

## BNDM

Starting the matching from the end enables long shifts.

- The Horspool algorithm bases the shift on a **single character**.
- The Boyer–Moore algorithm uses the matching **suffix** and the mismatching character.
- Factor based algorithms continue matching until no pattern **factor** matches. This may require more comparisons but it enables longer shifts.

### Example 2.14:

Horspool shift

```
varmasti-aikaisen-ainainen
ainaisen-ainainen
ainaisen-ainaisen
```

Boyer–Moore shift

```
varmasti-aikaisen-ainainen
ainaisen-ainainen
ainaisen-ainainen
```

Factor shift

```
varmasti-aikaisen-ainainen
ainaisen-ainainen
ainaisen-ainainen
```

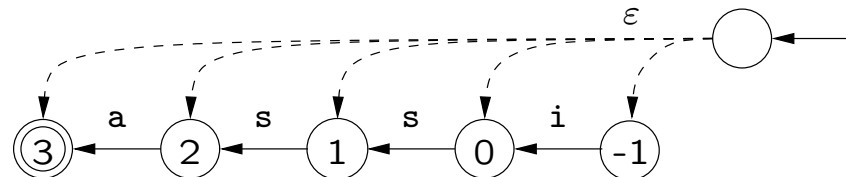
Factor based algorithms use an **automaton** that accepts **suffixes** of the **reverse pattern**  $P^R$  (or equivalently reverse prefixes of the pattern  $P$ ).

- BDM (Backward DAWG Matching) uses a **deterministic** automaton that accepts exactly the suffixes of  $P^R$ .

DAWG (Directed Acyclic Word Graph) is also known as suffix automaton.

- BNDM (Backward Nondeterministic DAWG Matching) simulates a **nondeterministic** automaton.

**Example 2.15:**  $P = \text{ass}i.$



- BOM (Backward Oracle Matching) uses a much **simpler deterministic** automaton that accepts all suffixes of  $P^R$  but may also accept some other strings. This can cause shorter shifts but not incorrect behaviour.

Suppose we are currently comparing  $P$  against  $T[j..j + m)$ . We use the automaton to scan the text backwards from  $T[j + m - 1]$ . When the automaton has scanned  $T[j + i..j + m)$ :

- If the automaton is in an accept state, then  $T[j + i..j + m)$  is a **prefix** of  $P$ .
  - $\Rightarrow$  If  $i = 0$ , we found an occurrence.
  - $\Rightarrow$  Otherwise, mark the prefix match by setting  $shift = i$ . This is the length of the shift that would achieve a matching alignment.
- If the automaton can still reach an accept state, then  $T[j + i..j + m)$  is a **factor** of  $P$ .
  - $\Rightarrow$  Continue scanning.
- When the automaton can no more reach an accept state:
  - $\Rightarrow$  Stop scanning and shift:  $j \leftarrow j + shift$ .

BNDM does a [bitparallel](#) simulation of the nondeterministic automaton, which is quite similar to Shift-And.

The state of the automaton is stored in a bitvector  $D$ . When the automaton has scanned  $T[j + i..j + m)$ :

- $D.k = 1$  if and only if there is a path from the initial state to state  $k$  with the string  $(T[j + i..j + m))^R$ , and thus  $T[j + i..j + m) = P[m - k - 1..2m - k - i - 1)$ .
- If  $D.(m - 1) = 1$ , then  $T[j + i..j + m)$  is a prefix of the pattern.
- If  $D = 0$ , then the automaton can no more reach an accept state.

Updating  $D$  uses precomputed bitvectors  $B[c]$ , for all  $c \in \Sigma$ :

- $B[c].i = 1$  if and only if  $P[m - 1 - i] = P^R[i] = c$ .

The update when reading  $T[j + i]$  is familiar:  $D = (D \ll 1) \& B[T[j + i]]$

- Note that there is no “| 1”. This is because  $D.(-1) = 0$  always after reading at least one character, so the shift brings the right bit to  $D.0$ .
- Before reading anything  $D.(-1) = 1$ . This exception is handled by starting the computation with the first shift already performed. Because of this, the shift is done at the end of the loop.

**Algorithm 2.16:** BNDM

Input: text  $T = T[0 \dots n]$ , pattern  $P = P[0 \dots m]$

Output: position of the first occurrence of  $P$  in  $T$

Preprocess:

- (1) **for**  $c \in \Sigma$  **do**  $B[c] \leftarrow 0$
- (2) **for**  $i \leftarrow 0$  **to**  $m - 1$  **do**  $B[P[m - 1 - i]] \leftarrow B[P[m - 1 - i]] + 2^i$

Search:

- (3)  $j \leftarrow 0$
- (4) **while**  $j + m \leq n$  **do**
- (5)      $i \leftarrow m$ ;  $shift \leftarrow m$
- (6)      $D \leftarrow 2^m - 1$                                  //  $D \leftarrow 1^m$
- (7)     **while**  $D \neq 0$  **do**
  - // Now  $T[j + i..j + m)$  is a pattern factor
  - (8)          $i \leftarrow i - 1$
  - (9)          $D \leftarrow D \& B[T[j + i]]$
  - (10)        **if**  $D \& 2^{m-1} \neq 0$  **then**
    - // Now  $T[j + i..j + m)$  is a pattern prefix
    - (11)           **if**  $i = 0$  **then return**  $j$
    - (12)           **else**  $shift \leftarrow i$
  - (13)          $D \leftarrow D \ll 1$
- (14)      $j \leftarrow j + shift$
- (15) **return**  $n$

**Example 2.17:**  $P = \text{assi}$ ,  $T = \text{apassi}$ .

$B[c], c \in \{\text{a,i,p,s}\}$	$D$ when scanning <b>apas</b> backwards
<u>a i p s</u>	<u>s a p a</u>
i 0 1 0 0	i 1 0 0 0
s 0 0 0 1	s 1 1 0 0
s 0 0 0 1	s 1 1 0 0
a 1 0 0 0	a 1 0 <u>1</u> 0 $\Rightarrow \text{shift} = 2$

$D$  when scanning **apassi** backwards

	<u>i s s a p a</u>
i	1 1 0 0 0
s	1 0 1 0 0
s	1 0 0 1 0
a	1 0 0 0 <u>1</u> $\Rightarrow \text{occurrence}$

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

- Preprocessing time is  $\mathcal{O}(\sigma + m)$ .
- In the worst case, the search time is  $\mathcal{O}(mn)$ .  
For example,  $P = \mathbf{a}^{m-1}\mathbf{b}$  and  $T = \mathbf{a}^n$ .
- In the best case, the search time is  $\mathcal{O}(n/m)$ .  
For example,  $P = \mathbf{b}^m$  and  $T = \mathbf{a}^n$ .
- In the average case, the search time is  $\mathcal{O}((n/m) + (n/m) \log_\sigma m)$ .  
This is optimal! It has been proven that any algorithm needs to inspect  $\Omega((n/m) + (n/m) \log_\sigma m)$  text characters on average.

When  $m > w$ , there are several options:

- Use multi-word bitvectors.
- Search for a pattern prefix of length  $w$  and check the rest when the prefix is found.
- Use BDM or BOM.

- The search time of BDM and BOM is  $\mathcal{O}((n/m) + (n/m) \log_{\sigma} m)$ , which is optimal on **average**. (BNDM is optimal only when  $m \leq w$ .)
- MP and KMP are optimal in the **worst case** but not in the average case.
- There are also algorithms that are optimal in **both** cases. They are based on similar techniques, but we will not describe them here.



## Crochemore

- Two of the exact string matching algorithms we have seen, the brute force algorithm and the Karp–Rabin algorithm, use only a **constant amount of extra space** in addition to the text and the pattern. All the other algorithms use some additional data structures whose size is proportional to the length of the pattern.
- On the other hand, the brute force algorithm and the Karp–Rabin algorithm have a worst case running time of  $\mathcal{O}(mn)$ .
- Remarkably, there exists algorithms with  $\mathcal{O}(n)$  worst case time that need only a constant amount of extra space. One of them is the **Crochemore algorithm**.
- We will only outline the main ideas of the algorithm here without detailed proofs.

All the linear time, constant extra space algorithms are based on the *periodicity* properties of the pattern.

**Definition 2.18:** Let  $S[0..m)$  be a string. An integer  $p \in [1..m]$  is a **period** of  $S$ , if  $S[i] = S[i + p]$  for all  $i \in [0..m - p)$ . The **smallest period** of  $S$  is denoted  $per(S)$ .  $S$  is  **$k$ -periodic** if  $m \geq k \cdot per(S)$ .

**Example 2.19:** The periods of  $S_1 = \text{aabaaabaa}$  are 4, 7, 8 and 9. The periods of  $S_2 = \text{abcabcabcabca}$  are 3, 6, 9, 12 and 13.  $S_2$  is 3-periodic but  $S_1$  is not.

There is a strong connection between periods and borders.

**Lemma 2.20:**  $p$  is a period of  $S[0..m)$  if and only if  $S$  has a proper border of length  $m - p$ .

**Proof.** Both conditions hold if and only if  $S[0..m - p) = S[p..m)$ .  $\square$

**Corollary 2.21:** The length of the longest proper border of  $S$  is  $m - per(S)$ .

The Crochemore algorithm resembles the Morris–Pratt algorithm at a high level:

- When the pattern  $P$  is aligned against a text factor  $T[j..j + m)$ , they compute the longest common prefix  $\ell = \text{lcp}(P, T[j..j + m))$  and report an occurrence if  $\ell = m$ . Otherwise, they shift the pattern forward.
- MP shifts the pattern forward by  $\ell - \text{fail}[\ell]$  positions. Recall that  $\text{fail}[\ell]$  in MP is the length of the longest proper border of  $P[0..\ell)$ . Thus the pattern shift by MP is  $\text{per}(P[0..\ell))$ .
- In the next lcp computation, MP skips the first  $\text{fail}[\ell] = \ell - \text{per}(P[0..\ell))$  characters (cf. lcp-comparison).
- Thus knowing  $\text{per}(P[0..\ell))$  is sufficient to emulate MP shift and skip.

Crochemore either does the same shift and skip as MP, or a shorter shift than MP and starts the lcp comparison from scratch:

- If  $P[0..\ell)$  is 3-periodic, then compute  $per(P[0..\ell))$  and do the MP shift and skip.
- If  $P[0..\ell)$  is not 3-periodic, then shift by  $\lfloor \ell/3 \rfloor + 1 \leq per(P[0..\ell))$  and start the lcp comparison from scratch.
- Note that the latter case is inoptimal but always safe: no occurrence is missed.

To find out if  $P[0..\ell)$  is 3-periodic and to compute  $per(P[0..\ell))$  if it is, Crochemore uses another combinatorial concept.

**Definition 2.22:** Let  $MS(S)$  denote the lexicographically maximal suffix of a string  $S$ . If  $S = MS(S)$ ,  $S$  is called self-maximal.

Period computation is easier for maximal suffixes and self-maximal strings than for arbitrary strings.

**Lemma 2.23:** Let  $S[0..m)$  be a self-maximal string and let  $p = \text{per}(S)$ . For any  $c \in \Sigma$ ,

$$\begin{array}{ll} MS(Sc) = Sc \text{ and } \text{per}(Sc) = p & \text{if } c = S[m - p] \\ MS(Sc) = Sc \text{ and } \text{per}(Sc) = m + 1 & \text{if } c < S[m - p] \\ MS(Sc) \neq Sc & \text{if } c > S[m - p] \end{array}$$

Furthermore, let  $r = m \bmod p$  and  $R = S[m - r..m)$ . Then  $R$  is self-maximal and

$$MS(Sc) = MS(Rc) \quad \text{if } c > S[m - p]$$

Crochemore's algorithm computes the maximal suffix and its period for  $P[0..\ell]$  incrementally using Lemma 2.23. The following algorithm updates the maximal suffix information when the match is extended by one character.

**Algorithm 2.24:** Update-MS( $P, \ell, s, p$ )

Input: a string  $P$  and integers  $\ell, s, p$  such that

$MS(P[0..\ell]) = P[s..\ell]$  and  $p = per(P[s..\ell])$ .

Output: a triple  $(\ell + 1, s', p')$  such that

$MS(P[0..\ell + 1]) = P[s'..\ell + 1]$  and  $p' = per(P[s'..\ell + 1])$ .

```

(1) if  $\ell = 0$  then return  $(1, 0, 1)$ 
(2)  $i \leftarrow \ell$ 
(3) while  $i < \ell + 1$  do
    //  $P[s..i]$  is self-maximal and  $p = per(P[s..i])$ 
(4)   if  $P[i] > P[i - p]$  then
(5)      $i \leftarrow i - ((i - s) \bmod p)$ 
(6)      $s \leftarrow i$ 
(7)      $p \leftarrow 1$ 
(8)   else if  $P[i] < P[i - p]$  then
(9)      $p \leftarrow i - s + 1$ 
(10)   $i \leftarrow i + 1$ 
(11) return  $(\ell + 1, s, p)$ 

```

As the final piece of the Crochemore algorithm, the following result shows how to use the maximal suffix information to obtain information about the periodicity of the full string.

**Lemma 2.25:** Let  $S[0..m)$  be a string and let  $S[s..m) = MS(S)$  and  $p = per(MS(S))$ .

- $S$  is 3-periodic if and only if  $p \leq m/3$  and  $S[0..s) = S[p..p + s)$ .
- If  $S$  is 3-periodic, then  $per(S) = p$ .

The algorithm is given on the next slide.

### Algorithm 2.26: Crochemore

Input: strings  $T[0..n)$  (text) and  $P[0..m)$  (pattern).

Output: position of the first occurrence of  $P$  in  $T$

```
(1)  $j \leftarrow \ell \leftarrow p \leftarrow s \leftarrow 0$ 
(2) while  $j + m \leq n$  do
(3)   while  $j + \ell < n$  and  $\ell < m$  and  $T[j + \ell] = P[\ell]$  do
(4)      $(\ell, s, p) \leftarrow \text{Update-MS}(P, \ell, s, p)$ 
      //  $\ell = \text{lcp}(P, T[j..j + m))$ 
(5)   if  $\ell = m$  then return  $j$ 
      //  $\text{MS}(P[0..\ell]) = P[s..\ell]$  and  $p = \text{per}(P[s..\ell])$ 
(6)   if  $p \leq \ell/3$  and  $P[0..s) = P[p..p + s)$  then
      //  $\text{per}(P[0..\ell]) = p$ 
(7)      $j \leftarrow j + p$ 
(8)      $\ell \leftarrow \ell - p$ 
(9)   else //  $\text{per}(P[0..\ell]) > \ell/3$ 
(10)     $j \leftarrow j + \lfloor \ell/3 \rfloor + 1$ 
(11)     $(\ell, s, p) \leftarrow (0, 0, 0)$ 
(12) return  $n$ 
```



For **general alphabet**:

- The time complexity is  $\mathcal{O}(n)$ .
- The algorithm uses only a **constant** number of integer variables in addition to the strings  $P$  and  $T$ .

Crochemore is not competitive in practice. However, there are situations, where the pattern can be very long and the space complexity is more important than speed.

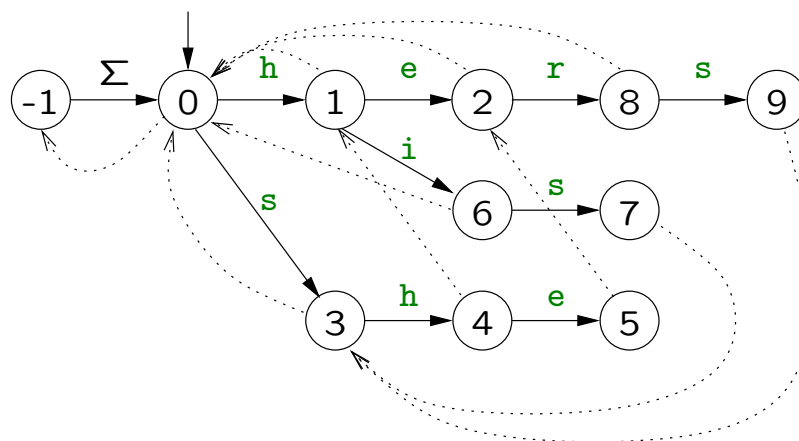
There are also other linear time, constant extra space algorithms. All of them are based on string periodicity in some way.

## Aho–Corasick

Given a text  $T$  and a set  $\mathcal{P} = \{P_1, P_2, \dots, P_k\}$  of patterns, the **multiple exact string matching** problem asks for the occurrences of all the patterns in the text. The Aho–Corasick algorithm is an extension of the Morris–Pratt algorithm for multiple exact string matching.

Aho–Corasick uses the trie  $trie(\mathcal{P})$  as an **automaton** and augments it with a failure function similar to the Morris–Pratt failure function.

**Example 2.27:** Aho–Corasick automaton for  $\mathcal{P} = \{\text{he}, \text{she}, \text{his}, \text{hers}\}$ .



Let  $S_v$  denote the string represented by a node  $v$  in the trie. The components of the AC automaton are:

- $root$  is the root and  $child()$  the child function of the trie.
- $fail(v) = u$  such that  $S_u$  is the **longest proper suffix** of  $S_v$  represented by any trie node  $u$ .
- $patterns(v)$  is the set of pattern indices  $i$  such that  $P_i$  is a **suffix** of  $S_v$ .

**Example 2.28:** For the automaton in Example 2.27,  $patterns(2) = \{1\}$  (**he**),  $patterns(5) = \{1, 2\}$  (**he**, **she**),  $patterns(7) = \{3\}$  (**his**),  $patterns(9) = \{4\}$  (**hers**), and  $patterns(v) = \emptyset$  for all other nodes  $v$ .

At each stage of the matching, the algorithm computes the node  $v$  such that  $S_v$  is the longest suffix of  $T[0..j]$  represented by any node.

**Algorithm 2.29:** Aho–Corasick

Input: text  $T$ , pattern set  $\mathcal{P} = \{P_1, P_2, \dots, P_k\}$ .

Output: all pairs  $(i, j)$  such that  $P_i$  occurs in  $T$  ending at  $j$ .

- (1)  $(root, child(), fail(), patterns()) \leftarrow \text{Construct-AC-Automaton}(\mathcal{P})$
- (2)  $v \leftarrow root$
- (3) **for**  $j \leftarrow 0$  **to**  $n - 1$  **do**
- (4)     **while**  $child(v, T[j]) = \perp$  **do**  $v \leftarrow fail(v)$
- (5)      $v \leftarrow child(v, T[j])$
- (6)     **for**  $i \in patterns(v)$  **do** output  $(i, j)$

The construction of the automaton is done in two phases: the trie construction and the failure links computation.

**Algorithm 2.30:** Construct-AC-Automaton

Input: pattern set  $\mathcal{P} = \{P_1, P_2, \dots, P_k\}$ .

Output: AC automaton:  $root$ ,  $child()$ ,  $fail()$  and  $patterns()$ .

- (1)  $(root, child(), patterns()) \leftarrow \text{Construct-AC-Trie}(\mathcal{P})$
- (2)  $(fail(), patterns()) \leftarrow \text{Compute-AC-Fail}(root, child(), patterns())$
- (3) **return**  $(root, child(), fail(), patterns())$

**Algorithm 2.31:** Construct-AC-Trie

Input: pattern set  $\mathcal{P} = \{P_1, P_2, \dots, P_k\}$ .

Output: AC trie:  $root$ ,  $child()$  and  $patterns()$ .

```
(1) Create new node  $root$ 
(2) for  $i \leftarrow 1$  to  $k$  do
(3)    $v \leftarrow root$ ;  $j \leftarrow 0$ 
(4)   while  $child(v, P_i[j]) \neq \perp$  do
(5)      $v \leftarrow child(v, P_i[j])$ ;  $j \leftarrow j + 1$ 
(6)   while  $j < |P_i|$  do
(7)     Create new node  $u$ 
(8)      $child(v, P_i[j]) \leftarrow u$ 
(9)      $v \leftarrow u$ ;  $j \leftarrow j + 1$ 
(10)   $patterns(v) \leftarrow \{i\}$ 
(11) return ( $root, child(), patterns()$ )
```

Lines (3)–(10) perform the standard trie insertion (Algorithm 1.2).

- Line (10) marks  $v$  as a representative of  $P_i$ .
- The creation of a new node  $v$  initializes  $patterns(v)$  to  $\emptyset$  (in addition to initializing  $child(v, c)$  to  $\perp$  for all  $c \in \Sigma$ ).

**Algorithm 2.32:** Compute-AC-Fail

Input: AC trie:  $root$ ,  $child()$  and  $patterns()$

Output: AC failure function  $fail()$  and updated  $patterns()$

- (1) Create new node  $fallback$
- (2) **for**  $c \in \Sigma$  **do**  $child(fallback, c) \leftarrow root$
- (3)  $fail(root) \leftarrow fallback$
- (4)  $queue \leftarrow \{root\}$
- (5) **while**  $queue \neq \emptyset$  **do**
- (6)      $u \leftarrow \text{popfront}(queue)$
- (7)     **for**  $c \in \Sigma$  such that  $child(u, c) \neq \perp$  **do**
- (8)          $v \leftarrow child(u, c)$
- (9)          $w \leftarrow fail(u)$
- (10)        **while**  $child(w, c) = \perp$  **do**  $w \leftarrow fail(w)$
- (11)         $fail(v) \leftarrow child(w, c)$
- (12)         $patterns(v) \leftarrow patterns(v) \cup patterns(fail(v))$
- (13)         $\text{pushback}(queue, v)$
- (14) **return** ( $fail()$ ,  $patterns()$ )

The algorithm does a **breath first traversal** of the trie. This ensures that correct values of  $fail()$  and  $patterns()$  are already computed when needed.

$fail(v)$  is correctly computed on lines (8)–(11):

- Let  $fail^*(v) = \{v, fail(v), fail(fail(v)), \dots, root\}$ . These nodes are exactly the trie nodes that represent suffixes of  $S_v$ .
- Let  $u = parent(v)$  and  $child(u, c) = v$ . Then  $S_v = S_u c$  and a string  $S$  is a suffix of  $S_u$  iff  $Sc$  is suffix of  $S_v$ . Thus for any node  $w$ 
  - If  $w \in fail^*(v) \setminus \{root\}$ , then  $parent(w) \in fail^*(u)$ .
  - If  $w \in fail^*(u)$  and  $child(w, c) \neq \perp$ , then  $child(w, c) \in fail^*(v)$ .
- Therefore,  $fail(v) = child(w, c)$ , where  $w$  is the first node in  $fail^*(u)$  other than  $u$  such that  $child(w, c) \neq \perp$ , or  $fail(v) = root$  if no such node  $w$  exists.

$patterns(v)$  is correctly computed on line (12):

$$\begin{aligned}
 patterns(v) &= \{i \mid P_i \text{ is a suffix of } S_v\} \\
 &= \{i \mid P_i = S_w \text{ and } w \in fail^*(v)\} \\
 &= \{i \mid P_i = S_v\} \cup patterns(fail(v))
 \end{aligned}$$

Assuming  $\sigma$  is constant:

- The search time is  $\mathcal{O}(n)$ .
- The space complexity is  $\mathcal{O}(m)$ , where  $m = \|\mathcal{P}\|$ .
  - The implementation of *patterns()* requires care (exercise).
- The preprocessing time is  $\mathcal{O}(m)$ , where  $m = \|\mathcal{P}\|$ .
  - The only non-trivial issue is the while-loop on line (10).
  - Let  $root, v_1, v_2, \dots, v_\ell$  be the nodes on the path from root to a node representing a pattern  $P_i$ . Let  $w_j = fail(v_j)$  for all  $j$ . Let  $depth(v)$  be the depth of a node  $v$  ( $depth(root) = 0$ ).
  - When processing  $v_j$  and computing  $w_j = fail(v_j)$ , we have  $depth(w_j) = depth(w_{j-1}) + 1$  before line (10) and  $depth(w_j) \leq depth(w_{j-1}) + 1 - t_j$  after line (10), where  $t_j$  is the number of rounds in the while-loop.
  - Thus, the total number of rounds in the while-loop when processing the nodes  $v_1, v_2, \dots, v_\ell$  is at most  $\ell = |P_i|$ , and thus over the whole algorithm at most  $\|\mathcal{P}\|$ .

The analysis when  $\sigma$  is not constant is left as an exercise.



## Summary: Exact String Matching

Exact string matching is a fundamental problem in stringology. We have seen several different algorithms for solving the problem.

The properties of the algorithms vary with respect to worst case time complexity, average case time complexity, type of alphabet (general/integer) and even space complexity.

The algorithms use a wide range of completely different techniques:

- There exists numerous algorithms for exact string matching but most of them use variations or combinations of the techniques we have seen (study groups).
- Many of the techniques can be adapted to other problems. All of the techniques have some uses in practice.