

# Projeto e Análise de Algoritmos

Processamento de Texto

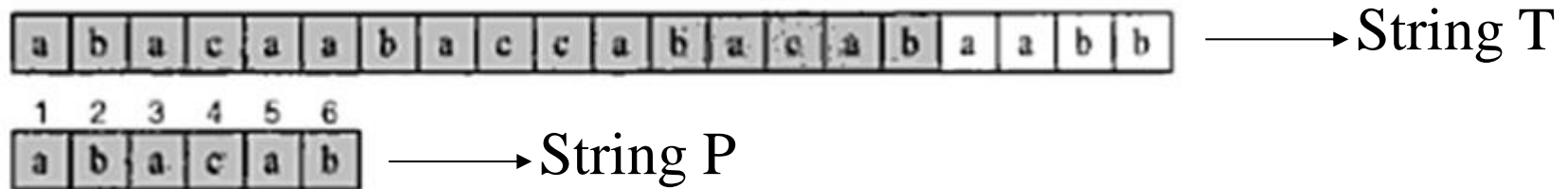
# Aplicações

- Pesquisas na Web (padrões)
- Aplicações linguísticas
- Cadeia de DNA
- Engenharia de software (controle de versões)

# Cadeias de Caracteres

- “CGTAAACTGCTTTAATCAAACGC”  
Cadeias de moléculas chamadas bases:  
adenina (A), guanina (G), citosina (C) e timina (T)
- “Botafogo é campeão brasileiro de 2011!”
- “<http://www.iprj.uerj.br>”
- “if ( (x==0) && (y!=3) )”

# Encontrando padrões em textos (método da força bruta)



- Objetivo: descobrir se o string P (padrão) está em T
- Força bruta: caminha elemento a elemento da esquerda para a direita
- Exercício: escrever o algoritmo

# Encontrando padrões em textos (força bruta)

Algoritmo ForcaBruta(T, P)

// T possui n caracteres e P possui m caracteres

Para  $i \leftarrow 0$  até  $n - m$  faça  $\longrightarrow$   $n - m + 1$  vezes

$j \leftarrow 0$

enquanto ( $j < m$  e  $T[i+j] == P[j]$ ) faça  $\longrightarrow$   $m$  vezes

$j \leftarrow j + 1$

se  $j == m$  então

retorne  $i$

$O(nm)$ . Se  $m = n/2$ ,  $O(n^2)$

Retorne “Não existe substring em T igual a P.”

# Encontrando padrões em textos (força bruta)

a b a c a a b a c c a b a c a b a a b b

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

7  
a b a c a b

8 9  
a b a c a b

10  
a b a c a b

27 comparações

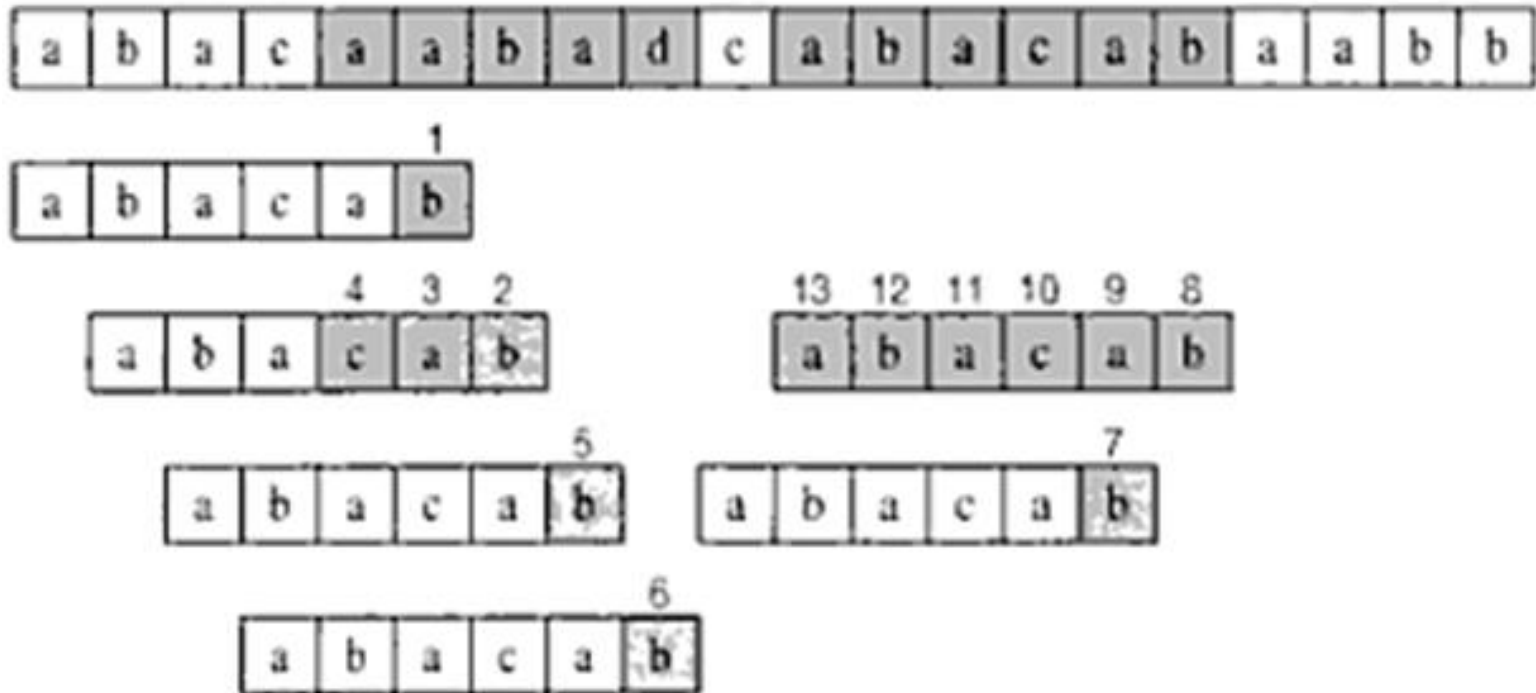
11 comparações

22 23 24 25 26 27  
a b a c a b

# Encontrando padrões em textos (algoritmo Boyer-Moore)

- Objetivo: melhorar o tempo de execução do algoritmo de força bruta
- Heurística do espelho: comparação de forma invertida
- Heurística do salto de caracteres
  - Imaginando  $T$  a sequência,  $P$  o padrão buscado e  $T[i] = c$ 
    - Se  $c$  não pertence a  $P$  então movemos  $P$  completamente para depois de  $T[i]$  (salto grande)
    - Se  $c$  pertence a  $P$ , movemos  $P$  para frente até que uma ocorrência do caracter  $c$  em  $P$  esteja alinhada com  $T[i]$

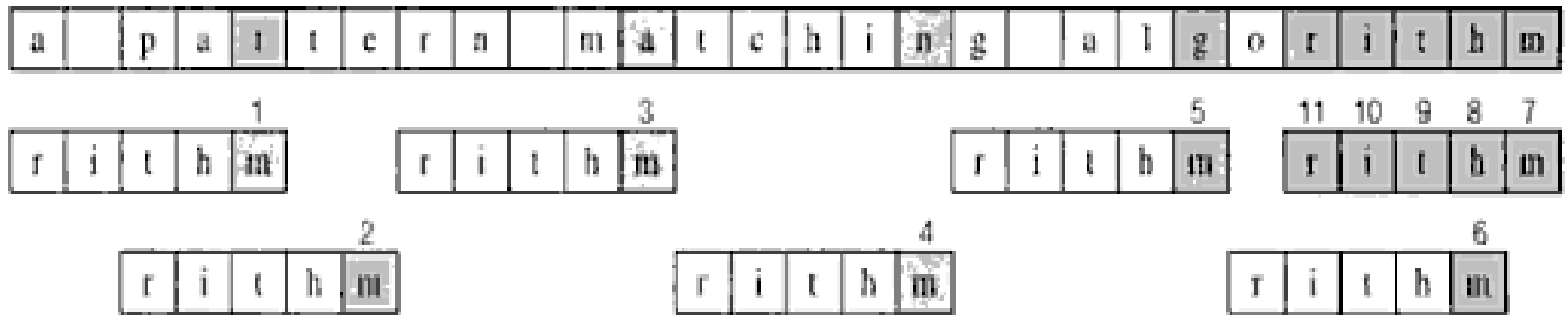
# Encontrando padrões em textos (algoritmo Boyer-Moore)



13 comparações



# Encontrando padrões em textos (algoritmo Boyer-Moore)



Exemplo em língua inglesa – Goodrich & Tamassia

# Encontrando padrões em textos (algoritmo Boyer-Moore)

- Criar função  $\text{last}(c)$  para tratar a heurística
- Se  $c$  está em  $P$ ,  $\text{last}(c)$  retorna o índice da última (mais a direita) ocorrência de  $c$  em  $P$
- Senão,  $\text{last}(c)$  retorna  $-1$

Exercício:  
Pensar e tentar  
fazer o algoritmo todo

A função  $\text{last}(c)$

$c$	a	b	c	d
$\text{last}(c)$	4	5	3	-1

a b a c a b

# Encontrando padrões em textos (algoritmo Boyer-Moore)

Algoritmo **BoyerMoore** (T, P)

// T possui n caracteres e P possui m caracteres

$i \leftarrow m - 1$

$j \leftarrow m - 1$

repita

se  $P[j] = T[i]$  então

se  $j = 0$  então

retorne  $i$  {achamos!}

senão

$i \leftarrow i - 1$

$j \leftarrow j - 1$

senão

$i \leftarrow i + m - \min(j, 1 + \text{last}(T[i]))$  {salto}

$j \leftarrow m - 1$

até  $i > n - 1$

retorne "Não existe *substring* em T igual a P."

Exercício:  
Fazer o chinês