

Titolo modulo: Approfondimenti su “String Matching” [P2_06]

Titolo unità didattica: algoritmo KMP di string matching [1-AT]

Algoritmi di string matching che riducono il numero dei confronti

Argomenti trattati:

- ✓ Il problema dello string matching
- ✓ Esempi di problemi di string matching
- ✓ Algoritmi di string matching
- ✓ Algoritmo KMP

Prerequisiti richiesti: algoritmo di ricerca diretta

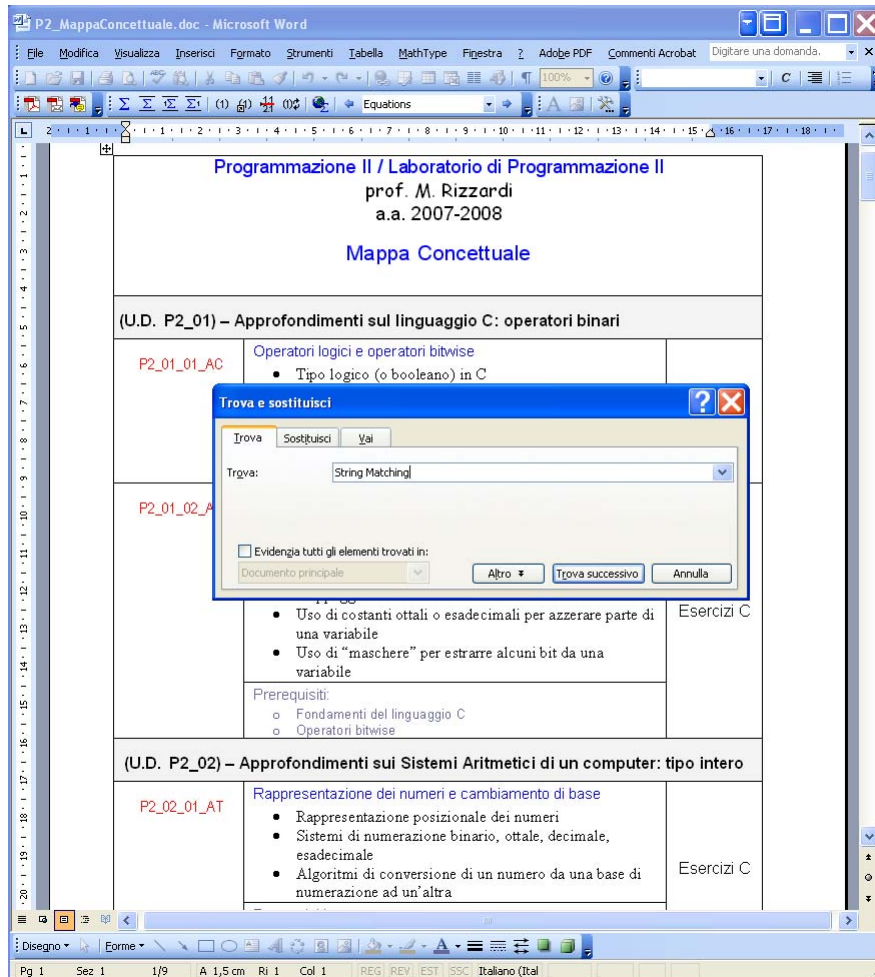
String Matching (**S.M.**):

ricerca di un **pattern P** in un **testo T**
(è un particolare problema di ricerca)

Il problema dello String Matching

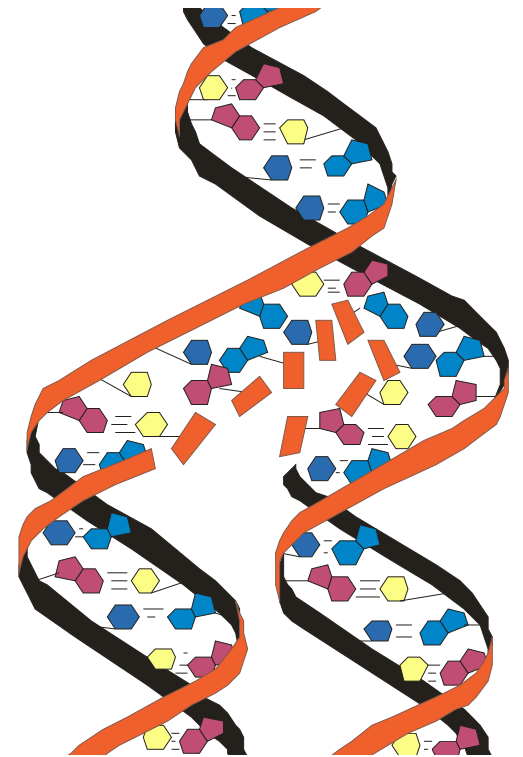
- **S.M. Decisionale:**
il pattern P è presente nel testo T?
- **S.M. Quantitativo:**
quante volte P si trova in T?
- **S.M. Enumerativo:**
in quali posizioni di T occorre P?

String matching (caso particolare di pattern matching)



È un problema che si presenta spesso nei programmi **editor** o **word processor**.

Gli stessi algoritmi di **string matching** sono anche usati per cercare particolari sottosequenze di **DNA**

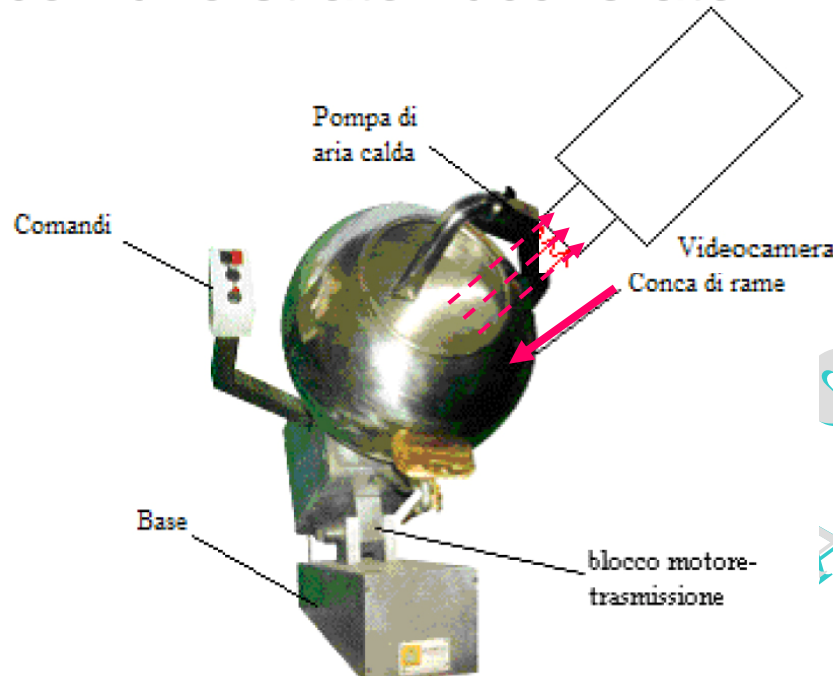


Altre applicazioni di Pattern/String matching

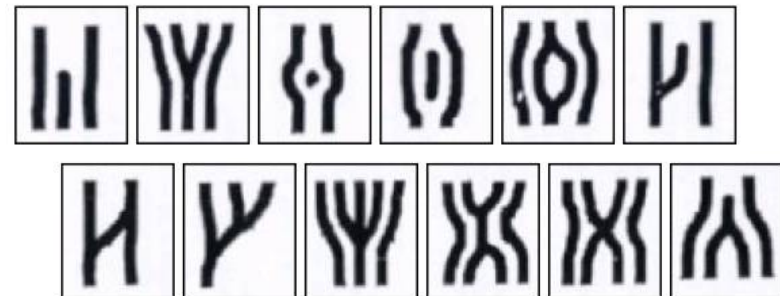
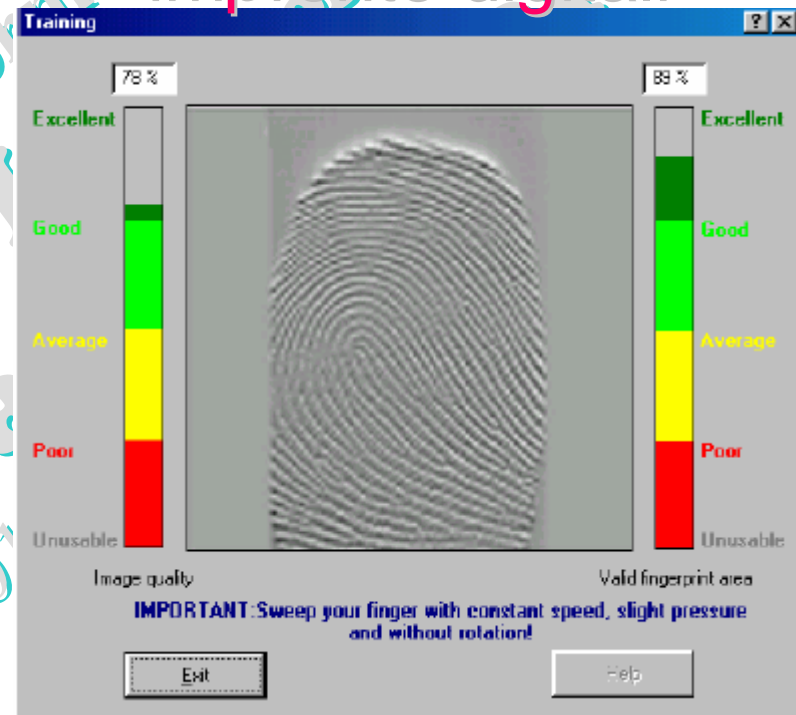
confettatura

=

rivestimento delle mandorle
con uno strato zuccherato



riconoscimento di
impronte digitali



minuzie

Algoritmi di string matching

❖ Naive (ricerca diretta);  già visto!

❖ Smart (automa a stata finiti);

❖ Rabin-Karp;

❖ Knuth-Morris-Pratt;

❖ Boyer-Moore.

algoritmi "veloci"

(prof. M. Rizzardi)

Algoritmo Naive (ricerca diretta): valutazione

- ❖ Complessità computazionale nel caso peggiore: $O(MN)$ dove N è la lunghezza del testo e M quella del pattern.
- ❖ Complessità computazionale nel caso medio: $O(N)$ se la distribuzione dei caratteri del testo e del pattern è uniforme e l'alfabeto contiene molti caratteri.

Vantaggi

- Semplice da implementare.
- È ragionevole usarlo per la prima occorrenza del pattern quando questa abbia una buona probabilità di trovarsi all'inizio del testo o quando il pattern sia lungo al più 3-4 caratteri.
- Il caso peggiore si verifica maggiormente con stringhe particolarmente ripetitive.

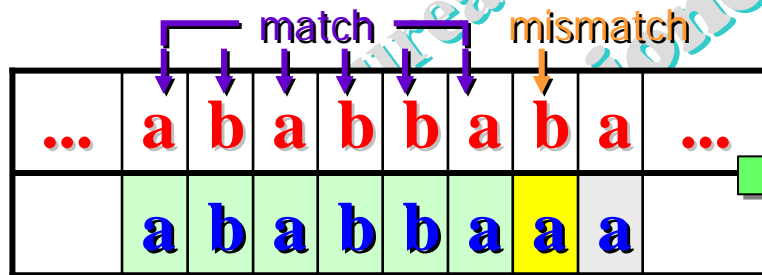
Svantaggi

- Inefficiente su sequenze lunghe del pattern (sequenze del DNA) in quanto l'algoritmo confronta ogni volta anche sottostringhe matched spostando il pattern sul testo soltanto di una posizione.

String matching: **Algoritmo KMP** (Knuth-Morris-Pratt)

Idea dell'algoritmo **KMP**:

quando il confronto fra il carattere del testo e del pattern fallisce (mismatch) dopo un certo numero di successi (match), invece che arretrare il puntatore sul pattern e sul **testo**, ignorando l'esito dei confronti già eseguiti, si possono sfruttare le conoscenze sul pattern evitando confronti il cui esito è già noto.



Si possono saltare dei passaggi

È richiesta la costruzione di una **tabella preliminare** a partire dai caratteri di **pattern**.

Algoritmo di ricerca più veloce in due fasi:

- (1) elaborazione del pattern;
- (2) scansione "veloce" del testo

algoritmo KMP: elaborazione del pattern

definizione di **prefisso** e **suffisso** di una stringa

STRINGA = "Buon giorno!"

"Buon" giorno!"

"Buon" giorno!"

"Buon" = $\text{pref}_4(\text{STRINGA})$
prefisso di lunghezza 4
di STRINGA

"giorno!" = $\text{suff}_7(\text{STRINGA})$
suffisso di lunghezza 7
di STRINGA

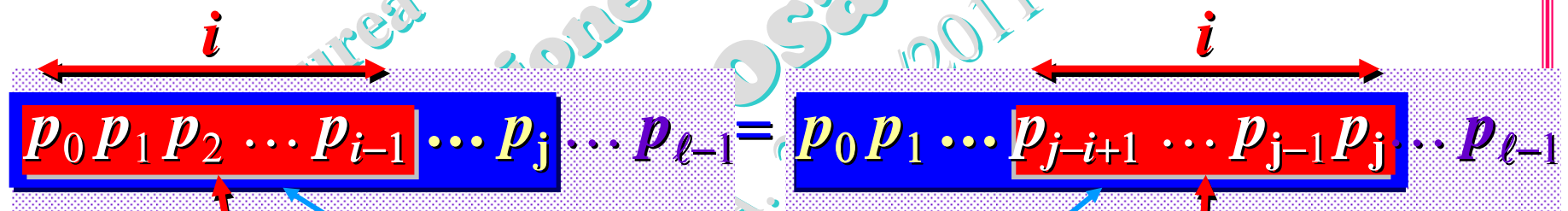
Si può calcolare il **prefisso** ed il **suffisso** anche di una sottostringa

Siano p_k i caratteri del **pattern** (di lunghezza totale ℓ):

$$\text{pattern} = "p_0 p_1 p_2 \dots p_{\ell-1}"$$

$\forall j$ sia $\text{sub}_j[\text{pattern}]$ la sottostringa di **pattern** di lunghezza $j+1$:

$$\text{sub}_j[\text{pattern}] = "p_0 p_1 p_2 \dots p_j" = \text{pattern}[0..j]$$



prefisso di lunghezza i
 $\text{pref}_i(\text{sub}_j[\text{pattern}])$

suffisso di lunghezza i
 $\text{suff}_i(\text{sub}_j[\text{pattern}])$

$\text{sub}_j[\text{pattern}] = \text{sottostringa di lunghezza } j+1 \text{ di pattern}$

Esempio: **prefisso** e **suffisso** di una sottostringa

STRINGA = "Buon giorno!"

$\text{sub}_6[\text{STRINGA}] = \text{"Buon gi"}$

"Buon"gi

Buon **gi**

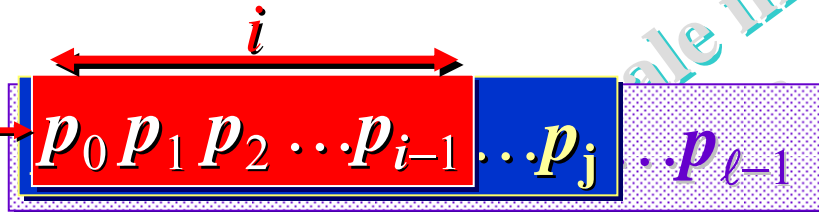
"Buon" = $\text{pref}_4(\text{sub}_6[\text{STRINGA}])$
prefisso di lunghezza 4
della sottostringa [6] di STRINGA

"_gi" = $\text{suff}_3(\text{sub}_6[\text{STRINGA}])$
suffisso di lunghezza 3
della sottostringa [6] di STRINGA

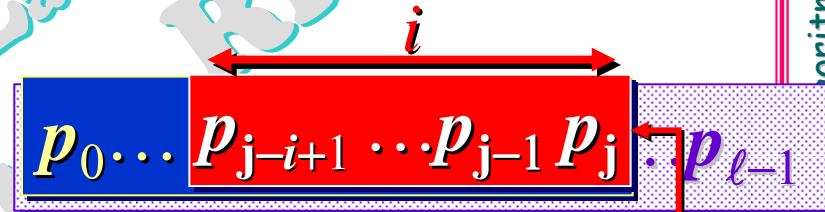
$p_0 p_1 p_2 \dots p_{i-1} \dots p_j \dots p_{\ell-1}$

sottostringa di lunghezza $j+1$ di pattern

definizione di **funzione insuccesso**



$\text{pref}_i(\text{sub}_j[\text{pattern}])$



$\text{suff}_i(\text{sub}_j[\text{pattern}])$

valore massimo di i ($0 < i \leq j$), se $\exists i$, tale che

$$\text{pref}_i(\text{sub}_j[\text{pattern}]) = \text{suff}_i(\text{sub}_j[\text{pattern}])$$

$$p_0 p_1 p_2 \dots p_{i-1} = p_{j-(i-1)} \dots p_{j-2} p_{j-1} p_j$$

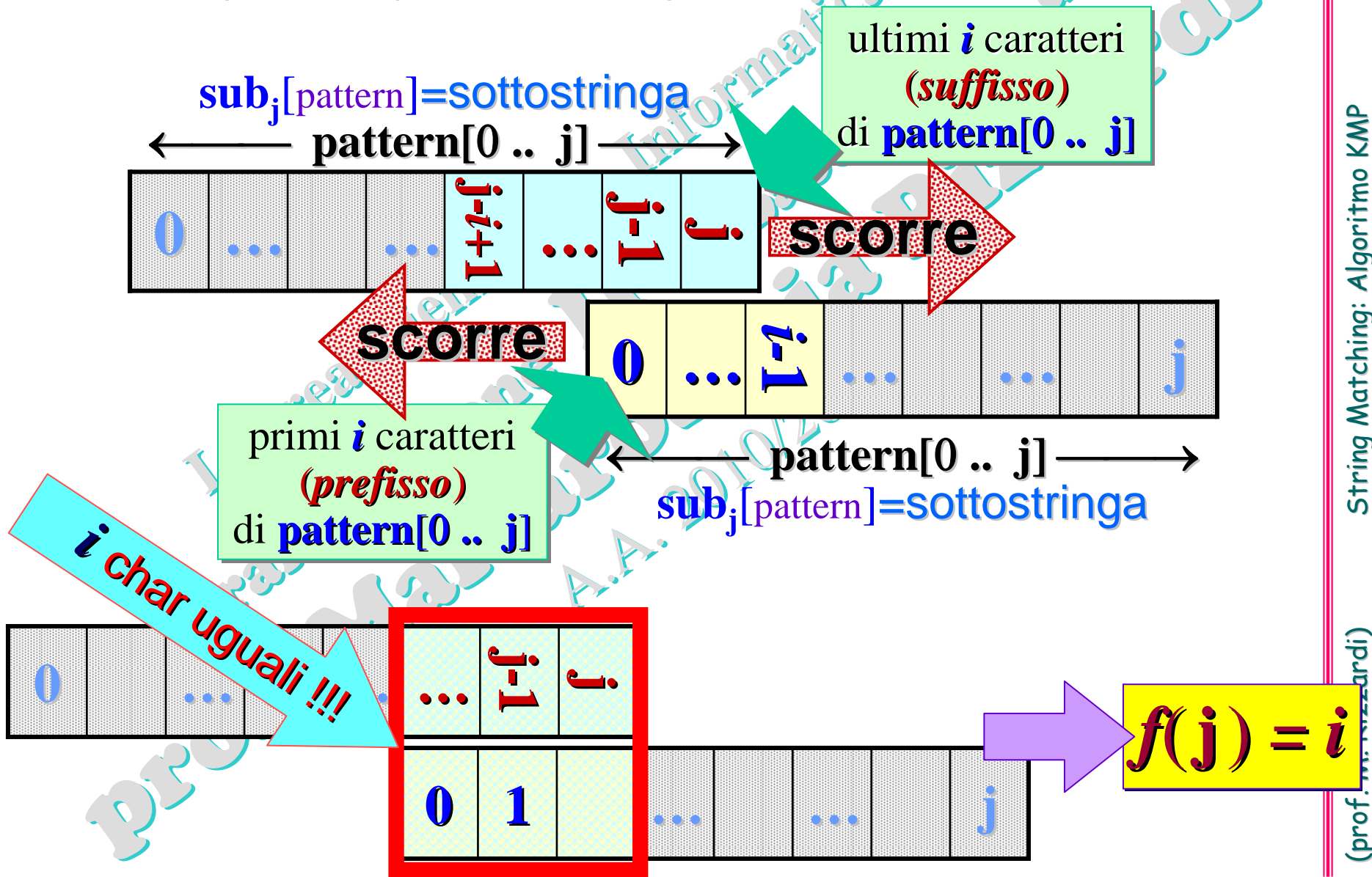
prefisso di lunghezza i suffisso di lunghezza i

$\forall j = 0, 1, \dots, \ell - 1$

$f(j) =$

0 altrimenti

Cerca nella sottostringa del pattern di lunghezza $j+1$ eventuali **prefissi** e **suffissi uguali**: sceglie quelli di **lunghezza massima**.



Esempio: pattern

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

$\text{sub}_3[\text{pattern}] = \text{pattern}[0..3]$

a	b	a	b	b	a	a	a
---	---	---	---	--------------	--------------	--------------	--------------

scorre

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

scorre

scorre

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

scorre

$\text{suff}_2(\text{sub}_3[\text{pattern}])$

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

a	b	a	b	b	a	a	a
---	---	---	---	---	---	---	---

$\text{pref}_2(\text{sub}_3[\text{pattern}])$

$f(3) = 2$

Esempio: pattern = “a b c a b c a c a b”
 “ $p_0 p_1 p_2 p_3 p_4 p_5 p_6 p_7 p_8 p_9$ ”

j	f(j)	costruzione tabella															
0	0	a	b	c	a	b	c	a	c	a	b						
			a	b	c	a	b	c	a	c	a	b					
1	0	a	b	c	a	b	c	a	c	a	b						
					a	b	c	a	b	c	a	c	a	b			
2	0	a	b	c	a	b	c	a	c	a	b						
						a	b	c	a	b	c	a	c	a	b		
3	1	a	b	c	a	b	c	a	c	a	b						
					a	b	c	a	b	c	a	c	a	b	a		
4	2	a	b	c	a	b	c	a	c	a	b						
					a	b	c	a	b	c	a	c	a	b	a		
5	3	a	b	c	a	b	c	a	c	a	b						
					a	b	c	a	b	c	a	c	a	b			
6	4	a	b	c	a	b	c	a	c	a	b						
					a	b	c	a	b	c	a	c	a	b			
7	0	a	b	c	a	b	c	a	c	a	b						
										a	b	c	a	b	c	a	c
8	1	a	b	c	a	b	c	a	c	a	b						
										a	b	c	a	b	c	a	c

legge

sottos

pre

suf

legenda

sottstringa

prefisso

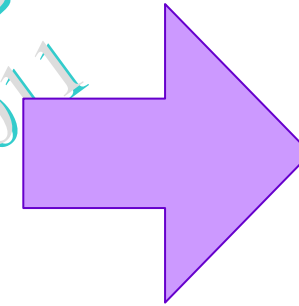
suffisso

Esempio: trovare il pattern **a b a b b a a a** nel seguente testo

a b b a a a b a b b a b a b b b a b a b b a a a b b b a b a

pattern = **a b a b b a a a**
 $p_0 p_1 p_2 p_3 p_4 p_5 p_6 p_7$

j	$f(j)$	costruzione tabella							
0	0	a	b	a	b	b	a	a	a
1	0	a	b	a	b	b	a	a	a
2	1	a	b	a	b	b	a	a	a
3	2	a	b	a	b	b	a	a	a
4	0	a	b	a	b	b	a	a	a
5	1	a	b	a	b	b	a	a	a
6	1	a	b	a	b	b	a	a	a
		p_0	p_1	p_2	p_3	p_4	p_5	p_6	p_7



j	$f(j)$
0	0
1	0
2	1
3	2
4	0
5	1
6	1

legenda

sottosstringa

prefisso

suffisso

algoritmo KMP: continuazione confronti

Se nel confronto dei caratteri t_k del **testo** con quelli p_h del **pattern** si verifica un "insuccesso" (mismatch) dopo j "successi" (match), cioè

$$\begin{array}{c} \text{j successi} \\ \longleftrightarrow \\ "t_{i-j} \dots t_{i-2} t_{i-1}" = "p_0 p_1 p_2 \dots p_{j-1}" \end{array} \text{ ma } \boxed{\begin{array}{c} \text{insuccesso} \\ t_i \neq p_j \end{array}}$$

allora la ricerca può ricominciare confrontando

next step {

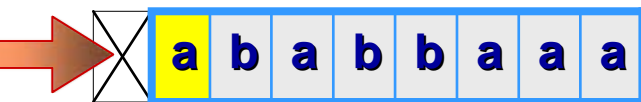
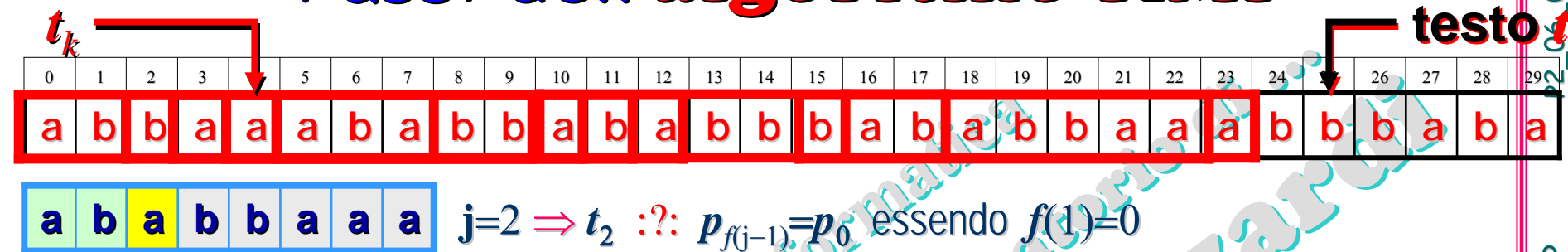
t_i	con	$p_{f(j-1)}$	se	$j \neq 0$
t_{i+1}	con	p_0	se	$j = 0$

Passi dell' **algoritmo KMP**

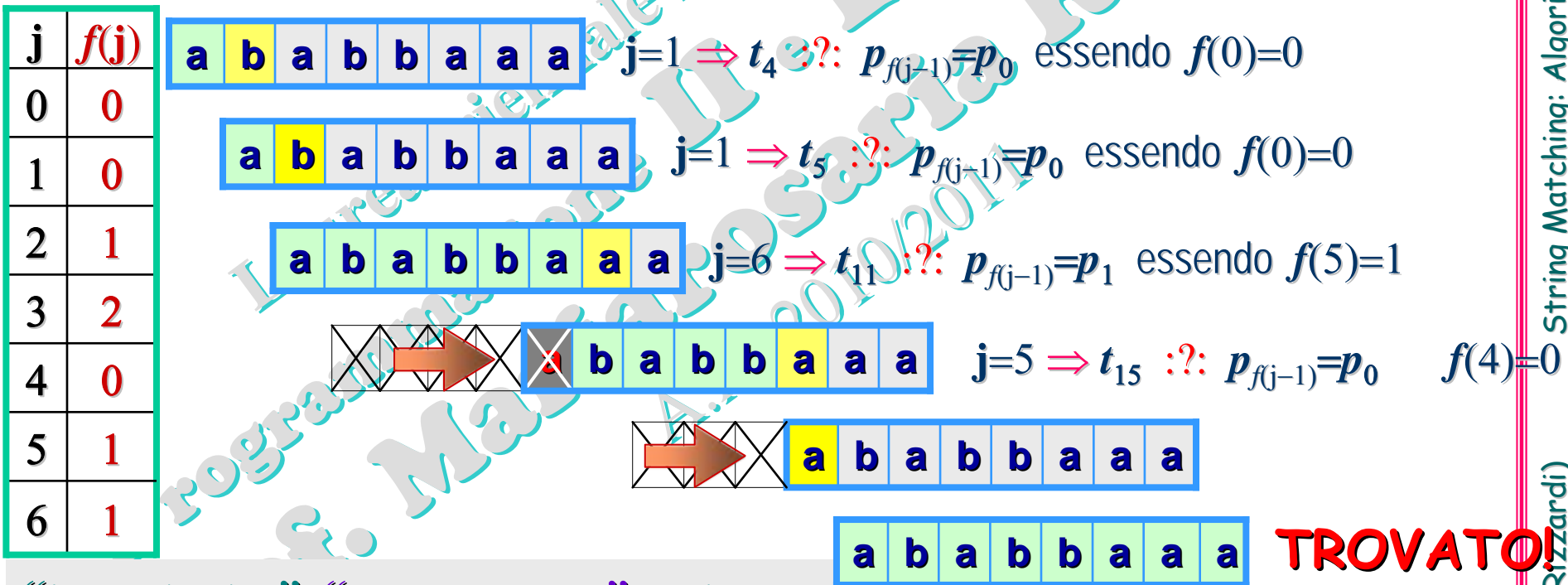
01.18

String Matching: Algoritmo KMP

(prof. M. Rizzardi)



Algoritmo efficiente!



" $t_{k-j} \dots t_{k-2} t_{k-1}$ " = " $p_0 p_1 p_2 \dots p_{j-1}$ " ma $t_k \neq p_j$

next step $\begin{cases} t_k & \text{con } p_{f(j-1)} & \text{se } j \neq 0 \\ t_{k+1} & \text{con } p_0 & \text{se } j = 0 \end{cases}$

Algoritmo kMP (Knuth-Morris-Pratt): valutazione

- ❖ Complessità computazionale nel caso peggiore: $O(MN)$.
- ❖ Complessità computazionale nel caso medio: $O(M+N)$.

Vantaggi

- Indipendenza dall'alfabeto e dalla distribuzione dei suoi caratteri.
- Velocità di esecuzione della ricerca.

Svantaggi

- Degrada all'algoritmo naive nel caso peggiore.

Esercizio

1

Determinare il numero totale di occorrenze di un pattern in un testo, usando per la ricerca l'*algoritmo KMP*.

[liv. 3]