

Fast and Linear-Time String Matching Algorithms Based on the Distances of q -Gram Occurrences

Satoshi Kobayashi, Diptarama Hendrian, Ryo Yoshinaka, Ayumi Shinohara

Graduate School of Information Sciences, Tohoku University, Japan

String matching problem

Input

Text T , Pattern P

Output

All positions i in T such that $T[i : i + |P| - 1] = P$

Example

	1	2	3	4	5	6	7	8	9	10	11	12	13
$T:$	a	b	a	a	b	a	b	b	a	b	b	a	b
$P:$	a	b	b	a									

Output : 6, 9

Naive solution : $O(nm)$

$n = |T|$ $m = |P|$

String matching problem

Input

Text T , Pattern P

Output

All positions i in T such that $T[i : i + |P| - 1] = P$

Example

	1	2	3	4	5	6	7	8	9	10	11	12	13
$T:$	a	b	a	a	b	a	b	b	a	b	b	a	b
$P:$	a	b	b	a									

Output : 6, 9

Naive solution : $O(nm)$

$n = |T|$ $m = |P|$

String matching problem

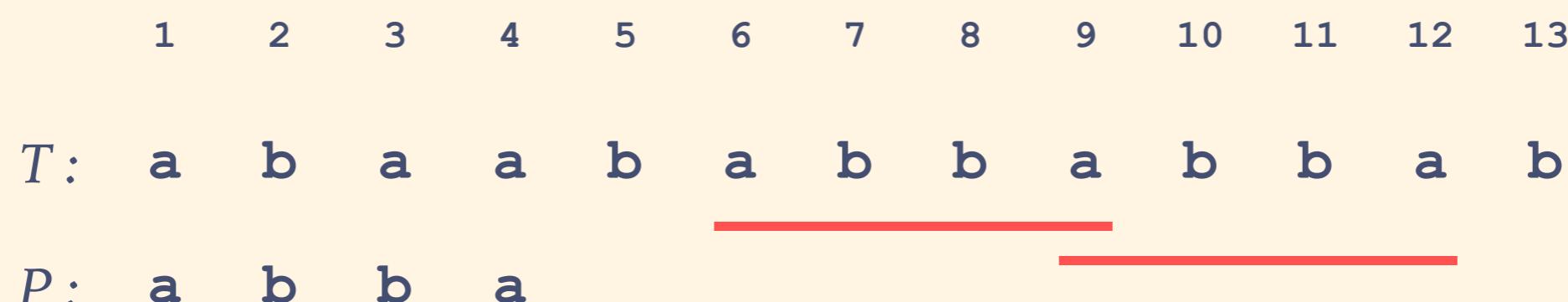
Input

Text T , Pattern P

Output

All positions i in T such that $T[i : i + |P| - 1] = P$

Example



Output : 6, 9

Naive solution : $O(nm)$

$n = |T|$ $m = |P|$

String matching algorithms

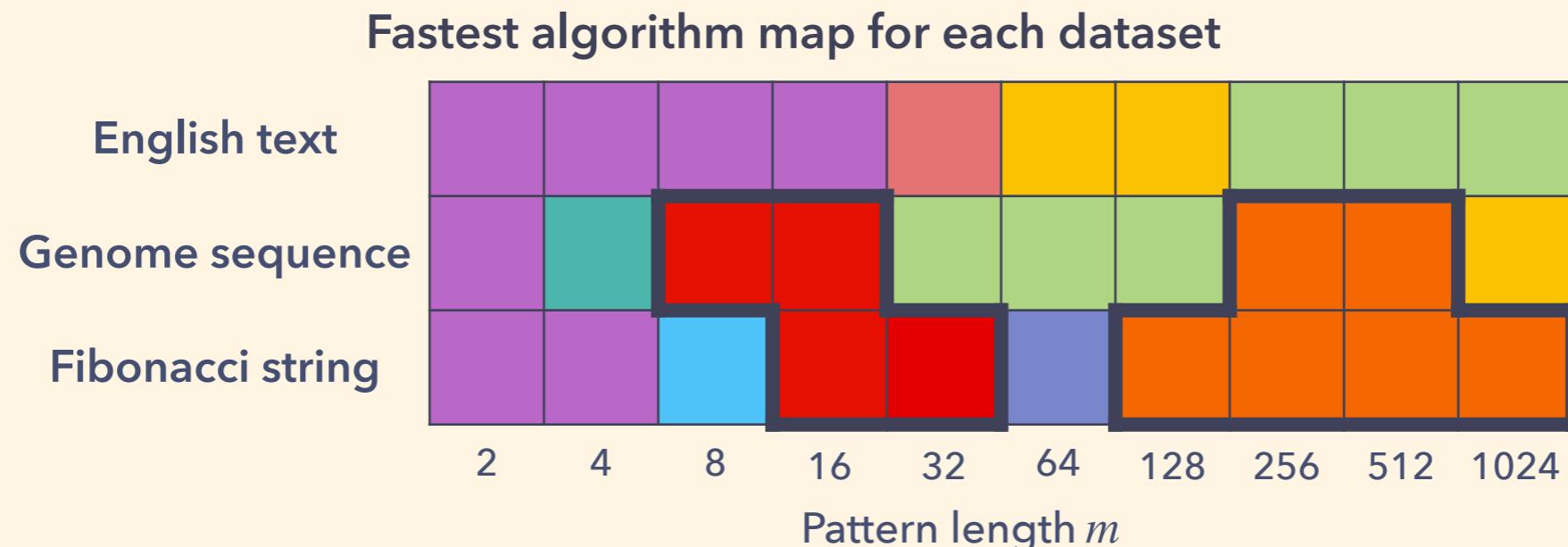
- Knuth-Morris-Pratt (KMP) algorithm [Knuth+, 1977]
 - Preprocessing time : $O(m)$
 - Searching time : $O(n)$
- Boyer-Moore algorithm [Boyer & Moore, 1977]
 - Preprocessing time : $O(m + \sigma)$
 - Searching time : $O(nm)$
 - **Runs fast in practice**

n : Text length
 m : Pattern length
 σ : Alphabet size

Our contributions

$n = |T|$ $m = |P|$ ω : Word length
 $\sigma = |\Sigma|$: Alphabet size q : q -gram

- Propose two string matching algorithms based on the distances of the q -gram occurrences
- Both algorithms work in linear time in the input string size



Comparing 15 powerful algorithms announced from 1977 to 2019 with the proposed algorithms

	Algorithm	Preprocess	Search
●	BNDM q [Navarro & Raffinot, 1998]	$O(m+\sigma)$	$O(nm \lceil m/\omega \rceil)$
●	SBNDM q [Holub & Durian, 2005]	$O(m+\sigma)$	$O(nm \lceil m/\omega \rceil)$
●	FJS [Franek+, 2005]	$O(m+\sigma)$	$O(n)$
●	HASH q [Leqroq, 2007]	$O(mq)$	$O(n(m+q))$
●	BSDM q [Faro & Leqroq, 2012]	$O(m)$	$O(nm)$

	Algorithm	Preprocess	Search
●	WFR q [Cantone+, 2017]	$O(m)$	$O(nm)$
●	LWFR q [Cantone+, 2019]	$O(m)$	$O(n)$
●	DIST q New	$O(mq)$	$O(nq)$
●	LDIST q New	$O(m)$	$O(n)$

Naive solution : $O(nm)$

Existing algorithms

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

$T:$ a b a b a b b b c a a c a c

$P:$ a b a b c

- Match
- ✗ Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

$T:$ a b a b a b b c a a c a c

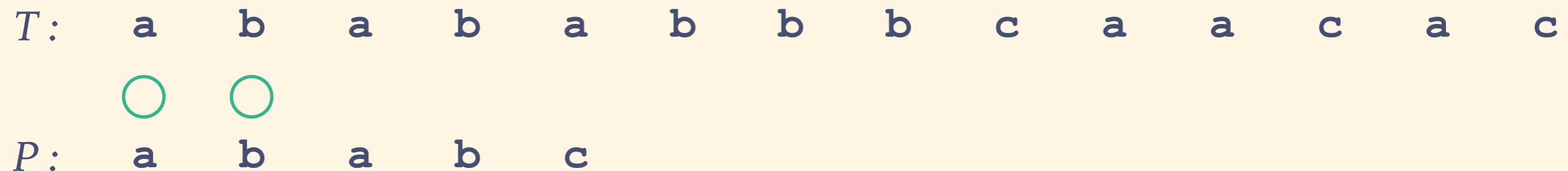
$P:$ a b a b c

- Match
- ✗ Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]



- Match
- ✗ Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

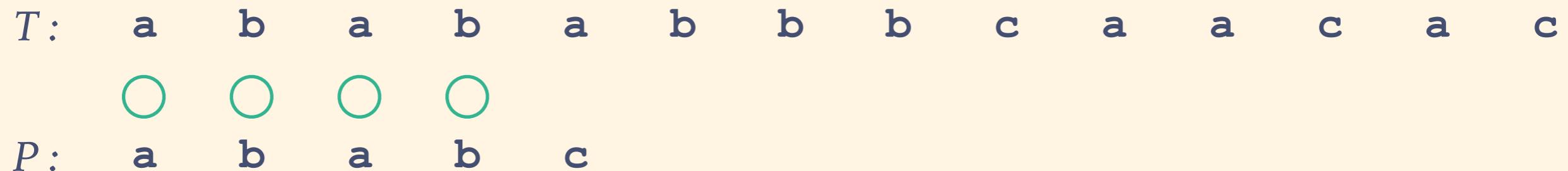


- Match
- ✗ Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

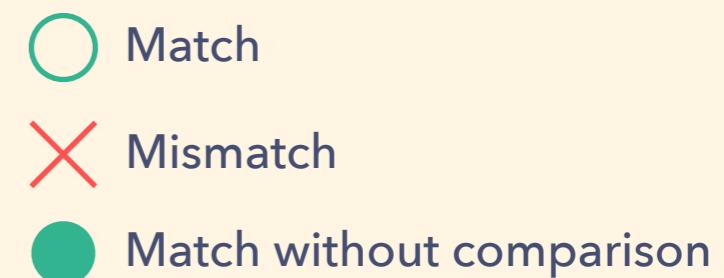
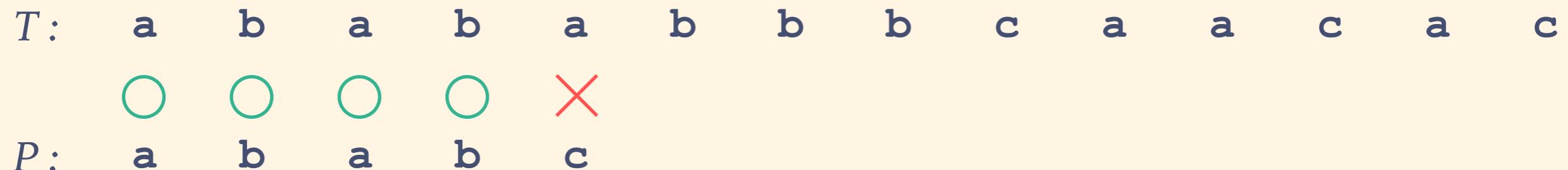


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

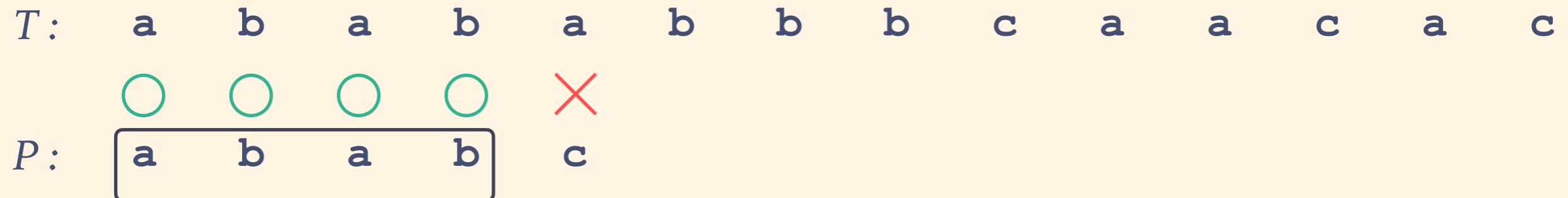
KMP algorithm [Knuth+, 1977]



Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

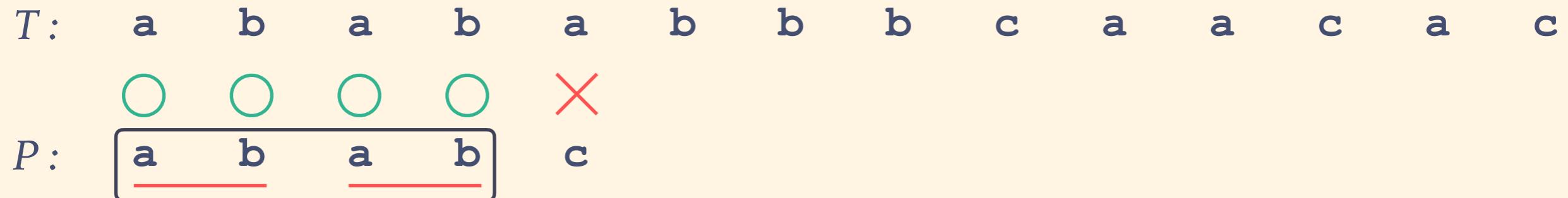


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

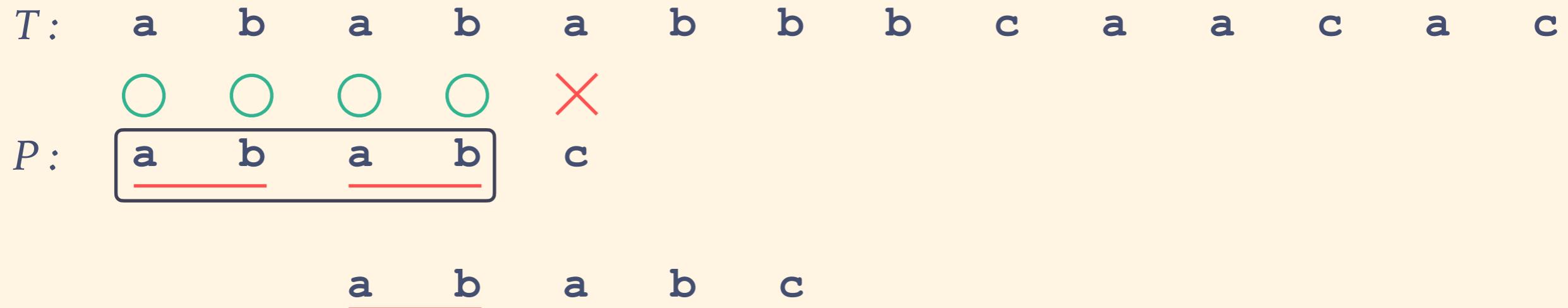


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

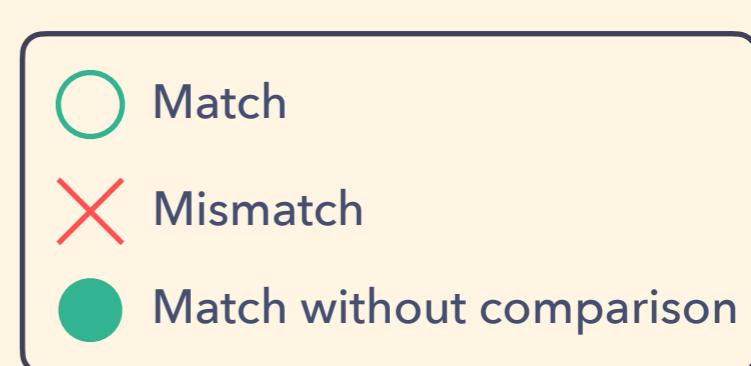
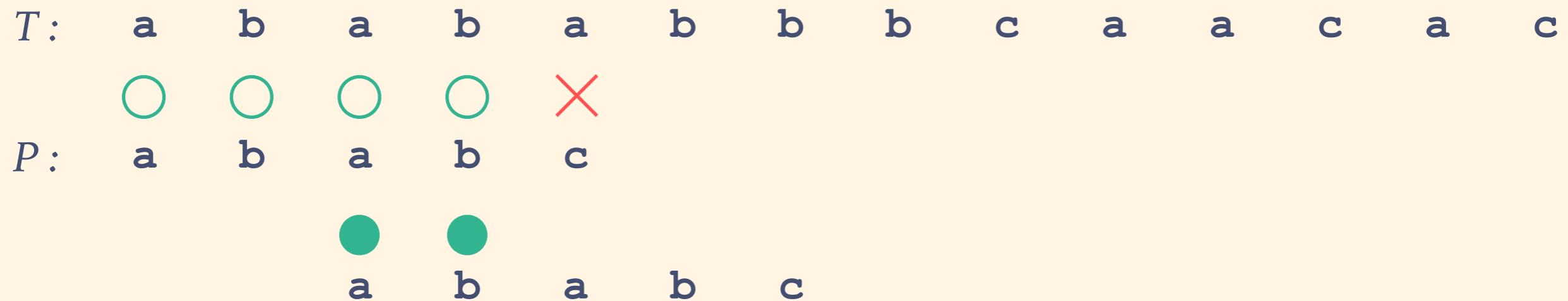


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

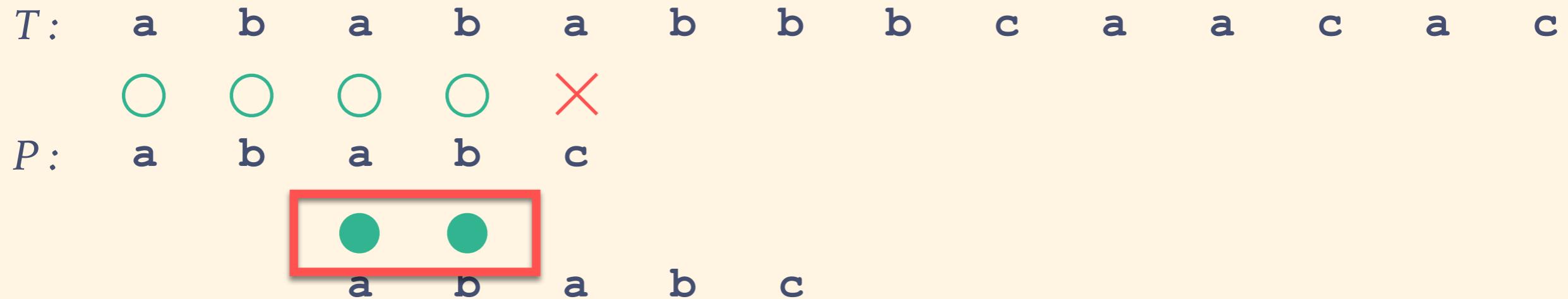
KMP algorithm [Knuth+, 1977]



Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

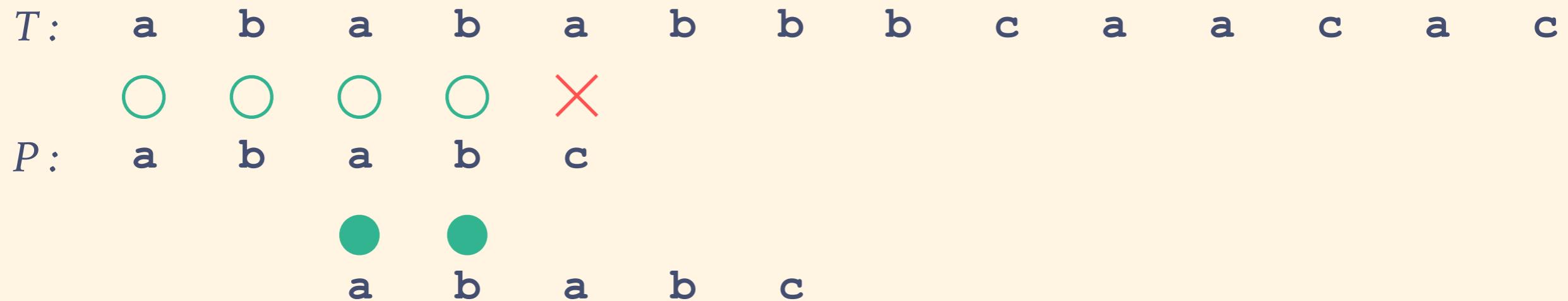


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]

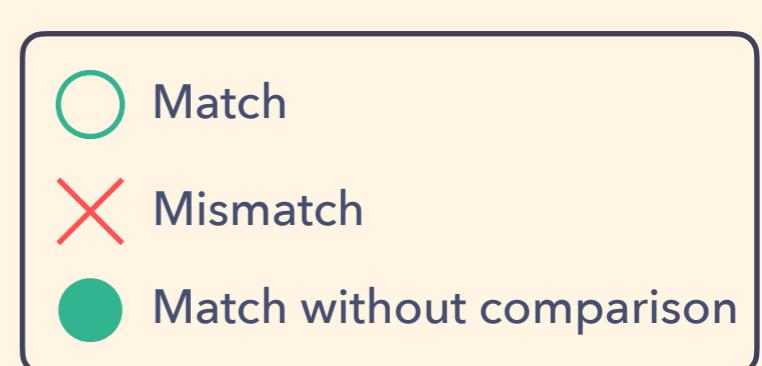
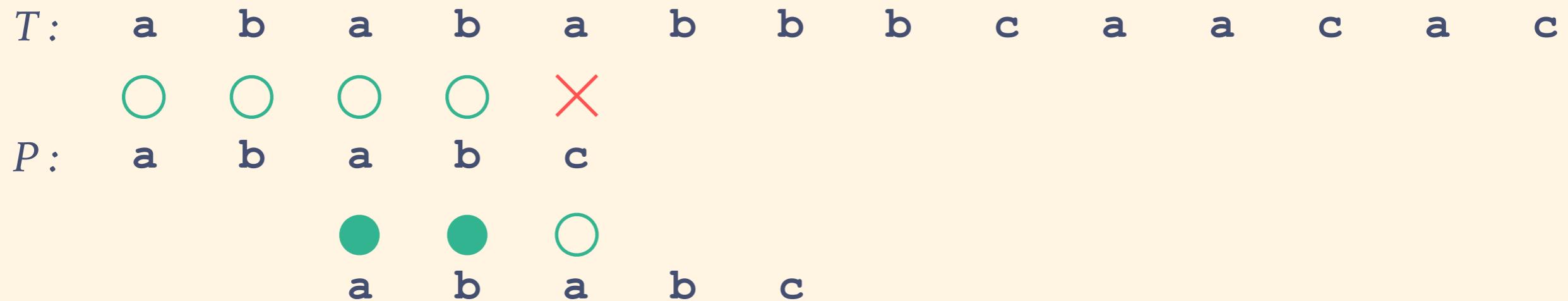


- Match
- Mismatch
- Match without comparison

Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

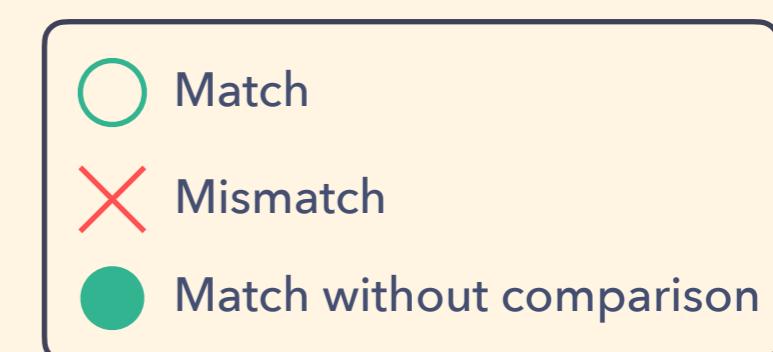
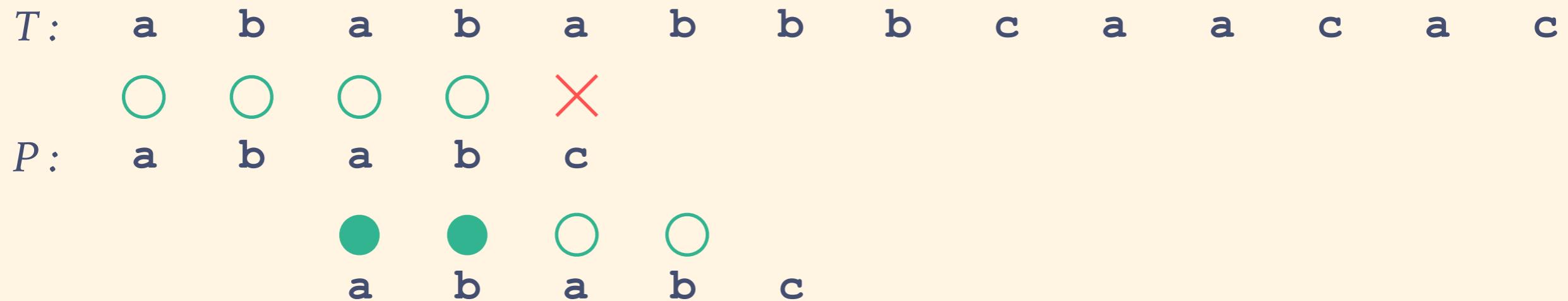
KMP algorithm [Knuth+, 1977]



Preprocessing time : $O(m)$ Searching time : $O(n)$

n : Text length
 m : Pattern length

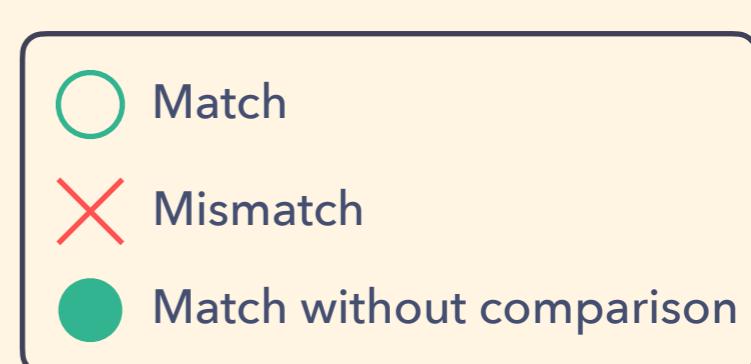
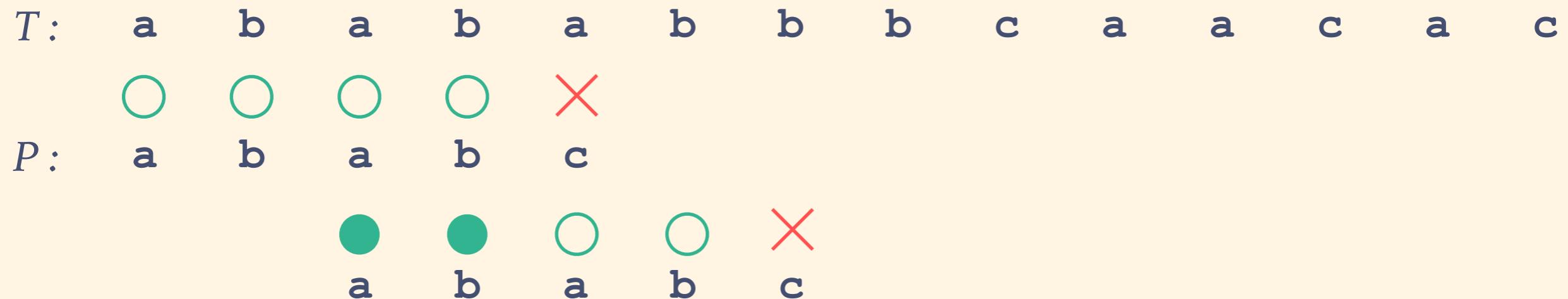
KMP algorithm [Knuth+, 1977]



Preprocessing time : $O(m)$ Searching time : $O(n)$

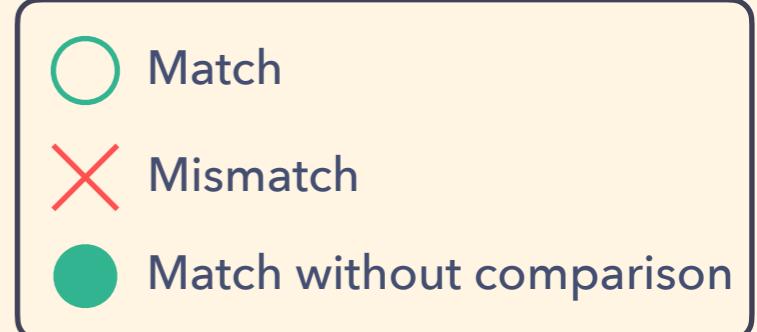
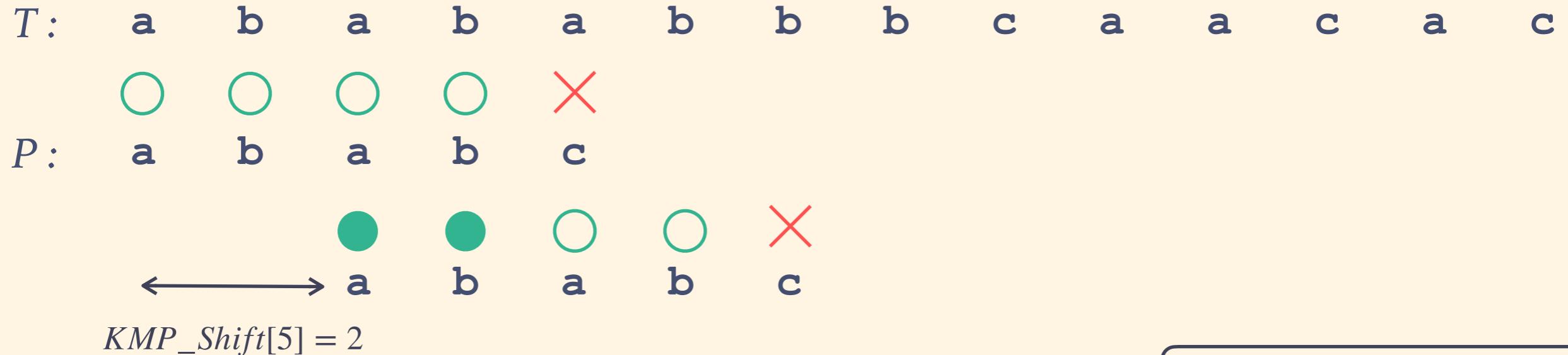
n : Text length
 m : Pattern length

KMP algorithm [Knuth+, 1977]



Preprocessing time : $O(m)$ Searching time : $O(n)$

KMP algorithm [Knuth+, 1977]



$Strong_Bord(j)$

Input : A mismatch position j in the pattern

Output : A maximum value $k(0 \leq k < j)$ that satisfies $P[1 : k] = P[j - k : j - 1]$ and $P[k + 1] \neq P[j]$
 (-1 if no such k exists)

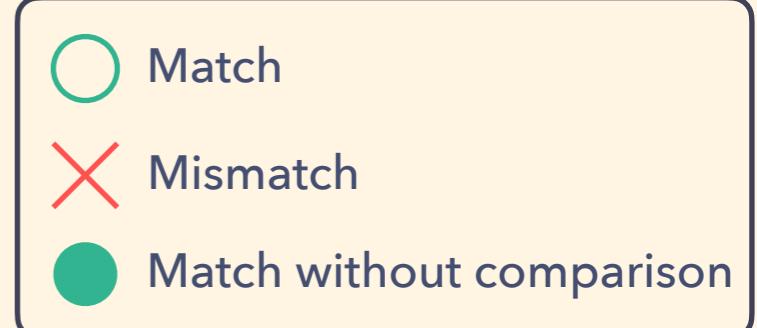
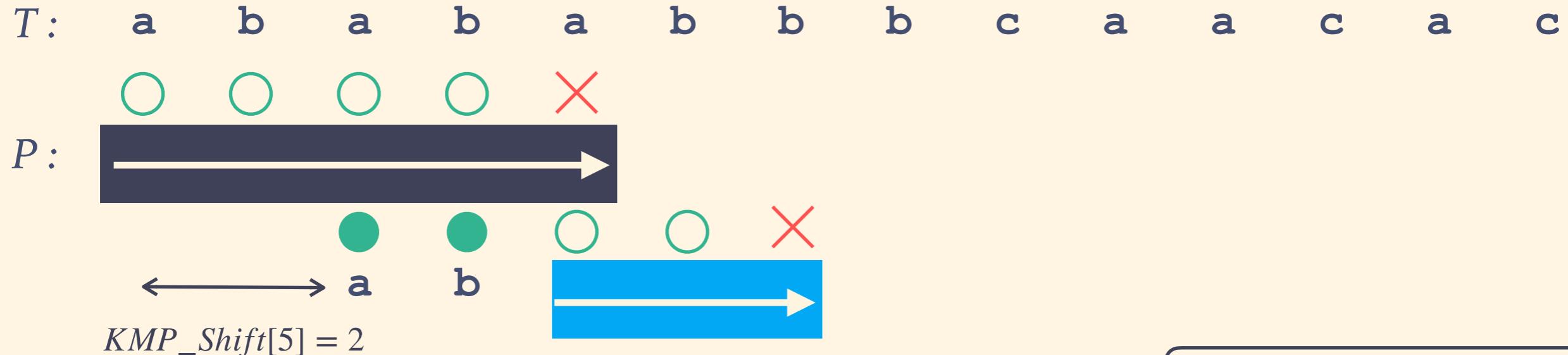
A shift amount when there is a mismatch in the j -th pattern

$$KMP_Shift[j] = j - Strong_Bord(j) - 1$$

j	1	2	3	4	5	6
P	a	b	a	b	c	
KMP_Shift	1	1	3	3	2	5

Preprocessing time : $O(m)$ Searching time : $O(n)$

KMP algorithm [Knuth+, 1977]



$Strong_Bord(j)$

Input : A mismatch position j in the pattern

Output : A maximum value $k(0 \leq k < j)$ that satisfies $P[1 : k] = P[j - k : j - 1]$ and $P[k + 1] \neq P[j]$
 (-1 if no such k exists)

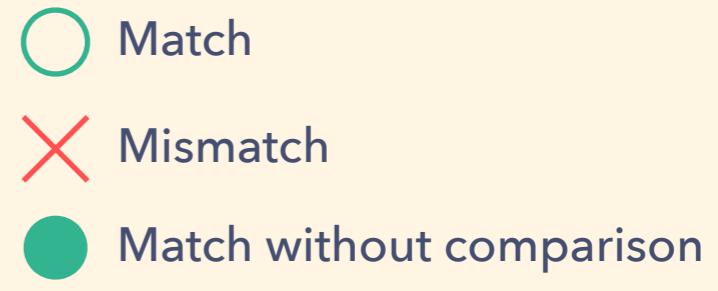
A shift amount when there is a mismatch in the j -th pattern

$$KMP_Shift[j] = j - Strong_Bord(j) - 1$$

j	1	2	3	4	5	6
P	a	b	a	b	c	
KMP_Shift	1	1	3	3	2	5

Preprocessing time : $O(m)$ Searching time : $O(n)$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

↑
 $m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

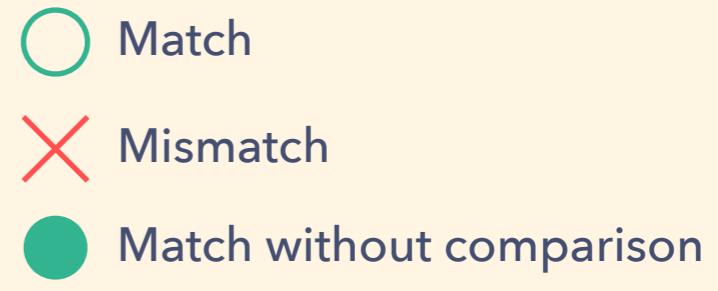
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

↑
 $m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

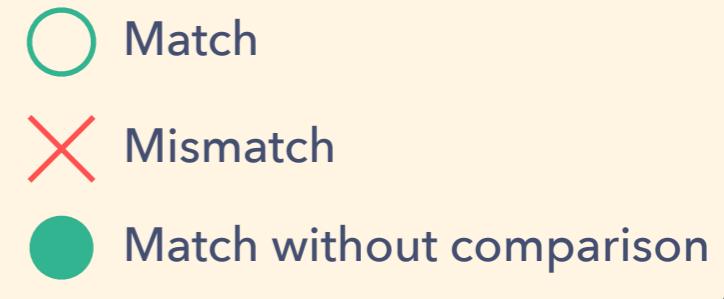
n : Text length

m : Pattern length

σ : Alphabet size

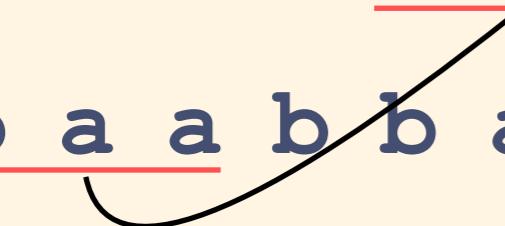
Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b



$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

↑
 $m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

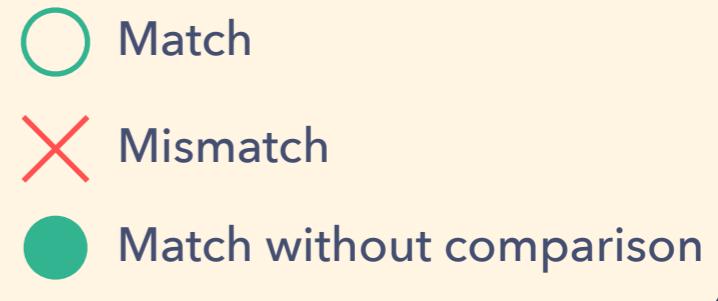
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

$\xleftarrow{\text{shift}[h(baa)] = 4}$ a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

$m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

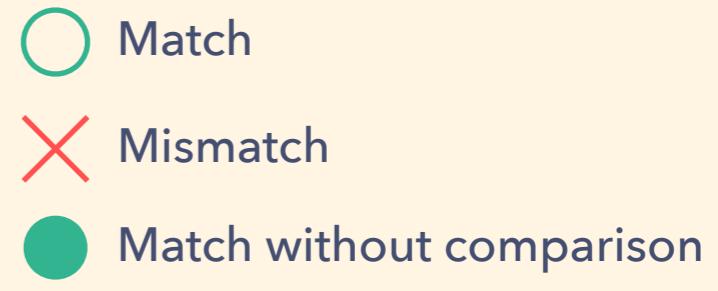
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

$$m - q + 1$$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

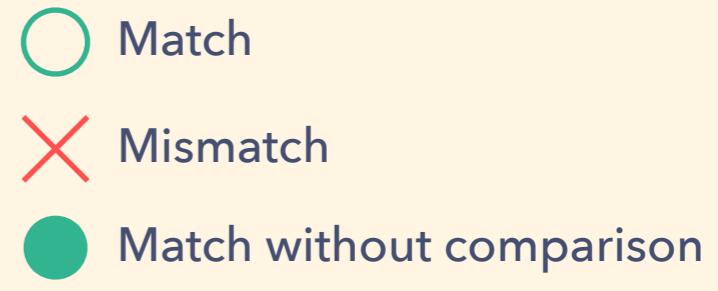
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

↑
 $m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

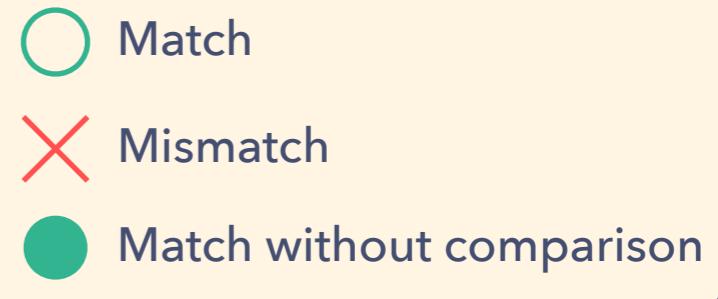
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b

a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

$m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

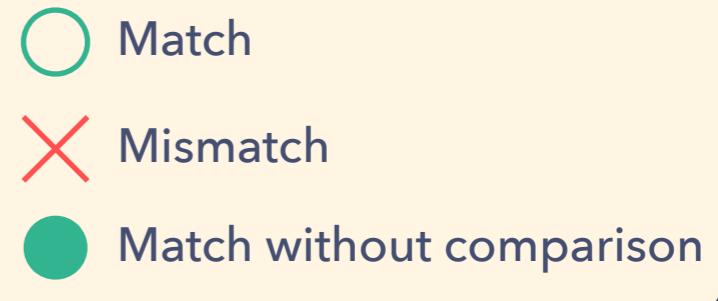
n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

HASH q algorithm [Leqroq, 2007]



T : a a a b a b a a b b a b a b b

P : a b a a b b a b
 ○ ○ ○ ○ ○ ○ ○
 a b a a b b a b

$P = a b a a b b a b$

x	$h(x)$	$Shift[h(x)]$
aba	681	5
baa	683	4
aab	680	3
abb	682	2
bba	685	1
bab	684	0
Others	-	6

$$shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

↑
 $m - q + 1$

- Determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (2^{q-1} \cdot x[1] + 2^{q-2} \cdot x[2] + \dots + 2 \cdot x[q-1] + x[q]) \bmod 2^8$$

x : String

(Treat characters as the ASCII code)

n : Text length

m : Pattern length

σ : Alphabet size

Preprocessing time : $O(mq)$ Searching time : $O(n(m + q))$

Proposed 1 DIST q algorithm

Idea of DIST q algorithm

Practically fast

Proposed

Linear time

$shift$ (HASH q) + q -gram distance array + KMP algorithm

- q -gram distance array

$$dist[i] = \min(\{ j \mid h(P[i - j - q + 1 : i - j]) = h(P[i - q + 1 : i]), q - 1 \leq j < i \} \cup \{i - q + 1\})$$

- A hash value is used to determine the equivalence of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16} \quad x : \text{String}$$

When $q = 3$									
i	1	2	3	4	5	6	7	8	9
P	a	b	a	a	b	<u>b</u>	a	a	a
$dist$	-	-	1	2	3	4	5	4	7

T: b b a b a b b a a b b b a b a a
P: a b a a b b a a a

Idea of DIST q algorithm

Practically fast

Proposed

Linear time

$shift$ (HASH q) + q -gram distance array + KMP algorithm

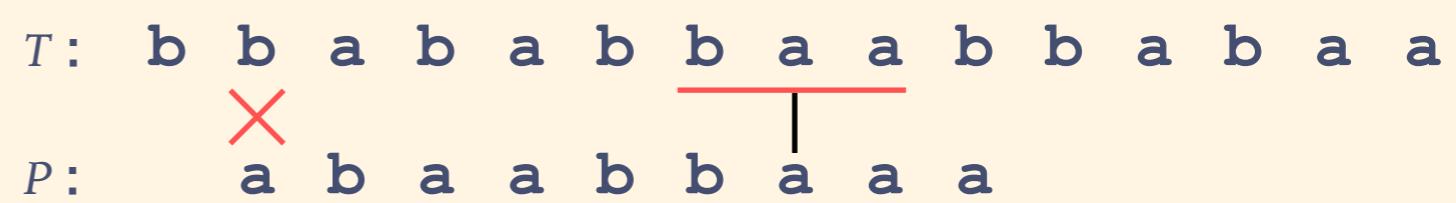
- q -gram distance array

$$dist[i] = \min(\{ j \mid h(P[i - j - q + 1 : i - j]) = h(P[i - q + 1 : i]), q - 1 \leq j < i \} \cup \{i - q + 1\})$$

- A hash value is used to determine the equivalence of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16} \quad x : \text{String}$$

When $q = 3$									
i	1	2	3	4	5	6	7	8	9
P	a	b	a	a	b	<u>b</u>	a	a	a
$dist$	-	-	1	2	3	4	5	4	7



Idea of DIST q algorithm

Practically fast

Proposed

Linear time

$shift$ (HASH q) + q -gram distance array + KMP algorithm

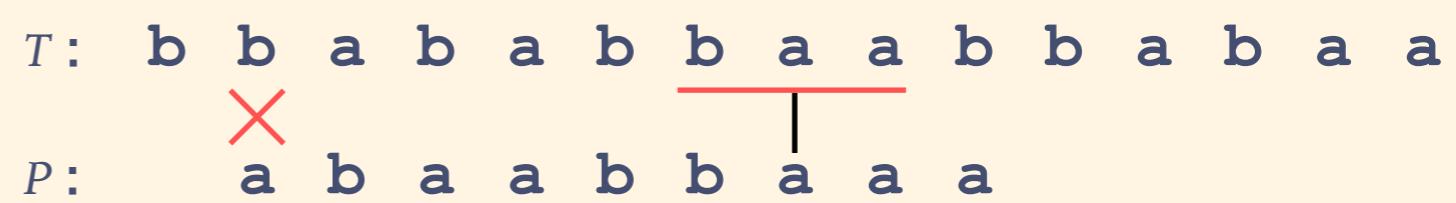
- q -gram distance array

$$dist[i] = \min(\{ j \mid h(P[i - j - q + 1 : i - j]) = h(P[i - q + 1 : i]), q - 1 \leq j < i \} \cup \{i - q + 1\})$$

- A hash value is used to determine the equivalence of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16} \quad x : \text{String}$$

When $q = 3$									
i	1	2	3	4	5	6	7	8	9
P	a	b	a	a	b	<u>b</u>	a	a	a
$dist$	-	-	1	2	3	4	5	4	7



Idea of DIST q algorithm

Practically fast

Proposed

Linear time

$shift$ (HASH q) + q -gram distance array + KMP algorithm

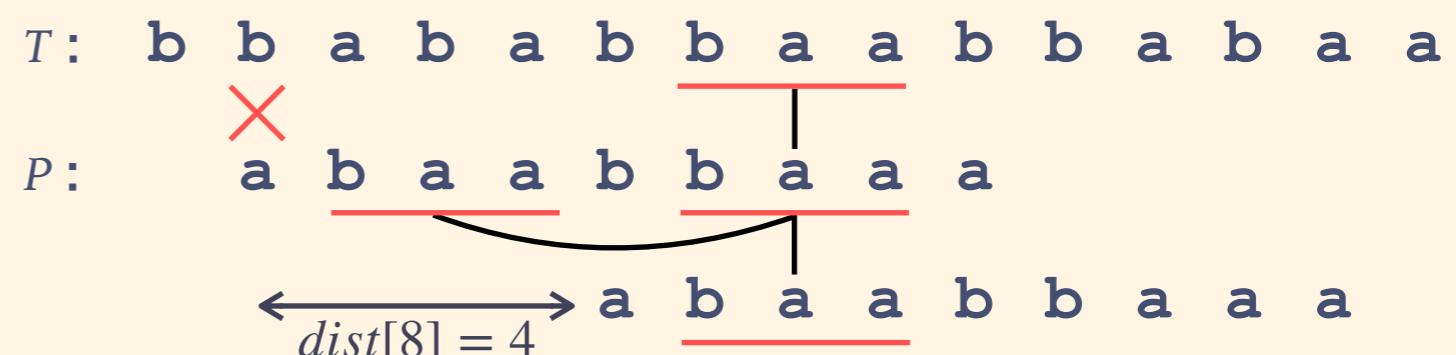
- q -gram distance array

$$dist[i] = \min(\{ j \mid h(P[i - j - q + 1 : i - j]) = h(P[i - q + 1 : i]), q - 1 \leq j < i \} \cup \{i - q + 1\})$$

- A hash value is used to determine the equivalence of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16} \quad x : \text{String}$$

When $q = 3$									
i	1	2	3	4	5	6	7	8	9
P	a	b	a	a	b	<u>b</u>	a	a	a
$dist$	-	-	1	2	3	4	5	4	7



Array HQ_Shift (Almost same as shift of HASH q)

$$HQ_shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

- Also determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16}$$

$P = a b a a b b a b$

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	5
baa	2053	4
aab	2038	3
abb	2042	2
bba	2057	1
bab	2054	0
Others	-	6

$T:$ **b b a a a b a b a a b a b b a**
 $P:$ **a b a a b b a b**

$m - q + 1$

Use this shift to align the q -gram in the pattern and the q -gram in the text which has the same hash value

Array HQ_Shift (Almost same as shift of HASH q)

$$HQ_shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

- Also determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16}$$

$P = a b a a b b a b$

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	5
baa	2053	4
aab	2038	3
abb	2042	2
bba	2057	1
bab	2054	0
Others	-	6

$T:$ **b b a a a b a b a a b a b b a**

$P:$ **a b a a b b a b**

$m - q + 1$

Use this shift to align the q -gram in the pattern and the q -gram in the text which has the same hash value

Array HQ_Shift (Almost same as shift of HASH q)

$$HQ_shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

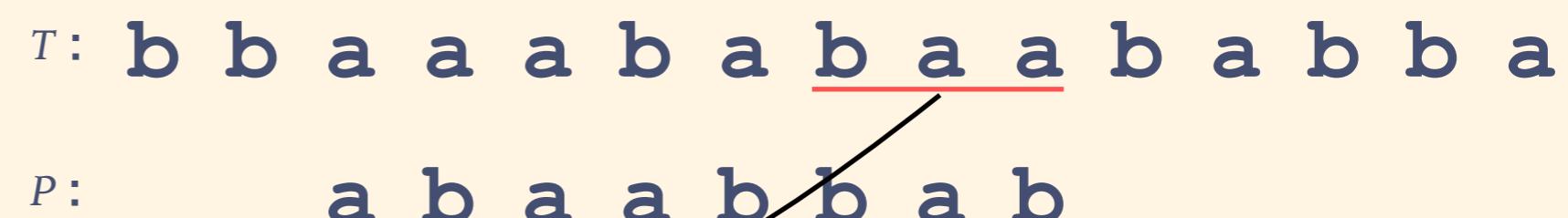
- Also determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16}$$

$P = a b a a b b a b$

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	5
baa	2053	4
aab	2038	3
abb	2042	2
bba	2057	1
bab	2054	0
Others	-	6

$m - q + 1$



Use this shift to align the q -gram in the pattern and the q -gram in the text which has the same hash value

Array HQ_Shift (Almost same as shift of HASH q)

$$HQ_shift[h(x)] = m - \max(\{j \mid h(P[j - q + 1 : j]) = h(x), q \leq j \leq m\} \cup \{q - 1\})$$

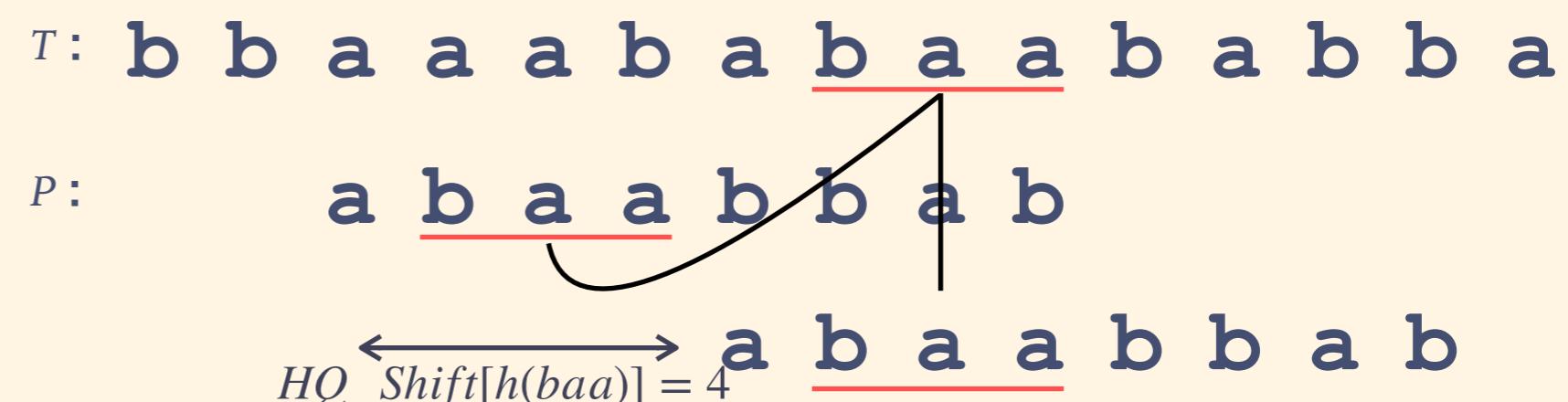
- Also determines the equivalence of q -grams using the hash value of q -grams

$$h(x) = (4^{q-1} \cdot x[1] + 4^{q-2} \cdot x[2] + \dots + 4 \cdot x[q-1] + x[q]) \bmod 2^{16}$$

$P = a b a a b b a b$

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	5
baa	2053	4
aab	2038	3
abb	2042	2
bba	2057	1
bab	2054	0
Others	-	6

$m - q + 1$



Use this shift to align the q -gram in the pattern and the q -gram in the text which has the same hash value

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

T : a b b a a b b a a b a b b a b b a a a b a a b a a b b a a a a

P : a b a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text

$T: a \ b \ b \ a \ a \ b \ \underline{b \ a \ a} \ b \ a \ b \ b \ a \ b \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

Match

Mismatch

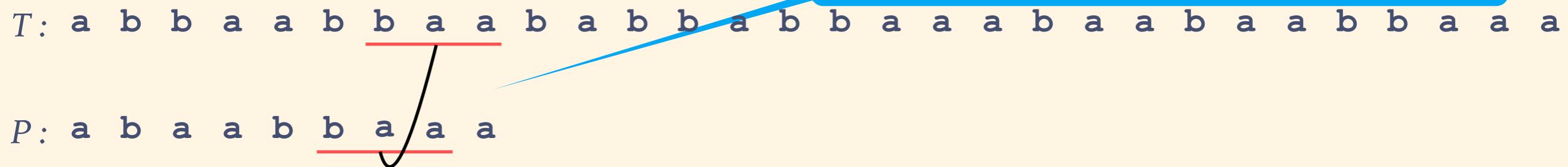
Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



- Match
- ✗ Mismatch
- Match without comparison

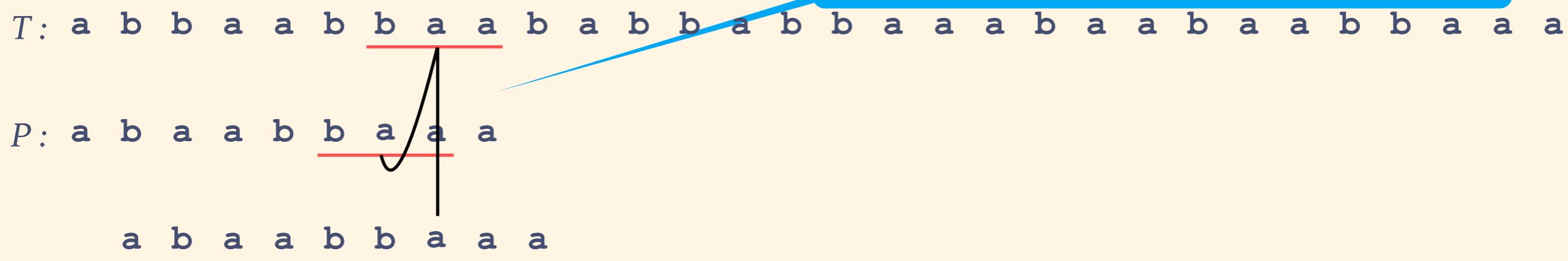
Searching

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4
abb	2042	3
bba	2054	2
aaa	2037	0
others	-	7

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



- Match
- ✗ Mismatch
- Match without comparison

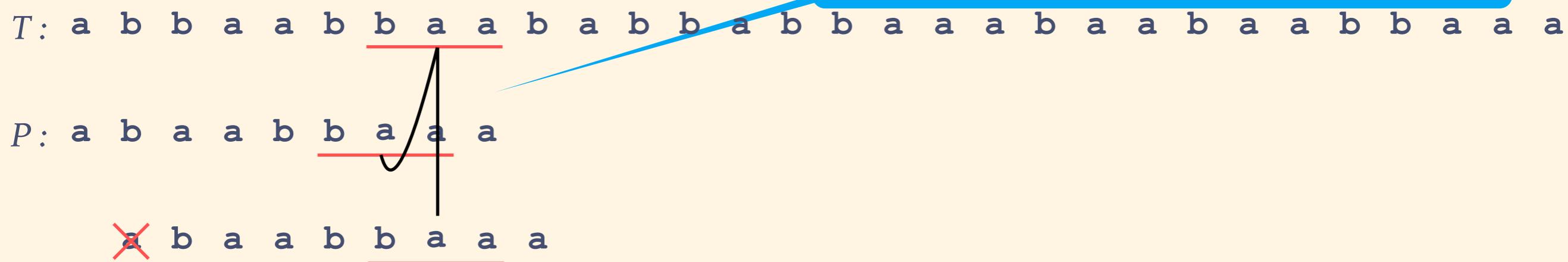
Searching

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4
abb	2042	3
bba	2054	2
aaa	2037	0
others	-	7

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



- Match
- ✗ Mismatch
- Match without comparison

Searching

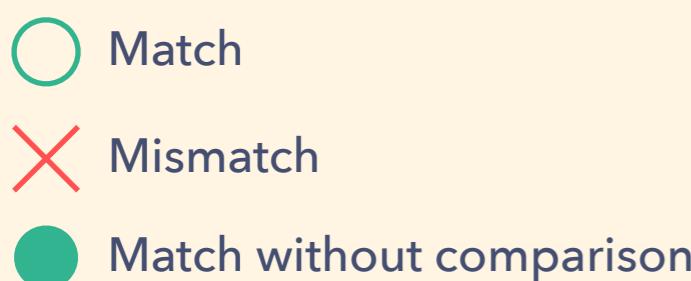
x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Shift the pattern using the distance array $dist$ if the first letter do not match

Alignment-Phase

Align “baa” by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



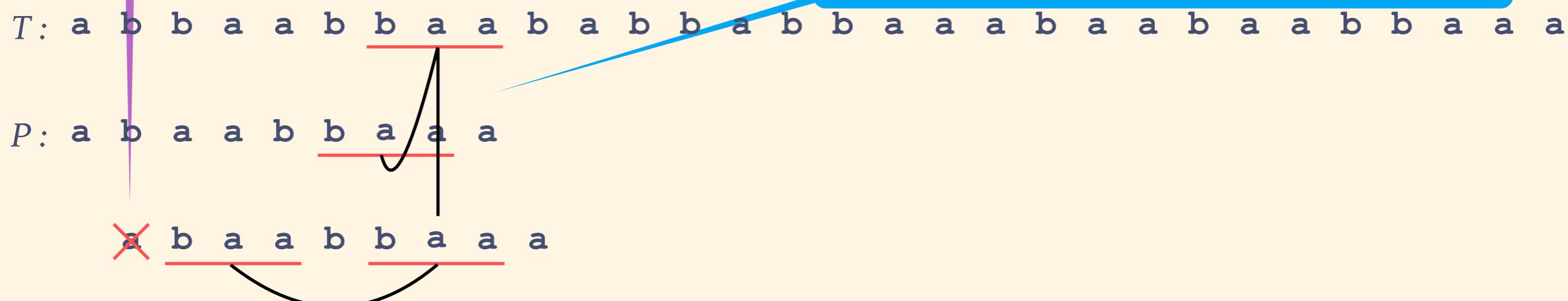
Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8

Shift the pattern using the distance array $dist$ if the first letter do not match

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



- Match
- ✗ Mismatch
- Match without comparison

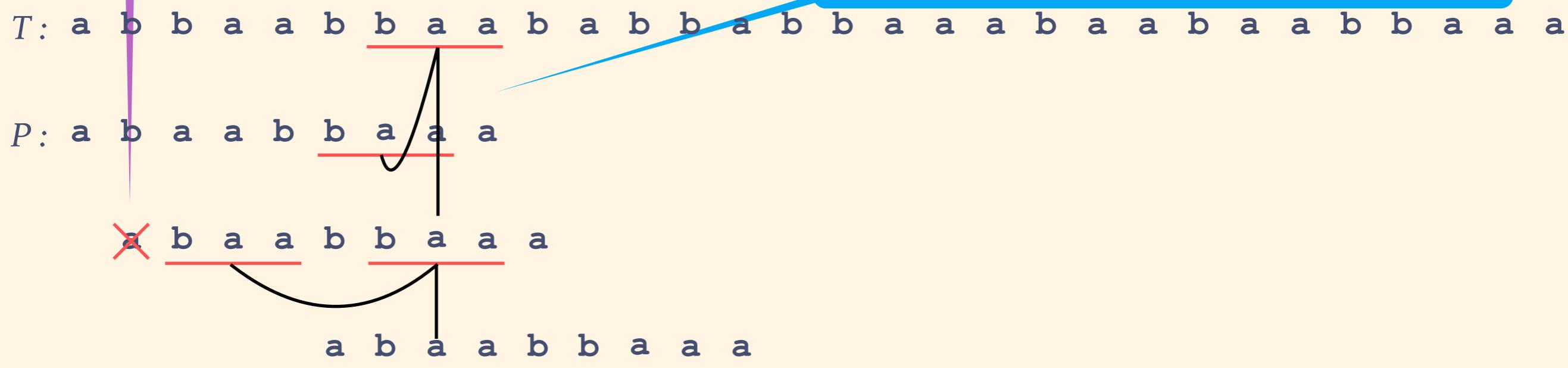
Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8

Shift the pattern using the distance array $dist$ if the first letter do not match

Alignment-Phase

Align "baa" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text



Match



Mismatch



Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "bba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

Match

Mismatch

Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "bba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

Match

Mismatch

Match without comparison

Searching

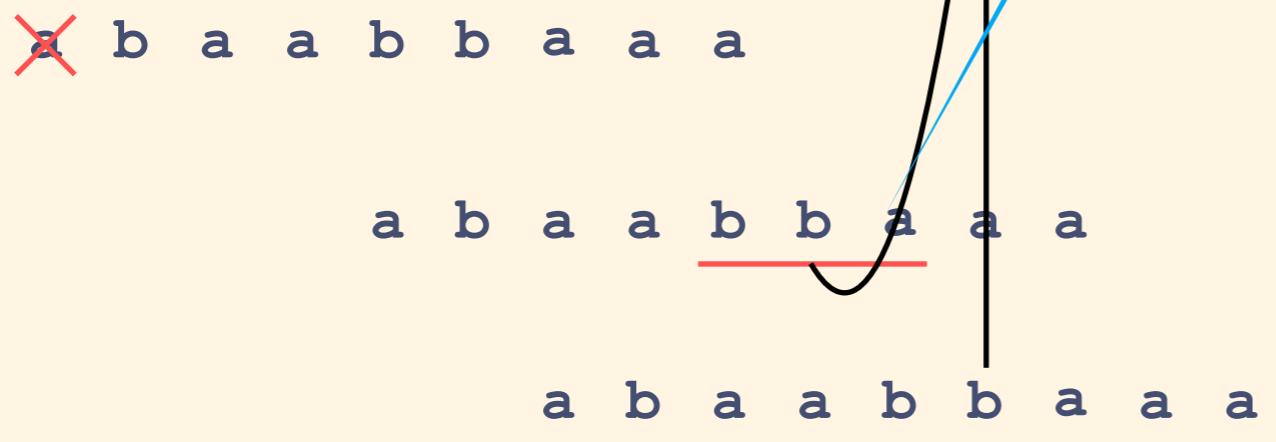
x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "bba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$



- Match
- ✗ Mismatch
- Match without comparison

Searching

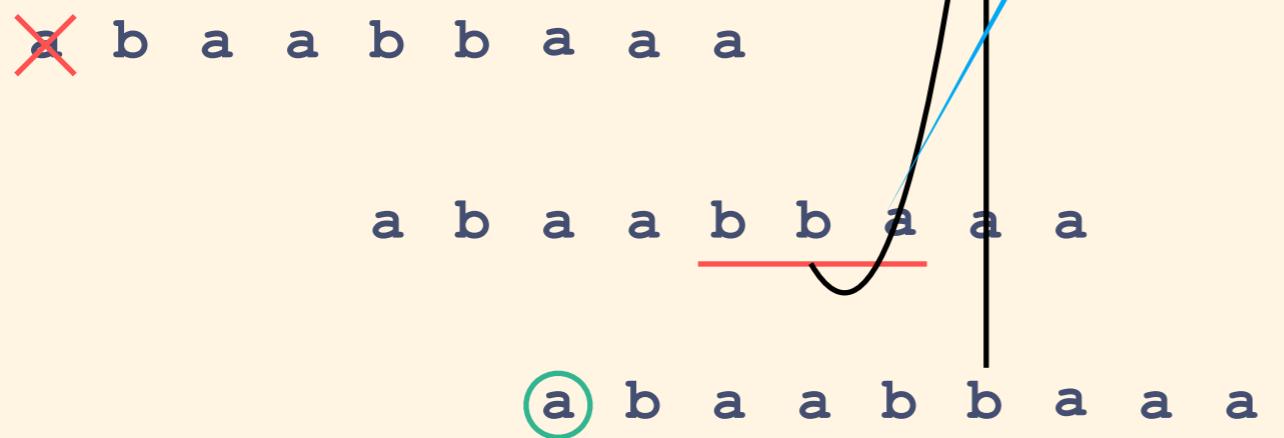
x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "bba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$



Match

Mismatch

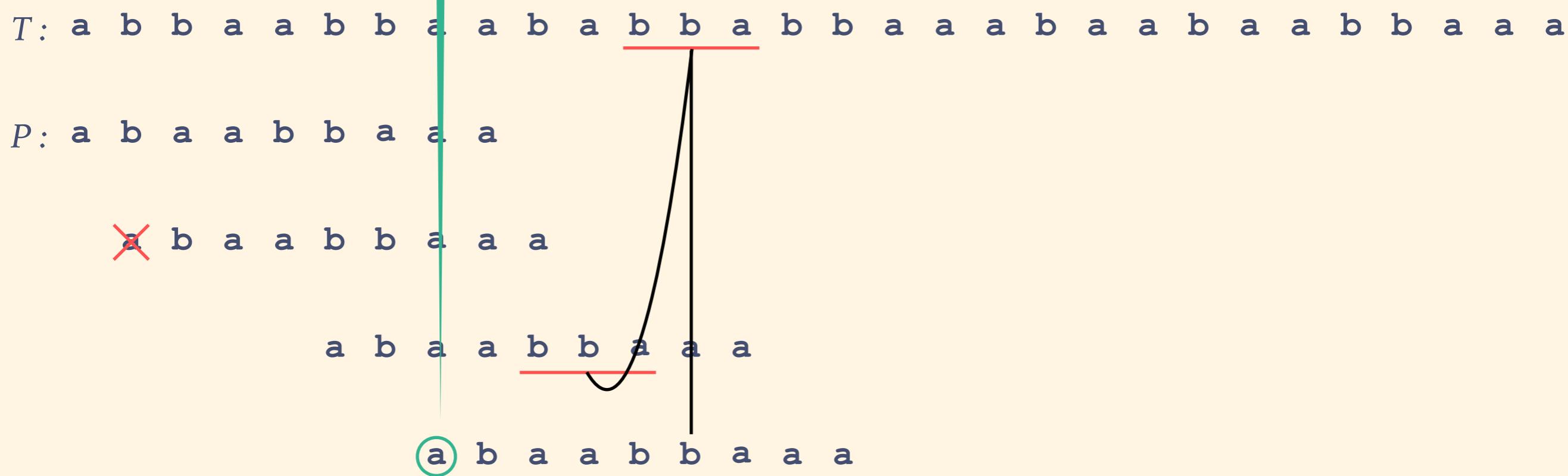
Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8

Comparison-Phase

compare $P[2:m]$ from left to right if the first letter match



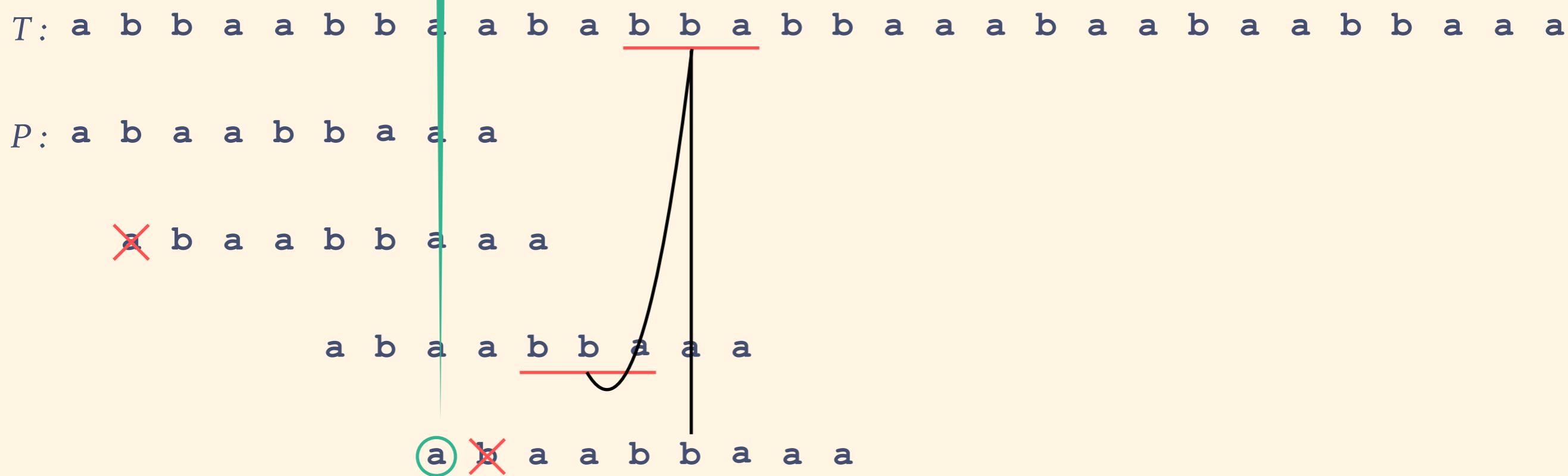
- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8

Comparison-Phase

compare $P[2:m]$ from left to right if the first letter match



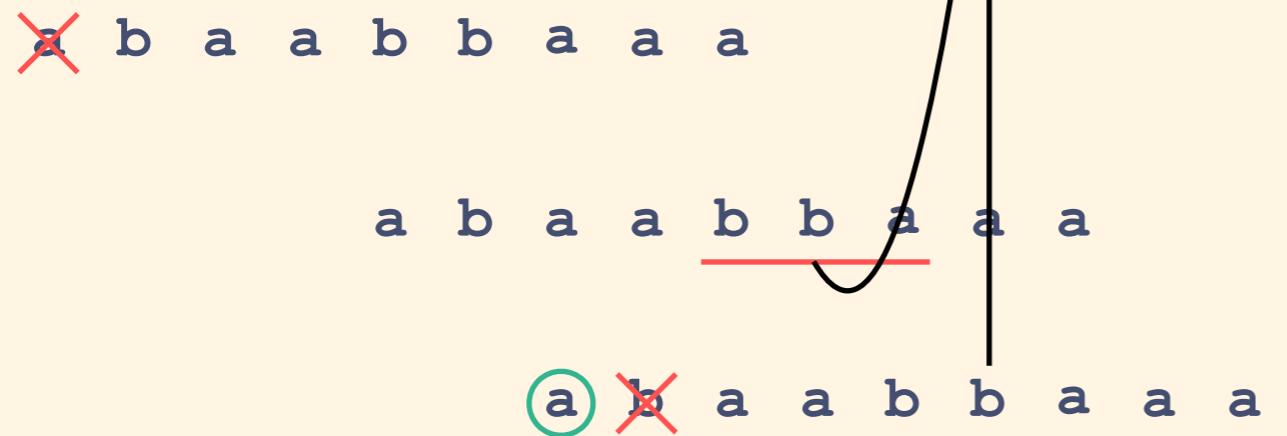
- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$



- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

Comparison-Phase

Select the one where the resumption position of the character comparison goes further to the right
 $KMP_Shift[2] = 1$

$dist[7] = 5$

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a

(a) ~~b~~ a a b b a a

Comparison-Phase

Select the one where the resumption position of the character comparison goes further to the right
 $KMP_Shift[2] = 1$

$dist[7] = 5$

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a

(a) ~~b~~ a a b b a a

a b a a b b a a

Comparison-Phase

Select the one where the resumption position of the character comparison goes further to the right
 $KMP_Shift[2] = 1$

$dist[7] = 5$

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
others	-	7											

Alignment-Phase

Align "aba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text character

$T: a b b a a b b a a b a b b a b b a a a a$

$P: a b a a b b a a a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

Match

Mismatch

Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4
abb	2042	3
bba	2054	2
aaa	2037	0
others	-	7

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	a
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Alignment-Phase

Align "aba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text character

T : a b b a a b b a a b a b b a b b a a b a a b b a a b b a a a

P : a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

Match

Mismatch

Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4
abb	2042	3
bba	2054	2
aaa	2037	0
others	-	7

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	a
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Alignment-Phase

Align "aba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text character

$T: a b b a a b b a a b a b b a b b a a b a a b b a a b b a a a$

$P: a b a a b b a a a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

a b a a b b a a a

Match

Mismatch

Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$
aba	2041	6
baa	2053	1
aab	2038	4
abb	2042	3
bba	2054	2
aaa	2037	0
others	-	7

j	1	2	3	4	5	6	7	8	9	10
P	a	b	a	a	b	b	a	a	a	a
$dist$	-	-	1	2	3	4	5	4	7	-
KMP_Shift	1	1	3	2	4	3	7	6	7	8

Alignment-Phase

Align "aba" by shifting with array HQ_Shift until $P[1]$ matches the corresponding text character

$T: a b b a a b b a a b a b b a b b a a b a a b a a b b a a a a$

$P: a b a a b b a a a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) b a a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) b a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

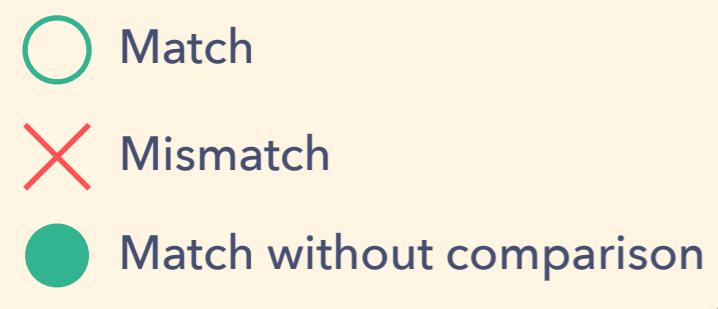
~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) a b b a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

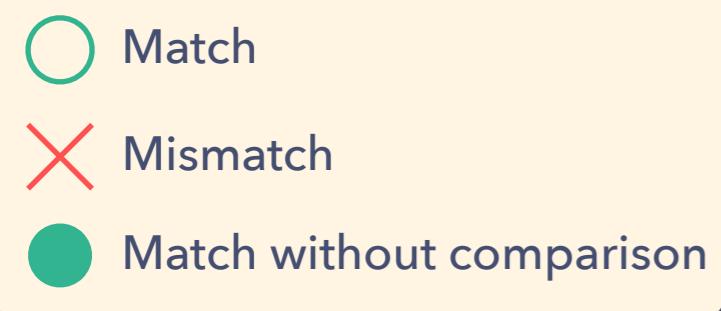
~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) b b a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

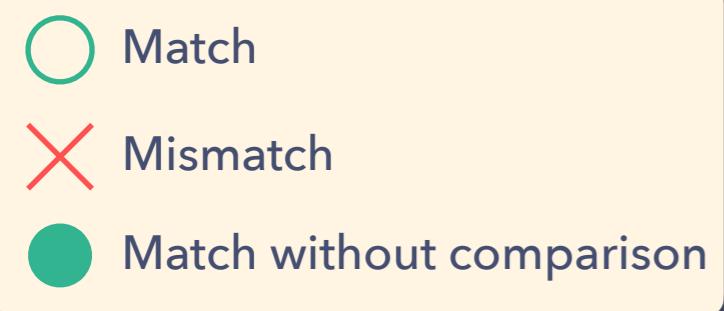
~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) b a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

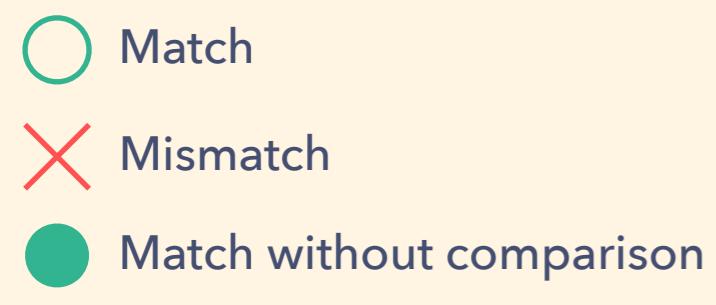
~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~a~~ a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T: a \ b \ b \ a \ a \ b \ b \ a \ a \ b \ a \ b \ b \ a \ a \ a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ a$

$P: a \ b \ a \ a \ b \ b \ a \ a \ a$

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a

a b a a b b a a

(a) (b) (a) (a) (b) ~~x~~ a a a

Comparison-Phase

Select the one where the resumption position of the character comparison goes further to the right
 $KMP_Shift[6] = 3$

$$dist[3] = 1$$

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

X b a a b b a a a
 a b a a b b a a a
 (a) X a a b b a a
 a b a a b b a a a

Comparison-Phase

Select the one where the resumption position of the character comparison goes further to the right

$$KMP_Shift[6] = 3$$

$$dist[3] = 1$$

(a) (b) (a) (a) (b) X a a a

() () a a b b a a a

- Match
- X Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () a a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

T : a b b a a b b a a b a b b a b b a a a b a a b b a a a a

P : a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) a b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) b b a a a

- Match
- ✗ Mismatch
- Match without comparison

Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

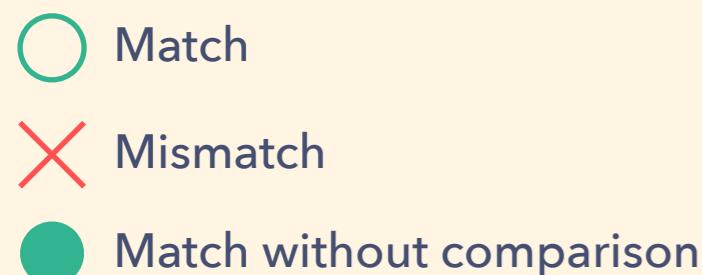
a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) (b) b a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

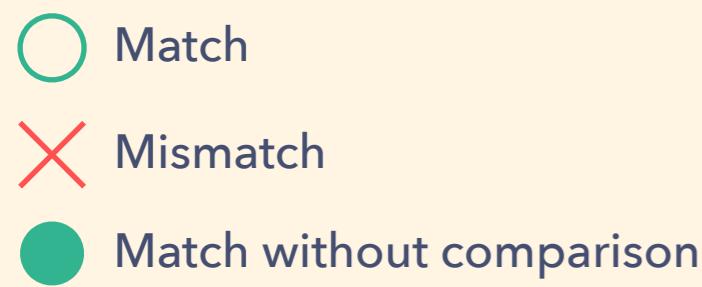
a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) (b) (b) a a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

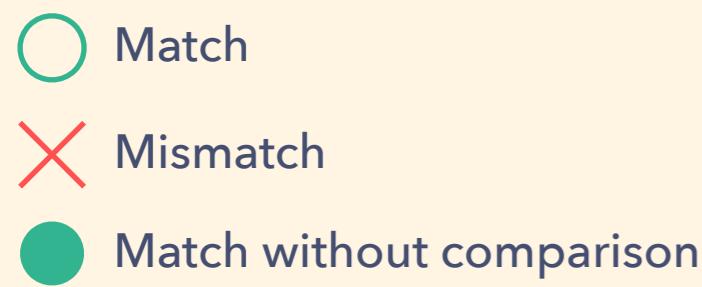
a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) (b) (b) (a) a a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

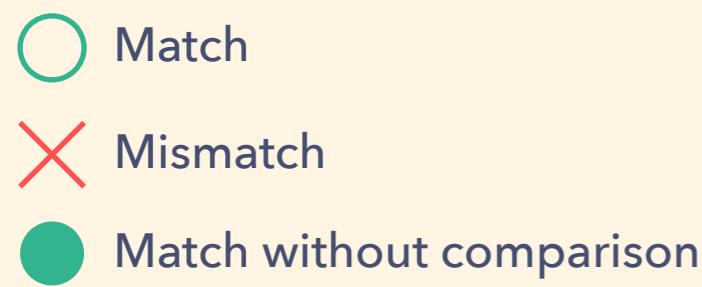
a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) (b) (b) (a) (a) a



Searching

x	$h(x)$	$HQ_Shift[h(x)]$	j	1	2	3	4	5	6	7	8	9	10
aba	2041	6	P	a	b	a	a	b	b	a	a	a	
baa	2053	1	$dist$	-	-	1	2	3	4	5	4	7	-
aab	2038	4	KMP_Shift	1	1	3	2	4	3	7	6	7	8
abb	2042	3											
bba	2054	2											
aaa	2037	0											
Others	-	7											

$T:$ a b b a a b b a a b a b b a b b a a a b a a b b a a a

$P:$ a b a a b b a a a

~~a~~ b a a b b a a a

a b a a b b a a a

(a) ~~b~~ a a b b a a a

a b a a b b a a a

(a) (b) (a) (a) (b) ~~x~~ a a a

() () (a) (a) (b) (b) (a) (a) (a)

- Match
- ✗ Mismatch
- Match without comparison

Time complexity of DIST q algorithm

- Preprocessing

- Array KMP_Shift : $O(m)$
- Array HQ_Shift : $O(mq)$
 - Calculate hash value that takes $O(q)$ time at $m - q + 1$ positions
- Array $dist$: $O(mq)$

$O(mq)$

- Searching

- Number of character comparisons : $O(n)$
- Calculate hash value at maximum $n - m + 1$ positions : $O(nq)$

$O(nq)$

$T :$	a	a	a	a	a	a	a	a	a
$P :$	b	a	a	a	a	a			
		x							
			b	a	a	a	a		
				x					
			b	a	a	a	a	a	
				x					
			b	a	a	a	a	a	a

Proposed 2 Linear-DIST_{*q*} (LDIST_{*q*}) algorithm

LDIST q algorithm

- Worst-case time complexity of the search phase of DISTq algorithm is $\Theta(nq)$

$T:$	a	a	a	a	a	a	a	a	a
$P:$	b	a	a	a	a	a			
		b	a	a	a	a	a		
		b	a	a	a	a	a		
		b	a	a	a	a	a	a	

- Since the hash function h used in DIST q algorithm is a rolling hash, when $h(T[i : j])$ has already been obtained, the hash value of $h(T[i + 1 : j + 1])$ can be computed in $O(1)$ time ($1 \leq i \leq j < |T|$)

$$h(T[i + 1 : j + 1]) = (4 \cdot (h(T[i : j]) - 16 \cdot T[i]) + T[j + 1]) \bmod 2^{16}$$

$$h(x[1 : 3]) = (16 \cdot x[1] + 4 \cdot x[2] + x[3]) \bmod 2^{16}$$

- Searching time can be reduced to $O(n)$

- by incrementally calculating the hash value of the q -gram using the previously calculated value of the other q -gram
- Preprocessing time is also reduced to $O(m)$

Experiments

Datasets

- English text
- Genome sequence
- Fibonacci string
- Texts with frequent pattern occurrences

- Implemented with C language
- Compiled with GCC9.2.0
- MacBook Pro (13-inch, 2018), macOS Catalina, Intel Core i7 2.7GHz quad core, 16GB memory

English text

- Use The King James version of the Bible as text
 - $n = 4017009$
 - $\sigma = 62$
- Patterns are randomly extracted from text

$n = |T|$ $m = |P|$ $\sigma = |\Sigma|$: Alphabet size

m	2	4	8	16	32	64	128	256	512	1024
BNDM q	140.48 ⁽²⁾	93.39 ⁽²⁾	70.84 ⁽⁴⁾	54.10 ⁽⁴⁾	<u>45.88⁽⁴⁾</u>	47.61 ⁽⁴⁾	47.64 ⁽⁴⁾	46.88 ⁽⁴⁾	47.40 ⁽⁶⁾	45.58 ⁽⁴⁾
SBNDM q	108.32⁽²⁾	73.50⁽²⁾	64.87⁽²⁾	50.73⁽⁴⁾	47.85 ⁽⁴⁾	48.45 ⁽⁴⁾	48.55 ⁽⁴⁾	47.56 ⁽⁴⁾	47.08 ⁽⁴⁾	45.98 ⁽⁴⁾
KBNDM	192.76	126.56	92.03	71.04	64.17	56.79	49.56	49.09	50.24	47.68
BSDM q	117.22 ⁽²⁾	79.96 ⁽²⁾	70.36 ⁽³⁾	57.72 ⁽⁴⁾	49.91 ⁽⁸⁾	48.00 ⁽⁸⁾	44.99 ⁽⁸⁾	43.62 ⁽⁸⁾	43.06 ⁽⁸⁾	42.70 ⁽⁸⁾
FJS	192.52	158.61	106.40	89.04	78.43	68.14	64.80	61.15	60.48	55.38
FJS+	196.88	155.86	101.77	77.74	67.40	60.31	54.77	52.44	51.62	47.74
HASH q	312.91 ⁽²⁾	218.95 ⁽²⁾	107.38 ⁽²⁾	70.85 ⁽²⁾	56.90 ⁽³⁾	49.33 ⁽⁶⁾	46.63 ⁽⁵⁾	45.27 ⁽⁸⁾	44.77 ⁽³⁾	42.98 ⁽⁷⁾
FS- w	129.50 ⁽⁶⁾	104.45 ⁽⁶⁾	71.57 ⁽⁶⁾	61.08 ⁽⁴⁾	54.19 ⁽⁶⁾	49.53 ⁽⁴⁾	48.95 ⁽⁶⁾	47.02 ⁽⁴⁾	46.96 ⁽⁶⁾	43.00 ⁽⁴⁾
IOM	193.37	148.51	104.36	86.22	75.09	69.08	63.63	60.62	58.83	51.76
WOM	199.23	153.81	112.45	86.61	75.26	62.80	60.54	62.74	58.91	55.65
SKIP q	161.38 ⁽²⁾	100.52 ⁽²⁾	72.33 ⁽³⁾	55.71 ⁽⁴⁾	49.06 ⁽⁴⁾	48.73 ⁽⁴⁾	49.87 ⁽⁴⁾	48.81 ⁽⁸⁾	45.75 ⁽²⁾	45.32 ⁽²⁾
WFR q	137.40 ⁽²⁾	85.22 ⁽²⁾	68.88 ⁽²⁾	52.85 ⁽⁵⁾	46.67 ⁽⁵⁾	<u>44.39⁽⁸⁾</u>	<u>42.09⁽⁸⁾</u>	41.66 ⁽⁵⁾	41.65 ⁽⁶⁾	40.53 ⁽²⁾
LWFR q	121.08 ⁽²⁾	85.47 ⁽²⁾	70.38 ⁽²⁾	53.83 ⁽³⁾	47.89 ⁽⁵⁾	44.49 ⁽⁵⁾	42.38 ⁽⁶⁾	<u>41.09⁽⁸⁾</u>	<u>41.23⁽⁸⁾</u>	<u>40.30⁽⁸⁾</u>
DIST q	115.61 ⁽²⁾	80.14 ⁽²⁾	65.84 ⁽³⁾	52.25 ⁽⁴⁾	48.13 ⁽⁴⁾	46.26 ⁽⁴⁾	43.32 ⁽⁴⁾	42.84 ⁽⁸⁾	42.98 ⁽⁴⁾	41.47 ⁽⁷⁾
LDIST q	229.62 ⁽²⁾	102.15 ⁽²⁾	69.70 ⁽³⁾	55.34 ⁽⁴⁾	49.46 ⁽⁵⁾	45.93 ⁽⁵⁾	44.37 ⁽³⁾	43.15 ⁽⁴⁾	42.21 ⁽⁵⁾	41.11 ⁽⁷⁾

algorithms that run in linear time in the input string size are marked with ★

Genome sequence

- Genome sequence of E. coli
 - $n = 4641652$
 - $\sigma = 4$
- Patterns are randomly extracted from text

$n = |T|$ $m = |P|$ $\sigma = |\Sigma|$: Alphabet size

m	2	4	8	16	32	64	128	256	512	1024
BNDM q	218.34 ⁽²⁾	174.55 ⁽²⁾	84.14 ⁽⁴⁾	64.89 ⁽⁶⁾	60.39 ⁽⁶⁾	61.07 ⁽⁶⁾	60.85 ⁽⁶⁾	62.17 ⁽⁶⁾	61.20 ⁽⁶⁾	60.16 ⁽⁶⁾
SBNDM q	179.90⁽²⁾	154.20 ⁽²⁾	80.94 ⁽⁴⁾	73.29 ⁽⁶⁾	57.10 ⁽⁸⁾	61.01 ⁽⁶⁾	61.74 ⁽⁶⁾	61.20 ⁽⁶⁾	61.67 ⁽⁶⁾	60.80 ⁽⁶⁾
KBNDM	311.78	201.99	150.15	113.84	83.23	67.83	75.65	75.39	76.58	74.72
BSDM q	195.38 ⁽²⁾	118.86⁽³⁾	84.23 ⁽⁵⁾	63.03 ⁽⁶⁾	61.26 ⁽⁶⁾	58.75 ⁽⁶⁾	57.50 ⁽⁷⁾	56.99 ⁽⁶⁾	57.01 ⁽⁶⁾	56.58 ⁽⁶⁾
FJS	407.02	353.60	311.96	279.13	308.42	297.00	266.12	317.79	317.34	296.19
FJS+	388.44	296.70	203.03	171.17	149.52	136.59	128.55	130.39	122.51	112.99
HASH q	571.44 ⁽²⁾	272.86 ⁽³⁾	126.00 ⁽³⁾	88.12 ⁽³⁾	68.22 ⁽³⁾	58.84 ⁽⁶⁾	55.09 ⁽⁶⁾	59.48 ⁽⁷⁾	57.34 ⁽⁷⁾	57.96 ⁽⁷⁾
FS- w	332.32 ⁽⁴⁾	245.99 ⁽⁴⁾	184.72 ⁽⁴⁾	158.72 ⁽⁴⁾	143.79 ⁽⁶⁾	125.05 ⁽⁴⁾	123.52 ⁽⁶⁾	117.90 ⁽⁶⁾	108.03 ⁽⁶⁾	100.72 ⁽⁶⁾
IOM	377.25	275.36	215.72	220.97	219.86	218.12	210.61	221.31	230.15	211.69
WOM	381.54	301.46	220.34	182.30	166.27	143.24	136.20	133.75	127.40	114.74
SKIP q	250.89 ⁽²⁾	136.18 ⁽³⁾	91.51 ⁽⁴⁾	63.96 ⁽⁶⁾	56.79 ⁽⁷⁾	53.09 ⁽⁷⁾	52.10 ⁽⁷⁾	57.12 ⁽⁷⁾	56.97 ⁽⁸⁾	58.00 ⁽⁶⁾
WFR q	219.50 ⁽²⁾	168.39 ⁽²⁾	88.86 ⁽⁴⁾	65.82 ⁽⁴⁾	57.19 ⁽⁸⁾	55.12 ⁽²⁾	51.77 ⁽³⁾	50.04 ⁽³⁾	55.02 ⁽⁶⁾	54.64⁽⁸⁾
LWFR q	216.10 ⁽²⁾	173.48 ⁽³⁾	88.71 ⁽⁴⁾	60.75 ⁽⁵⁾	53.84⁽⁵⁾	50.48⁽⁶⁾	49.65⁽⁸⁾	48.71 ⁽⁶⁾	54.90 ⁽⁶⁾	54.97 ⁽⁷⁾
DIST q	186.10 ⁽²⁾	125.44 ⁽³⁾	78.56⁽⁴⁾	60.48⁽⁵⁾	55.21 ⁽⁶⁾	52.05 ⁽⁷⁾	51.26 ⁽⁸⁾	50.44 ⁽⁸⁾	54.81 ⁽⁷⁾	55.58 ⁽⁸⁾
LDIST q	295.55 ⁽²⁾	181.99 ⁽³⁾	86.58 ⁽⁴⁾	65.29 ⁽⁶⁾	56.74 ⁽⁶⁾	52.31 ⁽⁶⁾	50.39 ⁽⁶⁾	48.70⁽⁷⁾	54.33⁽⁴⁾	55.22 ⁽⁷⁾

algorithms that run in linear time in the input string size are marked with ★

Fibonacci string

- Definition

$Fib_1 = b, \quad Fib_2 = a, \quad Fib_n = Fib_{n-1} \cdot Fib_{n-2}$ for $n > 2$

- Use Fib_{32} as text

- $n = 2178309, \sigma = 2$

- Patterns are randomly extracted from text

	$n = T \quad m = P \quad \sigma = \Sigma : \text{Alphabet size}$									
m	2	4	8	16	32	64	128	256	512	1024
BNDM q	343.25 ⁽²⁾	308.44 ⁽²⁾	283.26 ⁽⁴⁾	257.64 ⁽⁶⁾	233.94 ⁽⁶⁾	285.63 ⁽⁴⁾	284.37 ⁽⁴⁾	293.00 ⁽⁴⁾	307.82 ⁽⁴⁾	315.47 ⁽⁴⁾
SBNDM q	286.02⁽²⁾	292.15⁽²⁾	272.98 ⁽⁴⁾	276.35 ⁽⁶⁾	306.42 ⁽⁶⁾	372.03 ⁽⁶⁾	432.53 ⁽⁶⁾	493.20 ⁽⁶⁾	546.94 ⁽⁶⁾	602.09 ⁽⁶⁾
KBNDM	541.70	405.78	411.08	422.85	382.25	402.45	425.60	437.67	461.07	451.12
BSDM q	482.47 ⁽²⁾	500.43 ⁽³⁾	397.29 ⁽⁵⁾	362.52 ⁽⁸⁾	330.76 ⁽⁶⁾	736.89 ⁽¹⁾	766.57 ⁽¹⁾	782.98 ⁽¹⁾	790.26 ⁽¹⁾	508.80 ⁽³⁾
FJS	402.44	362.23	276.97	237.87	218.38	206.07	203.86	202.94	196.49	194.01
FJS+	456.20	396.04	335.70	319.93	295.64	300.36	296.48	295.37	288.56	289.00
HASH q	645.48 ⁽²⁾	406.70 ⁽²⁾	257.69⁽⁴⁾	251.91 ⁽⁷⁾	279.81 ⁽⁷⁾	344.99 ⁽⁷⁾	415.16 ⁽⁷⁾	470.71 ⁽⁷⁾	514.05 ⁽⁷⁾	579.57 ⁽⁷⁾
FS- w	383.23 ⁽¹⁾	396.23 ⁽¹⁾	347.72 ⁽¹⁾	289.15 ⁽¹⁾	253.36 ⁽¹⁾	246.82 ⁽¹⁾	248.73 ⁽¹⁾	235.66 ⁽¹⁾	235.35 ⁽¹⁾	230.61 ⁽¹⁾
IOM	381.92	414.42	453.54	497.84	543.93	641.13	751.42	839.92	899.19	1019.59
WOM	552.38	555.43	564.67	617.93	664.47	732.05	852.06	926.35	1036.31	1126.17
SKIP q	470.93 ⁽²⁾	394.09 ⁽²⁾	332.66 ⁽⁵⁾	336.16 ⁽⁸⁾	374.91 ⁽⁸⁾	464.23 ⁽³⁾	460.54 ⁽³⁾	450.18 ⁽³⁾	451.05 ⁽³⁾	464.13 ⁽³⁾
WFR q	442.32 ⁽²⁾	497.95 ⁽³⁾	528.45 ⁽³⁾	652.48 ⁽⁶⁾	2132.38 ⁽⁷⁾	3762.19 ⁽⁸⁾	6762.67 ⁽⁸⁾	12624.63 ⁽⁸⁾	24416.65 ⁽⁸⁾	48596.02 ⁽⁸⁾
LWFR q	552.64 ⁽²⁾	504.77 ⁽³⁾	428.34 ⁽⁵⁾	342.80 ⁽⁷⁾	297.07 ⁽⁷⁾	304.57 ⁽²⁾	274.47 ⁽⁶⁾	265.28 ⁽⁶⁾	258.06 ⁽⁶⁾	254.92 ⁽⁶⁾
DIST q	438.96 ⁽¹⁾	359.56 ⁽²⁾	293.80 ⁽²⁾	235.61⁽²⁾	213.07⁽⁷⁾	206.92 ⁽²⁾	201.18 ⁽⁸⁾	196.56 ⁽⁵⁾	194.86 ⁽⁴⁾	193.62 ⁽⁵⁾
LDIST q	565.13 ⁽²⁾	412.15 ⁽³⁾	300.98 ⁽⁴⁾	245.38 ⁽⁷⁾	215.89 ⁽⁸⁾	208.40 ⁽⁷⁾	200.45⁽⁴⁾	193.74⁽⁴⁾	190.85⁽³⁾	193.48⁽⁸⁾

algorithms that run in linear time in the input string size are marked with ★

In experiment of Fibonacci string

$$Fib_{10} =$$

abaababaabaababaababaabaababaabaababaabaababaabaababa

$$P = \mathbf{baaba}$$

$$Fib_{10} =$$

abaababaabaababaababaabaababaabaababaabaababaabaababaabaababa

$$P = \text{aabababaab}$$

- There are many repeating structures
 - The pattern is extracted from the text, so the number of pattern occurrences is very large

Hypothesize

Efficiency of proposed algorithms do not decrease when number of pattern occurrences is large

Texts with frequent pattern occurrences

- Generated by intentionally embedding a lot of randomly generated patterns without overlapping
 - $n = 4000000$
- Fixed pattern length

$m = 8, \sigma = 4$ occ : Pattern occurrences $n = |T|$ $m = |P|$ $\sigma = |\Sigma|$: Alphabet size

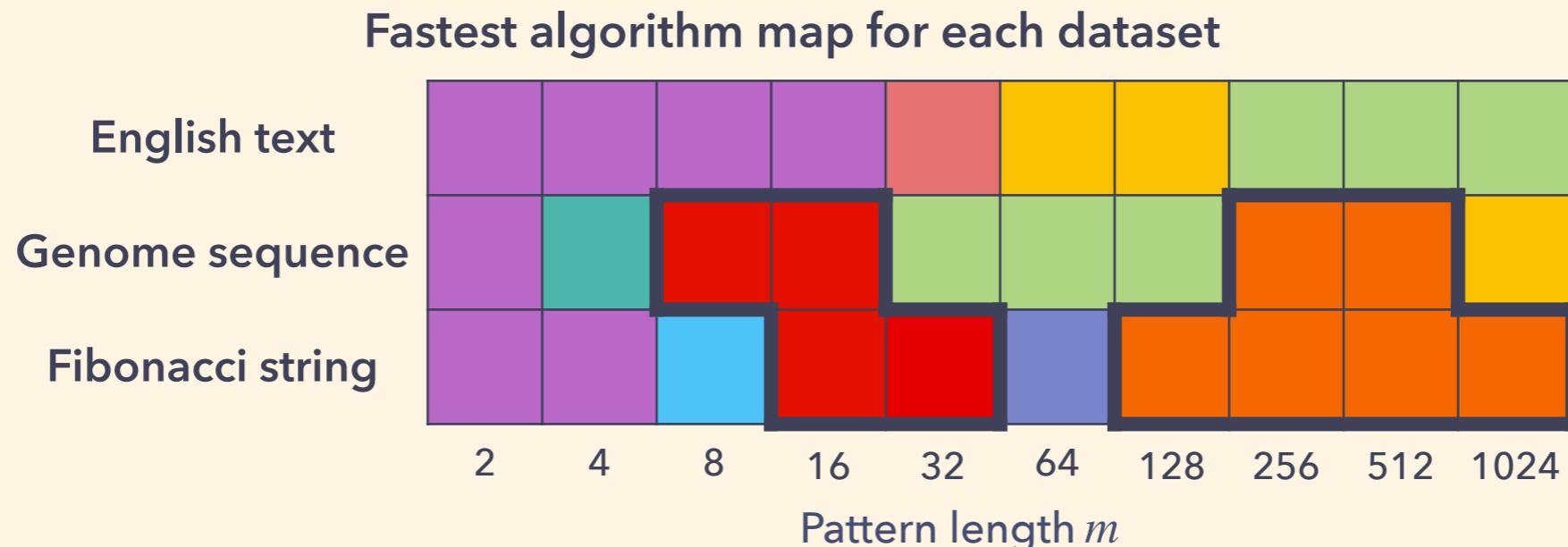
occ	0	128	256	512	1024	2048	4096	8192	16384	32768	65536	131072
BNDM q	73.53 ⁽⁴⁾	72.28 ⁽⁴⁾	72.87 ⁽⁴⁾	73.89 ⁽⁴⁾	74.39 ⁽⁴⁾	75.75 ⁽⁴⁾	77.60 ⁽⁴⁾	81.63 ⁽⁴⁾	82.67 ⁽⁴⁾	88.98 ⁽⁴⁾	108.56 ⁽⁴⁾	146.13 ⁽⁶⁾
SBNDM q	68.58⁽⁴⁾	68.92⁽⁴⁾	71.17 ⁽⁴⁾	77.26 ⁽⁴⁾	81.04 ⁽⁴⁾	80.76 ⁽⁴⁾	78.33 ⁽⁴⁾	82.82 ⁽⁴⁾	91.01 ⁽⁴⁾	96.83 ⁽⁴⁾	111.75 ⁽⁴⁾	140.21 ⁽⁴⁾
KBNDM	127.73	128.41	128.41	126.92	128.74	131.17	129.74	135.70	140.99	152.46	164.59	186.32
BSDM q	69.85 ⁽⁴⁾	69.04 ⁽⁴⁾	68.85⁽⁴⁾	69.19⁽⁴⁾	70.63⁽⁴⁾	72.00⁽⁴⁾	72.48⁽⁴⁾	76.25⁽⁴⁾	79.95⁽⁴⁾	89.91 ⁽⁴⁾	110.22 ⁽⁴⁾	143.91 ⁽⁴⁾
FJS	246.26	253.35	263.94	247.73	255.94	276.44	262.83	256.36	251.33	269.39	259.93	251.06
FJS+	174.38	173.86	180.46	173.05	178.65	182.98	177.56	181.17	182.48	190.71	197.69	199.60
HASH q	121.09 ⁽³⁾	121.36 ⁽³⁾	122.64 ⁽³⁾	122.50 ⁽³⁾	119.49 ⁽³⁾	123.98 ⁽³⁾	124.31 ⁽³⁾	124.93 ⁽³⁾	126.81 ⁽³⁾	128.94 ⁽³⁾	143.75 ⁽³⁾	153.08 ⁽⁴⁾
FS- w	154.89 ⁽⁴⁾	155.56 ⁽⁴⁾	165.18 ⁽⁴⁾	156.67 ⁽⁴⁾	159.03 ⁽⁴⁾	166.44 ⁽⁴⁾	165.62 ⁽⁴⁾	160.46 ⁽⁴⁾	167.46 ⁽⁴⁾	178.53 ⁽²⁾	185.86 ⁽²⁾	191.38 ⁽²⁾
IOM	185.31	181.07	195.53	187.98	194.18	200.99	195.98	195.89	197.67	202.61	214.90	209.52
WOM	192.59	192.28	207.57	189.21	196.02	203.86	196.98	199.79	203.59	215.79	219.94	230.59
SKIP q	79.32 ⁽⁴⁾	77.14 ⁽⁴⁾	79.19 ⁽⁴⁾	82.08 ⁽⁴⁾	81.82 ⁽⁴⁾	82.55 ⁽⁴⁾	83.83 ⁽⁴⁾	87.08 ⁽⁴⁾	89.82 ⁽⁴⁾	93.09 ⁽⁴⁾	106.41 ⁽⁴⁾	126.22 ⁽³⁾
WFR q	76.68 ⁽⁴⁾	76.38 ⁽⁴⁾	77.63 ⁽⁴⁾	83.21 ⁽⁴⁾	80.93 ⁽⁴⁾	81.42 ⁽⁴⁾	83.86 ⁽⁴⁾	88.55 ⁽⁴⁾	93.16 ⁽⁴⁾	106.07 ⁽⁴⁾	123.90 ⁽⁴⁾	164.77 ⁽⁴⁾
LWFR q	76.99 ⁽⁴⁾	76.33 ⁽⁴⁾	78.83 ⁽⁴⁾	77.57 ⁽⁴⁾	78.74 ⁽⁴⁾	76.46 ⁽⁴⁾	84.65 ⁽⁴⁾	89.87 ⁽⁴⁾	96.17 ⁽⁴⁾	108.67 ⁽⁴⁾	129.35 ⁽³⁾	168.70 ⁽³⁾
DIST q	69.67 ⁽⁴⁾	73.38 ⁽⁴⁾	74.35 ⁽⁴⁾	74.31 ⁽⁴⁾	75.12 ⁽⁴⁾	74.96 ⁽⁴⁾	77.21 ⁽⁴⁾	77.56 ⁽⁴⁾	80.09 ⁽⁴⁾	86.25⁽⁴⁾	101.12⁽⁴⁾	120.98⁽³⁾
LDIST q	75.82 ⁽⁴⁾	74.45 ⁽⁴⁾	74.89 ⁽⁴⁾	76.77 ⁽⁴⁾	74.62 ⁽⁴⁾	75.47 ⁽⁴⁾	77.10 ⁽⁴⁾	80.18 ⁽⁴⁾	83.40 ⁽⁴⁾	88.86 ⁽⁴⁾	103.73 ⁽⁴⁾	122.27 ⁽⁴⁾

algorithms that run in linear time in the input string size are marked with ★

Conclusion

$n = |T|$ $m = |P|$ ω : word length
 $\sigma = |\Sigma|$: alphabet size q : q -gram

- Proposed two string matching algorithms based on the distances of the q -gram occurrences
- Both algorithms run in linear time in the input string size



Comparing 15 powerful algorithms announced from 1977 to 2019 with the proposed algorithms

Algorithm	Preprocess	Search	Algorithm	Preprocess	Search	
BNDM q [Navarro & Raffinot, 1998]	$O(m+\sigma)$	$O(nm \lceil m/\omega \rceil)$	WFR q [Cantone+, 2017]	$O(m)$	$O(nm)$	
SBNDM q [Holub & Durian, 2005]	$O(m+\sigma)$	$O(nm \lceil m/\omega \rceil)$	LWFR q [Cantone+, 2019]	$O(m)$	$O(n)$	
FJS [Franek+, 2005]	$O(m+\sigma)$	$O(n)$	DIST q New	$O(mq)$	$O(nq)$	
HASH q [Leqroq, 2007]	$O(mq)$	$O(n(m+q))$	LDIST q New	$O(m)$	$O(n)$	
BSDM q [Faro & Leqroq, 2012]	$O(m)$	$O(nm)$	Naive solution : $O(nm)$			