

# ALGORITMI E STRUTTURE DATI

**Prof. Manuela Montangelo**

A.A. 2021/22

Algoritmi su stringhe:  
String matching

"E' vietata la copia e la riproduzione dei contenuti e immagini in qualsiasi forma.

E' inoltre vietata la redistribuzione e la pubblicazione dei contenuti e immagini non autorizzata espressamente dall'autore o dall'Università di Modena e Reggio Emilia."



**UNIMORE**  
UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

# Stringhe

**Alfabeto:** insieme di caratteri finito  $\Sigma = \{c_1, c_2, \dots, c_m\}$

**Stringa su alfabeto  $\Sigma$ :** sequenza lineare di elementi di  $\Sigma$ ,

$\alpha = s_0s_1s_2\dots s_n$ , tale che  $\forall i = 0, \dots, n \ s_i \in \Sigma$

**Insieme di stringhe definite su  $\Sigma$ :**  $\Sigma^*$  (stella di Kleene o monodie libero)  
insieme di tutte le stringhe di lunghezza finita sull'alfabeto  $\Sigma$

## Esempi:

- $\Sigma = \{0,1\}$ , un stringa  $\alpha = 1001010010$  è un numero binario,  
 $\Sigma^* =$  insieme di tutti i numeri binari (di qualunque lunghezza)
- $\Sigma = \{a, c, g, t\}$ , un stringa  $\alpha = acctgggtgca$  è una porzione di DNA
- $\Sigma = \{a, b, c, d, e, f, g, h, i, l, m, n, o, p, q, r, s, t, u, v, z\}$ , una stringa  
 $\alpha = algortimi$  è una parola contenente lettere dell'alfabeto italiano

# Stringhe

## DEFINIZIONI:

- Stringa vuota (non ha caratteri):  $\epsilon$
- Accesso carattere i-esimo di una stringa  $\alpha$ :  $\alpha[i]$  o  $\alpha_i$  (iniziamo a contare da indice zero)

$$\alpha = \text{ALGORITMI} \quad \alpha[3] = \text{O}$$

- Lunghezza di una stringa  $\alpha$ :  $|\alpha|$  = numero di caratteri di  $\alpha$  ( $|\epsilon| = 0$ )

$$\alpha = \text{ALGORITMI} \quad |\alpha| = 9$$

- $\beta$  è una sottostringa di  $\alpha$  se esistono due stringhe  $\gamma$  e  $\delta$  (eventualmente vuote) tale che  $\alpha = \gamma\beta\delta$   
( $\alpha[i..j] = \alpha[i]\alpha[i+1]\dots\alpha[j] = \alpha_{i..j}$  sottostringa di  $\alpha$  con i caratteri dall'i-esimo al j-esimo)

$$\alpha = \text{ALGORITMI}, \beta = \text{ORI} \longrightarrow \gamma = \text{ALG} \text{ e } \delta = \text{TMI}$$

- $\beta$  è un suffisso di  $\alpha$  se esiste una stringa  $\gamma$  (eventualmente vuota) tale che  $\alpha = \gamma\beta$

$$\alpha = \text{ALGORITMI}, \beta = \text{RITMI} \longrightarrow \gamma = \text{ALGO}$$

- $\beta$  è un prefisso di  $\alpha$  se esiste una stringa  $\delta$  (eventualmente vuota) tale che  $\alpha = \beta\delta$

$$\alpha = \text{ALGORITMI}, \beta = \text{AL} \longrightarrow \delta = \text{GORITMI}$$

# Stringhe

## OPERAZIONI SULLE STRINGHE:

- **Concatenazione:**

date su stringhe  $\alpha$  e  $\beta$  sullo stesso alfabeto  $\Sigma$ ,

la concatenazione è  $\gamma = \alpha\beta$  data dalla giustapposizione di  $\alpha$  e  $\beta$  ( $|\gamma| = |\alpha| + |\beta|$ )

$\alpha = \text{IPPO}$

$\beta = \text{POTAMO}$

$\gamma = \text{IPPOPOTAMO}$

- **Potenza:**

data una stringa  $\alpha$  su un alfabeto  $\Sigma$ ,

la potenza  $k$ -esima è la concatenazione di  $\alpha$  con se stessa  $k$  volte

$\alpha = \text{PIPO}$

$\alpha^3 = \text{PIPOPIPOPIPO}$

# String-matching

## OCCORRENZA:

- se  $\beta$  è una sottostringa di  $\alpha$ , allora  $\alpha$  contiene (almeno) un'occorrenza di

$\alpha = \text{ALGORITMI}$ ,  $\beta = \text{ORI}$   $\longrightarrow \gamma = \text{ALG}$  e  $\delta = \text{TMI}$

occorrenza di  $\beta$  in  $\alpha$

- se  $\alpha = \gamma\beta\delta$  e  $|\gamma| = k$ ,  $\beta$  occorre (ha un match) alla posizione  $k$  di  $\alpha$

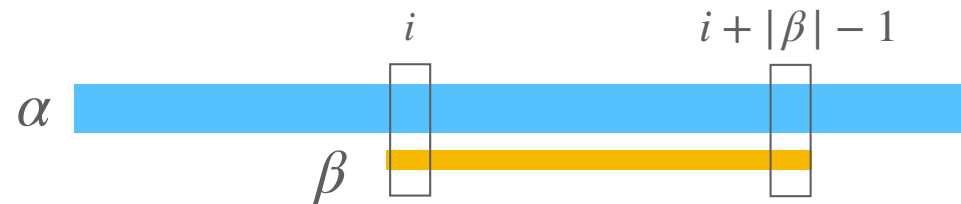
$\alpha = \text{ALGORITMI}$ ,  $\beta = \text{ORI}$   $\longrightarrow \gamma = \text{ALG}$  e  $\delta = \text{TMI}$

occorrenza di  $\beta$  in  $\alpha$  alla posizione 3

# String-matching

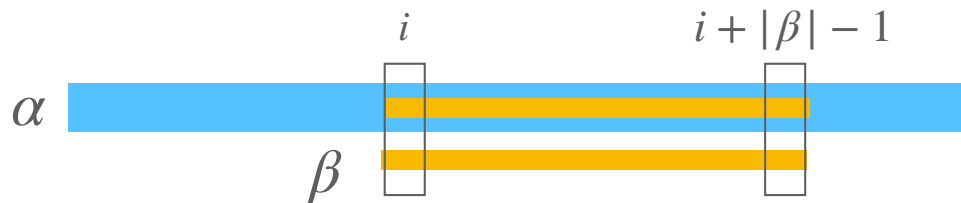
## ALLINEAMENTO:

L'allineamento di una stringa  $\beta$  alla posizione  $i$  della stringa  $\alpha$  è la sovrapposizione (ideale) di  $\beta$  con la sottostringa  $\alpha[i..i + |\beta| - 1]$



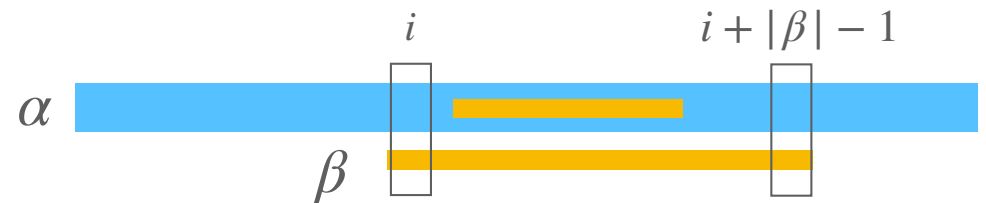
## ALLINEAMENTO con SUCCESSO (matching):

Esiste un'occorrenza della stringa  $\beta$  alla posizione  $i$  della stringa  $\alpha$



## ALLINEAMENTO FALLITO

Non esiste un'occorrenza della stringa  $\beta$  alla posizione  $i$  della stringa  $\alpha$



# String-matching

## PROBLEMA:

**INPUT:** Stringa  $T \in \Sigma^*$  di  $n$  caratteri  $\rightarrow$  TESTO

Stringa  $P \in \Sigma^*$  di  $m \leq n$  caratteri  $\rightarrow$  PATTERN

**OUTPUT:** le (eventuali) posizioni delle occorrenze di  $P$  in  $T$

## ESEMPIO:

$T = \text{perdirindina}$

$P = \text{din}$

Occorrenze di  $P$  in  $T$  alle posizioni 3 e 11

$\text{per}\text{din}\text{dirindina}$

# String-matching

## PROBLEMA:

**INPUT:** Stringa  $T \in \Sigma^*$  di  $n$  caratteri  $\rightarrow$  TESTO

Stringa  $P \in \Sigma^*$  di  $m \leq n$  caratteri  $\rightarrow$  PATTERN

**OUTPUT:** le (eventuali) posizioni delle occorrenze di  $P$  in  $T$

## Lower bound sul numero di confronti tra TESTO e PATTERN:

Ogni carattere del testo DEVE essere coinvolto in ALEMNO un confronto  
 $\rightarrow$  Lower bound  $\Omega(|T|)$

$$T = 1^n$$

$$P = 1$$

"saltando" anche solo un carattere di  $T$   
si "perde" un'occorrenza di  $P$  in  $T$



# String-matching

## PROBLEMA:

**INPUT:** Stringa  $T \in \Sigma^*$  di  $n$  caratteri  $\longrightarrow$  TESTO

Stringa  $P \in \Sigma^*$  di  $m \leq n$  caratteri  $\longrightarrow$  PATTERN

**OUTPUT:** le (eventuali) posizioni delle occorrenze di  $P$  in  $T$

## ALGORITMO a FORZA BRUTA

### IDEA:

proviamo tutti i possibili allineamenti del pattern nel testo  
e verifichiamo danno un successo o un fallimento

# String-matching

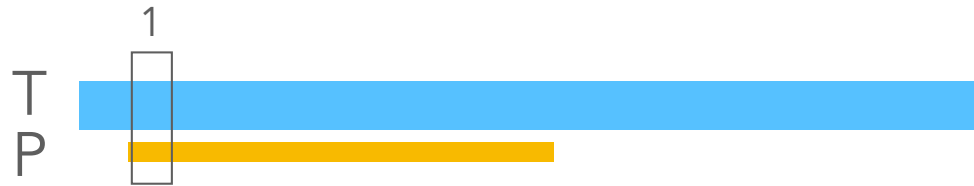
## ALGORITMO a FORZA BRUTA

### IDEA:

proviamo tutti i possibili allineamenti del pattern nel testo e verifichiamo danno un successo o un fallimento



1a possibilità —> indice 0 di T



2a possibilità —> indice 1 di T



3a possibilità —> indice 2 di T

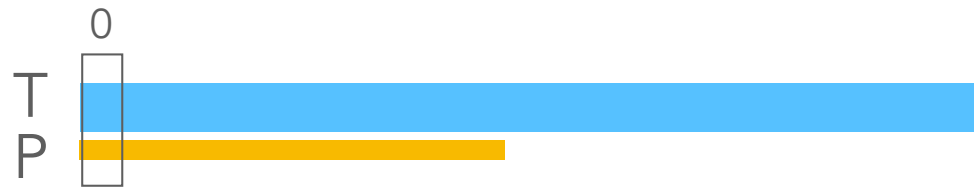
....

# String-matching

## ALGORITMO a FORZA BRUTA

### IDEA:

proviamo tutte le possibili collocazioni del pattern nel testo e verifichiamo se determinano un'occorrenza o no



1a possibilità  $\rightarrow$  indice 0 di T

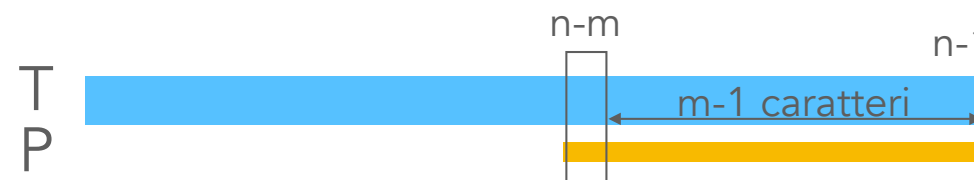


2a possibilità  $\rightarrow$  indice 1 di T



3a possibilità  $\rightarrow$  indice 2 di T

....



ultima possibilità  $\rightarrow$  indice  $n-m$  di T

# String-matching

## ALGORITMO a FORZA BRUTA

### IDEA:

proviamo tutte le possibili collocazioni del pattern nel testo e verifichiamo se determinano un'occorrenza o no

```
Brute-Force(T, P)
```

```
  n := |T|
```

```
  m := |P|
```

```
  for i = 0 to n - m do
```

```
    j := 0
```

```
    while (j < m AND T[i + j] = P[j]) do
```

```
      j := j + 1
```

```
    if j = m
```

```
      then print i
```

```
  print -1
```

Controlla se

$$P = T[i..i + m - 1]$$



# String-matching

## ALGORITMO a FORZA BRUTA

### IDEA:

proviamo tutte le possibili collocazioni del pattern nel testo e verifichiamo se determinano un'occorrenza o no

```
Brute-Force(T, P)  
  n := |T|  
  m := |P|  
  for i = 0 to n-m do  
    j := 0  
    while (j < m AND T[i+j] = P[j]) do  
      j := j + 1  
      O(m)  
    if j = m  
      then print i  
  print -1
```

### Costo computazionale

contiamo solo i confronti  
tra caratteri delle due stringhe

Il for viene eseguito  $n-m+1$  volte  
 $(n - m + 1) \cdot m = nm - m^2 + m$

se  $m \approx n \Rightarrow O(n)$

se  $m = n/c, c \in O(1) \Rightarrow O(n^2)$

se  $m \in O(1) \Rightarrow O(n)$

# String-matching

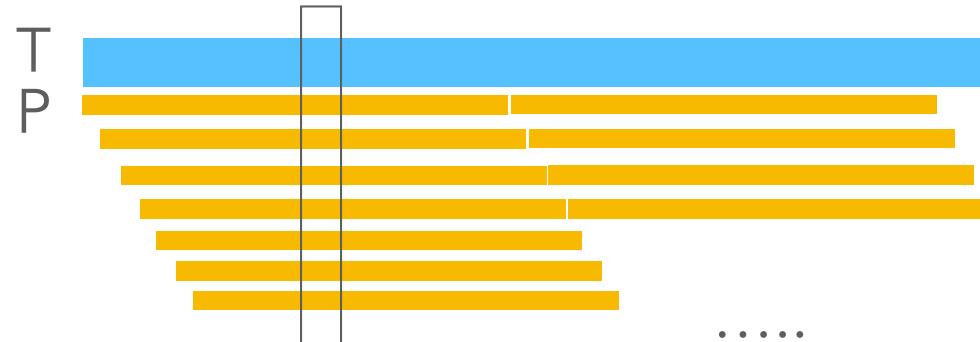
## ALGORITMO a FORZA BRUTA

Costo computazionale

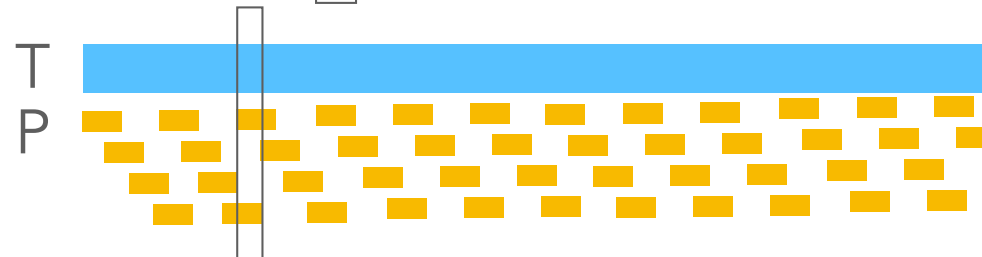
se  $m \approx n \Rightarrow O(n)$



se  $m = n/c, c \in O(1) \Rightarrow O(n^2)$



se  $m \in O(1) \Rightarrow O(n)$



# String-matching

L'algoritmo Brute-Force testa anche posizioni in cui è impossibile che inizi una nuova occorrenza (e lo sa!)

$i$

T	.... A A A B x x x x x .....
P	A A A A x x x x ...

è possibile che ci sia un'occorrenza  
del pattern nel testo in posizione  $i+1$ ?

$i$   $i+1$

T	.... A A A B x x x x x .....
P	A A A A x x x x ...

NO

# String-matching

L'algoritmo Brute-Force testa anche posizioni in cui è impossibile che inizi una nuova occorrenza (e lo sa!)

$i$

T	.... A A A B x x x x x .....
P	A A A A x x x x ...

qual è il primo indice (dopo  $i$ ) in cui è possibile che ci sia un'occorrenza del pattern nel testo?

$i$                        $i+4$

T	.... A A A B x x x x x .....
P	A A A A x x x x ...



# String-matching

L'algoritmo Brute-Force testa anche posizioni in cui è impossibile che inizi una nuova occorrenza (e lo sa!)

$i$

T	.... A B A B x x x x x .....
P	A B A A x x x x ...

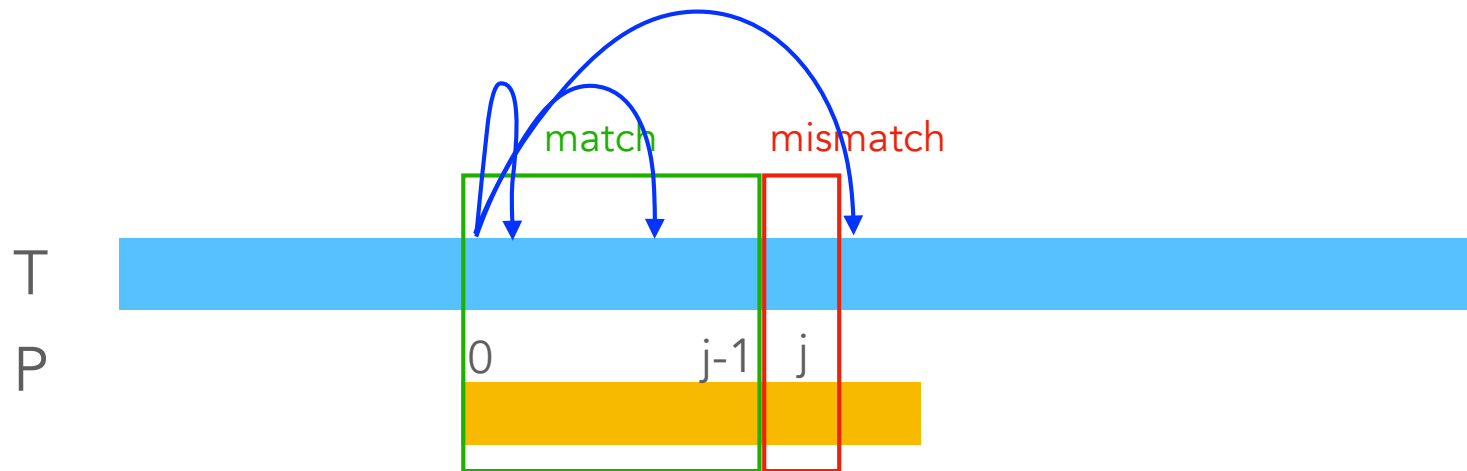
qual è il primo indice (dopo  $i$ ) in cui è possibile che ci sia un'occorrenza del pattern nel testo?

$i$      $i+2$

T	.... A B A B x x x x x .....
P	A B A A x x x x ...

# String-matching

Come determinare il prossimo allineamento?

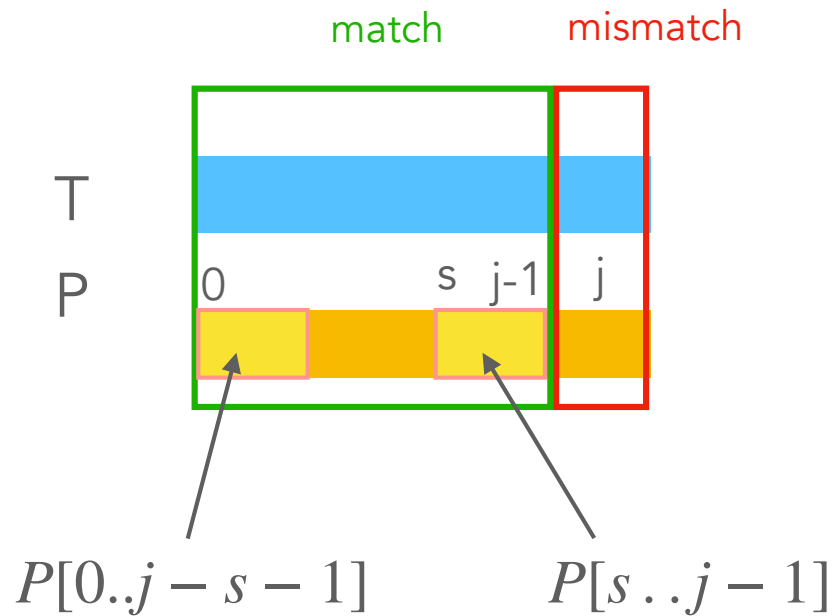


**shift:** #posizioni di cui spostarsi a destra nel testo per provare un nuovo allineamento

# String-matching

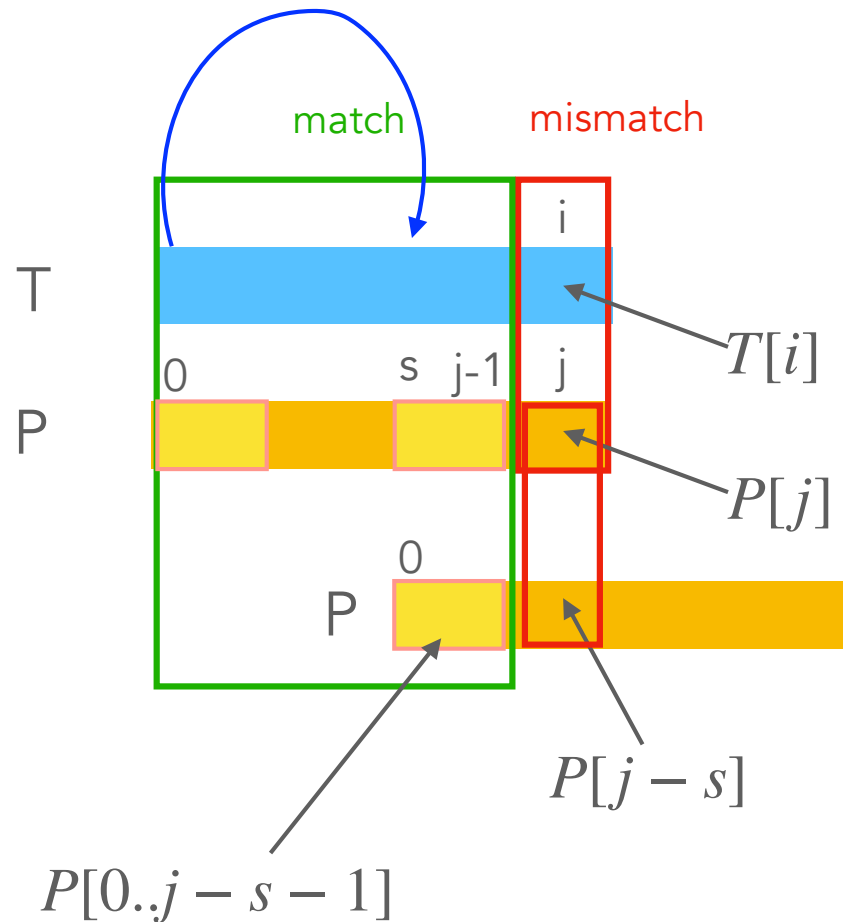
Come determinare il prossimo allineamento?

se  $P[0..j - s - 1] = P[s..j - 1]$ ?



# String-matching

Come determinare il prossimo allineamento?



se  $P[0..j-s-1] = P[s..j-1]$ ?

Vale la pena riprovare con uno  
**shift di  $s$**  posizioni?

dipende da  $P[j-s]$  e  $P[j]$

se  $P[j] = P[j-s]$  non è possibile  
perché  $P[j-s] \neq T[i]$

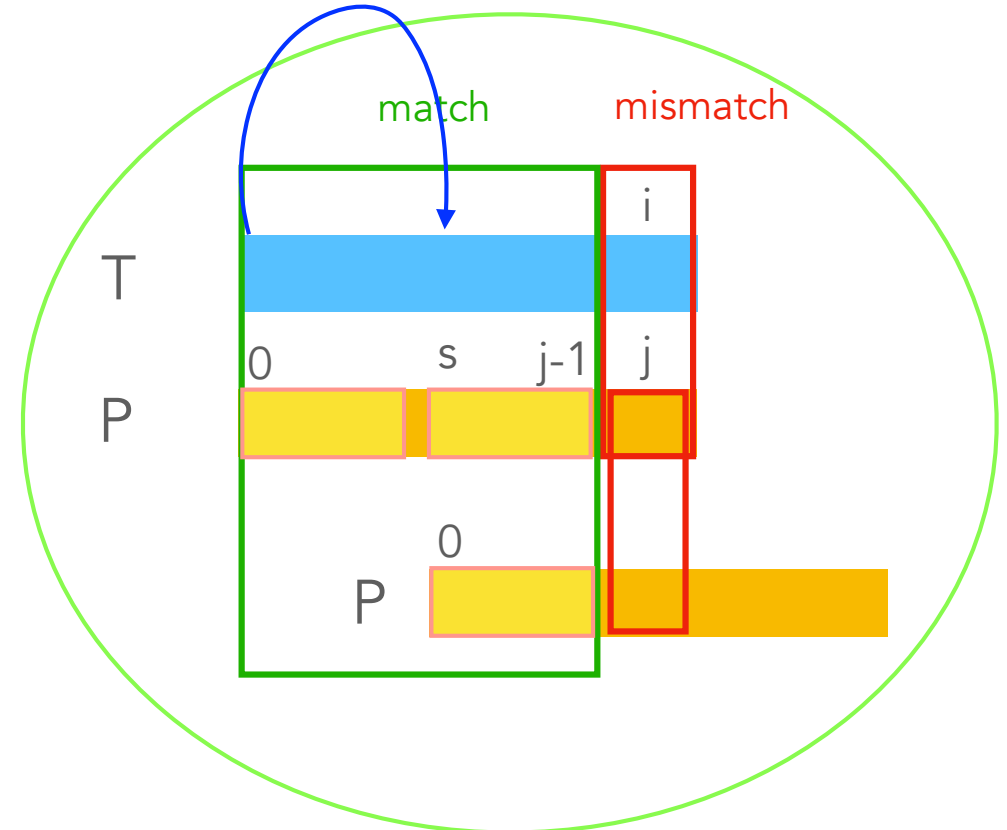
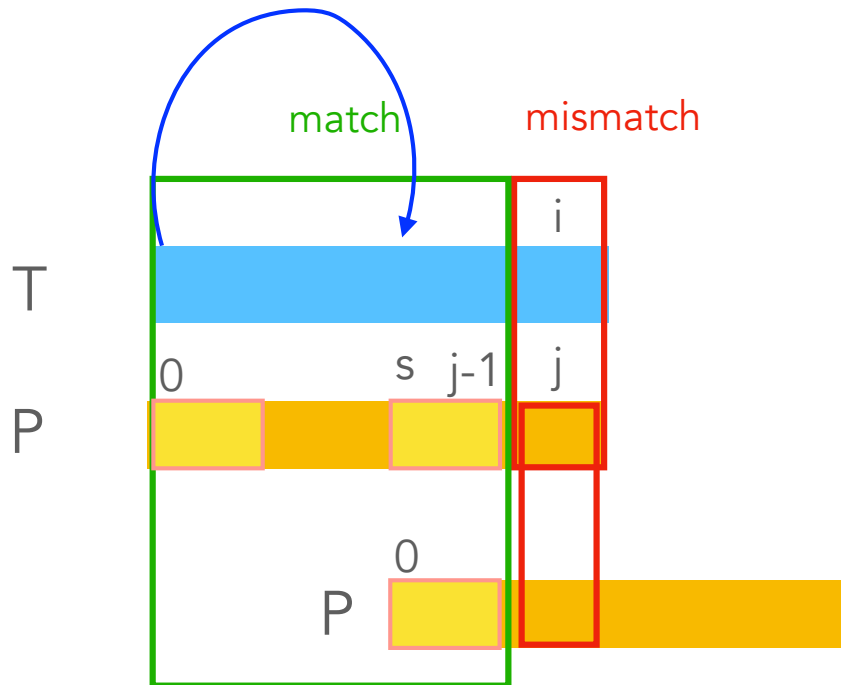
se  $P[j] \neq P[j-s]$  vale la pena provare

# String-matching

Come determinare il prossimo allineamento?

e se esiste **più di un**  $s$  per cui vale

$$P[0..j-s-1] = P[s..j-1] \text{ e } P[j] \neq P[j-s]?$$

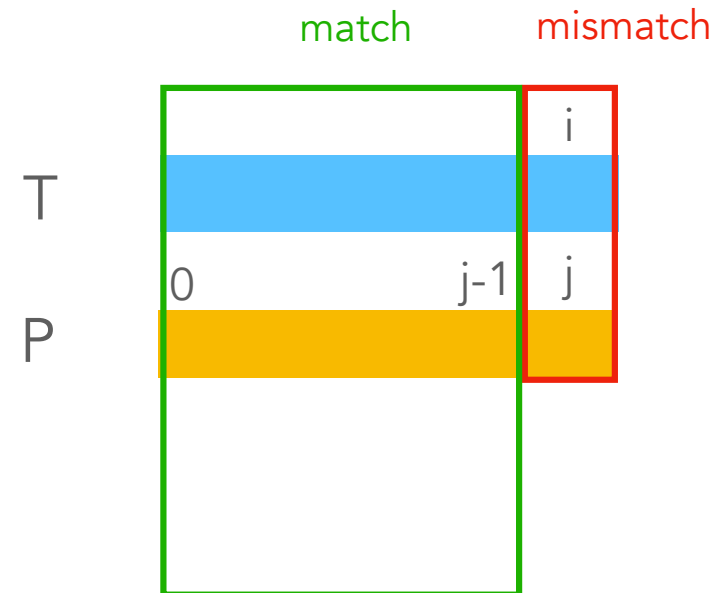
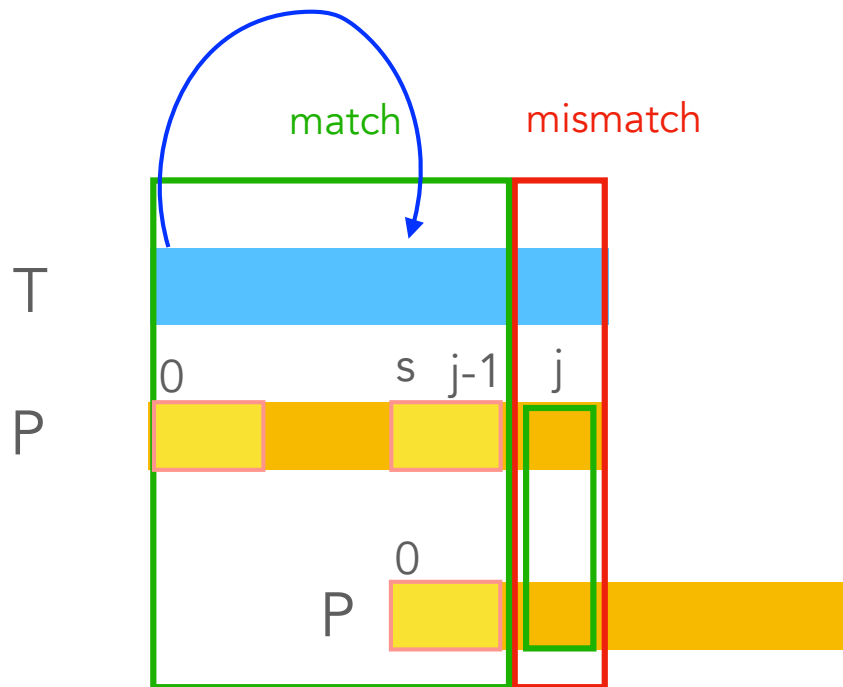


scegliamo l' $s$  più piccolo

# String-matching

Come determinare il prossimo allineamento?

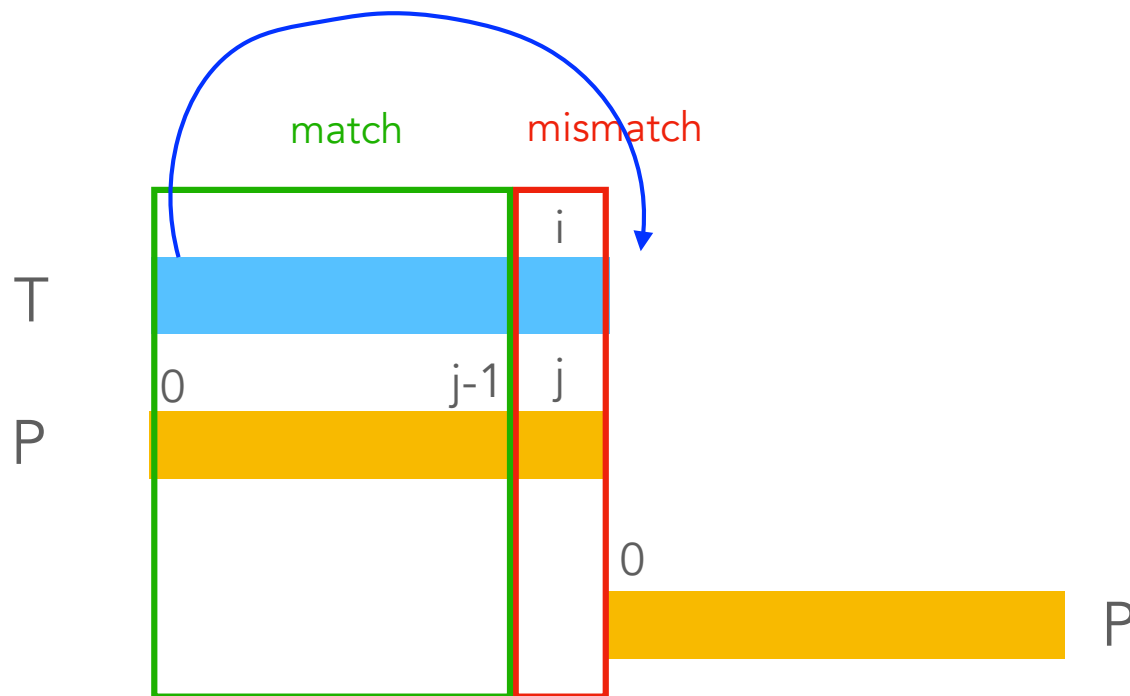
e se **non esiste** tale  $s$ ?



# String-matching

Come determinare il prossimo allineamento?

e se **non esiste** tale  $s$ ?

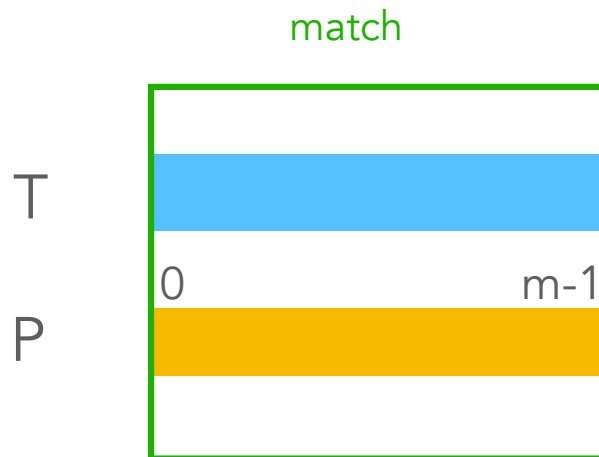


il prossimo tentativo va fatto con uno **shift di  $j + 1$**  posizioni

# String-matching

Come determinare il prossimo allineamento?

e se abbiamo trovato **un'occorrenza**  
(quindi non c'è mismatch)?



Si ragiona in modo analogo,  
senza dover controllare cosa  
succede nella posizione m-esima



# String-matching

Come determinare il prossimo allineamento?

Analizzando il pattern (e basta!)

Per ogni  $0 \leq j \leq |P|$  per ogni possibile posizione del mismatch,  
o se è stata trovata un'occorrenza del pattern

cerchiamo il più piccolo  $s$  tale che  $0 \leq s \leq j + 1$  e  $A_P(j, s) = \text{True}$   
minima lunghezza dello shift

$$A_P(j, s) = \begin{cases} \text{True} & \text{se } P[0..j-s-1] = P[s..j-1] \text{ AND } (P[j] \neq P[j-s] \text{ OR } j = |P|) \\ \text{True} & \text{se } s = j + 1 \\ \text{False} & \text{altrimenti} \end{cases}$$

# String-matching

Come determinare il prossimo allineamento?

Analizzando il pattern (e basta!)

Non sono coinvolti confronti con il testo

**Best-Prefix(P)**

$m := |P|$

**shift** := new array[0..m]

**for**  $j = 0$  **to**  $m$  **do** ←

$s := 1$

**while** **NOT**  $A_P(j, s)$  **do** ←

$s := s + 1$

**shift**[ $j$ ] :=  $s$  ←

**return** **shift**[ ]

per ogni possibile posizione del mismatch ( $j < m$ ),  
o se è stata trovata un'occorrenza del pattern ( $j = m$ )

ricerca dello shift di  
minima lunghezza

il **while** termina al più tardi  
quando  $s = j + 1$

quando non inizia una nuova iterazione del **while**,  
il valore di  $s$  è il primo per cui  $A_P(j, s) = \text{True}$

# String-matching

## ESEMPIO

**Best-Prefix(P)**

```
m := |P|  
shift := new array[0..m]  
for j = 0 to m do  
  s := 1  
  while NOT AP(j,s) do  
    s := s+1  
  shift[j] := s  
return shift[ ]
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	
2	di	
3	din	
4	dind	
5	dindi	
6	dindin	
7	dindina	

# String-matching

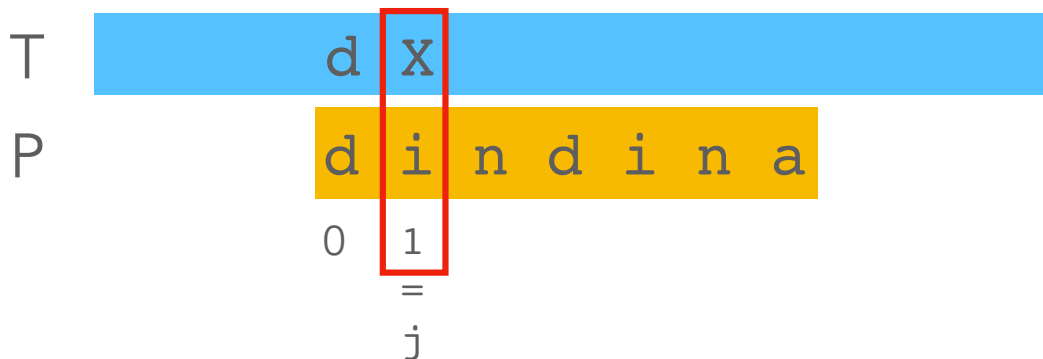
## ESEMPIO

### Best-Prefix(P)

```
m := |P|  
shift := new array[0..m]  
for j = 0 to m do  
  s := 1  
  while NOT AP(j,s) do  
    s := s+1  
  shift[j] := s  
return shift[ ]
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	
3	din	
4	dind	
5	dindi	
6	dindin	
7	dindina	



# String-matching

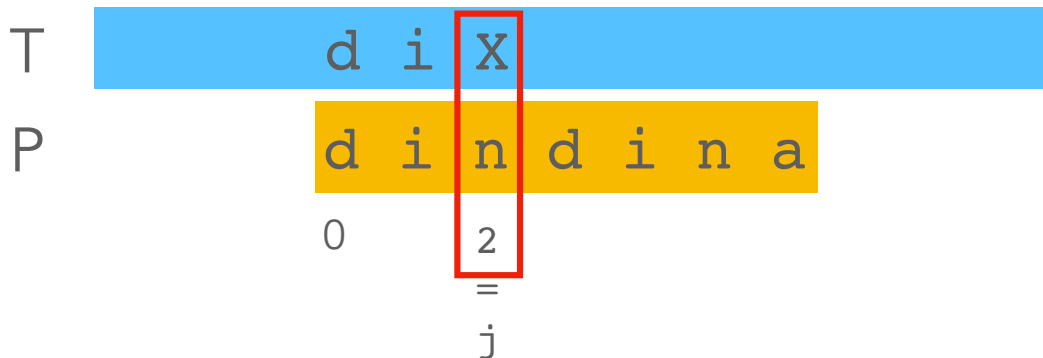
ESEMPIO

**Best-Prefix(P)**

```
m := |P|  
shift := new array[0..m]  
for j = 0 to m do  
  s := 1  
  while NOT AP(j,s) do  
    s := s+1  
  shift[j] := s  
return shift[ ]
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	
4	dind	
5	dindi	
6	dindin	
7	dindina	



# String-matching

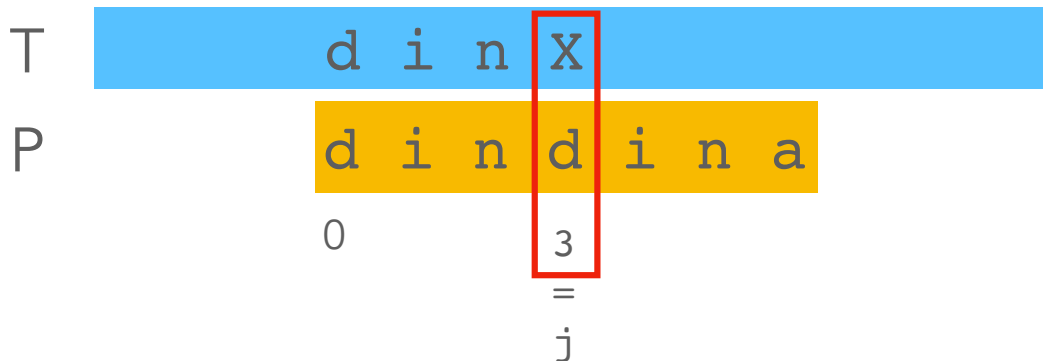
ESEMPIO

**Best-Prefix(P)**

```
m := |P|  
shift := new array[0..m]  
for j = 0 to m do  
  s := 1  
  while NOT AP(j,s) do  
    s := s+1  
  shift[j] := s  
return shift[ ]
```

P = dindina     m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	4
4	dind	
5	dindi	
6	dindin	
7	dindina	



# String-matching

ESEMPIO

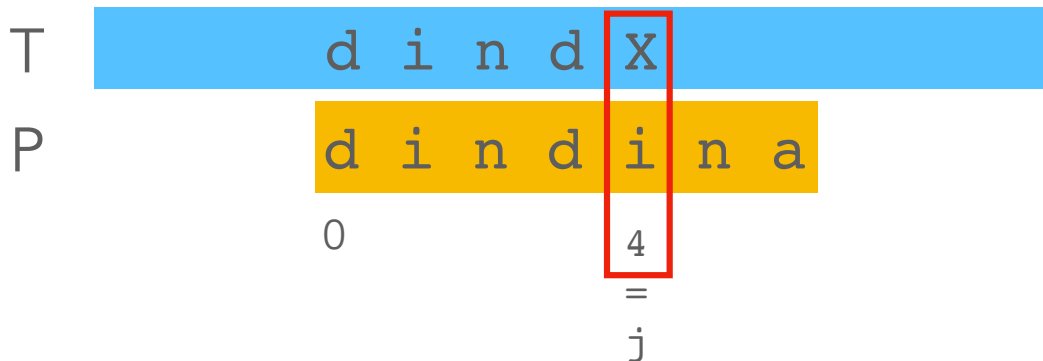
**Best-Prefix(P)**

```

m := |P|
shift := new array[0..m]
for j = 0 to m do
  s := 1
  while NOT AP(j,s) do
    s := s+1
  shift[j] := s
return shift[ ]
    
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	4
4	dind	4
5	dindi	
6	dindin	
7	dindina	



# String-matching

ESEMPIO

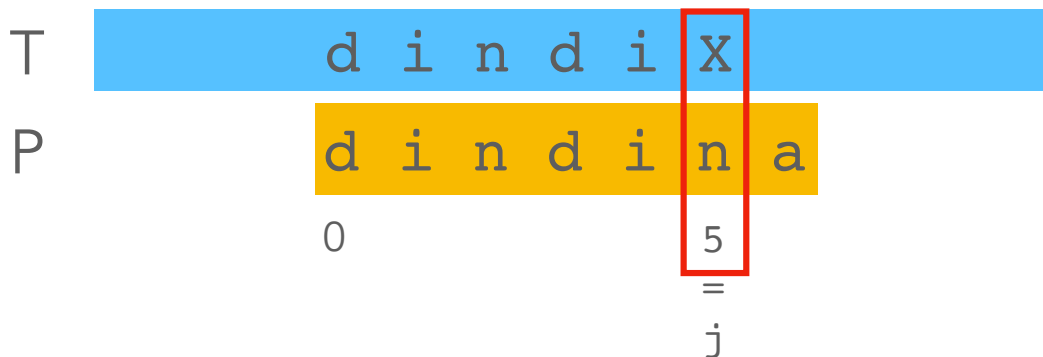
**Best-Prefix(P)**

```

m := |P|
shift := new array[0..m]
for j = 0 to m do
  s := 1
  while NOT AP(j,s) do
    s := s+1
  shift[j] := s
return shift[ ]
    
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	4
4	dind	4
5	dindi	5
6	dindin	
7	dindina	





# String-matching

ESEMPIO

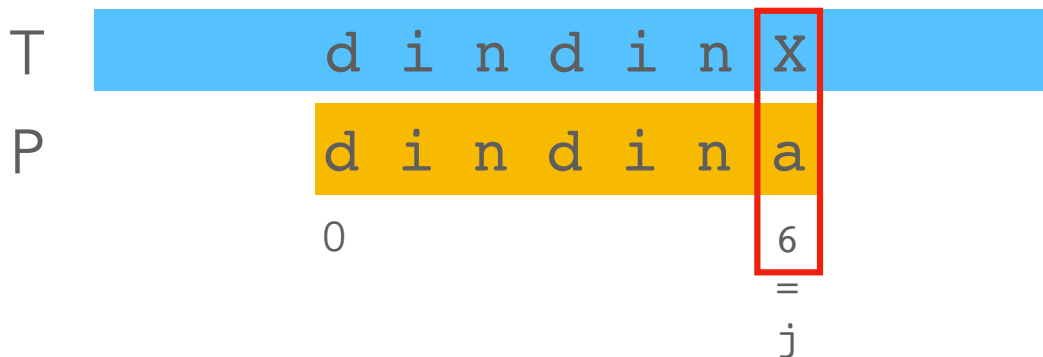
**Best-Prefix(P)**

```

m := |P|
shift := new array[0..m]
for j = 0 to m do
  s := 1
  while NOT AP(j,s) do
    s := s+1
  shift[j] := s
return shift[ ]
    
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	4
4	dind	4
5	dindi	5
6	dindin	3
7	dindina	



# String-matching

## ESEMPIO

**Best-Prefix(P)**

```
m := |P|  
shift := new array[0..m]  
for j = 0 to m do  
  s := 1  
  while NOT AP(j,s) do  
    s := s+1  
  shift[j] := s  
return shift[ ]
```

P = dindina      m = |P| = 7

j	P[0..j-1]	shift[j]
0		1
1	d	1
2	di	2
3	din	4
4	dind	4
5	dindi	5
6	dindin	3
7	dindina	7

T      d i n d i n a

P      d i n d i n a

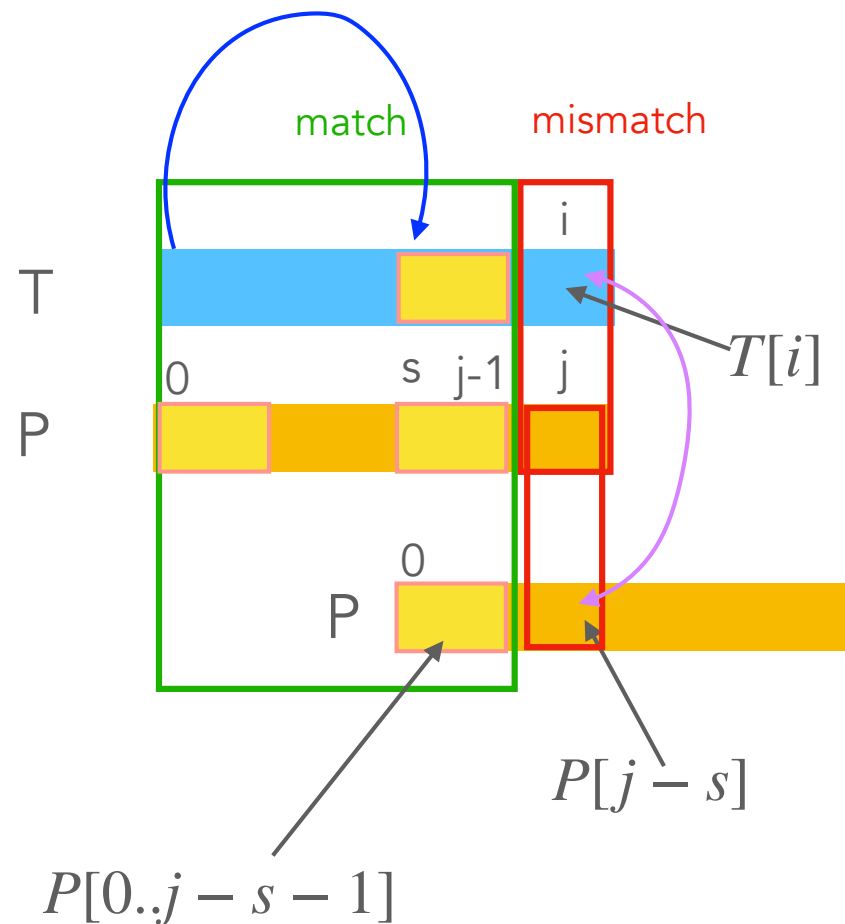
0

7  
=  
j

# String-matching

Sappiamo come determinare il prossimo tentativo di allineamento, qual'è il prossimo confronto da fare?

Evitiamo di ripetere confronti tra il testo e il pattern di cui conosciamo già l'esito



$$\text{shift} \leq j$$

Dall'allineamento precedente sappiamo che

$$P[0..j-s-1] = T[i-j+s..i-1]$$



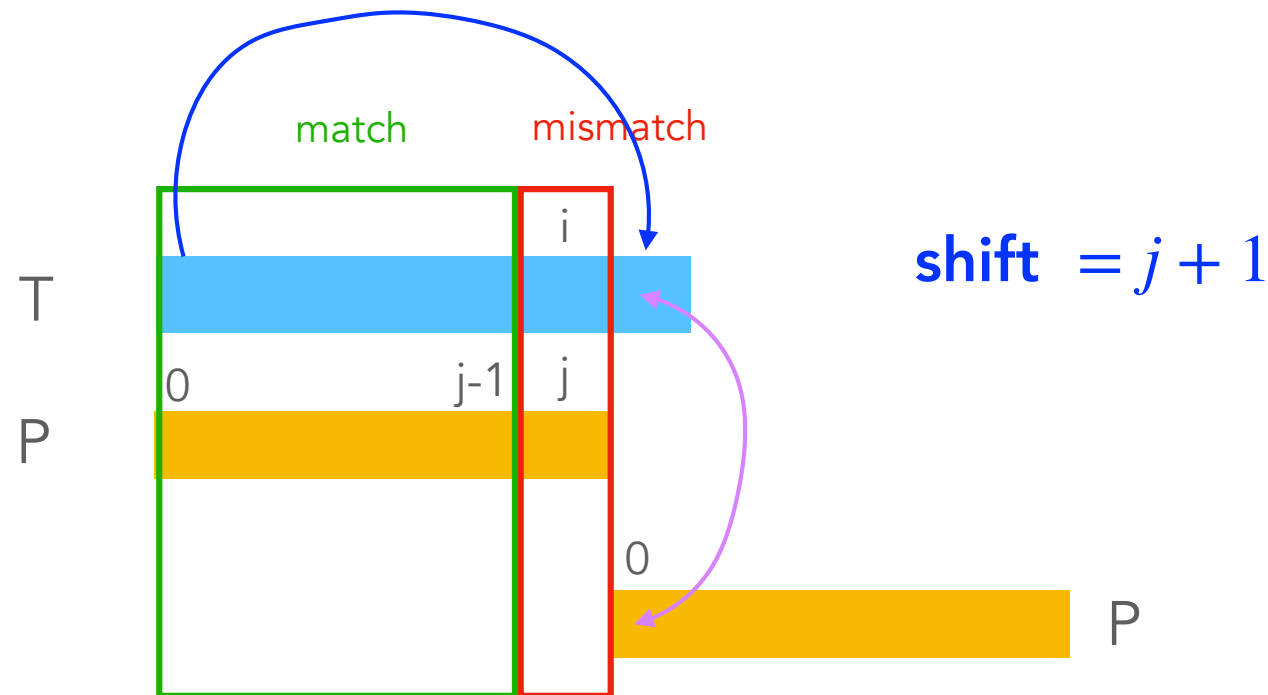
Il **prossimo** confronto interessante è tra

$$P[j-s] \text{ e } T[i]$$

# String-matching

Sappiamo come determinare il prossimo tentativo di allineamento

Evitiamo di ripetere confronti tra il testo e il pattern di cui conosciamo già l'esito



Il **prossimo** confronto interessante è tra  $P[0]$  e  $T[i + 1]$

# String-matching

Sappiamo come determinare il prossimo tentativo di allineamento

Evitiamo di ripetere confronti tra il testo e il pattern di cui conosciamo già l'esito

**shift = s**  $\leq j$        $P[j - s]$  e  $T[i]$       **prox** indice di P  $\Rightarrow j - s$

**shift = s = j + 1**  $P[0]$  e  $T[i + 1]$       **prox** indice di P  $\Rightarrow 0$

**prox** indice di P      **next** =  $\max\{0, j - s\}$

# String-matching

ALGORTIMO di KNUTH-MORRIS-PRATT (1975)

```
Knuth-Morris-Pratt(T, P)  
  n := |T|  
  m := |P|  
  shift := Best-Prefix(P)  
  for j = 0 to m do  
    next[j] := max{0, j - sift[j]}  
  i := 0 //indice che scorre T  
  j := 0 //indice che scorre P  
  while i ≤ m-n do  
    while (j < m AND P[j] = T[i+j]) do  
      j := j+1  
    if j = m then print i  
    i := i + shift[j]  
    j := next[j]  
  print -1
```

# String-matching

## ALGORTIMO di KNUTH-MORRIS-PRATT

```
Knuth-Morris-Pratt(T, P)  
  n := |T|  
  m := |P|  
  shift := Best-Prefix(P)  non coinvolge  
                                caratteri del testo!  
  for j = 0 to m do  
    next[j] := max{0, j - sift[j]}  
  i := 0 //indice che scorre T  
  j := 0 //indice che scorre P  
  while i ≤ m-n do  
    while (j < m AND P[j] = T[i+j]) do  
      j := j+1  
    if j = m then print i  
    i := i + shift[j]  
    j := next[j]  
  print -1
```

Costo computazionale

contiamo solo i confronti  
tra caratteri delle due stringhe

$$O(|T|)$$

ogni carattere del testo  
viene associato ad **AL PIÙ DUE**  
confronti con caratteri del  
pattern:

- un match
- il mismatch che occorre quando  
il pattern è allineato con il  
carattere

# ... nel mondo reale...

## Variazioni sul tema:

- String matching **esatto** VS string matching **approssimato**
- **Single pattern** string matching VS **multiple pattern** string matching

## Alcune applicazioni (... solo alcune di quelle più famose...):

- Correttore ortografico e operazione "cerca"
- Filtri anti-SPAM
- Antivirus
- Motori di ricerca
- Software anti-plagio
- Bioinformatica e sequenziamento DNA
- Indagini forensi digitali
- Information Retrieval
- ....