

**Algorithm 2.15:** BNDMInput: text  $T = T[0..n]$ , pattern  $P = P[0..m]$ 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  $n$ 

```

89

**Example 2.16:**  $P = \text{ass}i$ ,  $T = \text{apass}i$ .

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

 $D$  when scanning **ass***i* backwards

<u>a</u> s s i
i 0 0 0 1 1
s 0 0 1 0 1
s 0 1 0 0 1
a <u>1</u> 0 0 0 1 $\Rightarrow$ occurrence

90

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

91

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

92

**Karp–Rabin**The Karp–Rabin hash function (Definition 1.37) was originally developed for solving the exact string matching problem. The idea is to compute the hash values or **fingerprints**  $H(P)$  and  $H(T[j..j+m])$  for all  $j \in [0..n-m]$ .

- If  $H(P) \neq H(T[j..j+m])$ , then we must have  $P \neq T[j..j+m]$ .
- If  $H(P) = H(T[j..j+m])$ , the algorithm compares  $P$  and  $T[j..j+m]$  in brute force manner. If  $P \neq T[j..j+m]$ , this is a **false positive**.

The text factor fingerprints are computed in a **sliding window** fashion. The fingerprint for  $T[j+1..j+1+m] = \alpha T[j+m]$  is computed from the fingerprint for  $T[j..j+m] = T[j]\alpha$  in constant time using Lemma 1.38:

$$H(T[j+1..j+1+m]) = (H(T[j]\alpha) - H(T[j]) \cdot r^{m-1} \cdot r + H(T[j+m])) \bmod q$$

$$= (H(T[j..j+m]) - T[j] \cdot r^{m-1} \cdot r + T[j+m]) \bmod q.$$

A hash function that supports this kind of sliding window computation is known as a **rolling hash function**.

93

**Algorithm 2.17:** Karp–RabinInput: text  $T = T[0..n]$ , pattern  $P = P[0..m]$ Output: position of the first occurrence of  $P$  in  $T$ 

```

(1) Choose  $q$  and  $r$ ;  $s \leftarrow r^{m-1} \bmod q$ 
(2)  $hp \leftarrow 0$ ;  $ht \leftarrow 0$ 
(3) for  $i \leftarrow 0$  to  $m-1$  do  $hp \leftarrow (hp \cdot r + P[i]) \bmod q$  //  $hp = H(P)$ 
(4) for  $j \leftarrow 0$  to  $n-m$  do  $ht \leftarrow (ht \cdot r + T[j]) \bmod q$ 
(5) for  $j \leftarrow 0$  to  $n-m-1$  do
(6)   if  $hp = ht$  then if  $P = T[j..j+m]$  then return  $j$ 
(7)    $ht \leftarrow ((ht - T[j] \cdot s) \cdot r + T[j+m]) \bmod q$ 
(8) if  $hp = ht$  then if  $P = T[j..j+m]$  then return  $j$ 
(9) return  $n$ 

```

On an integer alphabet:

- The worst case time complexity is  $\mathcal{O}(mn)$ .
- The average case time complexity is  $\mathcal{O}(m+n)$ .

Karp–Rabin is not competitive in practice for a single pattern, but can be for multiple patterns (exercise).

94

**Crochemore**

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. In the next lcp computation, MP skips the first  $\text{fail}[\ell]$  characters (cf. lcp-comparison).
- Crochemore either does the same shift and skip as MP, or a shorter shift and starts the lcp comparison from scratch. Note that the latter case is inoptimal but always safe: no occurrence is missed.

Despite sometimes shorter shifts and less efficient lcp computation, Crochemore runs in **linear time**. More remarkably, it does so without any preprocessing and using only **constant extra space** in addition to  $P$  and  $T$ .We will only outline the main ideas of the algorithm without detailed proofs. Even then we will need some concepts from **combinatorics on words**, a branch of mathematics that studies combinatorial properties of strings.

95

**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  $\text{per}(S)$ .  $S$  is  **$k$ -periodic** if  $m/\text{per}(S) \geq k$ .**Example 2.19:** The periods of  $S_1 = \text{aabaabaa}$  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  $S[0..m-p] = S[p..m]$ . □**Corollary 2.21:** The length of the longest proper border of  $S$  is  $m - \text{per}(S)$ .

96

Recall that  $fail[\ell]$  in MP is the length of the longest proper border of  $P[0..\ell]$ . Thus the pattern shift by MP is  $\ell - fail[\ell] = per(P[0..\ell])$  and the lcp skip is  $fail[\ell] = \ell - per(P[0..\ell])$ . Thus knowing  $per(P[0..\ell])$  is sufficient to emulate MP shift and skip.

The Crochemore algorithm has two cases:

- 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.

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.

97

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

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

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

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

- (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  $P[i - p] > P[i]$  then
- (5)  $i \leftarrow i - ((i - s) \bmod p)$
- (6)  $s \leftarrow i$
- (7)  $p \leftarrow 1$
- (8) else if  $P[i - p] < P[i]$  then
- (9)  $p \leftarrow i - s + 1$
- (10)  $i \leftarrow i + 1$
- (11) return  $(\ell + 1, s, p)$

99

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

101

**Algorithm 2.28:** 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  $(i, j)$

Let  $S_v$  denote the string that node  $v$  represents.

- $root$  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  $S_v$ .

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

103

**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  $p = per(S)$ . For any  $a \in \Sigma$ ,

$$\begin{aligned} MS(Sa) &= Sa \text{ and } per(Sa) = p & \text{if } a = S[m - p] \\ MS(Sa) &= Sa \text{ and } per(Sa) = m + 1 & \text{if } a > S[m - p] \\ MS(Sa) &\neq Sa & \text{if } a < S[m - p] \end{aligned}$$

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

$$MS(Sa) = MS(Ra) \quad \text{if } a < S[m - p]$$

The proof is omitted.

98

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.

- Time complexity is  $\mathcal{O}(n)$ . (Proof omitted.)
- It uses only a constant number of integer variables in addition to the strings  $P$  and  $T$ .
- Works on ordered alphabet.

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.

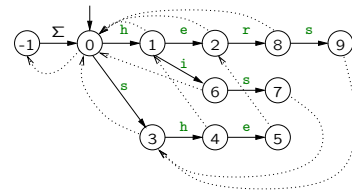
100

**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  $\mathcal{P} = \{\text{he, she, his, hers}\}$ .



101

102

**Algorithm 2.29:** Construct-AC-Automaton

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

Output: AC automaton:  $root, child(), fail()$  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)  $(fail(), patterns()) \leftarrow \text{Compute-AC-Fail}(root, child(), patterns())$
- (12) return  $(root, child(), fail(), patterns())$

Lines (3)–(10) form 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  $c \in \Sigma$ ).

104

**Algorithm 2.30:** Compute-AC-FailInput: 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)    $u \leftarrow popfront(queue)$ 
(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)     $pushback(queue, v)$ 
(14) return  $(fail(), patterns())$ 

```

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

105

Assuming  $\sigma$  is constant:

- The search time is  $\mathcal{O}(n)$ .
- The space complexity is  $\mathcal{O}(m)$ , where  $m = ||\mathcal{P}||$ .
  - 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.

107

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

- The nodes that represent suffixes of  $S_v$  that are exactly  $fail^*(v) = \{v, fail(v), fail(fail(v)), \dots, root\}$ .
- 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  $w \in fail^*(v)$ , then  $parent(fail(v)) \in fail^*(u)$ .
  - 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  $child(w, c) \neq \perp$ .

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

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

106

**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 (ordered/integer) and even space complexity.

The algorithms use a wide range of completely different techniques:

- There exists numerous algorithms for exact string matching but almost all them are based on these techniques.
- Many of the techniques can be adapted to other problems. All of the techniques have some uses in practice too.

108