

Example 3.3: $A = \text{ballad}, B = \text{handball}$

d		h	a	n	d	b	a	l	l
	0	1	2	3	4	5	6	7	8
b	1	1	2	3	4	4	5	6	7
a	2	2	1	2	3	4	4	5	6
l	3	3	2	2	3	4	5	4	5
l	4	4	3	3	3	4	5	5	4
a	5	5	4	4	4	4	4	5	5
d	6	6	5	5	4	5	5	5	6

$ed(A, B) = d_{mn} = d_{6,8} = 6$.

The recurrence gives directly a **dynamic programming** algorithm for computing the edit distance.

Algorithm 3.4: Edit distance
 Input: strings $A[1..m]$ and $B[1..n]$
 Output: $ed(A, B)$

- (1) for $i \leftarrow 0$ to m do $d_{i0} \leftarrow i$
- (2) for $j \leftarrow 1$ to n do $d_{0j} \leftarrow j$
- (3) for $j \leftarrow 1$ to n do
- (4) for $i \leftarrow 1$ to m do
- (5) $d_{ij} \leftarrow \min\{d_{i-1,j-1} + \delta(A[i], B[j]), d_{i-1,j} + 1, d_{i,j-1} + 1\}$
- (6) return d_{mn}

The time and space complexity is $\mathcal{O}(mn)$.

Algorithm 3.5: Edit distance in $\mathcal{O}(m)$ space

Input: strings $A[1..m]$ and $B[1..n]$

Output: $ed(A, B)$

- (1) for $i \leftarrow 0$ to m do $C[i] \leftarrow i$
- (2) for $j \leftarrow 1$ to n do
- (3) $c \leftarrow C[0]; C[0] \leftarrow j$
- (4) for $i \leftarrow 1$ to m do
- (5) $d \leftarrow \min\{c + \delta(A[i], B[j]), C[i-1] + 1, C[i] + 1\}$
- (6) $c \leftarrow C[i]$
- (7) $C[i] \leftarrow d$
- (8) return $C[m]$

- Note that because $ed(A, B) = ed(B, A)$ (exercise), we can assume that $m \leq n$.

Example 3.6: $A = \text{ballad}, B = \text{handball}$

d		h	a	n	d	b	a	l	l
	0	1	2	3	4	5	6	7	8
b	1	1	2	3	4	4	5	6	7
a	2	2	1	2	3	4	4	5	6
l	3	3	2	2	3	4	5	4	5
l	4	4	3	3	3	4	5	5	4
a	5	5	4	4	4	4	4	5	5
d	6	6	5	5	4	5	5	5	6

There are 7 paths from $(0,0)$ to $(6,8)$ corresponding to 7 different optimal edit sequences and alignments, including the following three:

```

IIIIINNDD    SNISSNIS    SNSSINSI
----ballad   ba-lla-d     ball-ad-
handball--   handball     handball
    
```

Proof of Theorem 3.2. We use induction with respect to $i + j$. For brevity, write $A_i = A[1..i]$ and $B_j = B[1..j]$.

Basis:

$$\begin{aligned}
 d_{00} &= 0 = ed(\epsilon, \epsilon) \\
 d_{i0} &= i = ed(A_i, \epsilon) \quad (i \text{ deletions}) \\
 d_{0j} &= j = ed(\epsilon, B_j) \quad (j \text{ insertions})
 \end{aligned}$$

Induction step: We show that the claim holds for d_{ij} , $1 \leq i \leq m, 1 \leq j \leq n$. By induction assumption, $d_{pq} = ed(A_p, B_q)$ when $p + q < i + j$.

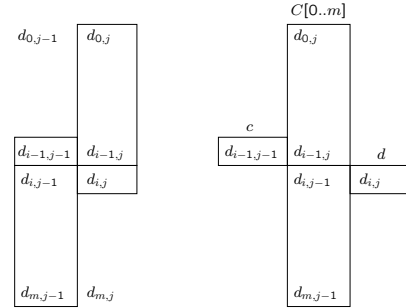
Let E_{ij} be an optimal edit sequence with the cost $ed(A_i, B_j)$. We have three cases depending on what the last operation symbol in E_{ij} is:

- N or S: $E_{ij} = E_{i-1,j-1}N$ or $E_{ij} = E_{i-1,j-1}S$ and
 $ed(A_i, B_j) = ed(A_{i-1}, B_{j-1}) + \delta(A[i], B[j]) = d_{i-1,j-1} + \delta(A[i], B[j])$.
- I: $E_{ij} = E_{i,j-1}I$ and $ed(A_i, B_j) = ed(A_i, B_{j-1}) + 1 = d_{i,j-1} + 1$.
- D: $E_{ij} = E_{i-1,j}D$ and $ed(A_i, B_j) = ed(A_{i-1}, B_j) + 1 = d_{i-1,j} + 1$.

One of the cases above is always true, and since the edit sequence is optimal, it must be one with the minimum cost, which agrees with the definition of d_{ij} . \square

The space complexity can be reduced by noticing that each column of the matrix (d_{ij}) depends **only on the previous column**. We do not need to store older columns.

A more careful look reveals that, when computing d_{ij} , we only need to store the bottom part of column $j-1$ and the already computed top part of column j . We store these in an array $C[0..m]$ and variables c and d as shown below:



It is also possible to find optimal edit sequences and alignments from the matrix d_{ij} .

An edit graph is a directed graph, where the nodes are the cells of the edit distance matrix, and the edges are as follows:

- If $A[i] = B[j]$ and $d_{ij} = d_{i-1,j-1}$, there is an edge $(i-1, j-1) \rightarrow (i, j)$ labelled with N.
- If $A[i] \neq B[j]$ and $d_{ij} = d_{i-1,j-1} + 1$, there is an edge $(i-1, j-1) \rightarrow (i, j)$ labelled with S.
- If $d_{ij} = d_{i,j-1} + 1$, there is an edge $(i, j-1) \rightarrow (i, j)$ labelled with I.
- If $d_{ij} = d_{i-1,j} + 1$, there is an edge $(i-1, j) \rightarrow (i, j)$ labelled with D.

Any path from $(0,0)$ to (m,n) is labelled with an optimal edit sequence.

Approximate String Matching

Now we are ready to tackle the main problem of this part: **approximate string matching**.

Problem 3.7: Given a text $T[1..n]$, a pattern $P[1..m]$ and an integer $k \geq 0$, report all positions $j \in [1..n]$ such that $ed(P, T(j-\ell..j)) \leq k$ for some $\ell \geq 0$.

The factor $T(j-\ell..j)$ is called an **approximate occurrence** of P .

There can be multiple occurrences of different lengths ending at the same position j , but usually it is enough to report just the end positions. We ask for the end position rather than the start position because that is more natural for the algorithms.

Define the values g_{ij} with the recurrence:

$$g_{0j} = 0, 0 \leq j \leq n,$$

$$g_{i0} = i, 1 \leq i \leq m, \text{ and}$$

$$g_{ij} = \min \begin{cases} g_{i-1,j-1} + \delta(P[i], T[j]) \\ g_{i-1,j} + 1 \\ g_{i,j-1} + 1 \end{cases} \quad 1 \leq i \leq m, 1 \leq j \leq n.$$

Theorem 3.8: For all $0 \leq i \leq m, 0 \leq j \leq n$:

$$g_{ij} = \min\{ed(P[1..i], T(j-\ell..j)) \mid 0 \leq \ell \leq j\}.$$

In particular, j is an ending position of an approximate occurrence if and only if $g_{mj} \leq k$.

Proof. We use induction with respect to $i + j$.

Basis:

$$g_{00} = 0 = ed(\epsilon, \epsilon)$$

$$g_{0j} = 0 = ed(\epsilon, \epsilon) = ed(\epsilon, T(j-0..j)) \quad (\text{min at } \ell = 0)$$

$$g_{i0} = i = ed(P[1..i], \epsilon) = ed(P[1..i], T(0-0..0)) \quad (0 \leq \ell \leq j = 0)$$

Induction step: Essentially the same as in the proof of Theorem 3.2.

121

122

Example 3.9: $P = \text{match}, T = \text{remachine}, k = 1$

g	r	e	m	a	c	h	i	n	e
m	0	0	0	0	0	0	0	0	0
a	1	1	1	0	1	1	1	1	1
t	2	2	2	1	0	1	2	2	2
c	3	3	3	2	1	1	2	3	3
h	4	4	4	3	2	1	2	3	4
e	5	5	5	4	3	2	1	2	3

One occurrence ending at position 6.

Algorithm 3.10: Approximate string matching

Input: text $T[1..n]$, pattern $P[1..m]$, and integer k

Output: end positions of all approximate occurrences of P

- (1) for $i \leftarrow 0$ to m do $g_{i0} \leftarrow i$
- (2) for $j \leftarrow 1$ to n do $g_{0j} \leftarrow 0$
- (3) for $j \leftarrow 1$ to n do
- (4) for $i \leftarrow 1$ to m do
- (5) $g_{ij} \leftarrow \min\{g_{i-1,j-1} + \delta(A[i], B[j]), g_{i-1,j} + 1, g_{i,j-1} + 1\}$
- (6) if $g_{mj} \leq k$ then output j

- Time and space complexity is $\mathcal{O}(mn)$ on ordered alphabet.
- The space complexity can be reduced to $\mathcal{O}(m)$ by storing only one column as in Algorithm 3.5.

123

124

Ukkonen's Cut-off Heuristic

We can speed up the algorithm using the diagonal monotonicity of the matrix (g_{ij}) :

A diagonal d , $-m \leq d \leq n$, consists of the cells g_{ij} with $j - i = d$. Every diagonal in (g_{ij}) is monotonically non-decreasing.

Example 3.11: Diagonals -3 and 2.

g	r	e	m	a	c	h	i	n	e
m	0	0	0	0	0	0	0	0	0
a	1	1	1	0	1	1	1	1	1
t	2	2	2	1	0	1	2	2	2
c	3	3	3	2	1	1	2	3	3
h	4	4	4	3	2	1	2	3	4
e	5	5	5	4	3	2	1	2	3

125

126

We can reduce computation using diagonal monotonicity:

- Whenever the value on a diagonal d grows larger than k , we can discard d from consideration, because we are only interested in values at most k on the row m .
- We keep track of the smallest undiscarded diagonal d . Each column is computed only up to diagonal d .

Example 3.13: $P = \text{match}, T = \text{remachine}, k = 1$

g	r	e	m	a	c	h	i	n	e
m	0	0	0	0	0	0	0	0	0
a	1	1	1	0	1	1	1	1	1
t		2	2	1	0	1	2	2	2
c					1	1	2	3	
h						1	2	3	
e							1	2	

127

128

More specifically, we have the following property.

Lemma 3.12: For every $i \in [1..m]$ and every $j \in [1..n]$, $g_{ij} = g_{i-1,j-1}$ or $g_{ij} = g_{i-1,j-1} + 1$.

Proof. By definition, $g_{ij} \leq g_{i-1,j-1} + \delta(P[i], T[j]) \leq g_{i-1,j-1} + 1$. We show that $g_{ij} \geq g_{i-1,j-1}$ by induction on $i + j$.

The induction assumption is that $g_{pq} \geq g_{p-1,q-1}$ when $p \in [1..m]$, $q \in [1..n]$ and $p + q < i + j$. At least one of the following holds:

1. $g_{ij} = g_{i-1,j-1} + \delta(P[i], T[j])$. Then $g_{ij} \geq g_{i-1,j-1}$.
2. $g_{ij} = g_{i-1,j} + 1$ and $i > 1$. Then

$$g_{ij} = g_{i-1,j} + 1 \stackrel{\text{ind. assump.}}{\geq} g_{i-2,j-1} + 1 + 1 \stackrel{\text{definition}}{\geq} g_{i-1,j-1}$$

3. $g_{ij} = g_{i,j-1} + 1$ and $j > 1$. Then

$$g_{ij} = g_{i,j-1} + 1 \stackrel{\text{ind. assump.}}{\geq} g_{i-1,j-2} + 1 + 1 \stackrel{\text{definition}}{\geq} g_{i-1,j-1}$$

4. $g_{ij} = g_{i-1,j} + 1$ and $i = 1$. Then $g_{ij} = 0 + 1 > 0 = g_{i-1,j-1}$.

5. $g_{ij} = g_{i,j-1} + 1$ and $j = 1$. Then $g_{ij} = i + 1 = (i-1) + 2 = g_{i-1,j-1} + 2$, which cannot be true. Thus this case can never happen. \square

The position of the smallest undiscarded diagonal on the current column is kept in a variable top .

Algorithm 3.14: Ukkonen's cut-off algorithm

Input: text $T[1..n]$, pattern $P[1..m]$, and integer k

Output: end positions of all approximate occurrences of P

- (1) for $i \leftarrow 0$ to $\min(k, m)$ do $g_{i0} \leftarrow i$
- (2) for $j \leftarrow 1$ to n do $g_{0j} \leftarrow 0$
- (3) $top \leftarrow \min(k + 1, m)$
- (4) for $j \leftarrow 1$ to n do
- (5) for $i \leftarrow 1$ to top do
- (6) $g_{ij} \leftarrow \min\{g_{i-1,j-1} + \delta(A[i], B[j]), g_{i-1,j} + 1, g_{i,j-1} + 1\}$
- (7) while $g_{top,j} > k$ do $top \leftarrow top - 1$
- (8) if $top = m$ then output j
- (9) else $top \leftarrow top + 1$

127

128

Myers' Bitparallel Algorithm

Another way to speed up the computation is bitparallelism.

Instead of the matrix (g_{ij}) , we store **differences** between adjacent cells:

$$\text{Vertical delta: } \Delta v_{ij} = g_{ij} - g_{i-1,j}$$

$$\text{Horizontal delta: } \Delta h_{ij} = g_{ij} - g_{i,j-1}$$

$$\text{Diagonal delta: } \Delta d_{ij} = g_{ij} - g_{i-1,j-1}$$

Because $g_{i0} = i$ ja $g_{0j} = 0$,

$$\begin{aligned} g_{ij} &= \Delta v_{1j} + \Delta v_{2j} + \dots + \Delta v_{ij} \\ &= i + \Delta h_{i1} + \Delta h_{i2} + \dots + \Delta h_{ij} \end{aligned}$$

Because of diagonal monotonicity, $\Delta d_{ij} \in \{0, 1\}$ and it can be stored in one bit. By the following result, Δh_{ij} and Δv_{ij} can be stored in two bits.

Lemma 3.15: $\Delta h_{ij}, \Delta v_{ij} \in \{-1, 0, 1\}$ for every i, j that they are defined for.

The proof is left as an exercise.

129

130

The time complexity is proportional to the computed area in the matrix (g_{ij}) .

- The worst case time complexity is still $\mathcal{O}(mn)$ on **ordered alphabet**.
- The **average case** time complexity is $\mathcal{O}(kn)$. The proof is not trivial.

There are many other algorithms based on diagonal monotonicity. Some of them achieve $\mathcal{O}(kn)$ **worst case** time complexity.

Example 3.16: '-' means -1, '=' means 0 and '+' means +1

	r	e	m	a	c	h	i	n	e
0	= 0	= 0	= 0	= 0	= 0	= 0	= 0	= 0	= 0
m	+	+	+	+	+	+	+	+	+
1	= 1	= 1	= 1	= 0	+ 1	= 1	= 1	= 1	= 1
a	+	+	+	+	+	= +	= +	+	+
2	= 2	= 2	= 2	= 1	= 0	+ 1	+ 2	= 2	= 2
t	+	+	+	+	+	= +	= +	= +	+
3	= 3	= 3	= 3	= 2	= 1	= 1	+ 2	+ 3	= 3
c	+	+	+	+	+	= +	= +	= +	+
4	= 4	= 4	= 4	= 3	= 2	= 1	+ 2	+ 3	+ 4
h	+	+	+	+	+	= +	= +	= +	= +
5	= 5	= 5	= 5	= 4	= 3	= 2	= 1	+ 2	+ 3

131

132

The computation rule is defined by the following result.

Lemma 3.17: If $Eq = 1$ or $\Delta v^{\text{in}} = -1$ or $\Delta h^{\text{in}} = -1$, then $\Delta d = 0$, $\Delta v^{\text{out}} = -\Delta h^{\text{in}}$ and $\Delta h^{\text{out}} = -\Delta v^{\text{in}}$. Otherwise $\Delta d = 1$, $\Delta v^{\text{out}} = 1 - \Delta h^{\text{in}}$ and $\Delta h^{\text{out}} = 1 - \Delta v^{\text{in}}$.

Proof. We can write the recurrence for g_{ij} as

$$\begin{aligned} g_{ij} &= \min\{g_{i-1,j-1} + \delta(P[i], T[j]), g_{i,j-1} + 1, g_{i-1,j} + 1\} \\ &= g_{i-1,j-1} + \min\{1 - Eq, \Delta v^{\text{in}} + 1, \Delta h^{\text{in}} + 1\}. \end{aligned}$$

Then $\Delta d = g_{ij} - g_{i-1,j-1} = \min\{1 - Eq, \Delta v^{\text{in}} + 1, \Delta h^{\text{in}} + 1\}$ which is 0 if $Eq = 1$ or $\Delta v^{\text{in}} = -1$ or $\Delta h^{\text{in}} = -1$ and 1 otherwise.

Clearly $\Delta d = \Delta v^{\text{in}} + \Delta h^{\text{out}} = \Delta h^{\text{in}} + \Delta v^{\text{out}}$. Thus $\Delta v^{\text{out}} = \Delta d - \Delta h^{\text{in}}$ and $\Delta h^{\text{out}} = \Delta d - \Delta v^{\text{in}}$. \square

133

134

Now the computation rules can be expressed as follows.

Lemma 3.18: $Pv^{\text{out}} = Mh^{\text{in}} \vee \neg(Xv \vee Ph^{\text{in}})$ $Mv^{\text{out}} = Ph^{\text{in}} \wedge Xv$
 $Ph^{\text{out}} = Mv^{\text{in}} \vee \neg(Xh \vee Pv^{\text{in}})$ $Mh^{\text{out}} = Pv^{\text{in}} \wedge Xh$

where $Xv = Eq \vee Mv^{\text{in}}$ and $Xh = Eq \vee Mh^{\text{in}}$.

Proof. We show the claim for Pv and Mv only. Ph and Mh are symmetrical.

By Lemma 3.17,

$$Pv^{\text{out}} = (\neg \Delta d \wedge Mh^{\text{in}}) \vee (\Delta d \wedge \neg Ph^{\text{in}})$$

$$Mv^{\text{out}} = (\neg \Delta d \wedge Ph^{\text{in}}) \vee (\Delta d \wedge 0) = \neg \Delta d \wedge Ph^{\text{in}}$$

Because $\Delta d = \neg(Eq \vee Mv^{\text{in}} \vee Mh^{\text{in}}) = \neg(Xv \vee Mh^{\text{in}}) = \neg Xv \wedge \neg Mh^{\text{in}}$,

$$\begin{aligned} Pv^{\text{out}} &= ((Xv \vee Mh^{\text{in}}) \wedge Mh^{\text{in}}) \vee (\neg Xv \wedge \neg Mh^{\text{in}} \wedge \neg Ph^{\text{in}}) \\ &= Mh^{\text{in}} \vee \neg(Xv \vee Mh^{\text{in}} \vee Ph^{\text{in}}) \\ &= Mh^{\text{in}} \vee \neg(Xv \vee Ph^{\text{in}}) \end{aligned}$$

$$Mv^{\text{out}} = (Xv \vee Mh^{\text{in}}) \wedge Ph^{\text{in}} = Xv \wedge Ph^{\text{in}}$$

All the steps above use just basic laws of Boolean algebra except the last step, where we use the fact that Mh^{in} and Ph^{in} cannot be 1 simultaneously. \square

135

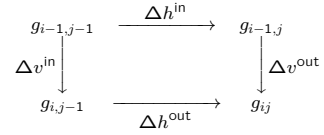
136

In the standard computation of a cell:

- **Input** is $g_{i-1,j}$, $g_{i-1,j-1}$, $g_{i,j-1}$ and $\delta(P[i], T[j])$.
- **Output** is g_{ij} .

In the corresponding bitparallel computation:

- **Input** is $\Delta v^{\text{in}} = \Delta v_{i,j-1}$, $\Delta h^{\text{in}} = \Delta h_{i,j-1}$ and $Eq_{ij} = 1 - \delta(P[i], T[j])$.
- **Output** is $\Delta v^{\text{out}} = \Delta v_{i,j}$ and $\Delta h^{\text{out}} = \Delta h_{i,j}$.



The algorithm does not compute the Δd values but they are useful in the proofs.

To enable bitparallel operation, we need two changes:

- The Δv and Δh values are "trits" not bits. We encode each of them with two bits as follows:

$$\begin{aligned} Pv &= \begin{cases} 1 & \text{if } \Delta v = +1 \\ 0 & \text{otherwise} \end{cases} & Mv &= \begin{cases} 1 & \text{if } \Delta v = -1 \\ 0 & \text{otherwise} \end{cases} \\ Ph &= \begin{cases} 1 & \text{if } \Delta h = +1 \\ 0 & \text{otherwise} \end{cases} & Mh &= \begin{cases} 1 & \text{if } \Delta h = -1 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Then

$$\begin{aligned} \Delta v &= Pv - Mv \\ \Delta h &= Ph - Mh \end{aligned}$$

- We replace arithmetic operations (+, -, min) with Boolean (logical) operations (\wedge , \vee , \neg).

According to Lemma 3.18, the bit representation of the matrix can be computed as follows.

```

for i ← 1 to m do
  Pvi0 ← 1; Mvi0 ← 0
for j ← 1 to n do
  Ph0j ← 0; Mh0j ← 0
for i ← 1 to m do
  for j ← 1 to n do
    Xhij ← Eqij ∨ Mhi-1,j
    Phij ← Mvi,j-1 ∨ ¬(Xhij ∨ Pvi,j-1)
    Mhij ← Pvi,j-1 ∧ Xhij
  for i ← 1 to m do
    Xvij ← Eqij ∨ Mvi,j-1
    Pvij ← Mhi-1,j ∨ ¬(Xvij ∨ Phi-1,j)
    Mvij ← Phi-1,j ∧ Xvij

```

This is not yet bitparallel though.

To obtain a bitparallel algorithm, the columns Pv_{*j} , Mv_{*j} , Xv_{*j} , Ph_{*j} , Mh_{*j} , Xh_{*j} and Eq_{*j} are stored in bitvectors.

Now the second inner loop can be replaced with the code

```
Xv_{*j} ← Eq_{*j} ∨ Mv_{*j-1}
Pv_{*j} ← (Mh_{*j} << 1) ∨ ¬(Xv_{*j} ∨ (Ph_{*j} << 1))
Mv_{*j} ← (Ph_{*j} << 1) ∧ Xv_{*j}
```

A similar attempt with the for first inner loop leads to a problem:

```
Xh_{*j} ← Eq_{*j} ∨ (Mh_{*j} << 1)
Ph_{*j} ← Mv_{*j-1} ∨ ¬(Xh_{*j} ∨ Pv_{*j-1})
Mh_{*j} ← Pv_{*j-1} ∧ Xh_{*j}
```

Now the vector Mh_{*j} is used in computing Xh_{*j} before Mh_{*j} itself is computed! Changing the order does not help, because Xh_{*j} is needed to compute Mh_{*j} .

To get out of this dependency loop, we compute Xh_{*j} without Mh_{*j} using only Eq_{*j} and Pv_{*j-1} which are already available when we compute Xh_{*j} .

137

At first sight, we cannot use Lemma 3.19 to compute even a single bit in constant time, not to mention a whole vector Xh_{*j} . However, it can be done, but we need more bit operations:

- Let \vee denote the xor-operation: $0 \vee 1 = 1 \vee 0 = 1$ and $0 \vee 0 = 1 \vee 1 = 0$.
- A bitvector is interpreted as an integer and we use **addition** as a bit operation. The **carry** mechanism in addition plays a key role. For example $0001 + 0111 = 1000$.

In the following, for a bitvector B , we will write

$$B = B[1..m] = B[m]B[m-1] \dots B[1]$$

The reverse order of the bits reflects the interpretation as an integer.

139

- b) The following calculation shows that $Y[i] = 1$ in this case:

$$\begin{array}{r} \\ \\ \\ \\ \\ Y = (((((E \wedge P) + P) \vee P) \vee E)[\dots i]) = 11 \dots 11 \end{array}$$

where \mathbf{b} is the unknown bit $P[i]$, \mathbf{c} is the possible carry bit coming from the summation of bits $1 \dots, \ell - 1$, and $\bar{\mathbf{b}}$ and $\bar{\mathbf{c}}$ are their negations.

- c) Because for all bitvectors B , $0 \wedge B = 0$ ja $0 + B = B$, we get $Y = (((0 \wedge P) + P) \vee P) \vee 0 = (P \vee P) \vee 0 = 0$.
- d) Consider the calculation in case b). A key point there is that the carry bit in the summation travels from position ℓ to i and produces $\bar{\mathbf{b}}$ to position i . The difference in this case is that at least one bit $P[k]$, $\ell \leq k < i$, is zero, which stops the carry at position k . Thus $((E \wedge P) + P)[i] = \mathbf{b}$ and $Y[i] = 0$.

□

141

On an **integer alphabet**, when $m \leq w$:

- Pattern preprocessing time is $\mathcal{O}(m + \sigma)$.
- Search time is $\mathcal{O}(n)$.

When $m > w$, we can store each bit vector in $\lceil m/w \rceil$ machine words:

- The worst case search time is $\mathcal{O}(n \lceil m/w \rceil)$.
- Using Ukkonen's cut-off heuristic, it is possible reduce the **average case** search time to $\mathcal{O}(n \lceil k/w \rceil)$.

143

Lemma 3.19: $Xh_{ij} = \exists \ell \in [1, i] : Eq_{\ell j} \wedge (\forall x \in [\ell, i-1] : Pv_{x,j-1})$.

Proof. We use induction on i .

Basis $i = 1$: The right-hand side reduces to Eq_{1j} , because $\ell = 1$. By Lemma 3.18, $Xh_{1j} = Eq_{1j} \vee Mh_{0j}$, which is Eq_{1j} because $Mh_{0j} = 0$ for all j .

Induction step: The induction assumption is that $Xh_{i-1,j}$ is as claimed. Now we have

$$\begin{aligned} & \exists \ell \in [1, i] : Eq_{\ell j} \wedge (\forall x \in [\ell, i-1] : Pv_{x,j-1}) \\ &= Eq_{ij} \vee \exists \ell \in [1, i-1] : Eq_{\ell j} \wedge (\forall x \in [\ell, i-1] : Pv_{x,j-1}) \\ &= Eq_{ij} \vee (Pv_{i-1,j-1} \wedge \exists \ell \in [1, i-1] : Eq_{\ell j} \wedge (\forall x \in [\ell, i-2] : Pv_{x,j-1})) \\ &= Eq_{ij} \vee (Pv_{i-1,j-1} \wedge Xh_{i-1,j}) \quad (\text{ind. assum.}) \\ &= Eq_{ij} \vee Mh_{i-1,j} \quad (\text{Lemma 3.18}) \\ &= Xh_{ij} \quad (\text{Lemma 3.18}) \end{aligned}$$

□

138

Lemma 3.20: Denote $X = Xh_{*j}$, $E = Eq_{*j}$, $P = Pv_{*j-1}$ and let $Y = (((E \wedge P) + P) \vee P) \vee E$. Then $X = Y$.

Proof. By Lemma 3.19, $X[i] = 1$ iff and only if

- $E[i] = 1$ or
- $\exists \ell \in [1, i] : E[\ell \dots i] = 00 \dots 01 \wedge P[\ell \dots i-1] = 11 \dots 1$.

and $X[i] = 0$ iff and only if

- $E_{1..i} = 00 \dots 0$ or
- $\exists \ell \in [1, i] : E[\ell \dots i] = 00 \dots 01 \wedge P[\ell \dots i-1] \neq 11 \dots 1$.

We prove that $Y[i] = X[i]$ in all of these cases:

- The definition of Y ends with " $\vee E$ " which ensures that $Y[i] = 1$ in this case.

As a final detail, we compute the bottom row values g_{mj} using the equalities $g_{m0} = m$ ja $g_{mj} = g_{m,j-1} + \Delta h_{mj}$.

Algorithm 3.21: Myers' bitparallel algorithm

Input: text $T[1..n]$, pattern $P[1..m]$, and integer k

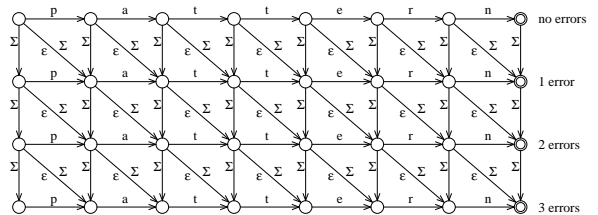
Output: end positions of all approximate occurrences of P

- for $c \in \Sigma$ do $B[c] \leftarrow 0^m$
- for $i \leftarrow 1$ to m do $B[P[i]][i] = 1$
- $Pv \leftarrow 1^m$; $Mv \leftarrow 0$; $g \leftarrow m$
- for $j \leftarrow 1$ to n do
 - $Eq \leftarrow B[T[j]]$
 - $Xh \leftarrow (((Eq \wedge Pv) + Pv) \vee Pv) \vee Eq$
 - $Ph \leftarrow Mv \vee \neg(Xh \vee Pv)$
 - $Mh \leftarrow Pv \wedge Xh$
 - $Xv \leftarrow Eq \vee Mv$
 - $Pv \leftarrow (Mh << 1) \vee \neg(Xv \vee (Ph << 1))$
 - $Mv \leftarrow (Ph << 1) \wedge Xv$
 - $g \leftarrow g + Ph[m] - Mh[m]$
 - if $g \leq k$ then output j

142

There are also algorithms based on bitparallel simulation of a nondeterministic automaton.

Example 3.22: $P = \text{pattern}$, $k = 3$



- The algorithm of Wu and Manber uses a bit vector for each row. It can be seen as an extension of Shift-And. The search time complexity is $\mathcal{O}(kn \lceil m/w \rceil)$.
- The algorithm of Baeza-Yates and Navarro uses a bit vector for each diagonal, packed into one long bitvector. The search time complexity is $\mathcal{O}(n \lceil km/w \rceil)$.

144