

3

Multiple string matching

3.1 Basic concepts

The single string matching problem may be extended in a natural way to search simultaneously for a set of strings $P = \{p^1, p^2, \dots, p^r\}$, where each p^i is a string $p^i = p_1^i p_2^i \dots p_{m_i}^i$ over a finite character set Σ . Denote by $|P|$ the sum of the lengths of the strings in P , more formally $|P| = \sum_{i=1}^r |p^i| = \sum_{i=1}^r m_i$. Let ℓ_{min} be the minimum length of a pattern in P and ℓ_{max} the maximum. As before, the search is done in a text $T = t_1 t_2 \dots t_n$.

Strings in P may be factors, prefixes, suffixes, or even the same as others. For example, if we search for the set $\{\text{ATATA}, \text{TATA}\}$ in a DNA sequence, each time we find an occurrence of **ATATA** we also find an occurrence of the second string. Hence, the total number of occurrences can be $r \times n$. To make the multistring matching problem precise, we consider that we are interested in reporting all pairs (i, j) such that $t_{j-|p^i|+1} \dots t_j$ is equal to p^i .

The simplest solution to this problem is to repeat r searches with one of the algorithms of Chapter 2. This leads to a total worst-case complexity of $O(|P|)$ for the preprocessing and $O(r \times n)$ for the search.

The worst-case search complexity can be reduced to $O(n + nocc)$, where $nocc$ is the total number of occurrences, by using some kind of *extension* of the search algorithms for a single pattern. The average complexity can also be improved, although it is difficult to think in terms of “average” complexity, since many parameters play a role in the running time of the algorithms. The most important parameters are the size of the alphabet, the number of patterns, the distribution of the lengths of the patterns (particularly the minimum size), and the memory available.

We again denote by θ an object that is not defined. For instance, when we write **While** $q \neq \theta$ **Do**, it means we iterate while q is defined.

Troughout this chapter, we will consider the example in Figure 3.1. We

will simultaneously search for the three strings “announce”, “annual”, and “annually”.

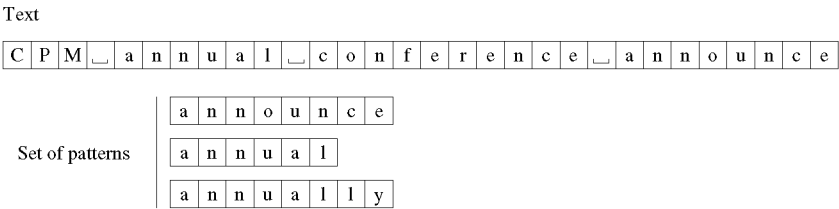


Fig. 3.1. Simultaneously searching three strings in our example text.

As with a single string, sets of natural language strings usually contain fewer repetitions than sets of DNA sequences. To show the tricky cases that could occur, we also show the behavior of our algorithms when searching for the set of strings **ATATATA**, **TATAT**, **ACGATAT** in the sequence **AGATACGATATATAC**.

The three approaches for searching a single string (Chapter 2) lead to several extensions for searching a set of strings. For each approach, there are usually many possible extensions, according to the way the set of patterns is managed and the way the shifts are obtained. The notion of a *search window* is not relevant for multiple string matching, which will become clear soon. We present in this chapter the empirically most efficient extensions, which are usually also the simplest.

Prefix searching (Figure 3.2) The search is done forward, reading the characters of the text one after another with an automaton built on the set P . For each position of the text, we compute through this automaton the longest suffix of the text read that is also a prefix of one of the strings of P . The most famous algorithm that uses this approach is **Aho-Corasick** [AC75].

Suffix searching (Figure 3.3) A position pos is slid along the text, from which we search backward for a suffix of any of the strings. As with a single pattern, we shift pos according to the next occurrence of the suffix read in P . This approach may avoid reading all the characters of the text.

Factor searching (Figure 3.4) A position pos is also slid along the text, from which we read backwards a factor of some prefix of size ℓ_{min} of the strings in P . It also may avoid reading all the characters of the text.

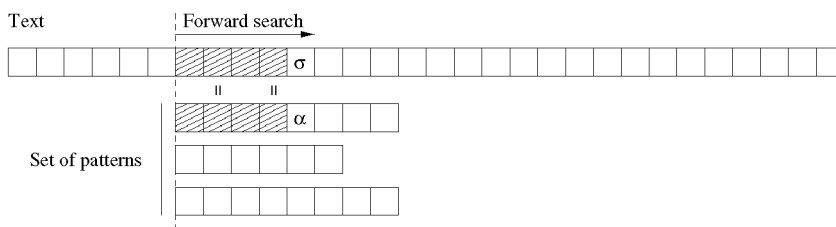


Fig. 3.2. First approach: We compute the longest prefix of a pattern in the set that is also a suffix of the text read. It requires reading all the characters of the text at least once.

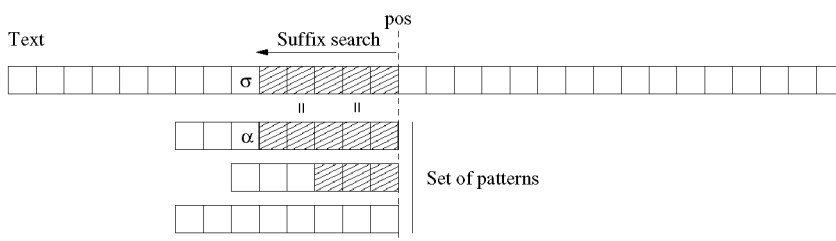


Fig. 3.3. Second approach: We search backwards for a suffix of one of the strings. It avoids, on average, reading all characters of the text.

Before describing these three approaches in depth, we introduce a basic data structure on a set of strings, called a *trie*. This structure is used by most of the classical multistring matching algorithms. The *trie* of the set $P = \{p^1, p^2, \dots, p^r\}$ is a rooted directed tree that represents the set P ; that is, every path starting from the root is labeled by one of the strings p^i , and, conversely, every string $p^i \in P$ labels a path from the root. Below, unless specified, paths start at the root. Every state q corresponding to an entire string is marked as *terminal*, and a function $F(q)$ points to a list of

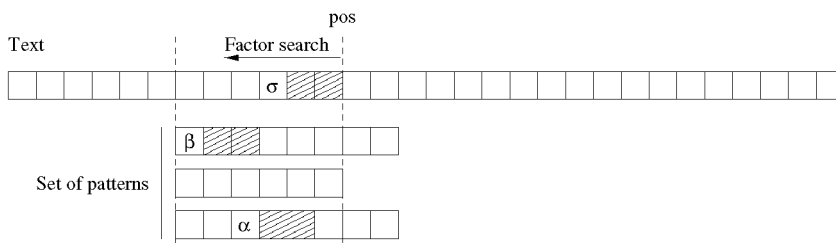


Fig. 3.4. Third approach: We search for a factor of any of the patterns in the current window.

all the numbers of the strings in P that correspond to q . We give the trie for $P = \{\text{announce, annual, annually}\}$ in Figure 3.5.

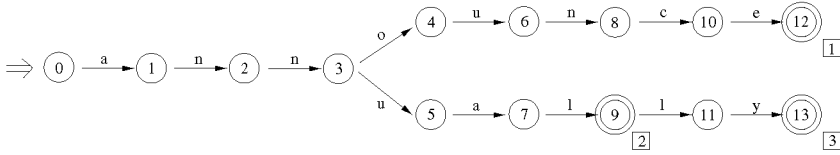


Fig. 3.5. Trie for $P = \{\text{announce, annual, annually}\}$. The function F of each terminal state is represented with squares. It indicates the identifier of the string in P .

We usually use the automata notation for describing a trie, since a trie is also a deterministic acyclic automaton recognizing the corresponding set of strings. The trie of a set P can be built in $O(|P|)$ time by inserting the strings p^i one by one into the tree, starting at the root, and building the corresponding transitions. Pseudo-code for the trie construction is given in Figure 3.6.

```

Trie( $P = \{p^1, p^2, \dots, p^r\}$ )
1.   Create an initial non terminal state 0
2.   For  $i \in 1 \dots r$  Do
3.        $Current \leftarrow$  initial state 0
4.        $j \leftarrow 1$ 
5.       While  $j \leq m_i$  AND  $\delta(Current, p_j^i) \neq \theta$  Do
6.            $Current \leftarrow \delta(Current, p_j^i)$ 
7.            $j \leftarrow j + 1$ 
8.       End of while
9.       While  $j \leq m_i$  Do
10.          Create a new non terminal state  $State$ 
11.           $\delta(Current, p_j^i) \leftarrow State$ 
12.           $Current \leftarrow State$ 
13.           $j \leftarrow j + 1$ 
14.       End of while
15.       If  $Current$  is already terminal Then  $F(Current) \leftarrow F(Current) \cup \{i\}$ 
16.       Else mark  $Current$  as terminal,  $F(Current) \leftarrow \{i\}$ 
17.   End of for

```

Fig. 3.6. Pseudo-code for the construction of a trie from a set of strings $P = \{p^1, p^2, \dots, p^r\}$. The strings are taken one by one and inserted into the tree.

The size of the trie depends on the implementation of the transitions. The simplest implementation is for each state q of the trie to code $\delta(q, *)$ in a table of size $|\Sigma|$. Then the total size of the trie of a set P is worst-case $|\Sigma| \times |P|$. This representation is usually used when the sizes of the set

of strings and of the alphabet are not too large. It has the advantage of passing through a transition in constant time $O(1)$ by performing an access to a table.

Since the total number of transitions is at most $|P|$, it is possible to code all the transitions in $O(|P|)$ space, independently of the size of the alphabet. However, the time to pass through a transition increases. If the transitions of each state are coded with a linked list, sorted or not, this time grows to $O(|\Sigma|)$ in the worst and average cases. It can be reduced to $O(\log |\Sigma|)$ by coding the transitions with balanced trees [CLR90], but this complicates the code.

We now describe in detail the three general approaches to search for a set of strings.

3.2 Prefix based approach

The extension of the prefix based approach leads to the **Multiple Shift-And** and **Aho-Corasick** algorithms. As with a single pattern (Section 2.2), we assume that we have read the text up to position i and that we know the length of the longest suffix of $t_1 \dots t_i$ that is a prefix of a pattern $p^k \in P$. The algorithmic problem is to calculate this length after reading the next character of the text.

In the single pattern case, there were two ways of finding this length. One was based on managing a bit array with bit-parallelism. For multiple pattern matching, this technique is only practical for very small patterns, because the total length $|P|$ has to be smaller than a few computer words. Nevertheless, this possibility is widely used for extended string matching (Chapter 4) and approximate string matching (Chapter 6). We call this algorithm **Multiple Shift-And**.

The solution, when the length of the set is too large to fit in computer words, is to find a mechanism that computes the size of the longest suffix of the text read that is also a prefix of one of the strings of P , in amortized constant time per character. This is what the **Aho-Corasick** algorithm does, with a linear time $O(|P|)$ preprocessing phase.

3.2.1 Multiple Shift-And algorithm

The bit-parallelism approach is only valuable when the set $P = \{p^1, \dots, p^r\}$ has a total length $|P|$ small enough to fit in a few computer words. For simplicity, we assume below that $|P|$ is smaller than w . The idea is to perform with bit-parallelism all the computations required by the **Shift-**

And algorithm (Section 2.2.2) for the r strings in the same computer word [BYG89b].

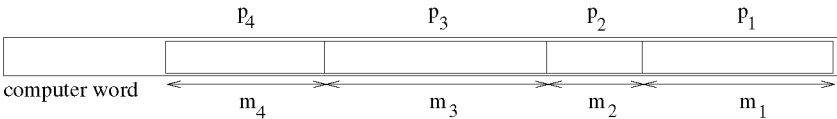


Fig. 3.7. **Multiple Shift-And** algorithm. The total size of the patterns has to fit in w .

We pack the strings together in the computer word, as in Figure 3.7. Then, for each new character of the text, we perform the computations for the strings of P like in the **Shift-And** algorithm. The initialization word DI is the concatenation of the initialization words for each string, that is,

$$DI \leftarrow 0^{m_r-1}1 \dots 0^{m_2-1}10^{m_1-1}1$$

Similarly, the final test is

$$DF \leftarrow 10^{m_r-1} \dots 10^{m_2-1}10^{m_1-1}$$

The main loop is the same as for the **Shift-And** algorithm. Pseudo-code is given in Figure 3.8.

The **Shift-Or** trick (Section 2.2.2) cannot be used here, since the shift “ $<<$ ” only introduces a zero to the right, and we need a zero in each position that begins a new string of P in the computer word.

Example using English We search for the set of strings $P = \{\text{announce, annual, annually}\}$ in the text “annual.announce”.

Table B

a	00010001010001000000001
c	00000000000000001000000
e	000000000000000010000000
l	01100000100000000000000
n	00000110000011000100110
o	000000000000000000001000
u	00001000001000000100000
y	100000000000000000000000
*	000000000000000000000000

$DI = 0000000100000100000001$
 $DF = 1000000010000010000000$
 $D = 0000000000000000000000$

1. Reading a

$B[a]$	00000001000001000000001
$D =$	00010001010001000000001
	00000001000001000000001

2. Reading n

$B[n]$	00000011000011000000011
$D =$	0000011000011000100110
	00000010000010000000010

3. Reading n

$B[n]$	00000101000101000000101
$D =$	0000011000011000100110
	00000100000100000000100

Multiple Shift-And($P = \{p^1, p^2, \dots, p^r\}$, $T = t_1 t_2 \dots t_n$)

1. **Preprocessing**
2. **For** $c \in \Sigma$ **Do** $B[c] \leftarrow 0^{|P|}$
3. $\ell \leftarrow 0$
4. **For** $k \in 1 \dots r$ **Do**
5. **For** $j \in 1 \dots m_k$ **Do** $B[p_j^k] \leftarrow B[p_j^k] \mid 0^{|P|-\ell-j} 10^{\ell+j-1}$
6. $\ell \leftarrow \ell + m_k$
7. **End of for**
8. $DI \leftarrow 0^{m_r-1} 1 \dots 0^{m_2-1} 10^{m_1-1} 1$
9. $DF \leftarrow 10^{m_r-1} \dots 10^{m_2-1} 10^{m_1-1}$
10. **Searching**
11. $D \leftarrow 0^{|P|}$
12. **For** $pos \in 1 \dots n$ **Do**
13. $D \leftarrow ((D < 1) \mid DI) \ \& \ B[t_{pos}]$
14. **If** $D \ \& \ DF \neq 0^{|P|}$ **Then**
15. Check which patterns match
16. Report the corresponding occurrences ending in pos
17. **End of if**
18. **End of for**

Fig. 3.8. **Multiple Shift-And** algorithm. The total length of the patterns $|P|$ has to be less than w . We let $m_k = |p^k|$.

4. Reading u	8. Reading a
$B[u]$ 00001001001001000001001	$B[a]$ 00000001000001000000001
$D =$ 00001000001000000000000	$D =$ 00000001000001000000001
5. Reading a	9. Reading n
$B[a]$ 00010001010001000000001	$B[n]$ 00000011000011000000011
$D =$ 00010001010001000000001	$D =$ 00000010000010000000010
6. Reading l	10. Reading n
$B[l]$ 00100011100011000000011	$B[n]$ 00000101000101000000101
$D =$ 00100000100000000000000	$D =$ 00000100000100000000100
$D \ \& \ DF \neq 0^{ P }$, we check the patterns that match, and we mark an occurrence of annual .	11. Reading o
7. Reading _	$B[o]$ 00001001001001000001001
$B[_]$ 00000000000000000000000	$D =$ 000000000000000000001000
$D =$ 00000000000000000000000	12. Reading u
	$B[u]$ 0000000100000100010001
	$D =$ 000000000000000000001000

13. Reading n

$$\begin{array}{r} 0000000100000100100001 \\ B[n] \ 0000011000011000100110 \\ \hline D = 0000000000000000100000 \end{array}$$

14. Reading c

$$\begin{array}{r} 0000000100000101000001 \\ B[c] \ 0000000000000000100000 \\ \hline D = 0000000000000000100000 \end{array}$$

15. Reading e

$$\begin{array}{r} 0000000100000110000001 \\ B[c] \ 0000000000000001000000 \\ \hline D = 0000000000000001000000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$, we check the patterns that match, and we mark an occurrence of **announce**.

Example using DNA We search for the set of strings $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ in the text **AGATACGATATATAC**.

Table B

A	0101001010101010101
C	00000100000000000000
G	00001000000000000000
T	1010000101010101010
*	00000000000000000000

$$\begin{array}{l} DI = 0000001000010000001 \\ DF = 1000000100001000000 \\ D = 0000000000000000000 \end{array}$$

1. Reading A

$$\begin{array}{r} 0000001000010000001 \\ B[A] \ 0101001010101010101 \\ \hline D = 0000001000000000001 \end{array}$$

2. Reading G

$$\begin{array}{r} 0000011000010000011 \\ B[G] \ 0000100000000000000 \\ \hline D = 0000000000000000000 \end{array}$$

3. Reading A

$$\begin{array}{r} 0000001000010000001 \\ B[A] \ 0101001010101010101 \\ \hline D = 0000001000000000001 \end{array}$$

4. Reading T

$$\begin{array}{r} 0000011000010000011 \\ B[T] \ 1010000101010101010 \\ \hline D = 0000000000010000010 \end{array}$$

5. Reading A

$$\begin{array}{r} 0000001000110000101 \\ B[A] \ 0101001010101010101 \\ \hline D = 0000001000100000101 \end{array}$$

6. Reading C

$$\begin{array}{r} 0000011001010001011 \\ B[C] \ 0000010000000000000 \\ \hline D = 0000010000000000000 \end{array}$$

7. Reading G

$$\begin{array}{r} 0000101000010000001 \\ B[G] \ 0000100000000000000 \\ \hline D = 0000100000000000000 \end{array}$$

8. Reading A

$$\begin{array}{r} 0001001000010000001 \\ B[A] \ 0101001010101010101 \\ \hline D = 0001001000000000001 \end{array}$$

9. Reading T

$$\begin{array}{r} 0010011000010000011 \\ B[T] \ 1010000101010101010 \\ \hline D = 0010000000010000010 \end{array}$$

10. Reading A

$$\begin{array}{r} 0100001000110000101 \\ B[A] \ 0101001010101010101 \\ \hline D = 0100001000100000101 \end{array}$$

11. Reading T

$$\begin{array}{r} 1000011001010001011 \\ B[T] \ 1010000101010101010 \\ \hline D = 1000000001010001010 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$, we check the patterns that match, and we mark an occurrence of **ACGATAT**.

12. Reading A

$$\begin{array}{r} 00000010101110010101 \\ B[A] \quad 0101001010101010101 \\ \hline D = 0000001010100010101 \end{array}$$

13. Reading T

$$\begin{array}{r} 0000011101010101011 \\ B[T] \quad 1010000101010101010 \\ \hline D = 0000000101010101010 \end{array}$$

$D \& DF \neq 0^{|P|}$, we check the patterns that match, and we mark an occurrence of TATAT.

14. Reading A

$$\begin{array}{r} 0000001010111010101 \\ B[A] \quad 0101001010101010101 \\ \hline D = 0000001010101010101 \end{array}$$

$D \& DF \neq 0^{|P|}$, we check the patterns that match, and we mark an occurrence of ATATATA.

15. Reading C

$$\begin{array}{r} 0000011101010101011 \\ B[C] \quad 0000010000000000000 \\ \hline D = 0000000000000000000 \end{array}$$

3.2.2 Basic Aho-Corasick algorithm

The algorithm of Aho and Corasick [AC75] is an extension of the **Knuth-Morris-Pratt** algorithm (Section 2.2.1) for a set of patterns.

The algorithm uses a special automaton, called the *Aho-Corasick automaton*, built on P . It is the trie of P augmented with a “supply function” S_{AC} .

Formally, we denote by q a state of the trie of P , and by $L(q)$ the label of the path from the initial state to q . Then $S_{AC}(q)$ is defined, except for the initial state, as the state reached when the automaton reads the longest suffix of $L(q)$ that is also a prefix of some $p^i \in P$. This is a kind of extension of a *border* (Section 2.2.1) to a set of strings. The supply state of the initial state is set to θ . A *supply link* goes from each state q to $S_{AC}(q)$, and a *supply path* is a chain of supply links.

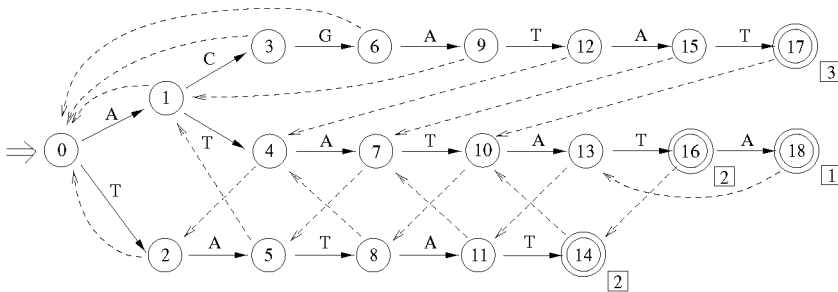


Fig. 3.9. Aho-Corasick automaton for the set $\{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$. The dashed links represent the state-to-state supply function S_{AC} . Double-circled states are terminal.

The Aho-Corasick automaton for the set $\{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ is shown in Figure 3.9. On this automaton, for instance, $L(15) = \text{ACGATA}$, its longest suffix that is also a prefix of one of the patterns is **ATA**, which leads

to state 7, and hence $S_{AC}(15) = 7$. The terminal states are those of the trie that correspond to an entire pattern, and also all states whose supply paths go down through another terminal state on their way to the root. In Figure 3.9, for instance, state 16 is terminal because $S_{AC}(16)$ is terminal.

We assume that a prefix $t_1 t_2 \dots t_i$ of the text has already been read, and that the longest suffix of $t_1 \dots t_i$ that is also a prefix of one of the patterns leads to a state *Current* in the Aho-Corasick automaton. We denote this longest suffix $v = L(\textit{Current})$. We want to read t_{i+1} and compute for $t_1 \dots t_i t_{i+1}$ the new longest suffix u . There are two cases.

1. If there exists an outgoing transition from *Current* to another state f in the trie labeled by t_{i+1} , then the new *Current* state becomes f , and $u = L(f) = ut_{i+1}$ is the new longest prefix of one of the patterns that is a suffix of $t_1 \dots t_{i+1}$.
2. If not (i.e., we fail reading t_{i+1} in the tree), we go down the supply path of q until either
 - (a) we find a state on the path followed by t_{i+1} . In this case, the current state becomes the arrival state f by the transition t_{i+1} , and $u = L(f)$.
 - (b) we reach θ , which means that the longest suffix u we search for is the empty string ε , and we move to the initial state.

Pseudo-code for the search algorithm is given in Figure 3.10. The complexity of the search phase is simple to evaluate, if we observe that we cannot go down more supply links than text characters we read. The number of supply links crossed through is then bounded by n , and the number of transitions used (real transitions plus supply links) is bounded by $2n$. The number of character comparisons depends on how the transitions of the automaton are implemented. The complexity is $O(n + \textit{nocc})$ if they are coded with a table, and $O(n \log |\Sigma| + \textit{nocc})$ with balanced trees.

To construct the Aho-Corasick automaton we begin by building the trie of the set of strings P with the algorithm in Figure 3.6. The states of the Aho-Corasick automaton are those of the trie. The initial state is the same and the terminal states of the trie are also terminal. We build the supply function S_{AC} on this trie in transversal order, which is the order we numbered the states in Figure 3.9.

We assume that we have computed the supply function of all the states before state *Current* in transversal order. We consider the parent *Parent* of *Current* in the trie, leading to *Current* by σ , that is, $\textit{Current} = \delta_{AC}(\textit{Parent}, \sigma)$. The supply state $S_{AC}(\textit{Parent})$ has already been computed. We search

```

Aho-Corasick( $P = \{p^1, p^2, \dots, p^r\}, T = t_1 t_2 \dots t_n$ )
1.  Preprocessing
2.     $AC \leftarrow \text{Build\_AC}(P)$ 
3.  Searching
4.     $Current \leftarrow$  Initial state of the automaton  $AC$ 
5.    For  $pos \in 1 \dots n$  Do
6.      While  $\delta_{AC}(Current, t_{pos}) = \theta$  AND  $S_{AC}(Current) \neq \theta$  Do
7.         $Current \leftarrow S_{AC}(Current)$ 
8.      End of while
9.      If  $\delta_{AC}(Current, t_{pos}) \neq \theta$  Then
10.         $Current \leftarrow \delta_{AC}(Current, t_{pos})$ 
11.      Else  $Current \leftarrow$  initial state of  $AC$ 
12.      End of if
13.      If  $Current$  is terminal Then
14.        Mark all the occurrences  $(F(Current), pos)$ 
15.      End of if
16.    End of for

```

Fig. 3.10. **Aho-Corasick** algorithm to search for a set $P = \{p^1, p^2, \dots, p^r\}$ of strings. It uses the Aho-Corasick automaton to compute at each text character t_{pos} the longest prefix of any pattern p^k that is also a suffix of the text read $t_1 \dots t_{pos}$.

for the state where u ends, u being the longest suffix of $v = L(Current)$ that labels a path in the trie. The string v has the form $v'\sigma$. If there exists such a nonempty string u , since it is a suffix of v , it must be of the form $u = u'\sigma$. In that case, u' is a suffix of v' that is the label of a path in the trie.

If $S_{AC}(Parent)$ has an outgoing transition by σ to a state h , then $w = L(S_{AC}(Parent))$ is the longest suffix of v' that is the label of a path, and $w\sigma$ is also a label of a path in the trie. Consequently, it is the longest suffix u of $v = v'\sigma$ that we are searching for, and $S_{AC}(Current)$ has to be set to h .

If $S_{AC}(Parent)$ does not have an outgoing transition by σ , we consider $S_{AC}(S_{AC}(Parent))$ and so on. We repeat the operation, until either we find a state on the supply path that has an outgoing transition by σ , or we find θ , which means that u is the empty string ε and $S_{AC}(Current)$ has to be set to the initial state.

The mechanism is similar to the **Aho-Corasick** search algorithm itself. Its pseudo-code is given in Figure 3.11. Complexity is evaluated with the observation we made for the whole algorithm: We do not go down more supply links than the total number of real transitions, which is bounded by $O(|P|)$. So the number of total transitions used (real transitions plus supply links) is bounded by $2 \times |P|$. Like for the search phase, the complexity in terms of comparisons of characters depends on how the transitions of the

automaton are implemented. It is $O(|P|)$ if they are coded with a table, and $O(|P| \log |\Sigma|)$ with balanced trees.

```
Build_AC( $P = \{p^1, p^2, \dots, p^r\}$ )
1.   $AC\_trie \leftarrow \mathbf{Trie}(P)$ 
     $\delta_{AC}$  is its transition function
2.   $Initial\_state \leftarrow \text{root of } AC\_trie$ 
3.   $S_{AC}(Initial\_state) \leftarrow \theta$ 
4.  For  $Current$  in transversal order Do
5.     $Parent \leftarrow \text{parent of } Current \text{ in } AC\_trie$ 
6.     $\sigma \leftarrow \text{label of the transition from } Parent \text{ to } Current$ 
7.     $Down \leftarrow S_{AC}(Parent)$ 
8.    While  $Down \neq \theta$  AND  $\delta_{AC}(Down, \sigma) = \theta$  Do
9.       $Down \leftarrow S_{AC}(Down)$ 
10.   End of while
11.   If  $Down \neq \theta$  Then
12.      $S_{AC}(Current) \leftarrow \delta_{AC}(Down, \sigma)$ 
13.     If  $S_{AC}(Current)$  is terminal Then
14.       Mark  $Current$  as terminal
15.        $F(Current) \leftarrow F(Current) \cup F(S_{AC}(Current))$ 
16.     End of if
17.   Else  $S_{AC}(Current) \leftarrow Initial\_state$ 
18.   End of if
19. End of for
```

Fig. 3.11. Construction of the Aho-Corasick automaton. The state *Current* goes in transversal order through the trie *AC_trie* built on *P*. The state *Down* goes down the supply links from the parent of *Current*, looking for an outgoing transition labeled with the same character as between *Current* and its parent. $F(Current)$ is initialized as empty when *Current* is first marked as terminal.

Example using English We search for the set of strings $P = \{\text{announce, annual, annually}\}$ in the text “annual_announce”. The Aho-Corasick automaton built on *P* is shown in Figure 3.12.

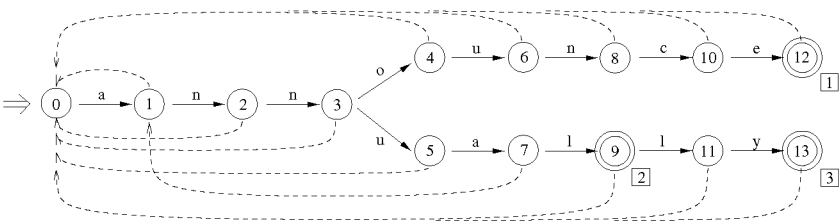


Fig. 3.12. Aho-Corasick automaton of our example set $P = \{\text{announce, annual, annually}\}$. Double-circled states are terminal.

$Current \leftarrow 0$.

- | | |
|---|---|
| 1. Reading a
$Current \leftarrow 1 = \delta(0, a)$ | 8. Reading a
$Current \leftarrow 1 = \delta(0, a)$ |
| 2. Reading n
$Current \leftarrow 2 = \delta(1, n)$ | 9. Reading n
$Current \leftarrow 2 = \delta(1, n)$ |
| 3. Reading n
$Current \leftarrow 3 = \delta(2, n)$ | 10. Reading n
$Current \leftarrow 3 = \delta(2, n)$ |
| 4. Reading u
$Current \leftarrow 5 = \delta(3, u)$ | 11. Reading o
$Current \leftarrow 4 = \delta(3, o)$ |
| 5. Reading a
$Current \leftarrow 7 = \delta(5, a)$ | 12. Reading u
$Current \leftarrow 6 = \delta(4, u)$ |
| 6. Reading l
$Current \leftarrow 9 = \delta(7, a)$.
The state 9 is terminal; we mark an occurrence of $F(9) \rightarrow \text{annual}$. | 13. Reading n
$Current \leftarrow 8 = \delta(6, n)$ |
| 7. Reading -
$\delta(9, -) = \theta$. We jump to $0 = S_{AC}(9)$.
$\delta(0, -) = \theta$, we jump to $\theta = S_{AC}(0)$.
We continue the search from the initial state 0, $Current \leftarrow 0$. | 14. Reading c
$Current \leftarrow 10 = \delta(8, c)$ |
| | 15. Reading e
$Current \leftarrow 12 = \delta(10, e)$.
The state 12 is terminal; we mark an occurrence of $F(12) \rightarrow \text{announce}$. |

Example using DNA We search for the set of strings $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ in the text **AGATACGATATATAC**. We again use the Aho-Corasick automaton built on P already shown in Figure 3.9.

$Current \leftarrow 0$.

- | | |
|---|--|
| 1. Reading A
$Current \leftarrow 1 = \delta(0, A)$ | 5. Reading A
$Current \leftarrow 7 = \delta(4, A)$ |
| 2. Reading G
$\delta(1, G) = \theta$. We jump to $0 = S_{AC}(1)$.
$\delta(0, G) = \theta$; we jump to $\theta = S_{AC}(0)$.
We continue the search from the initial state 0, $Current \leftarrow 0$. | 6. Reading C
$\delta(7, C) = \theta$. We jump to $5 = S_{AC}(7)$.
$\delta(5, C) = \theta$; we jump to $1 = S_{AC}(7)$.
$\delta(1, C) = 3$, $Current \leftarrow 3$. |
| 3. Reading A
$Current \leftarrow 1 = \delta(0, A)$ | 7. Reading G
$Current \leftarrow 6 = \delta(3, G)$ |
| 4. Reading T
$Current \leftarrow 4 = \delta(1, T)$ | 8. Reading A
$Current \leftarrow 9 = \delta(6, A)$ |

9. Reading T
 $Current \leftarrow 12 = \delta(9, T)$
10. Reading A
 $Current \leftarrow 15 = \delta(12, A)$
11. Reading T
 $Current \leftarrow 17 = \delta(12, T)$. The state 17 is terminal; we mark an occurrence of $F(17) \rightarrow \text{ACGATAT}$.
12. Reading A
 $\delta(17, A) = \theta$. We jump to $10 = S_{AC}(17)$. $\delta(10, A) = 13$, $Current \leftarrow 13$.
13. Reading T
 $Current \leftarrow 16 = \delta(13, T)$. The state 16 is terminal; we mark an occurrence of $F(16) \rightarrow \text{TATAT}$.
14. Reading A
 $Current \leftarrow 18 = \delta(16, A)$. The state 18 is terminal; we mark an occurrence of $F(18) \rightarrow \text{ATATATA}$.
15. Reading C
 $\delta(18, C) = \theta$. We jump to $13 = S_{AC}(18)$. $\delta(13, C) = \theta$; we jump to $11 = S_{AC}(13)$. $\delta(11, C) = \theta$; we jump to $7 = S_{AC}(11)$. $\delta(7, C) = \theta$; we jump to $1 = S_{AC}(7)$. $\delta(1, C) = 3$, $Current \leftarrow 3$.

3.2.3 Advanced Aho-Corasick algorithm

The above algorithm permits a powerful variant. The idea is to precompute all the transitions simulated by the supply function. We then obtain a complete automaton (all the states have an outgoing transition by every character of the alphabet) that we name the *extended Aho-Corasick automaton*.

This completion can be computed using the supply function. We first complete the outgoing transitions of the initial state with a loop, which means $\delta(0, \sigma) \leftarrow 0$ for each new letter σ . Now, let $Current$ be a state of the automaton taken in transversal order. We compute the missing outgoing transitions of $Current$ by using the formula $\delta(Current, \sigma) = \delta(S_{AC}(Current), \sigma)$ for each new letter σ .

The drawback to this automaton is the large amount of memory space it requires. It is $O(|P| \times |\Sigma|)$ independently of the way the transitions are implemented. This construction is useful for relatively small sets and alphabets. A trade-off that is often used is to compute the new transitions *on the fly* if there is memory left. This was done in the first version of the well-known Unix application *Grep*.

3.3 Suffix based approach

The experimental results of Chapter 2 show that the suffix based approach is usually faster than the prefix based one. So it is natural to try to extend

the suffix based approach to sets of patterns. The first attempt was that of Commentz-Walter in 1979 [CW79]. It is a direct extension of the **Boyer-Moore** algorithm. The **Horspool** algorithm has also been extended, but it is much less powerful for multiple string matching than for single patterns. A stronger extension is the **Wu-Manber** algorithm, which is practical, simple, and efficient.

3.3.1 Commentz-Walter idea

The **Commentz-Walter** [CW79] algorithm is a “natural” extension of the **Boyer-Moore** algorithm (Section 2.3.1). This algorithm is never faster in practice than **Aho-Corasick** or other algorithms presented below. However, it is historically important because it was the first expected sublinear multistring matching algorithm, and it was implemented in the second version of the Unix application *Grep*. Currently, this algorithm does not have a real case of application, and we just present the idea it is based on.

The **Commentz-Walter** algorithm represents $P = \{p^1, \dots, p^r\}$ using a trie of the reverse patterns $P^{rv} = \{(p^1)^{rv}, \dots, (p^r)^{rv}\}$ inside which the text is read. A position pos is slid along the text, beginning at position ℓ_{min} so as not to skip a possible occurrence. For each such new position, we read backwards the longest suffix u of $t_1 \dots t_{pos}$ that is also a suffix of one of the patterns. If we find an occurrence, we mark it. Then, we shift the position of the search to the right, using the three functions d_1, d_2, d_3 of the **Boyer-Moore** algorithm extended to a set of strings. The first two functions are computed for each state of the trie, and to shift we consider them at the last state q we crossed when reading the longest suffix u .

- $d_1(q)$ is the minimal shift such that $u = L(q)$ matches a factor of some $p^j \in P$.
- $d_2(q)$ is the minimal shift such that a suffix of $u = L(q)$ matches a prefix of some $p^j \in P$.

The last function $d_3[\alpha, k]$ is computed for each character α of the alphabet for positions $0 \leq k < \ell_{max}$. It is the minimal shift such that α read at position $pos - k$ matches another character of some $p^j \in P$.

For a visual idea of what these three functions do, the reader may refer to Figures 2.8, 2.9, and 2.10 of Chapter 2, where the three corresponding functions of the **Boyer-Moore** algorithm are shown.

We combine these three functions to compute a shift. Suppose that we read backwards k characters of the text from a position pos and this led to

state q . The shift $s[q, pos, k]$ is then computed with the following formula:

$$s[q, pos, k] = \min \left\{ \begin{array}{l} \max(d_1[q], d_3[t_{pos-k}, k]) \\ d_2[q] \end{array} \right.$$

The formula is the direct extension of the computation of the window shifts in the **Boyer-Moore** algorithm. As $d_2 \leq \ell_{min}$, the longest shift is bounded by ℓ_{min} , which is a necessary condition to avoid skipping an occurrence when shifting the position pos .

The **Commentz-Walter** algorithm is worst-case time $O(n \times \ell_{max})$ but sublinear on average if the number of patterns is not too large. The computation of the three functions d_1 , d_2 , and d_3 can be done in $O(|P|)$ time.

3.3.2 Set Horspool algorithm

The **Horspool** algorithm, similarly to **Boyer-Moore**, is directly extensible to a set of patterns. The new algorithm, which we call **Set Horspool**, can also be considered as a simplification of **Commentz-Walter**.

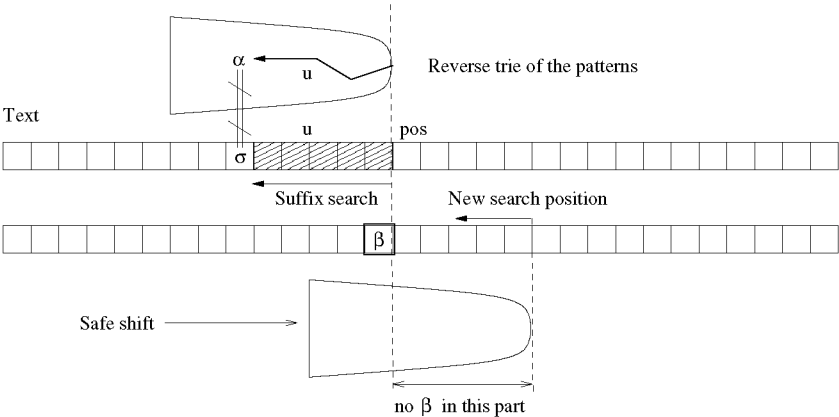


Fig. 3.13. **Horspool** algorithm for a set of patterns. The set is shifted according to the last character of the search window.

The general scheme is shown in Figure 3.13. We start reading the text backwards from a position pos initialized to ℓ_{min} to avoid skipping any occurrence. We read these characters in the trie built on the reverse patterns. If we reach a terminal state, we mark an occurrence. When we fail reading the text, we shift the position pos using the first character read (β in the figure). We shift until β is aligned with another β in the trie. If such a β does not exist, we simply shift by ℓ_{min} characters.

The **Set Horspool** algorithm is $O(n \times \ell_{max})$ time in the worst case. In

Set Horspool ($P = \{p^1, p^2, \dots, p^r\}$, $T = t_1 t_2 \dots t_n$)

1. **Preprocessing**
2. $HO \leftarrow \text{Trie}(P^{rv} = \{(p^1)^{rv}, \dots, (p^r)^{rv}\})$
 δ_{HO} is its transition function
3. **For** $c \in \Sigma$ **Do** $d[c] \leftarrow \ell_{min}$
4. **For** $j \in 1 \dots r$ **Do**
5. **For** $k \in 1 \dots m_j - 1$ **Do** $d[p_k^j] \leftarrow \min(d[p_k^j], m_j - k)$
6. **End of for**
7. **Searching**
8. $pos \leftarrow \ell_{min}$
9. **While** $pos \leq n$ **Do**
10. $j \leftarrow 0$, $Current \leftarrow$ initial state of HO
11. **While** $pos - j > 0$ AND $\delta_{HO}(t_{pos-j}, Current) \neq \theta$ **Do**
12. **If** $Current$ is terminal **Then**
13. Mark all the occurrences ($F(Current), pos$)
14. **End of if**
15. $Current \leftarrow \delta_{HO}(t_{pos-j}, Current)$
16. $j \leftarrow j + 1$
17. **End of while**
18. $pos \leftarrow pos + d[t_{pos}]$
19. **End of while**

Fig. 3.14. **Horspool** algorithm for a set of patterns. The shift is obtained with the first character t_{pos} read.

general, it is only efficient for a very small number of patterns on a relatively large alphabet.

Example using English We search for the set $P = \{\text{announce, annual, annually}\}$ in the text “CPM_{annual} conference announce”. The trie of the reverse patterns is shown in Figure 3.15

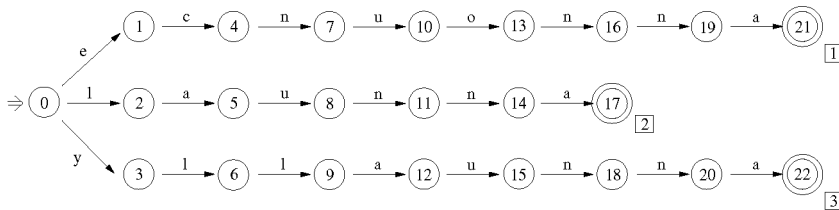


Fig. 3.15. Trie for the reverse set of $P = \{\text{announce, annual, annually}\}$, $P^{rv} = \{\text{ecnuonna, launna, yllaunna}\}$. Double-circled states are terminal.

$\ell_{min} = 6,$

$d = \begin{Bmatrix} \text{a} & \text{c} & \text{e} & \text{l} & \text{n} & \text{o} & \text{u} & \text{y} & * \\ 1 & 1 & 6 & 1 & 2 & 4 & 2 & 6 & 6 \end{Bmatrix}$

1. CPM_a [n] nual_conference_announce
n $\notin \{e, l, y\}$, $d[n] = 2$
2. CPM_ann [u] al_conference_announce
u $\notin \{e, l, y\}$, $d[u] = 2$
3. CPM_annua [l] _conference_announce
We read in the trie l, a, u, n, n, a. We reach the terminal state 17 and mark an occurrence of $F(17) \rightarrow \text{annual}$.
We re-use the first character read, $d[l] = 1$.
4. CPM_annual [-] conference_announce
- $\notin \{e, l, y\}$, $d[-] = 6$

5. CPM_annual_confe [r] ence_announce
r $\notin \{e, l, y\}$, $d[r] = 6$
6. CPM_annual_conference_ [a] nnonce
a $\notin \{e, l, y\}$, $d[a] = 1$
7. CPM_annual_conference_a [n] nnonce
n $\notin \{e, l, y\}$, $d[n] = 2$
8. CPM_annual_conference_ann [o] unce
o $\notin \{e, l, y\}$, $d[o] = 4$
9. CPM_annual_conference_announc [e]
We read in the trie e, c, n, u, o, n, n, a. We reach the terminal state 21 and mark an occurrence of $F(21) \rightarrow \text{announce}$.

Example using DNA We search for the set of strings $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ in the text **AGATACGATATATAC**. The trie of the reverse patterns is shown in Figure 3.16

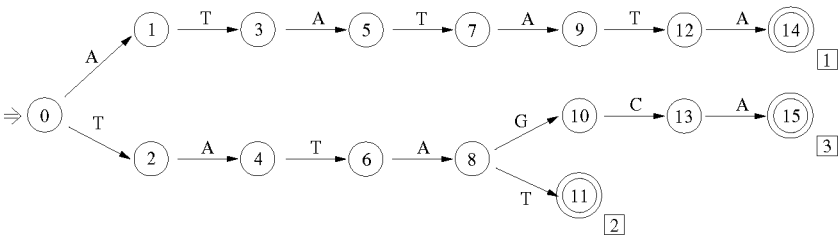


Fig. 3.16. Trie for the reverse set of $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$, $P^{rv} = \{\text{ATATATA}, \text{TATAT}, \text{TATAGCA}\}$. Double-circled states are terminal.

$\ell_{min} = 5,$

$\delta = \begin{Bmatrix} \text{A} & \text{C} & \text{G} & \text{T} & * \\ 1 & 5 & 4 & 1 & 5 \end{Bmatrix}$

1. AGAT [A] CGATATATAC
We read in the trie A, T, A, and we fail on the next G. We re-use the last character of the window, $d[A] = 1$.

2. AGATA [C] GATATATAC
C $\notin \{A, T\}$, $d[C] = 5$.
3. AGATACGATA [T] ATAC
We read in the trie T, A, T, A, G, C, A. We reach the terminal state 15 and mark an occurrence of $F(15) \rightarrow \text{ACGATAT}$.
We re-use the last character of the window, $d[T] = 1$.

4. AGATACGATAT A TAC

We read in the trie A, T, A, T, A, and we fail on the next G.

We re-use the first character read,
 $d[A] = 1$.

5. AGATACGATATA T AC

We read in the trie T, A, T, A, T. We reach the terminal state 11 and mark an occurrence of $F(11) \rightarrow \text{TATAT}$.

We re-use the first character read,
 $d[T] = 1$.

6. AGATACGATATAT A C

We read in the trie A, T, A, T, A, T, A. We reach the terminal state 14 and mark an occurrence of $F(14) \rightarrow \text{ATATATA}$.

We re-use the first character read,
 $d[A] = 1$.

7. AGATACGATATATA C

$C \notin \{A, T\}$, $d[C] = 5$. Then $pos > n$ and we stop the search.

3.3.3 Wu-Manber algorithm

The poor performance of the extension of **Horspool** to search a set of patterns is a direct consequence of the fact that the lengths of the shifts are usually decreasing, due to the high probability of finding each character of the alphabet in one of the strings.

The algorithm of Wu and Manber [WM94] bypasses this obstacle by reading blocks of characters, which reduces the probability that each block appears in one of the patterns. We consider blocks of length B . The difficulty is that there could be $|\Sigma|^B$ different blocks, requiring too much memory if B becomes large.

Wu and Manber overcome this problem by hashing all the possible blocks using a function h_1 into a limited size table *SHIFT*. Two blocks B_1 and B_2 can be associated with the same position in *SHIFT*. If we consider that for each new position we are reading a block Bl instead of the last character of the **Horspool** algorithm, then the shift given by Bl , $SHIFT(h_1(Bl))$, must be safe. To guarantee this, we save in $SHIFT(j)$ the minimum of the shifts of the blocks Bl such that $j = h_1(Bl)$. More precisely, the table *SHIFT* is built in the following way:

- If a block Bl does not appear in any string of P , we can safely shift $\ell_{min} - B + 1$ characters to the right. Hence we initialize the table by placing $\ell_{min} - B + 1$ everywhere.
- If Bl appears in one of the strings of P , we find its rightmost occurrence ending in j in a string p^i , and set the value of $SHIFT(h_1(Bl))$ to $m_i - j$. To compute all the values of the table *SHIFT*, we consider separately each $p^i = p_1^i \dots p_{m_i}^i$. For each block $B = p_{j-B+1}^i \dots p_j^i$, we find its corresponding cell $h_1(B)$ in *SHIFT*, and we place in $SHIFT(h_1(B))$ the minimum between the previous value and $m_i - j$.

```

Wu-Manber( $P = \{p^1, p^2, \dots, p^r\}$ ,  $T = t_1 t_2 \dots t_n$ )
1.   Preprocessing
2.       Computation of  $B$ 
3.       Construction of the hash tables SHIFT and HASH
4.   Searching
5.        $pos \leftarrow \ellmin$ 
6.       While  $pos \leq n$  Do
7.            $i \leftarrow h_1(t_{pos-B+1} \dots t_{pos})$ 
8.           If SHIFT[ $i$ ] = 0 Then
9.                $list \leftarrow HASH[h_2(t_{pos-B+1} \dots t_{pos})]$ 
10.              Verify all the patterns in list one by one against the text
11.               $pos \leftarrow pos + 1$ 
12.           Else  $pos \leftarrow pos + SHIFT[i]$ 
13.           End of if
14.       End of while

```

Fig. 3.17. **Wu-Manber** algorithm for searching a set of strings.

The size B varies with ℓmin , with the number of patterns, and with the size of the alphabet. Wu and Manber show that the value $B = \log_{|\Sigma|}(2 \times \ellmin \times r)$ yields the best experimental results. The size of the table *SHIFT* can also vary with the memory space available.

We can shift the search position along the text as long as the value of the shift is strictly positive. When the shift is zero, the text to the left of the search position may be one of the pattern strings. In this case Wu and Manber use a new hash table *HASH*. Each position *HASH*(j) contains a list of all the strings whose last block is hashed to j by a second hash function h_2 . This table permits us to select from P a subset of strings whose last block maps the block Bl read in the text.

For the search, similarly to the **Set Horspool** algorithm, we slide a position pos along the text, reading backwards a block Bl of B characters. The position pos is initialized to ℓmin . If $j = SHIFT(h_1(Bl)) > 0$, then we shift the window to $pos + j$ and continue the search. Otherwise, $SHIFT(h_1(Bl)) = 0$ and we select a set of strings using *HASH* that we compare to the text. Pseudo-code is given in Figure 3.17.

The original description of the algorithm [WM94] is quite fuzzy. Nothing is given in the article that permits you to calculate the best size of the tables *SHIFT* and *HASH*. Likewise, the hash functions are not specified. All these parameters affect the complexity. In practice, well parametrized, this algorithm is very fast. It is implemented in *Agrep* (Section 7.1.2).

We now present our two running examples. The **Wu-Manber** algorithm uses many hash functions and tables that are difficult to represent. We have chosen some that do not correspond to a real example, but permit us to show the interesting cases. We let $B = 2$ in the following tables.

Example using English We search for the set $P = \{\text{announce, annual, annually}\}$ in the text “CPM_annual_conference_announce”.

$$SHIFT[B] = \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline \text{string} & ll & no & ou & an & un & nc & ua & al & ly & nn & nu & ce & * \\ \hline \text{shift} & 1 & 3 & 4 & 1 & 0 & 0 & 2 & 0 & 5 \\ \hline \end{array}$$

$$HASH[B] = \begin{array}{|c|c|c|c|} \hline \text{string} & ce & ly & ua & al & * \\ \hline \text{string number in } P & 3, 1 & 2 & \emptyset \\ \hline \end{array}$$

1. CPM_ an nual_conference_announce
 $SHIFT[an] = 4$.
2. CPM_annu al _conference_announce
 $SHIFT[al] = 0$. $L = HASH[al] = \{2\}$.
3. CPM_annua l_ conference_announce
 $SHIFT[l_] = 5$.
4. CPM_annual_con fe rence_announce
 $SHIFT[fe] = 5$.
5. CPM_annual_conferen ce _announce
 $SHIFT[ce] = 0$. $L = HASH[ce] = \{3, 1\}$.
6. CPM_annual_conferenc e_ announce
 $SHIFT[e_] = 5$.
7. CPM_annual_conference_ann ou nce
 $SHIFT[ou] = 3$.
8. CPM_annual_conference_announ ce _
 $SHIFT[ce] = 0$. $L = HASH[ce] = \{3, 1\}$.

We compare p^2 against the text and mark its occurrence. We then shift the search position by 1.

We compare p^1 and p^3 against the text. The test succeeds for the string p^1 . Hence, we mark its occurrence.

We compare p^1 and p^3 against the text. No string matches. We shift the search position by 1.

Example using DNA We search the set of strings $P = \{\text{ATATATA, TATAT, ACGATAT}\}$ in the text AGATACGATATATAC.

$$SHIFT[B] = \begin{array}{|c|c|c|c|c|c|} \hline \text{string} & GA & TA & AT & CG & GA & AC & * \\ \hline \text{shift} & 0 & 0 & 3 & 4 & 4 \\ \hline \end{array}$$

$$HASH[B] = \begin{array}{|c|c|c|c|} \hline \text{string} & TA & AT & * \\ \hline \text{string number in } P & 1 & 2, 3 & \emptyset \\ \hline \end{array}$$

1. AGA TA CGATATATAC
 $SHIFT[TA] = 0$. $L = HASH[TA] = \{1\}$.

We compare p^1 against the text. The test fails. We shift the search position by 1.

2. AGAT AC GATATATAC
 $SHIFT[AC] = 4$.

3. AGATACGA TA TATAC
 $SHIFT[TA] = 0$. $L = HASH[TA] = \{1\}$.

We compare p^1 against the text. The test fails. We shift the search position by 1.

4. AGATACGAT AT ATAC
 $SHIFT[AT] = 0$. $L = HASH[AT] = \{2, 3\}$.

We compare p^2 and p^3 against the text. The string p^3 matches. We mark its occurrence. We shift the search position by 1.

5. AGATACGATA TA TAC
 $SHIFT[TA] = 0$. $L = HASH[TA] = \{1\}$.

We compare p^1 against the text. The test fails. We shift the search position by 1.

6. AGATACGATAT AT AC
 $SHIFT[AT] = 0$. $L = HASH[AT] = \{2, 3\}$.

We compare p^2 and p^3 against the text. The string p^2 matches. We mark its occurrence. We shift the search position by 1.

7. AGATACGATATA TA C
 $SHIFT[TA] = 0$. $L = HASH[TA] = \{1\}$.

We compare p^1 against the text. The string p^1 matches. We mark its occurrence. We shift the search position by 1.

8. AGATACGATATAT AC
 $SHIFT[AC] = 4$.

3.4 Factor based approach

The general factor based approach can be extended directly to a set of strings. We search backwards for the longest suffix u of the text that is also a factor of one of the strings in P . If we fail on a letter σ , then σu is not a factor in any of the strings; thus no string of P can overlap σu .

There are, however, two technical difficulties to overcome. The first problem is to shift the set of patterns safely to avoid skipping an occurrence; the second difficulty is to recognize the factors of a set of strings.

The first two factor based algorithms were the **Dawg-Match** [CCG⁺93, CCG⁺99] and the **MultiBDM** [CR94, Raf97]. They were developed with the aim of obtaining fast algorithms on average with good worst-case complexity. Indeed, they are all worst-case linear in the size of the text. But they are inherently complicated, and in practice their performance is poor.

As we aim to present the simplest and most efficient algorithms, we will not describe these two.

The two algorithms left that use the factor based approach are **Set Backward Dawg Matching (SBDM)** and **Set Backward Oracle Matching (SBOM)** [AR99]. They extend **BDM** and **BOM**, respectively (Chapter 2). We present the **SBDM** idea, on which **SBOM** is based.

Bit-parallelism is only valuable for a small set of strings. However, similarly to the **Multiple Shift-And** (Section 3.2.1), it permits efficient extended string matching (Chapter 4) and also approximate matching (Chapter 6). We present it first.

3.4.1 Multiple BNDM algorithm

The use of bit-parallelism to search a set of strings $P = \{p^1, \dots, p^r\}$ is efficient for sets such that $r \times \ell_{min}$ fits in a few computer words [NR00]. For simplicity we assume below that $r \times \ell_{min} \leq w$.

To perform longer shifts, we keep only the prefixes of size ℓ_{min} of the patterns. If we match a prefix, we directly verify the entire string against the text.

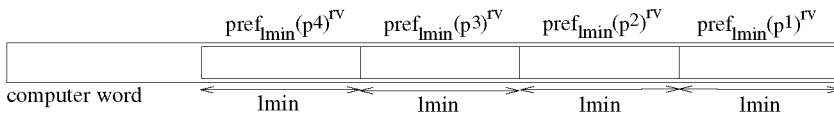


Fig. 3.18. **Multiple BNDM** algorithm. The total $r \times \ell_{min}$ has to fit in w . The notation $pref_{\ell_{min}}(p^i)$ denotes the prefix of size ℓ_{min} of p^i .

The prefixes are packed together as in Figure 3.18 and the search is similar to **BNDM** (Section 2.4.2), with the search performed for all the prefixes at the same time. The only difference is that we need to clear some bits after a shift. The mask CL in Figure 3.19 does that. It prevents the bits used to search for p^i from interfering with those used for p^{i+1} . The variable *last* is still used, but in this case it represents the position where a prefix of one of the strings begins. Pseudo-code of the whole algorithm is shown in Figure 3.19.

Example using English We search the text “CPM_{annual} conference announcement” for the set of strings $P = \{\text{announce, annual, annually}\}$.

Multiple BNDM ($p = p_1 p_2 \dots p_m$, $T = t_1 t_2 \dots t_n$)

```

1.  Preprocessing
2.    For  $c \in \Sigma$  Do  $B[c] \leftarrow 0^{|P|}$ 
3.     $\ell \leftarrow 0$ 
4.    For  $k \in 1 \dots r$  Do
5.       $\ell \leftarrow \ell + \ell_{min}$ 
6.      For  $j \in 1 \dots \ell_{min}$  Do  $B[p_j^k] \leftarrow B[p_j^k] \mid 0^{|P|-\ell+j-1} 10^{\ell-j}$ 
7.    End of for
8.     $CL \leftarrow (1^{\ell_{min}-1} 0)^r$ 
9.     $DF \leftarrow (10^{\ell_{min}-1})^r$ 
10. Searching
11.    $pos \leftarrow 0$ 
12.   While  $pos \leq n - m$  Do
13.      $j \leftarrow \ell_{min}$ ,  $last \leftarrow \ell_{min}$ 
14.      $D = 1^{|P|}$ 
15.     While  $D \neq 0^{|P|}$  Do
16.        $D \leftarrow D \& B[t_{pos+j}]$ 
17.        $j \leftarrow j - 1$ 
18.       If  $D \& DF \neq 0^{|P|}$  Then
19.         If  $j > 0$  Then  $last \leftarrow j$ 
20.       Else /* at least one prefix matches */
21.         Check which prefixes of length  $\ell_{min}$  match
22.          $p^i$  needs to be checked if
23.            $D \& 0^{|P|-\ell_{min} \times i} 10^{\ell_{min}-1} 0^{\ell_{min} \times (i-1)} \neq 0^{|P|}$ 
24.         Verify the corresponding string(s) against the text
25.         Report the occurrence(s) at  $pos + 1$ 
26.       End of if
27.     End of while
28.      $D \leftarrow (D << 1) \& CL$  /* Shifting and cleaning */
29.   End of while

```

Fig. 3.19. Bit-parallel code for the **Multiple BNDM** algorithm.

$$B = \begin{cases} \text{a} & 100010100010100000 \\ \text{l} & 000001000001000000 \\ \text{n} & 011000011000011001 \\ \text{o} & 0000000000000000100 \\ \text{u} & 000100000100000010 \\ \text{*} & 000000000000000000 \end{cases}$$

$DF = 100000100000100000$
 $CL = 111110111110111110$

$$\begin{array}{r} \begin{array}{r} 111111111111111111 \\ B[n] \ 011000011000011001 \\ D = 011000011000011001 \end{array} \\ \begin{array}{r} 110000110000110010 \\ B[a] \ 100010100010100000 \\ D = 100000100000100000 \end{array} \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

1. CPM_an nual_conference_announce
 $last \leftarrow 6$.
 $D \leftarrow 111111111111111111$

$$\begin{array}{r} \begin{array}{r} 000000000000000000 \\ B[_] \ 000000000000000000 \\ D = 000000000000000000 \end{array} \end{array}$$

2. CPM_ annual _conference_announce
 $last \leftarrow 6$.
 $D \leftarrow 11111111111111111111$

$$\begin{array}{r} 11111111111111111111 \\ B[l] \ 000001000001000000 \\ \hline D = 000001000001000000 \end{array}$$

$$\begin{array}{r} 000010000010000000 \\ B[a] \ 100010100010100000 \\ \hline D = 000010000010000000 \end{array}$$

$$\begin{array}{r} 000100000100000000 \\ B[u] \ 000100000100000010 \\ \hline D = 000100000100000000 \end{array}$$

$$\begin{array}{r} 001000001000000000 \\ B[n] \ 011000011000011001 \\ \hline D = 001000001000000000 \end{array}$$

$$\begin{array}{r} 010000010000000000 \\ B[n] \ 011000011000011001 \\ \hline D = 010000010000000000 \end{array}$$

$$\begin{array}{r} 100000100000000000 \\ B[a] \ 100010100010100000 \\ \hline D = 100000100000000000 \end{array}$$

D & $DF \neq 0^{|P|}$ and $j = 0$, so we check the patterns “annual” and “annually” against the text and mark the occurrence of “annual”.

3. CPM_annual_ confe _rence_announce
 $last \leftarrow 6$.
 $D \leftarrow 11111111111111111111$

$$\begin{array}{r} 11111111111111111111 \\ B[e] \ 000000000000000000 \\ \hline D = 000000000000000000 \end{array}$$

4. CPM_annual_confe rence_ announce
 $last \leftarrow 6$.
 $D \leftarrow 11111111111111111111$

$$\begin{array}{r} 11111111111111111111 \\ B[_] \ 000000000000000000 \\ \hline D = 000000000000000000 \end{array}$$

5. CPM_annual_conference_ announ ce
 $last \leftarrow 6$.
 $D \leftarrow 11111111111111111111$

$$\begin{array}{r} 11111111111111111111 \\ B[n] \ 011000011000011001 \\ \hline D = 011000011000011001 \end{array}$$

$$\begin{array}{r} 110000110000110010 \\ B[u] \ 000100000100000010 \\ \hline D = 000000000000000010 \end{array}$$

$$\begin{array}{r} 0000000000000000100 \\ B[o] \ 0000000000000000100 \\ \hline D = 0000000000000000100 \end{array}$$

$$\begin{array}{r} 00000000000000001000 \\ B[n] \ 011000011000011001 \\ \hline D = 00000000000000001000 \end{array}$$

$$\begin{array}{r} 00000000000000001000 \\ B[n] \ 011000011000011001 \\ \hline D = 00000000000000001000 \end{array}$$

$$\begin{array}{r} 000000000000000010000 \\ B[a] \ 100010100010100000 \\ \hline D = 000000000000000010000 \end{array}$$

D & $DF \neq 0^{|P|}$ and $j = 0$, so we check the string “announce” against the text and mark its occurrence.

The next shift of the search window gives $pos > n - \ell_{min}$ and the search stops.

Example using DNA We search for the set of strings $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ in the text **AGATACGATATATAC**.

$$B = \begin{Bmatrix} \begin{array}{|c|c|} \hline \text{A} & 101010101010010 \\ \hline \text{C} & 000000000001000 \\ \hline \text{G} & 000000000000100 \\ \hline \text{T} & 010101010100001 \\ \hline * & 000000000000000 \\ \hline \end{array} \end{Bmatrix}$$

$$\begin{array}{lcl} DF & = & 100001000010000 \\ CL & = & 111101111011110 \end{array}$$

1. AGATA CGATATATAC $last \leftarrow 5.$ $D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[A] \ 101010101010010 \\ \hline D = 101010101010010 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

$$\begin{array}{r} 010101010000100 \\ B[T] \ 010101010100001 \\ \hline D = 010101010000000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 101000100000000 \\ B[A] \ 101010101010010 \\ \hline D = 101000100000000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 2$

$$\begin{array}{r} 010001000000000 \\ B[G] \ 000000000000100 \\ \hline D = 000000000000000 \end{array}$$

2. AG ATACG ATATATAC $last \leftarrow 5.$ $D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[G] \ 000000000000100 \\ \hline D = 000000000000100 \end{array}$$

$$\begin{array}{r} 000000000001000 \\ B[C] \ 000000000001000 \\ \hline D = 000000000001000 \end{array}$$

$$\begin{array}{r} 000000000010000 \\ B[A] \ 101010101010010 \\ \hline D = 000000000010000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 2$

$$\begin{array}{r} 000000000000000 \\ B[T] \ 010101010100001 \\ \hline D = 000000000000000 \end{array}$$

3. AGAT ACGAT ATATAC $last \leftarrow 5.$ $D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[T] \ 010101010100001 \\ \hline D = 010101010100001 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

$$\begin{array}{r} 101000101000010 \\ B[A] \ 101010101010010 \\ \hline D = 101000101000010 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 010001010000100 \\ B[G] \ 000000000000100 \\ \hline D = 000000000000100 \end{array}$$

$$\begin{array}{r} 000000000001000 \\ B[C] \ 000000000001000 \\ \hline D = 000000000001000 \end{array}$$

$$\begin{array}{r} 000000000010000 \\ B[A] \ 101010101010010 \\ \hline D = 000000000010000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j = 0$, so we check
the pattern ACGATAT against the text
and mark its occurrence.

4. AGATACG ATATA TAC $last \leftarrow 5.$ $D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[A] \ 101010101010010 \\ \hline D = 101010101010010 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

$$\begin{array}{r} 010101010000100 \\ B[T] \ 010101010100001 \\ \hline D = 010101010000000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 101000100000000 \\ B[A] \ 101010101010010 \\ \hline D = 101000100000000 \end{array}$$

$D \ \& \ DF \neq 0^{|P|}$ and $j > 0$, then

$last \leftarrow 2$

$$\begin{array}{r} 0100010000000000 \\ B[T] 010101010100001 \\ \hline D = 0100010000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 1$

$$\begin{array}{r} 1000000000000000 \\ B[A] 101010101010010 \\ \hline D = 1000000000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j = 0$, so we check
the string ATATATA against the text and
mark its occurrence.

5. AGATACGA TATAT AC

$last \leftarrow 5$.

$D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[T] 010101010100001 \\ \hline D = 010101010100001 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

$$\begin{array}{r} 101000101000010 \\ B[A] 101010101010010 \\ \hline D = 101000101000010 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 010001010000100 \\ B[T] 010101010100001 \\ \hline D = 010001010000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 2$

$$\begin{array}{r} 1000001000000000 \\ B[A] 101010101010010 \\ \hline D = 1000001000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 1$

$$\begin{array}{r} 0000010000000000 \\ B[T] 010101010100001 \\ \hline D = 0000010000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j = 0$, so we check
the string TATAT against the text and
mark its occurrence.

6. AGATACGAT ATATA C

$last \leftarrow 5$.

$D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[A] 101010101010010 \\ \hline D = 101010101010010 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 4$

$$\begin{array}{r} 010101010000100 \\ B[T] 010101010100001 \\ \hline D = 010101010000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 1010001000000000 \\ B[A] 101010101010010 \\ \hline D = 1010001000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 2$

$$\begin{array}{r} 0100010000000000 \\ B[T] 010101010100001 \\ \hline D = 0100010000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 1$

$$\begin{array}{r} 1000000000000000 \\ B[A] 101010101010010 \\ \hline D = 1000000000000000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j = 0$, so we check
the string ATATATA against the text,
but we fail to recognize an occurrence.

7. AGATACGATA TATAC

$last \leftarrow 5$.

$D \leftarrow 1111111111111111$

$$\begin{array}{r} 1111111111111111 \\ B[C] 000000000001000 \\ \hline D = 000000000001000 \end{array}$$

$$\begin{array}{r} 000000000010000 \\ B[A] 101010101010010 \\ \hline D = 000000000010000 \end{array}$$

$D \& DF \neq 0^{|P|}$ and $j > 0$, then
 $last \leftarrow 3$

$$\begin{array}{r} 0000000000000000 \\ B[T] 010101010100001 \\ \hline D = 0000000000000000 \end{array}$$

3.4.2 Set Backward Dawg Matching idea

The **SBDM** algorithm uses a suffix automaton to recognize backwards the factors in a window of size ℓ_{min} that is shifted along the text.

3.4.2.1 Suffix automaton for a set of strings

The suffix automaton for a set of strings [BBE⁺87] is an automaton that recognizes the suffixes of the set of strings P it is built on. Let γ be the number of states of the trie built on P ($\gamma \leq |P| + 1$). Then the number of states of the suffix automaton is at least γ and at most 2γ . It is also $O(\gamma)$ in its number of transitions.

The construction algorithm is an extension of the construction for a single string (Section 2.4.1), but this time the resulting automaton is not necessarily minimal. The construction is linear in the size of P , but it is complex and slow.

3.4.2.2 Search algorithm

The suffix automaton is built in $O(r \times \ell_{min})$ time on $P_{\ell_{min}}^{rv}$, the set of reverse prefixes of length ℓ_{min} of the strings in P [BBE⁺87]. The search is done through a window of size ℓ_{min} , which we slide along the text. In this window, we read backwards the longest suffix that is also a factor of one of the prefixes of length ℓ_{min} of the strings in P . Two cases may occur.

- (i) We fail to recognize a factor, that is, we reach a letter σ that does not correspond to a transition in the suffix automaton of $P_{\ell_{min}}^{rv}$. No other prefix of a string can overlap the part of the window read. We therefore shift the window so that its new starting position corresponds to the position next to σ .
- (ii) We reach the beginning of the window in a state q of the suffix automaton. This means that we recognized a prefix $L(q)$ of a string in $F(q)$ (Section 3.1). We then verify a possible occurrence by comparing each string in $F(q)$ against the text. We finally move the search window by 1 and start the search again.

The worst-case complexity of **SBDM** is $O(n \times |P|)$, which is very high. However, for reasonable numbers of strings on a not too small alphabet, this algorithm is sublinear on average. The practical limit of this algorithm is the construction of the suffix automaton. For large sets of strings, it is too slow to be amortized by the time saved on the search phase. Moreover, the memory the suffix automaton requires quickly becomes too large as the set increases. We do not describe this algorithm in depth, nor give a pseudo-code, because **SBOM** uses the same approach but overcomes the bottleneck

of the suffix automaton with a lighter and simpler data structure. **SBOM** is faster than **SDBM** in all cases.

Note that the algorithm can be improved using the variable *last*, as was done for **BDM** (Section 2.4.1).

3.4.3 Set Backward Oracle Matching algorithm

The **Set Backward Oracle Matching algorithm (SBOM)** [AR99] uses a factor oracle of the set of strings. The factor oracle of P recognizes at least all the factors of the strings in P . The search algorithm is similar to **SDBM**. We slide a window of size ℓ_{min} along the text, reading backward a suffix of the window in the factor oracle. If we fail on a letter σ , we can safely shift the window past σ . If not, we reach the beginning of the window and verify a subset of P against the text.

3.4.3.1 Factor oracle of a set of strings

The factor oracle construction on a set of strings resembles the Aho-Corasick automaton construction. The only difference appears when going down the supply path looking for an outgoing transition labeled by σ . In the Aho-Corasick automaton construction, if this transition does not exist, we just jump to the next state on the supply path (Section 3.2.2). In the factor oracle construction, we create in addition a transition labeled by σ from each state on the supply path to the state where the original transition leads.

More precisely, we begin by building the trie of the set of strings P with the algorithm given in Figure 3.6. The states of the factor oracle are those of the trie as well as the initial state I and the terminal states. Hence, the factor oracle has at most $|P| + 1$ states, including the initial one.

To build the “external transitions,” which are at most $|P|$, we associate to each state q a “supply state,” computed simultaneously with the new transitions in transversal order. The supply state of the initial state is set to θ .

To explain the construction, we assume that we have already computed the supply function of all the states before state *Current* in transversal order. We consider the parent *Parent* of *Current* in the trie, leading to *Current* by σ , that is, $Current = \delta_{OR}(Parent, \sigma)$. The supply state $S_{OR}(Parent)$ has already been computed, and we go down the supply function from state $S_{OR}(Parent)$. We use a variable *Down* initialized to $S_{OR}(Parent)$ and we repeat the following steps.

- ST_1 If $Down = \theta$, then $S_{OR}(Current) \leftarrow I$.
- ST_2 If $Down \neq \theta$ and there does not exist a transition from $Down$ labeled by σ , then build a transition from state $Down$ to state $Current$ by σ and return to step ST_1 with $Down \leftarrow S_{OR}(Down)$.
- ST_3 If $Down \neq \theta$ and there exists a transition from $Down$ labeled by σ leading to a state Im , then set $S_{OR}(Current) \leftarrow Im$ and stop processing state $Current$.

The resulting factor automaton recognizes at least all factors of P [AR99]. The construction algorithm is worst-case time $O(|P|)$. Its pseudo-code is given in Figure 3.20.

```

Build_Oracle_Multiple( $P = \{p^1, p^2, \dots, p^r\}$ )
1.   $OR\_trie \leftarrow \mathbf{Trie}(P)$ 
     $\delta_{OR}$  is its transition function
2.  Mark the states that correspond to an entire string  $p^i$  as terminal
3.   $I \leftarrow \text{root of } OR\_trie$ 
4.   $S_{OR}(I) \leftarrow \theta$ 
5.  For  $Current$  in transversal order Do
6.     $Parent \leftarrow \text{parent in } OR\_trie \text{ of } Current$ 
7.     $\sigma \leftarrow \text{label of the transition from } Parent \text{ to } Current$ 
8.     $Down \leftarrow S_{OR}(Parent)$ 
9.    While  $Down \neq \theta$  AND  $\delta_{OR}(Down, \sigma) = \theta$  Do
10.      $\delta_{OR}(Down, \sigma) \leftarrow Current$ 
11.      $Down \leftarrow S_{OR}(Down)$ 
12.    End of while
13.    If  $Down \neq \theta$  Then
14.      $S_{OR}(Current) \leftarrow \delta_{OR}(Down, \sigma)$ 
15.    Else  $S_{OR}(Current) \leftarrow I$ 
16.    End of if
17. End of for

```

Fig. 3.20. Construction of the factor oracle for a set $P = \{p^1, p^2, \dots, p^r\}$. The state $Current$ goes through the trie OR_trie built on P in transversal order. The state $Down$ goes down the supply links from the parent of $Current$ looking for an outgoing transition labeled with the same character as between $Current$ and its parent, creating it if it does not exist.

3.4.3.2 Search with the factor oracle

The factor oracle is built in $O(r \times \ell_{min})$ time on the reverse prefixes of length ℓ_{min} of the strings in P . The search is done through a window of size ℓ_{min} , which we slide along the text. In this window, we read backwards the longest suffix that labels a path from the initial state. Two cases may occur.

- (i) We fail to recognize a factor, that is, we reach a letter σ that does not correspond to a transition in the factor oracle of $P_{\ell min}^{rv}$. No other prefix of a string can overlap the part of the window read. We therefore shift the window so that its new starting position corresponds to the position next to σ .
- (ii) We reach the beginning of the window in a state q of the factor oracle. When using a suffix automaton in **SBDM**, we can be sure at this step that we recognized a prefix of one of the strings. However, the factor oracle accepts paths of size ℓmin ending in terminal states that do not correspond to any prefix. Hence, we have to verify first that we read the prefix $L(q)^{rv}$ and only if this is the case we verify a possible occurrence by comparing each string in $F(q)$ against the text. We finally move the search window by 1 and start the search again.

```

SBOM( $P = \{p^1, p^2, \dots, p^r\}, T = t_1 t_2 \dots t_n$ )
1.  Preprocessing
2.       $\ell min \leftarrow$  minimal length of strings in  $p^i \in P$ 
3.       $Or \leftarrow \mathbf{Build\_Oracle\_Multiple}(\{pref_{\ell min}(p^1)^{rv}, pref_{\ell min}(p^2)^{rv}$ 
                                      $\dots, pref_{\ell min}(p^r)^{rv}\})$ 
                                      $\delta_{Or}$  is its transition function
4.      For  $q$  state of  $Or$  Do  $F(q) \leftarrow \emptyset$ 
5.      For  $i \in 1 \dots r$  Do
6.           $F(q) \leftarrow F(q) \cup \{i\}$ , where  $q$  is the state reached by  $pref_{\ell min}(p^i)^{rv}$ 
7.      End of for
8.  Searching
9.       $pos \leftarrow 0$ 
10.     While  $pos \leq n - \ell min$  Do
11.          $Current \leftarrow$  initial state of  $Or$ 
12.          $j \leftarrow \ell min$ 
13.         While  $j \geq 1$  AND  $Current \neq \theta$  Do
14.              $Current \leftarrow \delta_{Or}(Current, t_{pos+j})$ 
15.              $j \leftarrow j - 1$ 
16.         End of while
17.         If  $Current \neq \theta$  AND  $j = 0$  AND  $T_{pos+1 \dots pos+\ell min} = L(Current)^{rv}$  Then
18.             Verify all the patterns in  $F(Current)$  one by one against the text
19.              $j \leftarrow 1$ 
20.         End of if
21.          $pos \leftarrow pos + j$ 
22.     End of while

```

Fig. 3.21. Pseudo-code for the **SBOM** algorithm. The notation $pref_{\ell min}(p^i)$ denotes the prefix of size ℓmin of the string p^i .

Pseudo-code for **SBOM** is given Figure 3.21. Its worst-case complexity is $O(n \times |P|)$, the same as **SBDM**. However, this algorithm is sublinear

on average. The construction of the factor oracle is fast and requires little memory, which permits using this algorithm to search large sets of strings on relatively small texts.

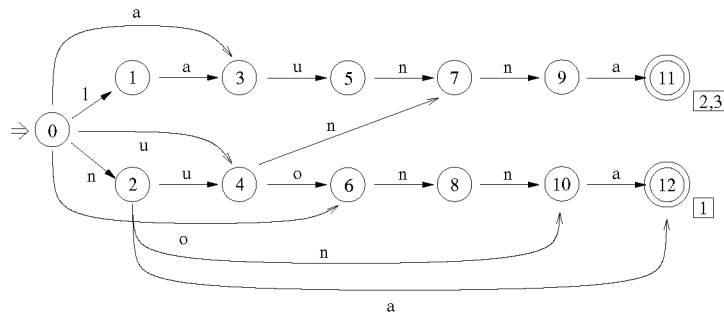


Fig. 3.22. Factor oracle for the reverse set of $P_{\ell min} = \{\text{announ}, \text{annual}\}$. Double-circled states are terminal.

Example using English We search the text “CPM_{annual} conference announce” for the set of strings $P = \{\text{announce}, \text{annual}, \text{annually}\}$. The factor oracle of the reverse set of $P_{\ell min} = \{\text{announ}, \text{annual}\}$ is shown in Figure 3.22.

1. CPM_{an} nual_conference_announce

We read n, a in the factor oracle. We fail on the next . We then shift the window after “.”.
2. CPM_ annual _conference_announce

We read l, a, u, n, n, a in the factor oracle. We reach the beginning of the window in state 11. We compare the strings $F(11) \rightarrow \text{annual}, \text{annually}$. We mark an occurrence of “annual”. We then shift the window by 1.
3. CPM_a nnual _conference_announce

We fail reading “.” in the oracle. We shift the window after “.”.
4. CPM_annual_ confer ence_announce

We fail reading r in the oracle. We shift the window after r.
5. CPM_annual_confer ence_a nnounce

We read a in the oracle, but we fail on the next “.”. We shift the window after “.”.
6. CPM_annual_conference_ announ ce

We read n, u, o, n, n, a in the factor oracle. We reach the beginning of the window in state 12. We compare the strings $F(12) \rightarrow \text{announce}$. We mark an occurrence of “announce”. We then shift the window by 1.
7. CPM_annual_conferencea_ nnounc e

We fail reading c in the oracle. We shift the window after c. Then $pos > n - \ell min$ and the search stops.

Example using DNA We search for the set of strings $P = \{\text{ATATATA}, \text{TATAT}, \text{ACGATAT}\}$ in the text AGATACGATATATAC . The factor oracle for the reverse set of $P^{\ell_{min}} = \{\text{ATATA}, \text{TATAT}, \text{ACGAT}\}$ is shown in Figure 3.23.

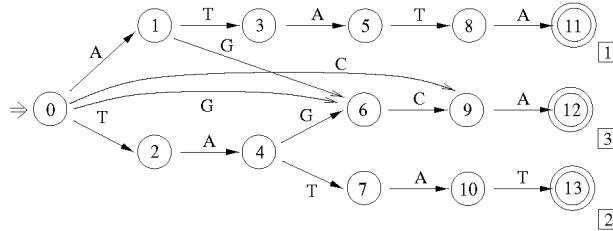


Fig. 3.23. Factor oracle for the reverse set of $P^{\ell_{min}} = \{\text{ATATA}, \text{TATAT}, \text{ACGAT}\}$. Double-circled states are terminal.

1. AGATA CGATATATAC

We read A, T, A in the factor oracle, and we fail on the next G. We then shift the window after G.

2. AG ATACG ATATATAC

We read G, C, A in the factor oracle, and we fail on the next T. We then shift the window after T.

3. AGAT ACGAT ATATAC

We read T, A, G, C, A in the factor oracle. We reach the beginning of the window in state 12. We compare the strings $F(12) \rightarrow \text{ACGATAT}$. We mark an occurrence and shift the window by 1.

4. AGATA CGATA TATAC

We read A, T, A in the factor oracle, and we fail on the next G. We then shift the window after G.

5. AGATACG ATATA TAC

We read A, T, A, T, A in the factor oracle. We reach the beginning of the window in state 11. We compare the string $F(11) \rightarrow \text{ATATATA}$. We mark an occurrence and shift the window by 1.

6. AGATACGA TATAT AC

We read T, A, T, A, T in the factor oracle. We reach the beginning of the window in state 13. We compare the string $F(13) \rightarrow \text{TATAT}$. We mark an occurrence and shift the window by 1.

7. AGATACGAT ATATA C

We read A, T, A, T, A in the factor oracle. We reach the beginning of the window in state 11. We compare the string $F(11) \rightarrow \text{ATATATA}$ and fail. We shift the window by 1.

8. AGATACGATA TATAC

We read C, A in the factor oracle and fail on the next T. We shift the window and the search stops since $pos > n - \ell_{min}$.

3.5 Experimental maps

We present in this section some maps of efficiency for the different multiple string matching algorithms, showing for all of them the zone in which they are most efficient in practice. The text of 10 megabytes is randomly built, as are the patterns. The experiments were performed on a $w = 32$ bits Ultra Sparc 1 running SunOs 5.6. The sets contain 5, 10, 100, and 1000 strings of the same length, varying from 5 to 100 in steps of 5. We tested all the algorithms presented. The **Wu-Manber** algorithm used in these experiments is the implementation found in *Agrep*. Its performance may vary, depending on the hash functions and the sizes of the tables used.

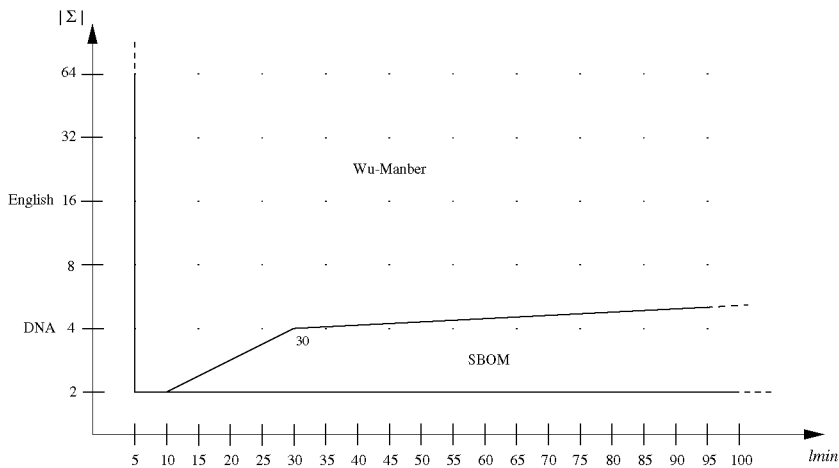


Fig. 3.24. Map of the most efficient algorithms when searching for 5 strings.

The maps are shown in Figures 3.24 to 3.27. The most efficient algorithms are just **Wu-Manber**, the advanced **Aho-Corasick**, and **SBOM**. As the set grows in size, **SBOM** becomes more and more attractive. The advanced **Aho-Corasick** also improves in comparison with the others for short strings, since it reads the text only once.

3.6 Other algorithms and references

Dynamic multiple string matching The algorithms presented in this chapter preprocess a fixed set of strings (a dictionary) in order to perform the search. However, if we need to modify the dictionary by adding or removing a string, we need to preprocess the new dictionary from scratch. The problem of searching for a set of strings in a text and allowing efficient modifications of the set is called *dynamic string matching*. It has been solved

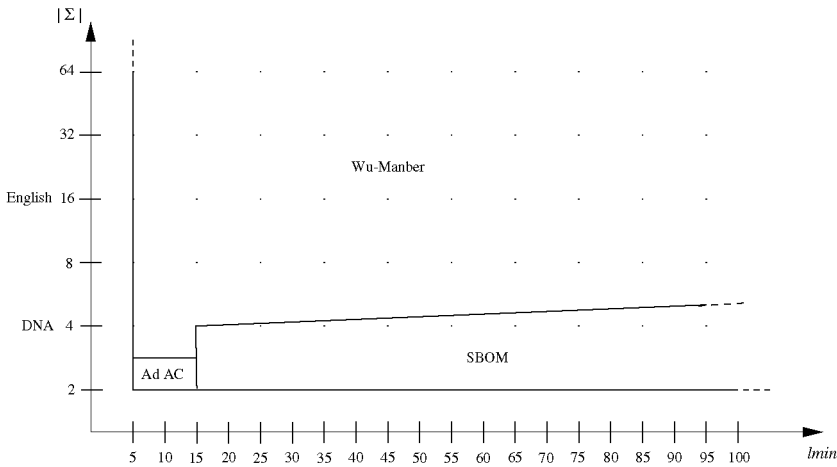


Fig. 3.25. Map of the most efficient algorithms when searching for 10 strings.

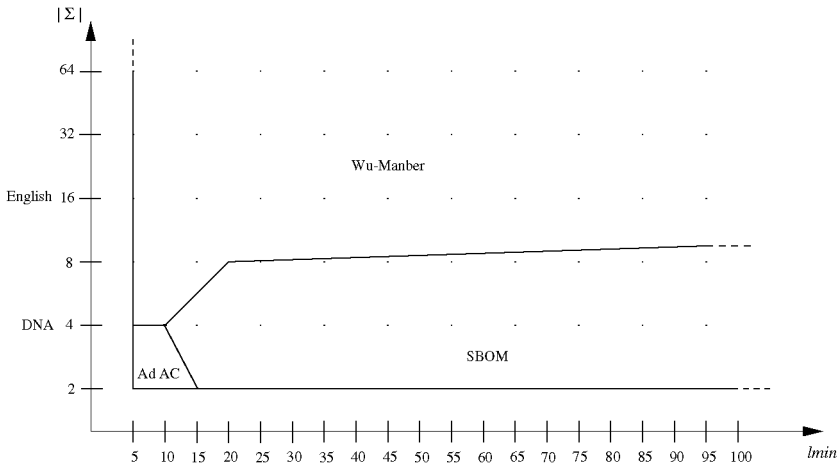


Fig. 3.26. Map of the most efficient algorithms when searching for 100 strings.

recently [SV96] with optimal worst-case complexities: (i) preprocessing of the set of strings in $O(|P|)$ time; (ii) adding or removing a string p in $O(|p|)$ time; and (iii) finding all occurrences of P in the text in $O(n + \text{nocc})$ time, where nocc is the number of occurrences of P in the text.

An application of dynamic string matching is in the matching of a set of strings with variable length “don’t cares” [KR95].

On the Commentz-Walter algorithm Several variations of the **Commentz-Walter** algorithm have been designed to limit its worst-case com-

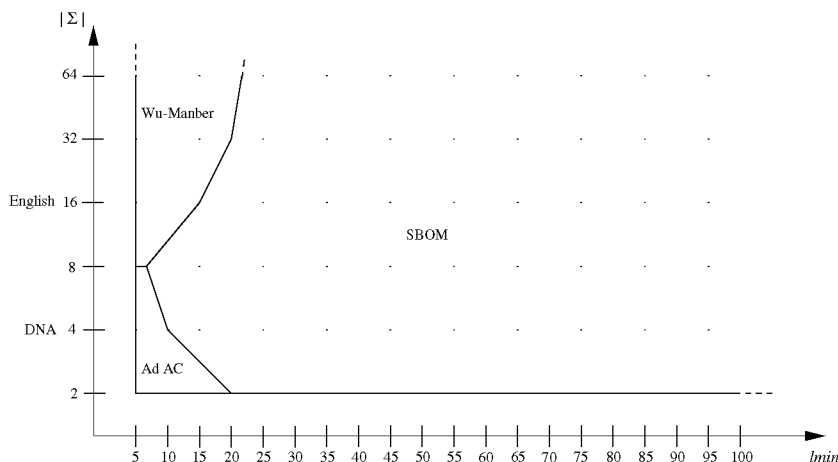


Fig. 3.27. Map of the most efficient algorithms when searching for 1000 strings.

plexity [Sri86] by using additional memory. These algorithms are, however, not efficient in practice.

On matching a set of strings on unbounded alphabets The problem of matching a set of strings of the same length m on an unbounded alphabet has been investigated [Bre95]. The resulting algorithm runs in $O((\log(|P|)/m + 1) \times n)$ comparisons after an $O(|P| \times m \times \log |A|)$ preprocessing time, where A is the alphabet on which the set P is built.