of having the straightforward entries (1,2,3,..7) as we match each character in P , q will hold this value and gets incremented each time.
- Initially $q = 0$ (points before the table, corresponds to i inside the above table).

From FA to KMP: Correct-matching



* It can be any $T[.]$, $T[1]$ is used for simplicity.

From FA to KMP: Correct-matching

				$q=1$ $q=2$							
											
				i	1	2	3	4	5	6	7
				$P[i]$	a	b	a	b	a	c	a
				$\pi[i]$	0	0	1	2	3	0	1

input

state	a	b	c
0	1	0	0
1	1	2	0
2	3	0	0
3	1	4	0
4	5	0	0
5	1	4	6
6	7	0	0
7	1	2	0

a

b

a

b

a

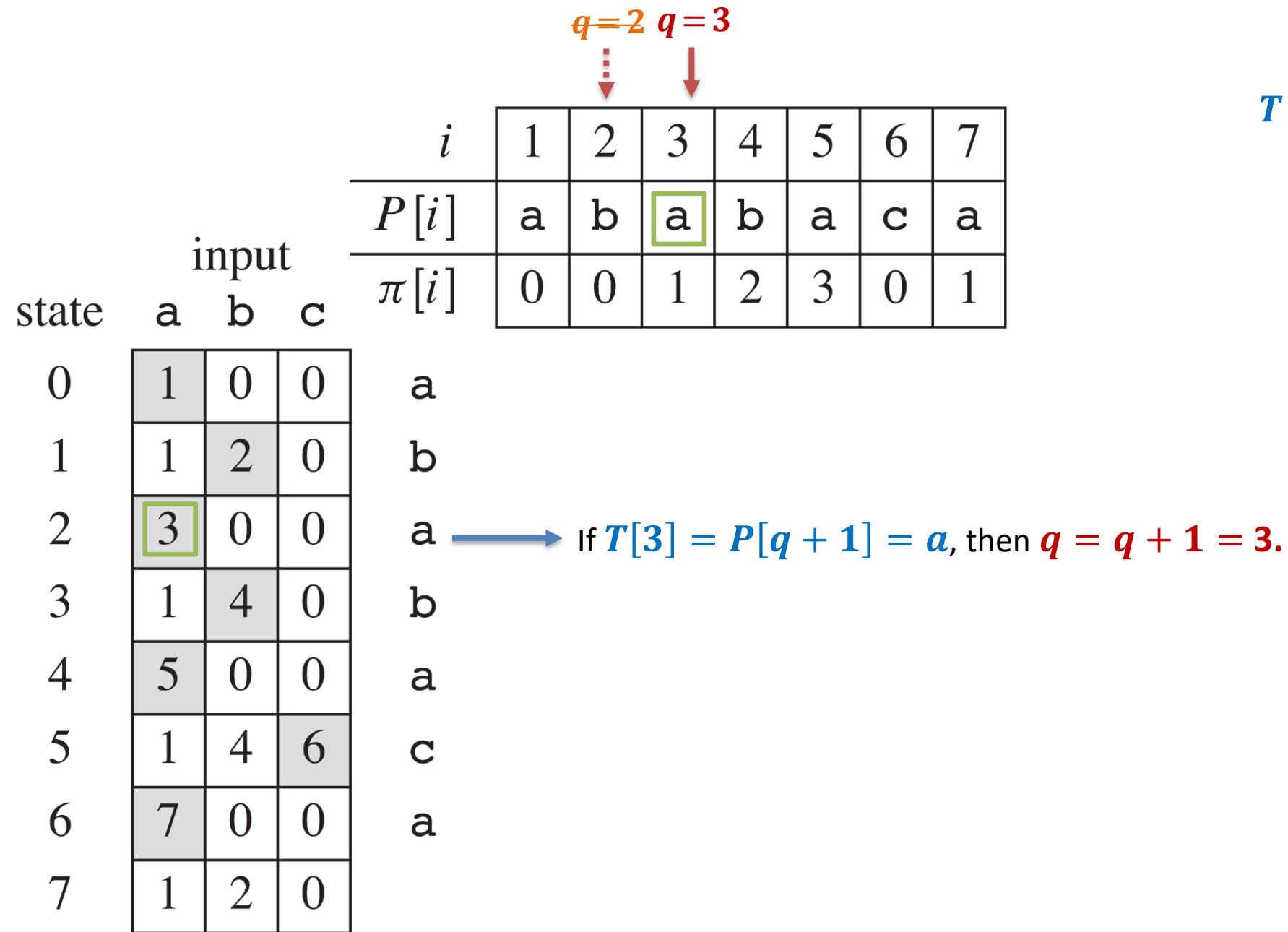
c

a

If $T[2] = P[q + 1] = b$, then $q = q + 1 = 2$.

$T = ab \dots$

From FA to KMP: Correct-matching



From FA to KMP: Correct-matching

Diagram illustrating the KMP algorithm's failure function calculation for the string "abacaba".

The top table shows the string P and its failure function π :

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

Red dashed arrow points to $i=3$ (where $P[3]=a$). Red solid arrow points to $i=4$ (where $P[4]=b$). A green box highlights $P[4]=b$.

The bottom table shows the state transitions for the automaton:



state	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

A blue arrow points from the state transition table to the text: "If $T[4] = P[q + 1] = b$, then".

$T = abab \dots$

→ If $T[4] = P[q + 1] = b$, then $q = q + 1 = 4$.

From FA to KMP: Correct-matching

				$q=4$ $q=5$							
											
				i	1	2	3	4	5	6	7
				$P[i]$	a	b	a	b	a	c	a
				$\pi[i]$	0	0	1	2	3	0	1
state	input										
	a	b	c								
0	1	0	0	a							
1	1	2	0	b							
2	3	0	0	a							
3	1	4	0	b							
4	5	0	0	a							
5	1	4	6	c							
6	7	0	0	a							
7	1	2	0								

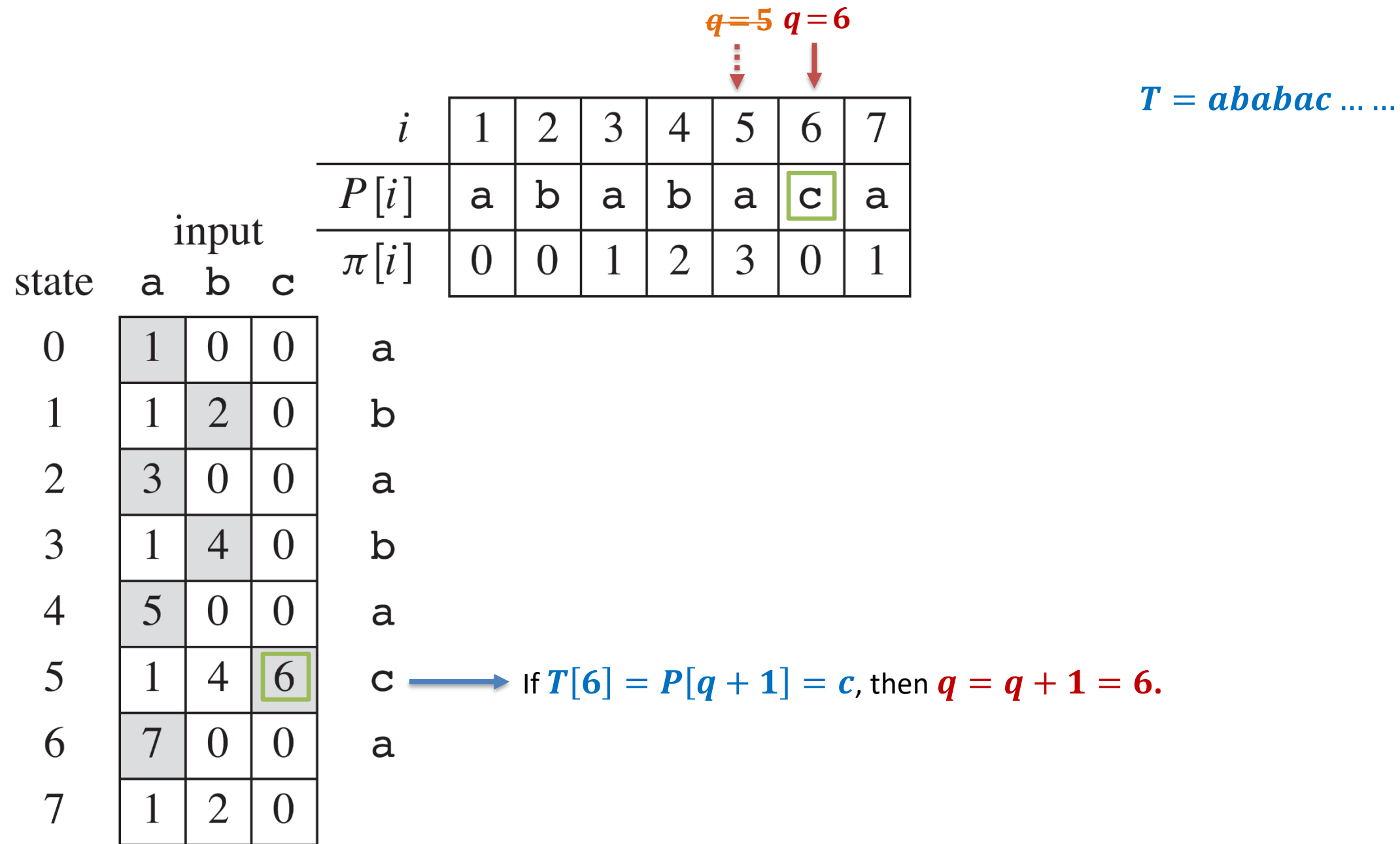
T

\longrightarrow If $T[5] = P[q + 1] = a$, then $q = q + 1 = 5$.

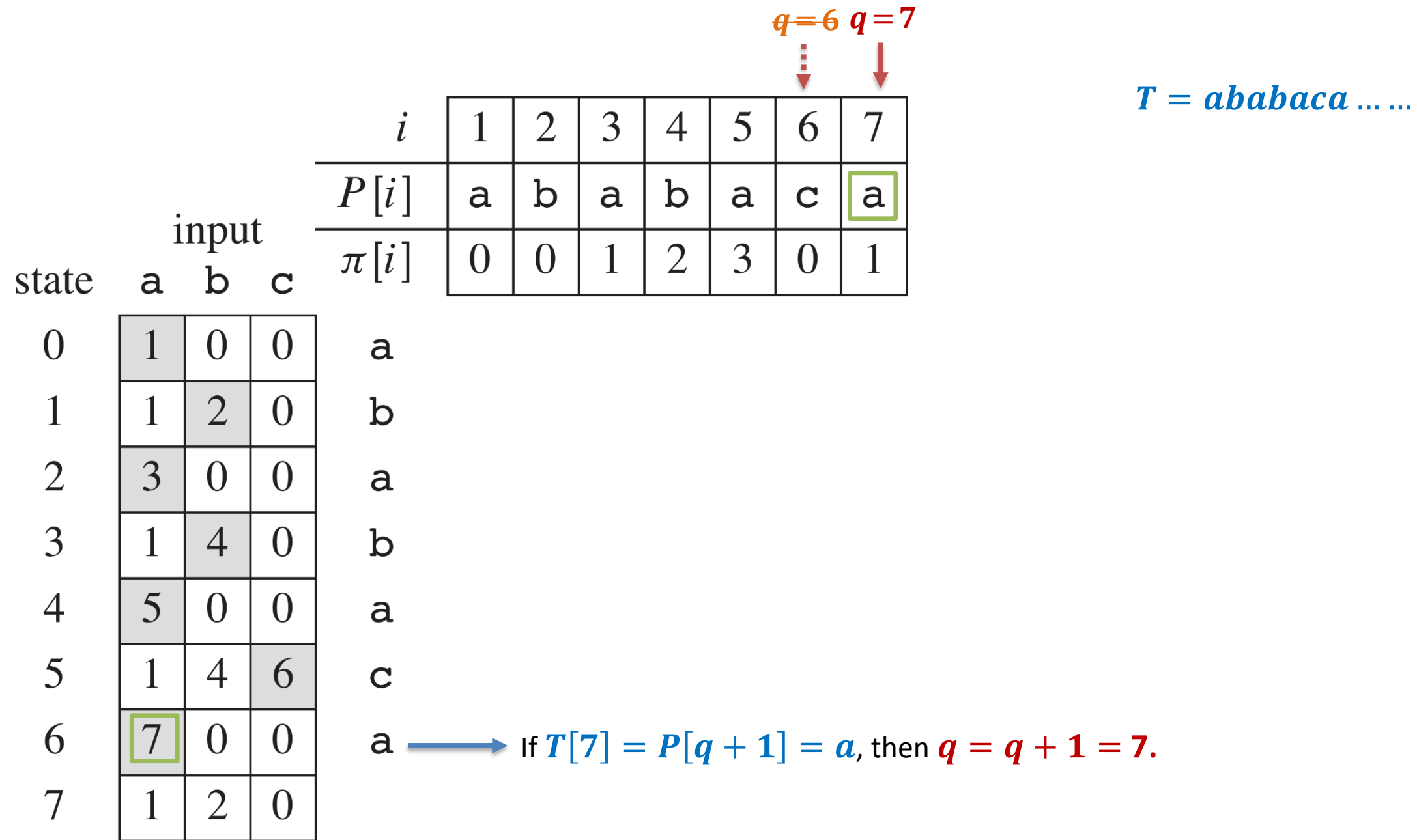
$T = ababa \dots$

\longrightarrow If $T[5] = P[q + 1] = a$, then $q = q + 1 = 5$.

From FA to KMP: Correct-matching



From FA to KMP: Correct-matching



From FA to KMP: Correct-matching

The diagram illustrates the KMP algorithm's failure function (pi) calculation and the corresponding automaton state transitions.

String and Failure Function (pi):

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

Red arrow: $q = 1$ (current state)
Orange arrow: $q' = 7$ (next state)

Automaton State Transitions:


state	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

After completing a match, we go to state $q' = 7$. So $q = \pi[q] = 1$. Means start at state 1.

$T = ababaca \dots$

After completing a match, we go back to the overlap between P and itself which is of size 1.
So $q = \pi[q] = 1$. Means start a **new** match with a .

From FA to KMP: Miss-Matching

$q = 0$


i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

state	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

\longrightarrow If $T[1] = b$ or $T[1] = c$ while $P[q + 1] = a$, then still $q = 0$.

From FA to KMP: Miss-matching

Diagram illustrating the step-by-step construction of a suffix array. The input string is `abcabcabc`. The permutation of indices is `[0, 2, 1, 5, 4, 7, 3, 6]`. The process involves sorting suffixes starting from index 1, where `b` is compared with `a` and `c` at index 2. The diagram shows the state of the suffix array after each step, with the current state being `[1, 2, 0, 3, 4, 7, 5, 6]`.

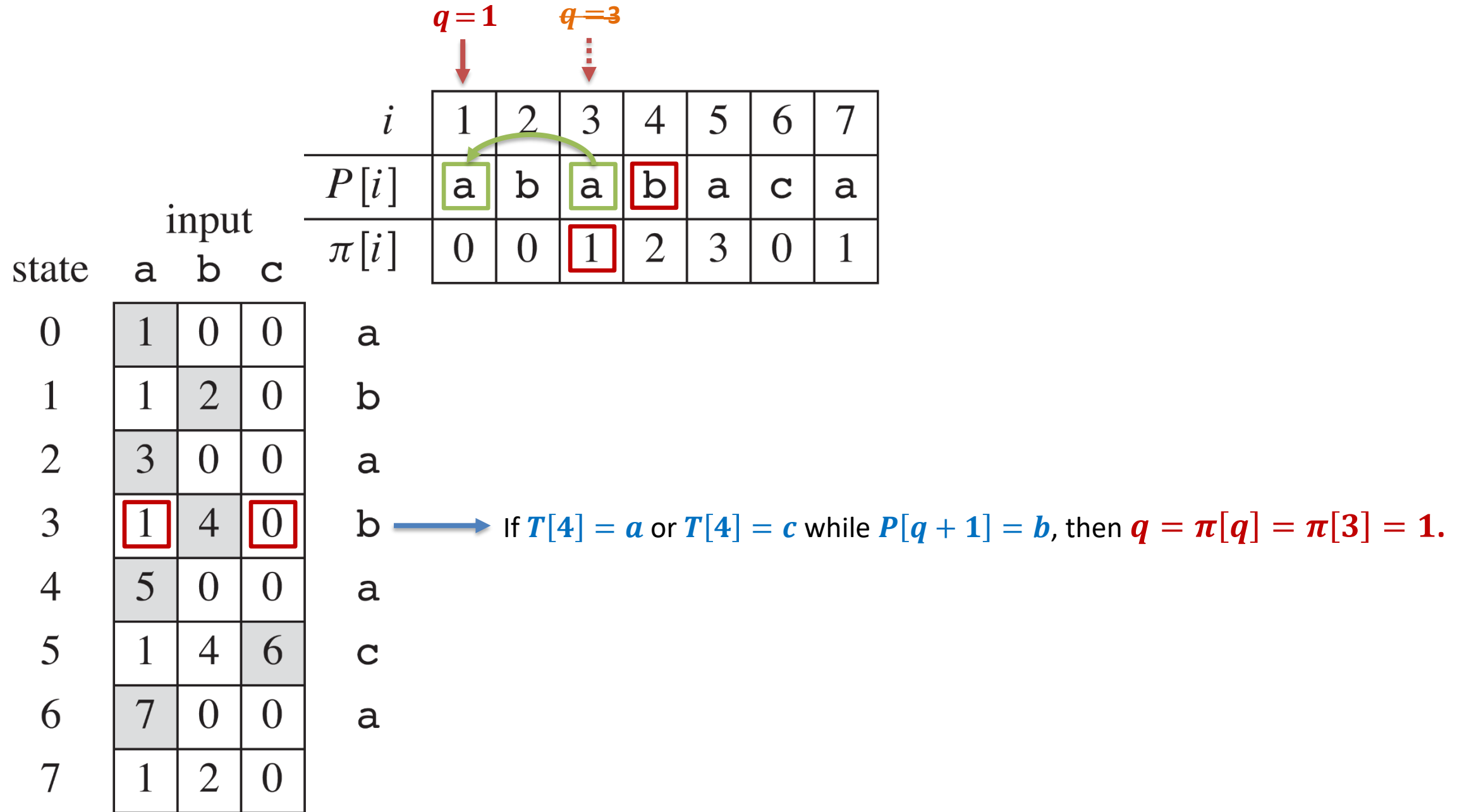
b \longrightarrow If $T[2] = a$ or $T[2] = c$ while $P[q + 1] = b$, then $q = \pi[q] = \pi[1] = 0$.

From FA to KMP: Miss-matching

				$q=0$	$q=2$						
				↓	⋮						
				i							
				1	2	3	4	5	6	7	
				$P[i]$	a	b	a	b	a	c	a
				$\pi[i]$	0	0	1	2	3	0	1
state	a	b	c								
0	1	0	0	a							
1	1	2	0	b							
2	3	0	0	a							
3	1	4	0	b							
4	5	0	0	a							
5	1	4	6	c							
6	7	0	0	a							
7	1	2	0								

If $T[3] = b$ or $T[3] = c$ while $P[q + 1] = a$, then $q = \pi[q] = \pi[2] = 0$.

From FA to KMP: Miss-matching



From FA to KMP: Miss-matching

				$q=0$	$q=1$	$q=3$							
				i	1	2	3	4	5	6	7		
				$P[i]$	a	b	a	b	a	c	a		
				$\pi[i]$	0	0	1	2	3	0	1		
state	a	b	c										
0	1	0	0	a									
1	1	2	0	b									
2	3	0	0	a									
3	1	4	0	b									
4	5	0	0	a									
5	1	4	6	c									
6	7	0	0	a									
7	1	2	0										

If $T[4] = a$ or $T[4] = c$ while
At $q = 1$, check again whether
again, so $q = \pi[q] = \pi[1]$

If $T[4] = a$ or $T[4] = c$ while $P[q + 1] = b$, then $q = \pi[q] = \pi[3] = 1$.
 At $q = 1$, check again whether a or c matches $P[q + 1] = b$ which is false again, so $q = \pi[q] = \pi[1] = 0$.

From FA to KMP: Miss-matching

Diagram illustrating the execution of the Rabin-Karp algorithm. The state transition table is shown below:

state	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

The pattern P and text T are shown as rows:

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

The diagram highlights the state transition from state 4 to state 5 when the input is 'a', and the corresponding values in the pattern and text.

→ If $T[5] = b$ or $T[5] = c$ while $P[q + 1] = a$, then $q = \pi[q] = \pi[4] = 2$.

From FA to KMP: Miss-matching

				$q=0$	$q=2$	$q=4$					
				\downarrow	\downarrow	\downarrow					
				i	1	2	3	4	5	6	7
				$P[i]$	a	b	a	b	a	c	a
				$\pi[i]$	0	0	1	2	3	0	1

From FA to KMP: Miss-matching

The diagram illustrates the KMP algorithm's failure function calculation. It shows a table for the failure function (labeled "state" in the image) and a table for the prefix function (labeled "P[i]" and "pi[i]" in the image).

Failure Function Table (state):

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

Prefix Function Table:

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

Annotations in the diagram:

- A green box highlights the 'a' at $P[3]$.
- A green arrow points from $P[3]$ to $P[2]$.
- A red box highlights the 'c' at $P[6]$.
- A red box highlights the '3' at $\pi[5]$.
- A blue arrow points from the text "If $T[6] = a$ or $T[6] = b$ while" to the 'c' column of the failure function table.

C \longrightarrow If $T[6] = a$ or $T[6] = b$ while $P[q + 1] = c$, then $q = \pi[q] = \pi[5] = 3$.

From FA to KMP: Miss-matching

				input						
state										
	a	b	c							
0	1	0	0	a						
1	1	2	0	b						
2	3	0	0	a						
3	1	4	0	b						
4	5	0	0	a						
5	1	4	6	c						
6	7	0	0	a						
7	1	2	0							

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

If $T[6] = a$ or $T[6] = b$ while $P[q + 1] = b$, then $q = \pi[q] = \pi[5] = 3$.
 At $q = 3$, check again whether a or b matches $P[q + 1] = b$:
 If $T[6] = b$ which is true, so continue normal matching with $q = 3$.

Knuth-Morris-Pratt algorithm

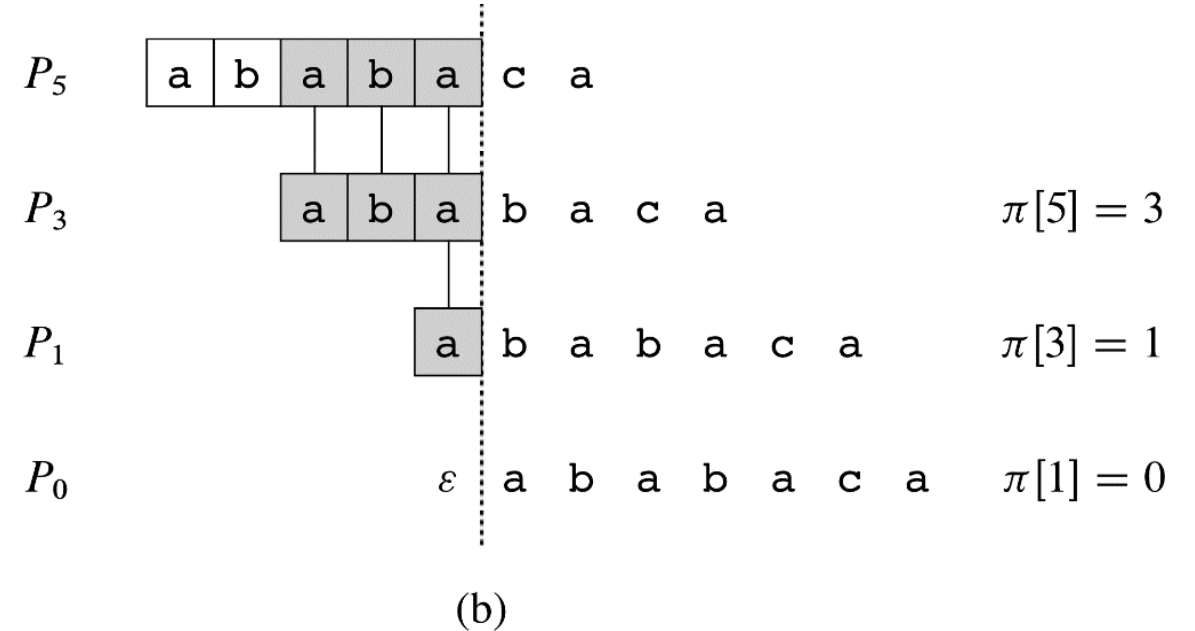
KMP-MATCHER(T, P)

```

1   $n = T.length$ 
2   $m = P.length$ 
3   $\pi = \text{COMPUTE-PREFIX-FUNCTION}(P)$ 
4   $q = 0$  // number of characters matched
5  for  $i = 1$  to  $n$  // scan the text from left to right
6      while  $q > 0$  and  $P[q + 1] \neq T[i]$ 
7           $q = \pi[q]$  // next character does not match
8      if  $P[q + 1] == T[i]$ 
9           $q = q + 1$  // next character matches
10     if  $q == m$  // is all of  $P$  matched?
11         print "Pattern occurs with shift"  $i - m$ 
12          $q = \pi[q]$  // look for the next match
    
```

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

(a)



KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in P that is equal to the prefix length of P .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0						

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in **P** that is equal to the prefix length of **P** .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0					

No suffix found to match
prefix inside this window

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in **P** that is equal to the prefix length of **P** .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1				

Suffix of length 1 matches
prefix inside this window

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in P that is equal to the prefix length of P .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2			

Suffix of length 2 matches
prefix inside this window

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in P that is equal to the prefix length of P .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3		

Suffix of length 3 matches
prefix inside this window

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in **P** that is equal to the prefix length of **P** .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	

No suffix found to match
prefix inside this window

KMP: π -table computation

- $\pi[1] = 0$ always
- For $q = 2..m$:
 - Find the longest suffix length in **P** that is equal to the prefix length of **P** .
 - Store it as $\pi[q]$.

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

Suffix of length 1 matches
prefix inside this window

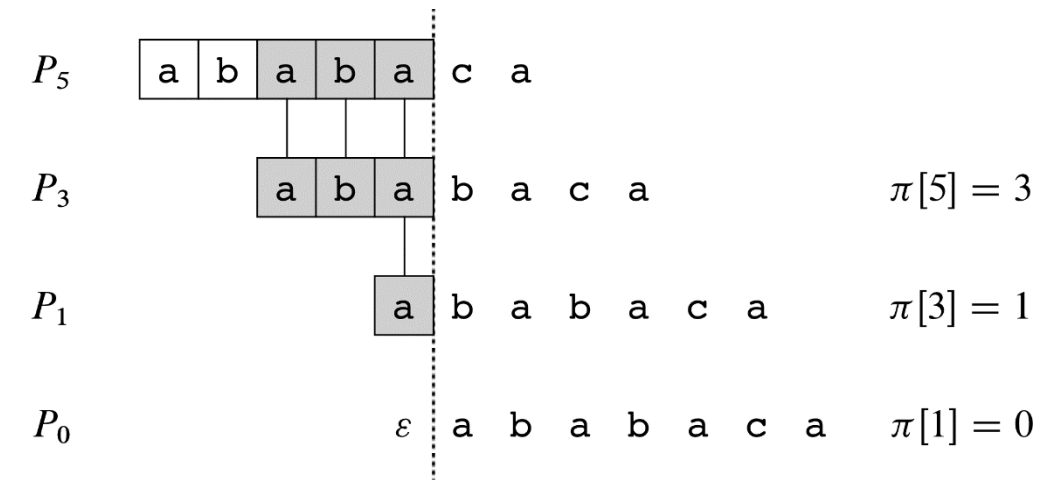
KMP: π -table computation algorithm

COMPUTE-PREFIX-FUNCTION(P)

```
1   $m = P.length$ 
2  let  $\pi[1..m]$  be a new array
3   $\pi[1] = 0$ 
4   $k = 0$ 
5  for  $q = 2$  to  $m$ 
6      while  $k > 0$  and  $P[k + 1] \neq P[q]$  // While next character does not match, keep backtracking
7           $k = \pi[k]$                        through the  $\pi$  table
8      if  $P[k + 1] == P[q]$  // If next character matches,
9           $k = k + 1$              increment the matching size  $k$ 
10      $\pi[q] = k$ 
11 return  $\pi$ 
```

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

(a)



(b)