

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) top ← min(k + 1, m)
(2) for i ← 0 to top do gi0 ← i
(3) for j ← 1 to n do g0j ← 0
(4) for j ← 1 to n do
(5)   for i ← 1 to top do
(6)     gij ← min{gi-1,j-1 + δ(A[i], B[j]), gi-1,j + 1, gi,j-1 + 1}
(7)   while gtop,j > k do top ← top - 1
(8)   if top = m then output j
(9)   else top ← top + 1; gtop,j ← k + 1
```

**Myers' Bitparallel Algorithm**

Another way to speed up the computation is bitparallelism.

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

Vertical delta:  $\Delta v_{ij} = g_{ij} - g_{i-1,j}$   
Horizontal delta:  $\Delta h_{ij} = g_{ij} - g_{i,j-1}$   
Diagonal delta:  $\Delta d_{ij} = g_{ij} - g_{i-1,j-1}$

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

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

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

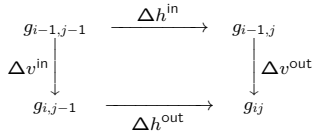
The proof is left as an exercise.

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



The algorithm does not compute the  $\Delta d$  values but they are useful in the proofs.

To enable bitparallel operation, we need two changes:

- The  $\Delta v$  and  $\Delta h$  values are "trits" not bits. We encode each of them with two bits as follows:

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

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

Then

$$\Delta v = Pv - Mv \\ \Delta h = Ph - Mh$$

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

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 |   |
|---|---------------------------------------|---|---|---|---|---|---|---|---|---|
| m | 0 = 0 = 0 = 0 = 0 = 0 = 0 = 0 = 0 = 0 | + | + | + | + | + | + | + | + | + |
| a | 1 = 1 = 1 = 0 + 1 = 1 = 1 = 1 = 1 = 1 | + | + | + | + | + | + | + | + | + |
| t | 2 = 2 = 2 = 2 = 2 = 2 = 2 = 2 = 2 = 2 | + | + | + | + | + | + | + | + | + |
| c | 3 = 3 = 3 = 2 = 1 = 1 + 2 + 3 = 3 = 3 | + | + | + | + | + | + | + | + | + |
| h | 4 = 4 = 4 = 3 = 2 = 1 + 2 + 3 + 4 = 4 | + | + | + | + | + | + | + | + | + |
|   | 5 = 5 = 5 = 4 = 3 = 2 = 1 + 2 + 3 + 4 | + | + | + | + | + | + | + | + | + |

The computation rule is defined by the following result.

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

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

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

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

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

Now the computation rules can be expressed as follows.

**Lemma 3.18:**  $Pv^{out} = Mh^{in} \vee \neg(Xv \vee Ph^{in})$      $Mv^{out} = Ph^{in} \wedge Xv$   
 $Ph^{out} = Mv^{in} \vee \neg(Xh \vee Pv^{in})$      $Mh^{out} = Pv^{in} \wedge Xh$   
where  $Xv = Eq \vee Mv^{in}$  and  $Xh = Eq \vee Mh^{in}$ .

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

By Lemma 3.17,

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

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

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

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.  $\square$

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
    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.

137

**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

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

□

139

**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

- $E_{1 \dots i} = 00 \dots 0$  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.

141

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  $P$

```

(1) for c ∈ Σ do B[c] ← 0m
(2) for i ← 1 to m do B[P[i]][i] = 1
(3) Pv ← 1m; Mv ← 0; g ← m
(4) for j ← 1 to n do
(5)   Eq ← B[T[j]]
(6)   Xh ← (((Eq ∧ Pv) + Pv) ∨ Pv) ∨ Eq
(7)   Ph ← Mv ∨ ¬(Xh ∨ Pv)
(8)   Mh ← Pv ∧ Xh
(9)   Xv ← Eq ∨ Mv
(10)  Pv ← (Mh << 1) ∨ ¬(Xv ∨ (Ph << 1))
(11)  Mv ← (Ph << 1) ∧ Xv
(12)  g ← g + Ph[m] - Mh[m]
(13)  if g ≤ k then output j

```

143

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  $Xh_{*j}$ .

138

At first sight, we cannot use Lemma 3.19 to compute even a single bit in constant time, let alone 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.

140

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

$$\begin{aligned}
& \begin{matrix} & i & \ell \\ E[\ell \dots i] & = & 00 \dots 01 \\ P[\ell \dots i] & = & \mathbf{b}1 \dots 11 \\ (E \wedge P)[\ell \dots i] & = & 00 \dots 01 \\ ((E \wedge P) + P)[\ell \dots i] & = & \mathbf{\bar{b}}0 \dots 0c \\ (((E \wedge P) + P) \vee P)[\ell \dots i] & = & 11 \dots 1\bar{c} \\ Y = (((E \wedge P) + P) \vee P) \vee E[\ell \dots i] & = & 11 \dots 11 \end{matrix}
\end{aligned}$$

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.

- Because for all bitvectors  $B$ ,  $0 \wedge B = 0$  ja  $0 + B = B$ , we get  $Y = (((0 \wedge P) + P) \vee P) \vee 0 = (P \vee P) \vee 0 = 0$ .
- Consider the calculation in case b). A key point there is that the carry bit in the summation travels from position  $\ell$  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  $Y[i] = (\mathbf{b} \vee \bar{\mathbf{b}}) \vee 0 = 0$ .

□

142

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

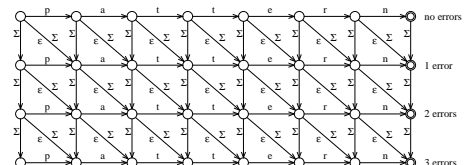
- Pattern preprocessing time is  $\mathcal{O}(m + \sigma)$ .
- Search time is  $\mathcal{O}(n)$ .

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  $\mathcal{O}(n \lceil k/w \rceil)$ .

There are also algorithms for bitparallel simulation of a nondeterministic automaton that recognizes the approximate occurrences of the pattern.

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



144

Another way to utilize Lemma 3.15 ( $\Delta h_{ij}, \Delta v_{ij} \in \{-1, 0, 1\}$ ) is to use precomputed tables to process multiple matrix cells at a time.

- There are at most  $3^m$  different columns. Thus there exists a deterministic automaton with  $3^m$  states and  $\sigma 3^m$  transitions that can find all approximate occurrences in  $\mathcal{O}(n)$  time. However, the space and construction time of the automaton can be too big to be practical.
- There is a super-alphabet algorithm that processes  $\mathcal{O}(\log_\sigma n)$  characters at a time and  $\mathcal{O}(\log_\sigma^2 n)$  matrix cells at a time using lookup tables of size  $\mathcal{O}(n)$ . This gives time complexity  $\mathcal{O}(mn/\log_\sigma^2 n)$ .
- A practical variant uses smaller lookup tables to compute multiple entries of a column at a time.

145

The following lemma shows the property used by the [Baeza-Yates–Perleberg algorithm](#) and proves that it satisfies the first condition.

**Lemma 3.23:** Let  $P_1 P_2 \dots P_{k+1} = P$  be a partitioning of the pattern  $P$  into  $k + 1$  nonempty factors. Any string  $S$  with  $ed(P, S) \leq k$  contains  $P_i$  as a factor for some  $i \in [1..k + 1]$ .

**Proof.** Each single symbol edit operation can change at most one of the pattern factors  $P_i$ . Thus any set of at most  $k$  edit operations leaves at least one of the factors untouched.  $\square$

147

Let us analyze the [average case](#) time complexity of the verification phase.

- The best pattern partitioning is as even as possible. Then each pattern factor has length at least  $r = \lfloor m/(k + 1) \rfloor$ .
- The expected number of exact occurrences of a random string of length  $r$  in a random text of length  $n$  is at most  $n/\sigma^r$ .
- The expected total verification time is at most
 
$$\mathcal{O}\left(\frac{m^2(k + 1)n}{\sigma^r}\right) \leq \mathcal{O}\left(\frac{m^3 n}{\sigma^r}\right).$$
 This is  $\mathcal{O}(n)$  if  $r \geq 3 \log_\sigma m$ .
- The condition  $r \geq 3 \log_\sigma m$  is satisfied when  $(k + 1) \leq m/(3 \log_\sigma m + 1)$ .

**Theorem 3.24:** The average case time complexity of the Baeza-Yates–Perleberg algorithm is  $\mathcal{O}(n)$  when  $k \leq m/(3 \log_\sigma m + 1) - 1$ .

149

## Summary: Approximate String Matching

We have seen two main types of algorithms for approximate string matching:

- Basic dynamic programming time complexity is  $\mathcal{O}(mn)$ . The time complexity can be improved to  $\mathcal{O}(kn)$  using diagonal monotonicity, and to  $\mathcal{O}(n \lceil m/w \rceil)$  using bitparallelism.
- Filtering algorithms can improve average case time complexity and are the fastest in practice when  $k$  is not too large.

Similar techniques can be useful for other variants of edit distance but not always straightforwardly.

151

## Baeza-Yates–Perleberg Filtering Algorithm

A [filtering algorithm](#) for approximate string matching searches the text for factors having some [property](#) that satisfies the following conditions:

1. Every approximate occurrence of the pattern has this property.
2. Strings having this property are reasonably rare.
3. Text factors having this property can be found quickly.

Each text factor with the property is a [potential occurrence](#), which is then [verified](#) for whether it is an actual approximate occurrence.

Filtering algorithms can achieve linear or even sublinear [average case time](#) complexity.

146

The algorithm has two phases:

**Filtration:** Search the text  $T$  for exact occurrences of the pattern factors  $P_i$ . Using the Aho–Corasick algorithm this takes  $\mathcal{O}(n)$  time for a constant alphabet.

**Verification:** An area of length  $\mathcal{O}(m)$  surrounding each potential occurrence found in the filtration phase is searched using the standard dynamic programming algorithm in  $\mathcal{O}(m^2)$  time.

The worst case time complexity is  $\mathcal{O}(m^2 n)$ , which can be reduced to  $\mathcal{O}(mn)$  by combining any overlapping areas to be searched.

148

Many variations of the algorithm have been suggested:

- The filtration can be done with a [different multiple exact string matching algorithm](#).
- The verification time can be reduced using a technique called [hierarchical verification](#).
- The pattern can be partitioned into [fewer than  \$k + 1\$  pieces](#), which are searched [allowing a small number of errors](#).

A lower bound on the average case time complexity is  $\Omega(n(k + \log_\sigma m)/m)$ , and there exists a filtering algorithm matching this bound.

150