

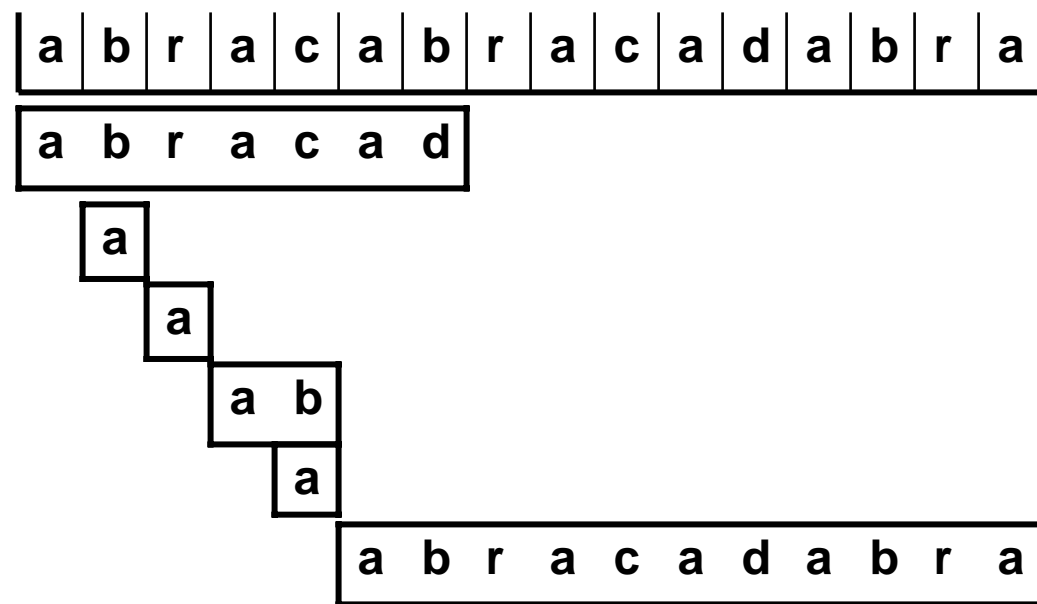
Text Searching

- Thus far, we've seen how to pre-process texts (stemming, "stopping" words, compressing, etc)
- We've also seen how to build, process, improve and assess the quality of queries and the returned answer sets
- We also saw how to effectively index the texts
- Now, how do we search text ? It goes beyond searching the indices, specially in the cases where the indices are not enough to answer queries (e.g., proximity queries)

Sequential Searching

- If no index is used this is the only option, however sometimes (e.g., when blocking is used) it is necessary even if an index exists
- The problem is “simple”
 - Given a pattern x of length m and a text y of length n ($n \gg m$) find all positions in y where x occurs
- There is much more work on this than we can cover, including many theoretical results. Thus we will discover some well known approaches

Brute Force



BRUTE_FORCE_MATCHER(x,m,y,n)

1. For p from 0 up to n-m
2. loop if $y[p..p+m-1]=x$
3. then report p

- The name says it all ...
- Starting from all possible initial positions (i.e., all positions), check whether the pattern could start at that position
- It takes $O(mn)$ time in the worst case:
 - For instance when $\alpha^{m-1}b$ is searched in α^n for any two symbols α, b .
- $O(n)$ in the average case – not that bad
- The most important thing is that it suggests the use of a sliding window over the text. The idea is to see whether we can see the pattern through the window

The Karp-Rabin Algorithm

- Hashing: a simple method for avoiding a quadratic number of symbol comparisons in most practical situations.
- At each position p of the window on the text, check if the part of the text delimited by the window $y[p \dots p + m - 1]$ "looks like" x .
- Use a hash function to check the resemblance
- The hash function should be highly discriminating for strings.
- The function should also have the following properties:
 - To be efficiently computable;
 - Easy computation of the value associated with the next part of the text from the value associated with the current part.
- If Σ is the alphabet of input symbols, the hash function h :

$$h(u) = \left(\sum_{i=0}^{|u|-1} u[i] \times d^{|u|-1-i} \right) \bmod q$$

where q and d are two constants. Then for each $v \in \Sigma^*$ for each symbols $a', a'' \in \Sigma$, $h(va'')$ is computed from $h(a'v)$ by the formula

$$h(va'') = \left(\left(h(a'v) - a' \times d^{|v|} \right) \times d + a'' \right) \bmod q$$

- Search for pattern x : compare the value $h(x)$ with the hash value associated each part of length m of text y .
- If the two values are equal: necessary to check whether this part of the text is equal to x or not by symbol comparisons

The Karp-Rabin algorithm (cont.)

KARP-RABIN-MATCHER(x, m, y, n)

```

1   $r \leftarrow 1$ 
2   $s \leftarrow x[0] \bmod q$ 
3   $t \leftarrow 0$ 
4  for  $i$  from 1 up to  $m - 1$ 
5      loop  $r \leftarrow (r \times d) \bmod q$ 
6           $s \leftarrow (s \times d + x[i]) \bmod q$ 
7           $t \leftarrow (t \times d + y[i - 1]) \bmod q$ 
8  for  $p$  from 0 up to  $n - m$ 
9      loop  $t \leftarrow (t \times d + y[p + m - 1]) \bmod q$ 
10         if  $t = s$  and  $y[p..p + m - 1] = x$ 
11             then report  $p$ 
12          $t \leftarrow ((c - 1) \times q + t - y[p] \times r) \bmod q$ 

```

- The value of symbols ranges from 0 to $c-1$
- The quantity $(c-1) \times q$ is added in line 8 to provide correct computations on positive integers
- Convenient values for d are powers of 2:
 - in this case, all the products by d can be computed as shifts on integers.
- The value of q is generally a large prime (such that the quantities $(q - 1) \times d + c - 1$ and $c \times q - 1$ do not cause overflows),

p	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$y[p]$	n	o	u	d	e	f	e	n	s	e	u	f	o	r	u	s	e	n	s	e
$h(y[p..p+4])$	8	8	6	28	9	18	28	26	22	12	17	24	16	0	1	9	—	—	—	—

- Searching for the pattern $x = \text{sense}$ in the text $y = \text{“no defense for sense”}$:
 - $c = 256$ (ASCII code), and the values of
 - q and d are set respectively to 31 and 2.
 - $h(x) = (115 \times 16 + 101 \times 8 + 110 \times 4 + 115 \times 2 + 101) \bmod 31 = 9$.
 - $h(y[4.. 8])=h(y[15.. 19])= h(x)$: two string-to-string comparisons against x are performed.

The Knuth-Morris-Pratt Algorithm

- The first algorithm with linear worst-case behaviour,
- On average it is not much faster than Brute-Force (BF) algorithm.
- Sliding window over the text: however, it does not try all window positions as BF does.
- It reuses information from previous checks.
- After the window is checked, a number of pattern letters were compared to the text window, and they all matched except possibly the last one compared.
- Hence, there is a *prefix* of the pattern that matched the text.
- The algorithm takes advantage of this information to avoid trying window positions which can be deduced not to match.
- The pattern is preprocessed in $O(m)$ time and space to compute a function called ψ .
- $\psi(i)$: the length of the longest proper prefix of $x_{0..i-1}$ which is also a suffix and the characters following prefix and suffix are different.

$$\psi[i] = \max \{k \mid (0 \leq k < i, x[i-k..i-1] = x[0..k-1], x[k] \neq x[i]) \text{ or } (k = -1)\}$$

where $i \in \{0, 1, \dots, m-1\}$ and

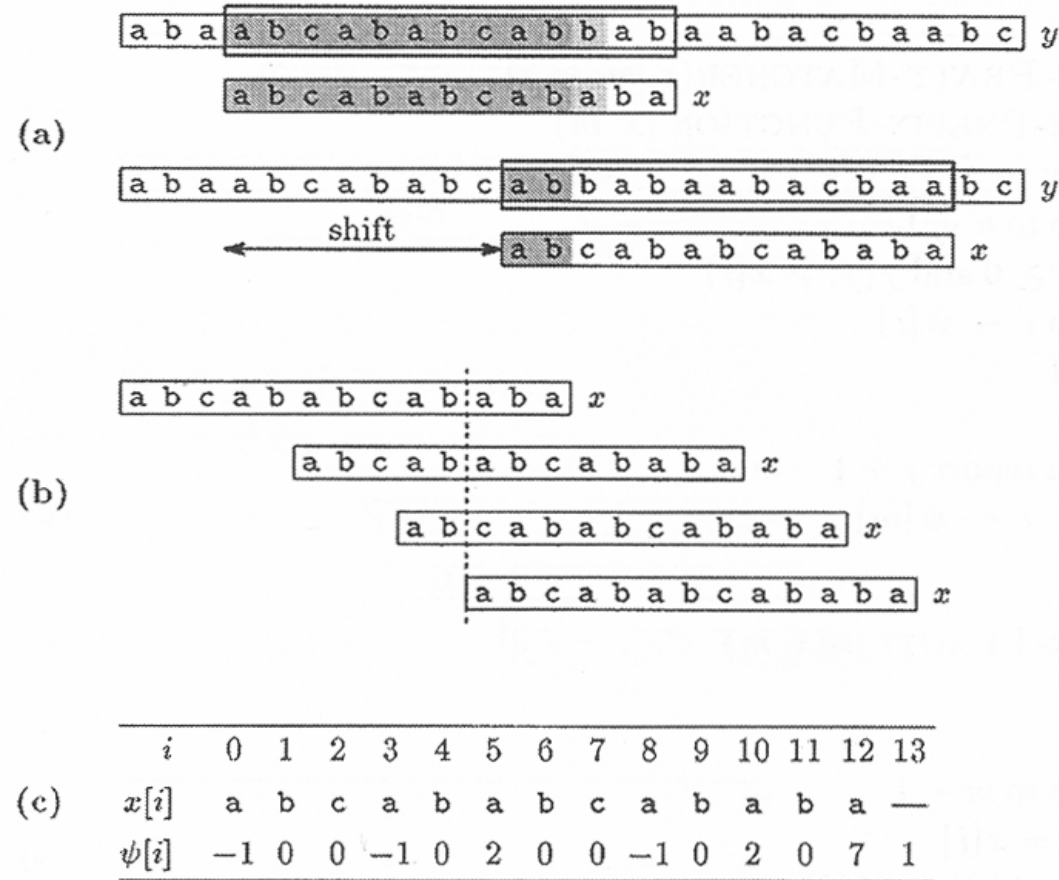
$$\psi[m] = \max \{k \mid 0 \leq k < m, x[m-k..m-1] = x[0..k-1]\}$$

- Hence $i - \psi[i]$ window positions can be safely skipped if the characters up to $i - 1$ matched, and the i -th did not:
 - Crucial observation: this shift depend only on the pattern, because if the text in the window matched up to position $i - 1$, then that text is equal to the pattern

The Knuth-Morris-Pratt Algorithm (Cont.)

```

KNUTH-MORRIS-PRATT-MATCHER( $x, m, y, n$ )
1   $\psi \leftarrow$  BETTER-PREFIX-FUNCTION( $x, m$ )
2   $i \leftarrow 0$ 
3  for  $j$  from 0 up to  $n - 1$ 
4    loop while  $i \geq 0$  and  $y[j] \neq x[i]$ 
5      loop  $i \leftarrow \psi[i]$ 
6     $i \leftarrow i + 1$ 
7     $i = m$ 
8      then report  $j + 1 - m$ 
9       $i \leftarrow \psi[m]$ 
    
```



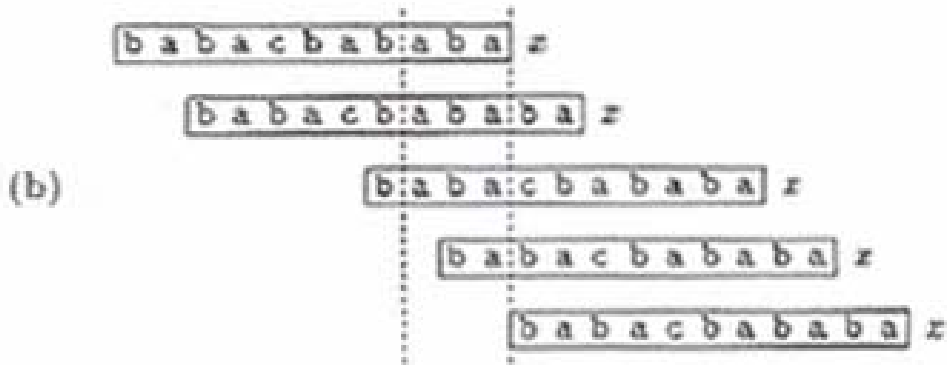
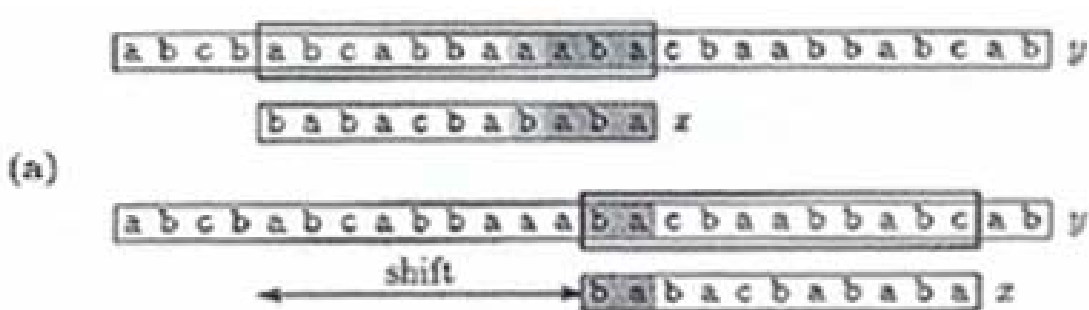
- Searching for the pattern $x = \text{abcababcbabab}$.
- (a) The window on the text y at position 3. A mismatch occurs at position 10 on x . The matching symbols are shown darkly shaded, and the current analyzed symbols lightly shaded. We should shift the window 8 positions to the right. The string-to-string comparison resumes at position 2 on the pattern.
- (b) The prefixes of x that are suffixes of $x[0..9] = \text{abcababcbab}$ are right-aligned along the discontinuous vertical line. String $x[0..4] = \text{abcab}$ is a suffix of $x[0..9]$, but is followed by symbol a which is identical to $x[10]$. String $x[0..1]$ is the expected prefix, since it is followed by symbol c .
- (c) The values of the function ψ for pattern x .

The Boyer-Moore Algorithm

- The most efficient string-matching algorithm in usual applications.
- A simplified version of it, or the entire algorithm, is often implemented in text editors for the "search" and "substitute" commands.
- The scan operation proceeds from right to left in the window on the text, instead of left to right as in the Knuth-Morris-Pratt algorithm.
- In case of a mismatch, the algorithm uses two functions to shift the window:
 - good-suffix shift function
 - bad-symbol shift function.
- Good-suffix shift function:
 - Let p be the current (left) position of the window on the text. Assume that a mismatch occurs between symbols $y[p + i]$ and $x[i]$ for some i , $0 \leq i \leq m-1$:
 - $y[p + i] \neq x[i]$ and $y[p + i + 1.. p + m - 1] = x [i + 1 .. m - 1]$.
 - Good-suffix shift: align $y [p + i + 1 .. p + m - 1]$ with its rightmost occurrence $x [k + 1 .. m - 1 - i + k]$ in x preceded by a symbol $x[k]$ different from $x[i]$ to avoid an immediate mismatch.
 - If no such occurrence exists, align the longest suffix of $y[p + i + 1 ... p + m - 1]$ with a matching prefix of x . The good-suffix shift function β is defined by

$$\beta[i] = \min\{i - k \mid (0 \leq k < i, x[k + 1...m-1-i + k] = x[i+1...m - 1], x[k] \neq x[i]) \\ \text{or } (i - m < k < 0, x = x [i - k...m - 1]x[m - i + k...m - 1]) \text{ or } (k = i - m)\}$$

The Boyer-Moore Algorithm



(c)

i	0	1	2	3	4	5	6	7	8	9	10
$x[i]$	b	a	b	a	c	b	a	b	a	b	a
$\beta[i]$	7	7	7	7	7	7	2	9	4	11	1

- The good-suffix shift in the Boyer-Moore algorithm when searching for the pattern $x = \text{babacbabababa}$.
- (a) The window on the text is at position 4. The string-to-string comparison, which proceeds from right to left, stops with a mismatch at position 7 on x .
- The window is shifted 9 positions to the right to avoid an immediate mismatch.
- (b) Indeed, the string $x[8..10] = \text{aba}$ is repeated three times in x , but is preceded each time by symbol $x[7] = \text{b}$. The expected matching substring in x is then the prefix ba of x . The substrings of x identical with aba and the prefixes of x ending with a suffix of aba are right-aligned along the rightmost discontinuous vertical line.
- (c) The values of the shift function β for pattern x .

The Boyer-Moore Algorithm (cont.)

- Bad-symbol shift function:
 - Text symbol $y[p + i]$ that causes a mismatch.
 - Assume that this symbol occurs in $x[0 \dots i - 1]$.
 - let k be the position of the rightmost occurrence of $y[p + i]$ in $x[0 \dots i - 1]$.
 - The window can be shifted $i - k$ positions to the right without missing an occurrence of x in y . Assume now that
 - If symbol $y[p + i]$ does not occur in $x[0 \dots i - 1]$, no occurrence of x in y can overlap the position $p + i$ on the text, and thus, the window can be shifted $i + 1$ positions to the right.
 - Let δ be the table indexed on alphabet Σ and $\{0, \dots, m - 1\}$, and defined for each symbol $a \in \Sigma$ and i by

$$\delta[a, i] = \min \{m\} \cup \{m - 1 - j \mid 0 \leq j < i, x[j] = a\}$$

Searching for the pattern $x = \text{babacbababa}$.

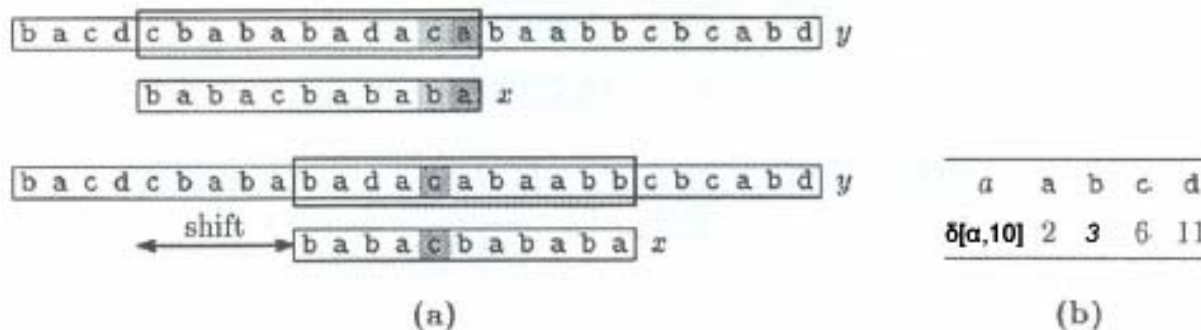
a) The window on the text at position 4.

Mismatch at position 9 on x .

$y[13] = c \rightarrow$ shift the window 5 positions to the right.

Notice that if the unexpected symbol were a or d , the applied shift would have been 1 and 10 respectively.

b) The values of the table δ for pattern X when alphabet Σ is reduced to $\{a, b, c, d\}$.

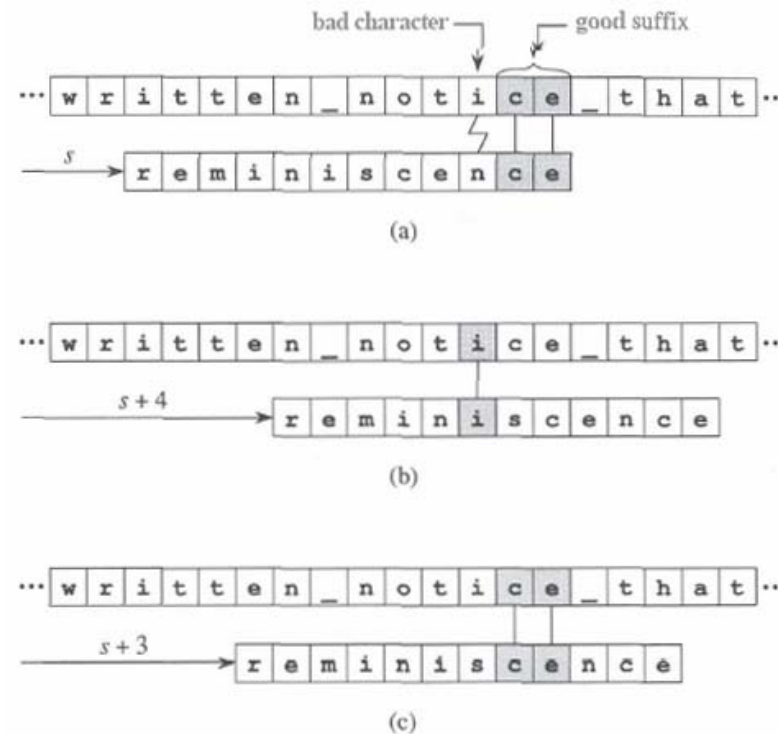


The Boyer-Moore Algorithm (cont.)

```
BOYER-MOORE-MATCHER( $x, m, y, n$ )
1   $\beta \leftarrow \text{GOOD\_SUFFIX\_FUNCTION}(x, m)$ 
2   $\delta \leftarrow \text{BAD\_SYMBOL\_FUNCTION}(x, m)$ 
3   $p \leftarrow 0$ 
4  while  $p \leq n - m$ 
5      loop  $i \leftarrow m - 1$ 
6          while  $i \geq 0$  and  $y[p + i] = x[i]$ 
7              loop  $i \leftarrow i - 1$ 
8          if  $i \geq 0$ 
9              then  $p \leftarrow p + \max\{\beta[i], \delta[y[p + i], i] + i - m + 1\}$ 
10             else report  $p$ 
11              $p \leftarrow p + \beta[0]$ 
```

- Pattern preprocessing time and space: $O(m + |\Sigma|)$ where $|\Sigma|$ is the size of the alphabet Σ
- $O(nm + |\Sigma|)$ worst-case run time, $O((n \log m)/m)$ on average
- The worst case is unlikely in English text
- Boyer-Moore's algorithm is significantly faster than the brute-force algorithm on English text

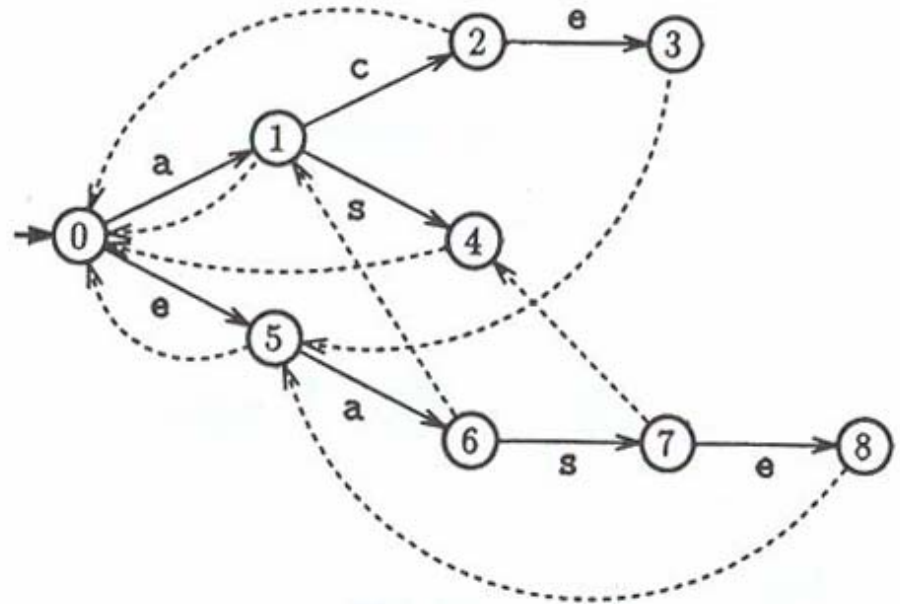
Another example of the Boyer-Moore heuristics.



- (a) Matching the pattern "reminiscence" against a text by comparing characters in a right-to-left manner. The shift s is invalid. Although a "good suffix" "ce" of the pattern matched correctly against the corresponding characters in the text (matching characters are shown shaded), the "bad character" "i", which didn't match the corresponding character "n" in the pattern, was discovered in the text.
- (b) The bad-symbol heuristic proposes moving the pattern to the right, if possible, by the amount that guarantees that the bad text character will match the rightmost occurrence of the bad character in the pattern. In this example, moving the pattern 4 positions to the right causes the bad text character "i" in the text to match the rightmost "i" in the pattern, at position 6. If the bad character doesn't occur in the pattern, then the pattern may be moved completely past the bad character in the text.
- (c) With the good suffix heuristic, the pattern is moved to the right by the least amount that guarantees that any pattern characters that align with the good suffix "ce" previously found in the text will match those suffix characters. In this example, moving the pattern 3 positions to the right satisfies this condition. Since the good suffix heuristic proposes a movement of 3 positions, which is smaller than the 4-position proposal of the bad-character heuristic, the Boyer-Moore algorithm increases the shift by 4.

The Aho-Corasick algorithm

Aho-Corasick trie example for the set 'ace', 'as' and 'ease' showing all failure transitions



- An extension of Knuth-Morris-Pratt in matching a set of patterns.
- The patterns are arranged in a trie-like data structure.
- Each trie node corresponds to matching a prefix of some pattern(s).
- There is also a set of *failure* transitions:
 - Those transitions go between nodes of the trie.
 - A transition from a node representing the prefix z to a node representing a prefix y :
 - y is the longest prefix in the set of patterns which is also a proper suffix of z .

The Aho-Corasick algorithm (cont.)

a	c	e	a	s	e	a	c	a	s
---	---	---	---	---	---	---	---	---	---

a	c	e
---	---	---

e	a	s	e
---	---	---	---

e	a
---	---

a	c
---	---

a	s
---	---

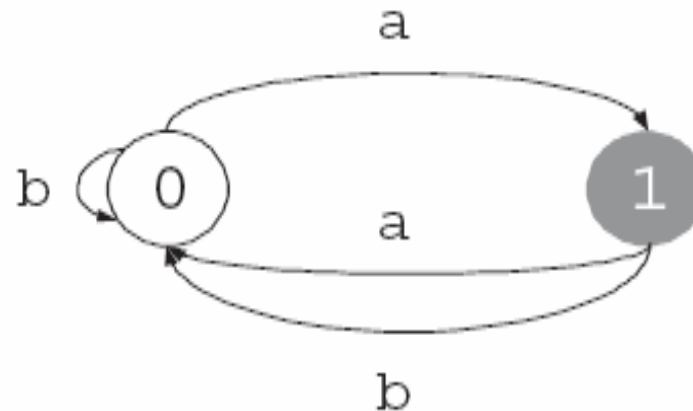
Finite Automaton

- **Definition 3** A finite automaton is a 5-tuple, $M = (Q, q_0, A, \Sigma, \delta)$, where
- Q is a finite set of states.
- $q_0 \in Q$ is the start state.
- $A \subseteq Q$ is the distinguished set of accepting states.
- Σ is a finite input alphabet.
- $\delta : Q \times \Sigma \rightarrow Q$ is the transition function of M .
- A finite automaton M induces a function $\varphi : \Sigma^* \rightarrow Q$ (called final-state function) such that
 - $\varphi(\varepsilon) = q_0$
 - $\varphi(wa) = \delta(\varphi(w), a)$ for $w \in \Sigma^*$, $a \in \Sigma$
- $\varphi(w)$ is the state M ends up after scanning the string w
- M accepts w if and only if $\varphi(w) \in A$

Example of finite automaton

- Given an input alphabet $\Sigma = \{a, b\}$ and state set $Q = \{0, 1\}$, where $q_0 = 0$ and $A = \{1\}$. The transition function is defined as follows:

δ	a	b
0	1	0
1	0	0



- The automaton accepts the string that end in an odd number of a's

Suffix Function

- **Definition.** Given pattern $P[1 \dots m]$, $\sigma : \Sigma^* \rightarrow \{0, 1, \dots, m\}$ as follows:

$$\sigma(x) = \max\{k \mid P_k \text{ is a suffix of } x\}$$

is called a suffix function corresponding to P .

- **Example.** For $P = ab$, we have $\sigma(\varepsilon) = 0$, $\sigma(ccaca) = 1$, $\sigma(ccab) = 2$.
- **Properties:**
 - If x is a suffix of y , then $\sigma(x) \leq \sigma(y)$.
 - $\sigma(x) = m$ iff P is a suffix of x .

String-matching Automaton

- The state set Q is $\{0, 1, \dots, m\}$, where
- $q_0 = 0$ and $A = \{m\}$.
- The transition function δ is defined as

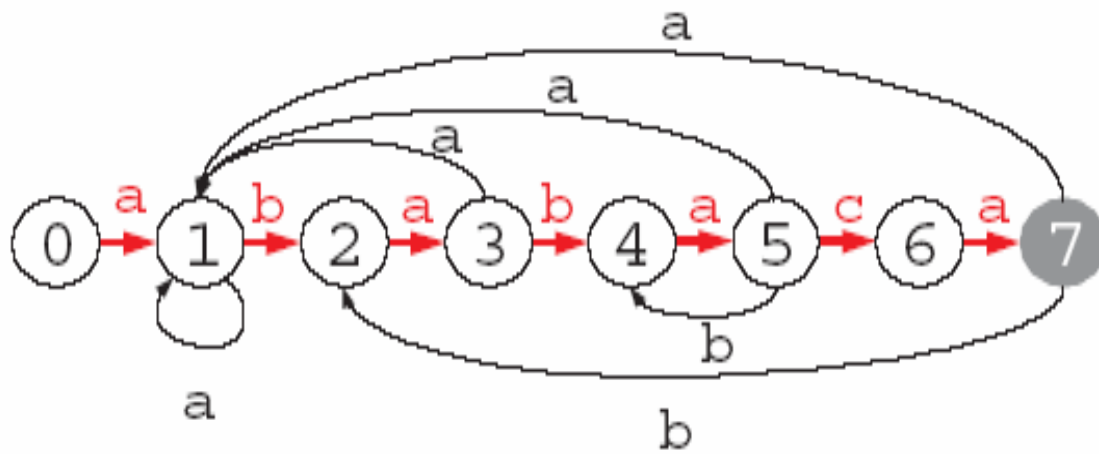
$$\delta(q, a) = \sigma(P_q a)$$

for any state q and character a .

- **Theorem.** If φ is the final-state function of a string-matching automaton for a given pattern P and $T[1 \dots n]$ is an input text for the automaton, then for $i = 0, 1, \dots, n$, we have

$$\varphi(T_i) = \sigma(T_i)$$

Example of String-matching Automaton



state	<i>a</i>	<i>b</i>	<i>c</i>	<i>P</i>
0	1	0	0	<i>a</i>
1	1	2	0	<i>b</i>
2	3	0	0	<i>a</i>
3	1	4	0	<i>b</i>
4	5	0	0	<i>a</i>
5	1	4	6	<i>c</i>
6	7	0	0	<i>a</i>
7	1	2	0	

<i>i</i>	—	1	2	3	4	5	6	7	8	9	10	11
<i>T</i> [<i>i</i>]	—	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>b</i>	<i>a</i>
$\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

Finite Automaton Algorithm

FINITE-AUTOMATON-MATCHER(T, δ, m)

```
1  $n \leftarrow \text{length}[T]$ 
2  $q \leftarrow 0$ 
3 for  $i \leftarrow 1$  to  $n$ 
4     do  $q \leftarrow \delta(q, T[i])$ 
5         if  $q = m$ 
6             then print "occur with shift"  $i - m$ 
```

COMPUTING-TRANSITION-FUNCTION(P, Σ)

```
1  $m \leftarrow \text{length}[P]$ 
2 for  $q \leftarrow 0$  to  $m$ 
3     do for each character  $a \in \Sigma$ 
4         do  $k \leftarrow \min(m + 1, q + 2)$ 
5             repeat  $k \leftarrow k - 1$ 
6                 until  $P_k$  is a suffix of  $P_q a$ 
7                  $\delta(q, a) \leftarrow k$ 
8 return  $\delta$ 
```

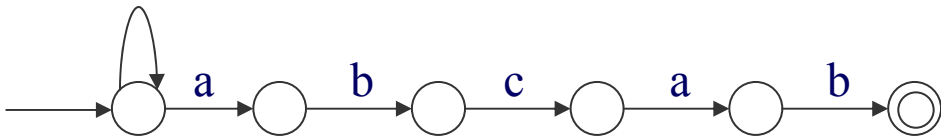
- The matching time is $\Theta(n)$
- The transition function δ can be computed in $O(m^3|\Sigma|)$ time

Shift-Or Algorithm

- Using bit-parallelism
 - simulate the operation of non-deterministic automaton that searches the pattern in the text.
- algorithm
 - builds a table B which for each character stores a bit mask $b_m \dots b_1$.
 - The mask in $B[c]$ has the i -th bit set to zero iff $p_i = c$.
 - The state of the search is kept in a machine word $D = d_m \dots d_1$
$$D' \leftarrow (D \ll 1) \mid B[T_j]$$
 - \ll : shifting all the bits in D one position to the left and setting the right most bit to zero
 - A match is reported whenever d_m is zero.

Shift Or Example(1)

text	d	a	b	c	a	b	c	a
pattern	a	b	c	a	b			



B[a]	0	1	1	0	1
B[b]	1	0	1	1	0
B[c]	1	1	0	1	1
B[*]	1	1	1	1	1

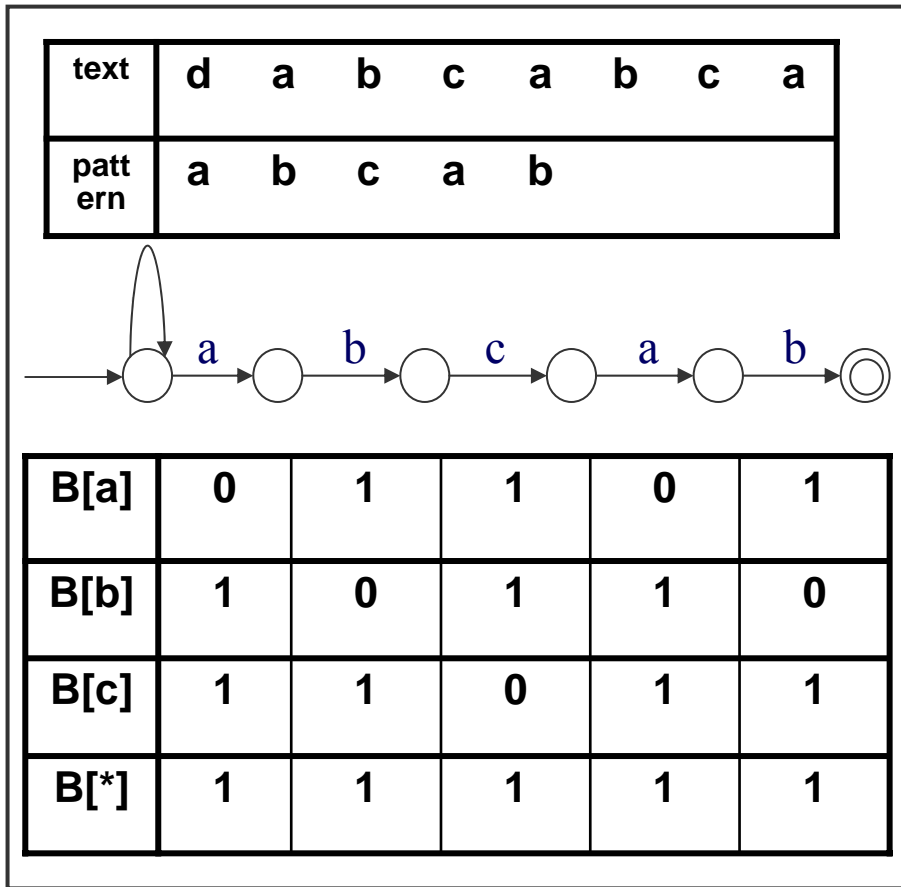
D<<1	1	1	1	1	0
B[d]	1	1	1	1	1
D'	1	1	1	1	1

D<<1	1	1	1	1	0
B[a]	1	0	1	1	0
D'	1	1	1	1	0

D<<1	1	1	1	0	0
B[b]	0	1	1	0	1
D'	1	1	1	0	1

D<<1	1	1	0	1	0
B[c]	1	1	0	1	1
D'	1	1	0	1	1

Shift-Or Example(2)



D<<1	1	0	1	1	0
B[a]	1	0	1	1	0
D'	1	0	1	1	0

D<<1	0	1	1	0	0
B[b]	0	1	1	0	1
D'	0	1	1	0	1

D<<1	1	1	0	1	0
B[c]	1	1	0	1	1
D'	1	1	0	1	1

D<<1	1	0	1	1	0
B[a]	1	0	1	1	0
D'	1	0	1	1	0

match

String Matching Allowing Errors

- Approximate string matching:
 - given a string $P=p_1p_2\dots p_m$ of length m , a string $T=t_1t_2\dots t_n$ of length n , match P to T with minimum number of errors (Levenshtein distance).
 - The classical solution to approximate string matching is based on dynamic programming:
 - Matrix $C[0..m, 0..n]$ where $C[i,j]$ represents the minimum number of errors needed to match $P_{1..i}$ to $T_{1..j}$

$$C[i, j] = \begin{cases} 0 & \text{if } i = 0, j = 0 \\ C[i-1, 0] + D(p_i) & \text{if } j = 0 \\ C[0, j-1] + I(t_j) & \text{if } i = 0 \\ \min\{C[i-1, j-1] + S(p_i, t_j), C[i-1, j] + D(p_i), C[i, j-1] + I(t_j)\} & \text{otherwise} \end{cases}$$

where $S(p_i, t_j)$ is the cost of substituting p_i with t_j , $D(p_i)$ the cost of deleting p_i and $I(t_j)$ the cost of inserting t_j .

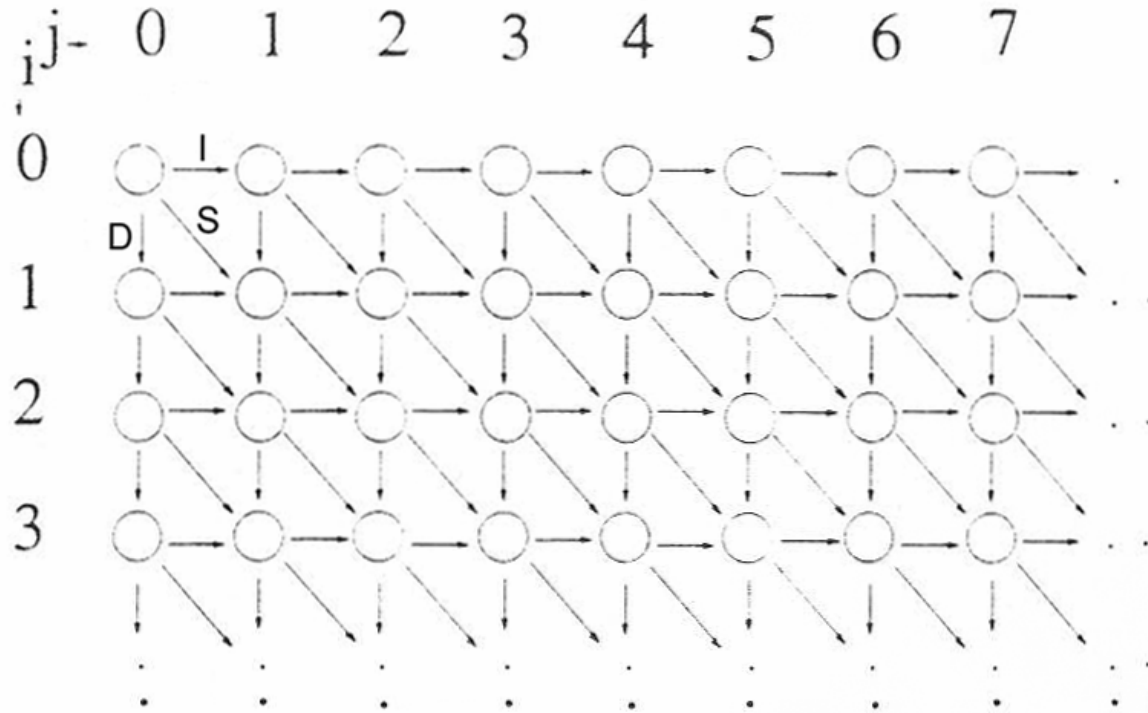
Common assumption: $D(p_i) = I(t_j) = 1$, $S(p_i, t_j) = \begin{cases} 1 & \text{if } p_i \neq t_j \\ 0 & \text{otherwise} \end{cases}$.

String Matching Allowing Errors (Cont.)

		s	u	r	g	e	r	y
	0	1	2	3	4	5	6	7
s	1	0	1	2	3	4	5	6
u	2	1	0	1	2	3	4	5
r	3	2	1	0	1	2	3	4
v	4	3	2	1	1	2	3	4
e	5	4	3	2	2	1	2	3
y	6	5	4	3	3	2	2	2

- The dynamic programming algorithm to compute the edit distance between “survey” and “surgery”. The bold entries show the path to the final result
- $O(mn)$ the running time of the algorithm
- Only $O(\min(m,n))$ the space requirement: only the previous column or row must be stored in order to compute the new column or row respectively.

String Matching Allowing Errors (cont.)



- The dependencies in the dynamic programming recurrence may be represented by a lattice graph
- The vertex in position (i,j) of the lattice graph represents entry (i,j) of the cost matrix.
- Each edge of the lattice graph is assigned a weight equal to the cost of the corresponding edit operation.
- The weights are obtained as follows:
 - The weight of an edge of type $\langle (i-1,j), (i,j) \rangle$ is $D(p_i)$
 - An edge of type $\langle (i-1,j-1), (i,j) \rangle$ has weight $S(p_i, t_j)$
 - The weight of an edge of type $\langle (i, j-1), (i,j) \rangle$ is $I(t_j)$
- $Cost(i,j)$ is the length of a shortest path from vertex $(0,0)$ to vertex (i,j)
- A match is reported at text positions j such that $C[m,j] \leq k$
- $O(mn)$ the execution time of the algorithm
- $O(m)$ space: only the previous column of the matrix is needed.

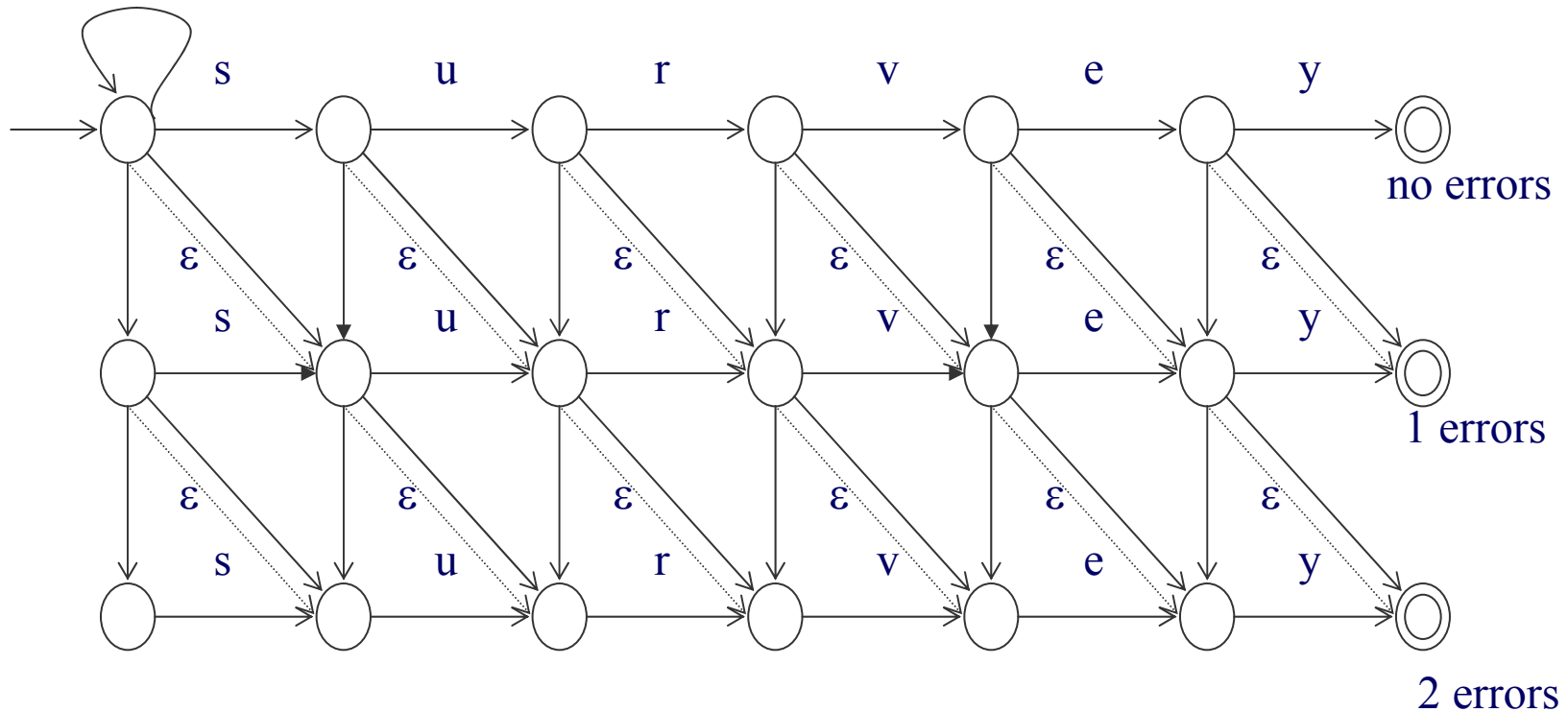
Text Searching allowing errors

- How to adapt the previous algorithm to search a short pattern $P[1..m]$ in a long text $T[1..n]$?
- Algorithm is basically the same:
 - Set $C[0,j]=0$ for all $j=0..n \rightarrow$ begin matching after the first j characters of text.
 - All other expressions are the same

		m	i	n	o	r	u	s	u	r	g	e	r	y
	0	0	0	0	0	0	0	0	0	0	0	0	0	0
s	1	1	1	1	1	1	1	0	1	1	1	1	1	1
u	2	2	2	2	2	2	2	1	0	1	2	2	2	2
r	3	3	3	3	3	2	3	2	1	0	1	2	2	3
v	4	4	4	4	4	3	4	3	2	1	1	2	3	3
e	5	5	5	5	5	4	5	4	3	2	2	1	2	3
y	6	6	6	6	6	5	6	5	4	3	3	2	2	2

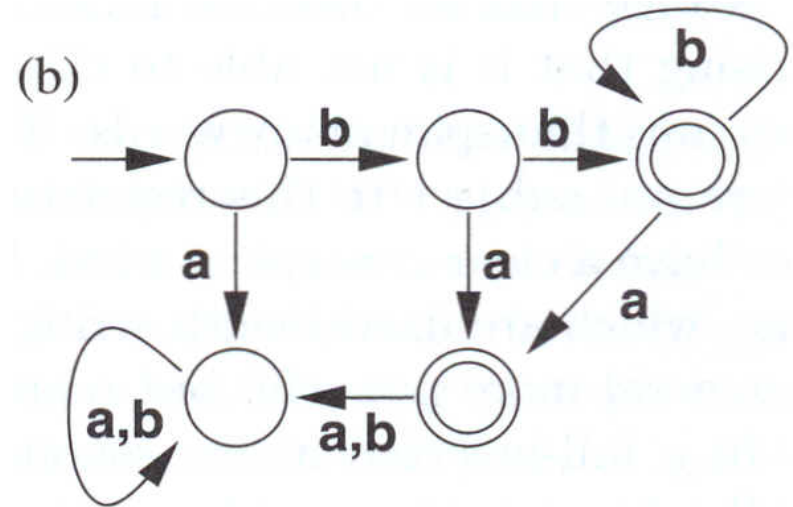
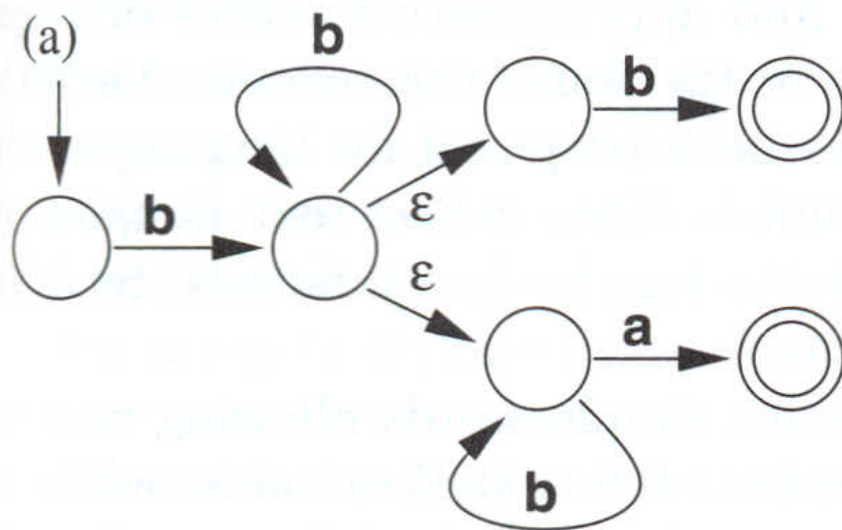
- The dynamic programming algorithm to search “survey” in the text “minor surgery” with two errors.
- Bold entries indicate matching text positions

Text Searching allowing errors (cont.)



- An non-deterministic finite automaton (NFA) for approximate string matching of the pattern 'survey' with two errors
- Unlabeled transitions match any character
- Horizontal arrows represent matching a character
- Vertical arrows represent insertions into the pattern
- Solid diagonal arrows represent replacements
- Dashed diagonal arrows represent deletions in the pattern

Regular Expressions



The non-deterministic (a) and deterministic (b) automata for the regular expression $bb^*(b \mid b^*a)$.

- $O(m)$ the size of the non-deterministic finite automaton, where m is the length of the regular expression
- $O(m2^m)$ the size and the construction time of the deterministic finite automaton
- $O(n)$ the time of searching any regular expression in text of n characters

Tradeoff of index space versus word searching time

