

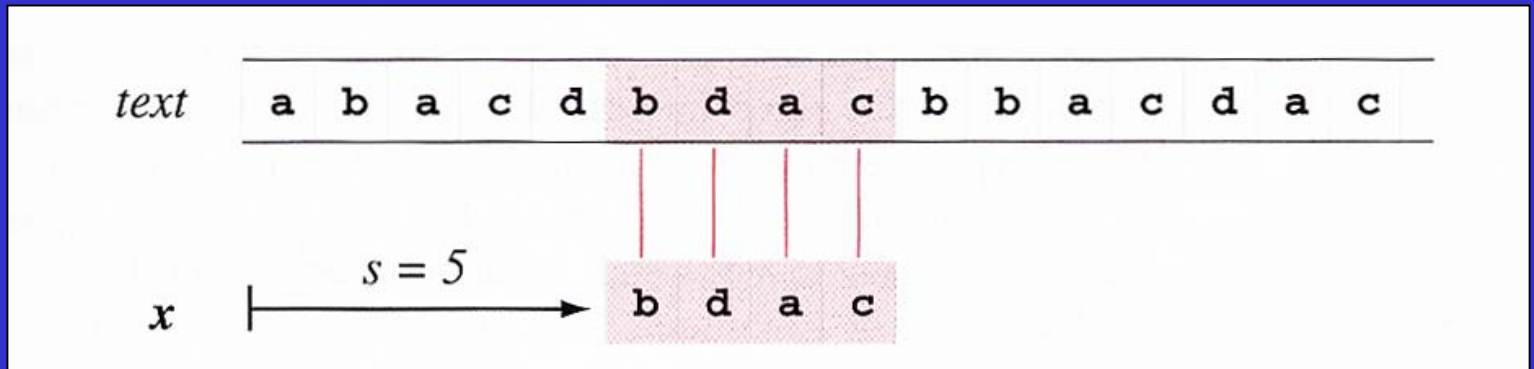
String Matching

Whether a candidate x is a factor of *text*

Definitions

- $\text{Length}[x] = |x|$
- $|\text{text}| \gg |x|$
- Alphabet A
 - Binary, $A = \{0,1\}$
 - Decimal, $A = \{0,1,2\dots,9\}$
 - English Letters, $A = \{a,b,c\dots,z\}$
 - DNA bases, $A = \{A,G,C,T\}$

Naive String Matching



■ **Algorithm 1.** (Naive String Matching)

```
1 begin initialize  $\mathcal{A}, x, text, n \leftarrow \text{length}[text], m \leftarrow \text{length}[x]$ 
2      $s \leftarrow 0$ 
3     while  $s \leq n - m$ 
4         if  $x[1 \dots m] = text[s + 1 \dots s + m]$ 
5             then print "pattern occurs at shift"  $s$ 
6                  $s \leftarrow s + 1$ 
7     return
8 end
```

Naïve String Matching

- $O(n-m+1)$ in the worst case: matches at very end
- Random x and *test* means algorithm is efficient
- Information from one candidate shift is not exploited in subsequent candidate shift

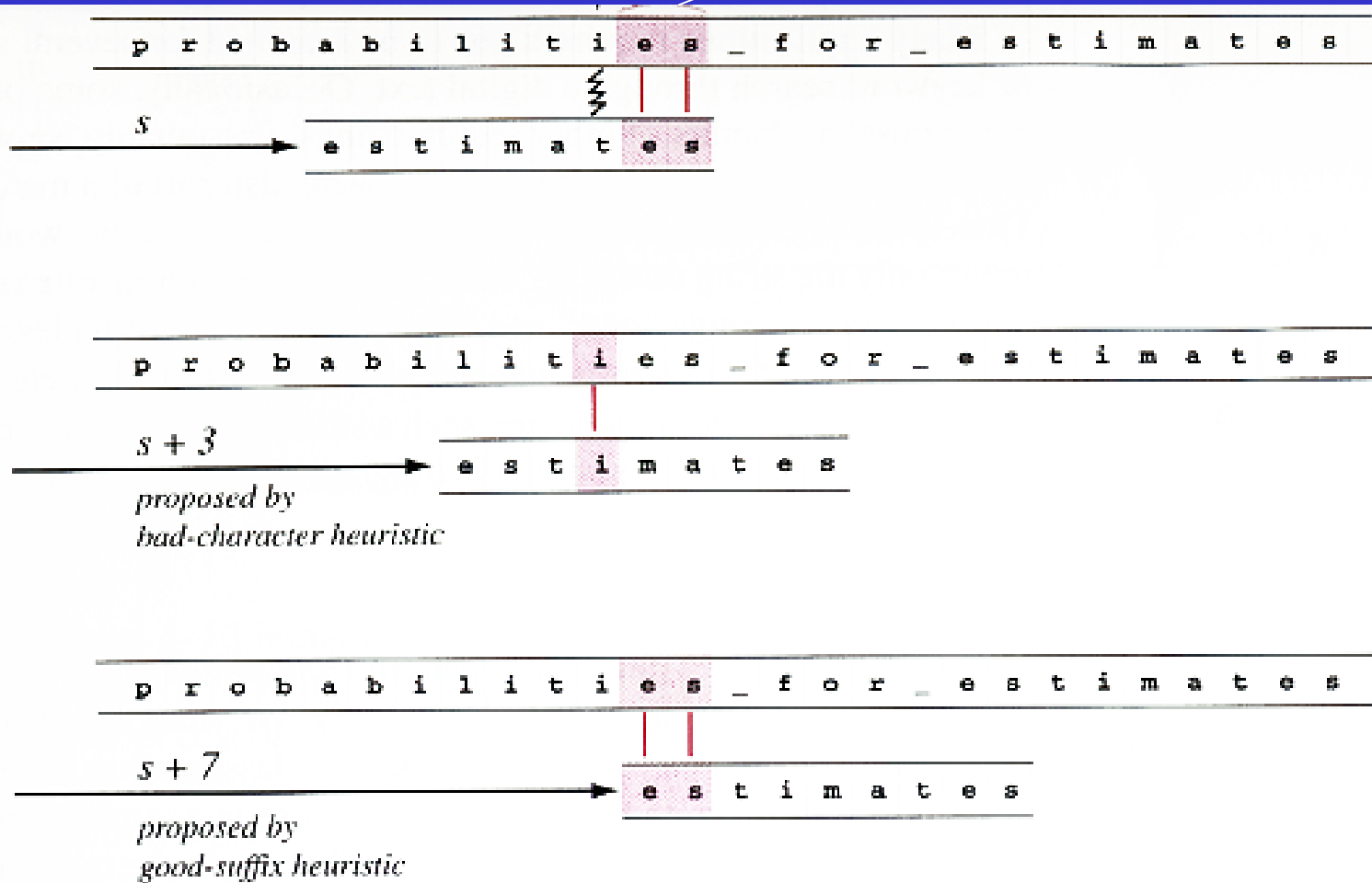
Boyer Moore Algorithm

- Resembles Naïve String Matching
- Except matching is done right-to-left
- Shifts are done in larger increments
- Good suffix heuristic
- Bad character heuristic

Boyer-Moore

Bad Character

Good Suffix



Boyer-Moore principle

- Good Suffix heuristic and Bad Character heuristic operate independently and in parallel
- After mismatch is detected, each proposes an amount by which s can be safely increased without missing a valid shift
- Larger of these shifts is selected

Boyer-Moore Implementation

- Good Suffix Function
 - $G(\mathbf{x})$ is a table that for each suffix gives the second rightmost occurrence in \mathbf{x}
 - For “estimates” s has value 2, es has value 1
- Last occurrence Function
 - Table containing every letter of alphabet and position of its rightmost occurrence in \mathbf{x}
 - For “estimates”, a=6, b=c=d=0, e=8, etc.

■ **Algorithm 2. (Boyer-Moore String Matching)**

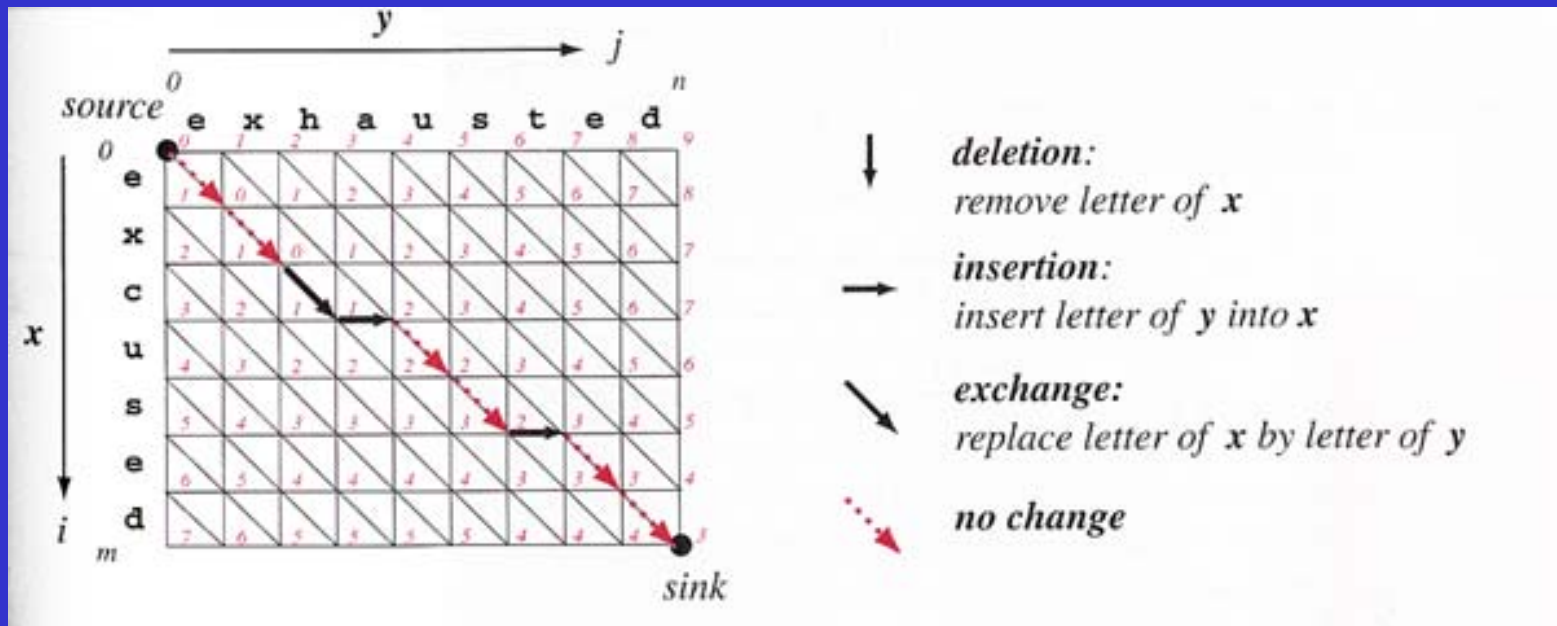
```
1 begin initialize  $\mathcal{A}, \mathbf{x}, \text{text}, n \leftarrow \text{length}[\text{text}], m \leftarrow \text{length}[\mathbf{x}]$ 
2      $\mathcal{F}(\mathbf{x}) \leftarrow$  last-occurrence function
3      $\mathcal{G}(\mathbf{x}) \leftarrow$  good-suffix function
4      $s \leftarrow 0$ 
5     while  $s \leq n - m$ 
6         do  $j \leftarrow m$ 
7         while  $j > 0$  and  $\mathbf{x}[j] = \text{text}[s + j]$ 
8             do  $j \leftarrow j - 1$ 
9             if  $j = 0$ 
10                 then print "pattern occurs at shift"  $s$ 
11                      $s \leftarrow s + \mathcal{G}(0)$ 
12                 else  $s \leftarrow s + \max[\mathcal{G}(j), j - \mathcal{F}(\text{text}[s + j])]$ 
13     return
14 end
```

Edit Distance

- Not clear whether “abbccc” is closer to “aabbcc” or to “abbcccb”
- Edit distance is the number of fundamental operations required to transform one string to the other

Cost Matrix

$C[i, j]$ = cost of transforming x to y



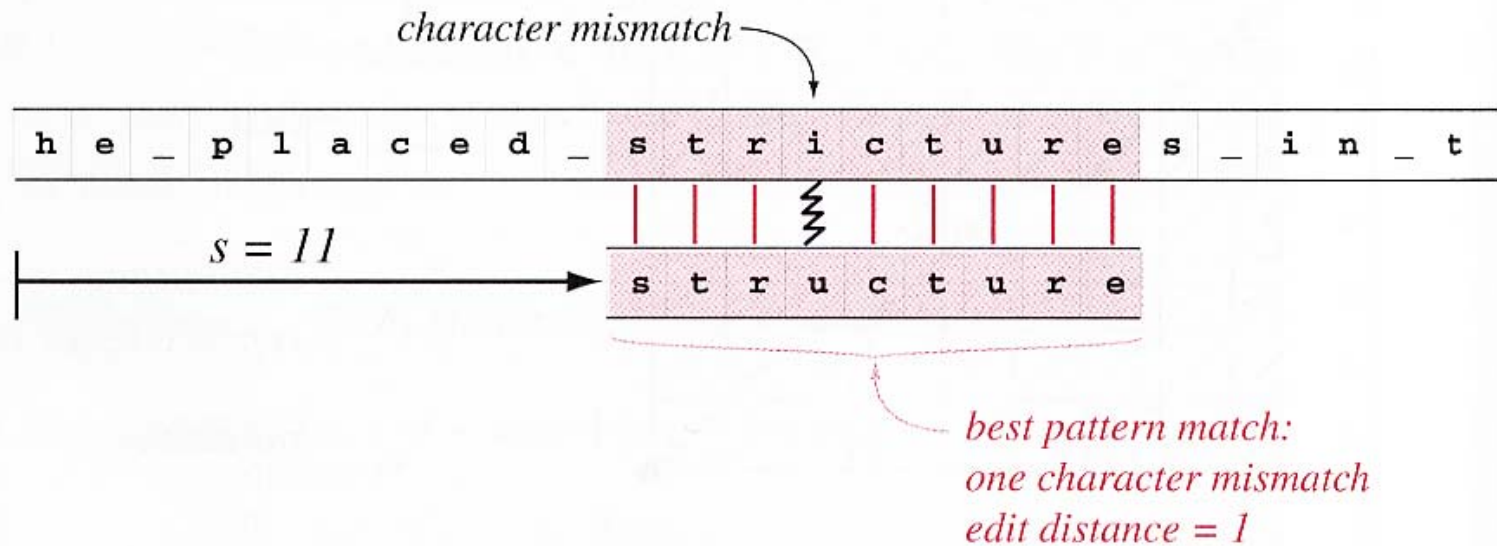
Edit Distance (excused, exhausted) = $C[i = 7, j = 9] = 3$

x	excused	source string	$C[0, 0] = 0$
	ex <u>h</u> used	substitute h for c	$C[3, 3] = 1$
	ex <u>h</u> aused	insert a	$C[3, 4] = 2$
	exhaust <u>e</u> d	insert t	$C[5, 7] = 3$
y	exhausted	target string	$C[7, 9] = 3$

■ Algorithm 3. (Edit Distance)

```
1 begin initialize  $\mathcal{A}$ ,  $x$ ,  $y$ ,  $m \leftarrow \text{length}[x]$ ,  $n \leftarrow \text{length}[y]$ 
2    $C[0, 0] \leftarrow 0$ 
3    $i \leftarrow 0$ 
4   do  $i \leftarrow i + 1$ 
5      $C[i, 0] \leftarrow i$ 
6   until  $i = m$ 
7   do  $j \leftarrow j + 1$ 
8      $C[0, j] \leftarrow j$ 
9   until  $j = n$ 
10   $i \leftarrow 0$ ;  $j \leftarrow 0$ 
11  do  $i \leftarrow i + 1$ 
12    do  $j \leftarrow j + 1$ 
13       $C[i, j] = \min \left[ \underbrace{C[i - 1, j] + 1}_{\text{insertion}}, \underbrace{C[i, j - 1] + 1}_{\text{deletion}}, \underbrace{C[i - 1, j - 1] + 1 - \delta(x[i], y[j])}_{\text{no change/exchange}} \right]$ 
14    until  $j = n$ 
15  until  $i = m$ 
16  return  $C[m, n]$ 
17 end
```

String Matching with Errors



String Matching with Don't Care

