

UNIVERSITY OF JOENSUU

DEPARTMENT OF COMPUTER SCIENCE

Report Series A

# **Efficient parameterized string matching**

Kimmo Fredriksson and Maxim Mozgovoy

A-2006-2

ACM F.2.2, H.3.3

ISBN 952-458-807-2

ISSN 0789-7316

# Sublinear parameterized single and multiple string matching

Kimmo Fredriksson\* and Maxim Mozgovoy  
Department of Computer Science  
University of Joensuu  
{kfredrik,mmozgo}@cs.joensuu.fi

## Abstract

We consider the following pattern matching problem. We have *pattern*  $P[0 \dots m-1]$  and *text* is  $T[0 \dots n-1]$ , where the symbols of  $P$  and  $T$  are taken from two disjoint finite alphabets  $\Sigma$  of size  $\sigma$  and  $\Lambda$  of size  $\lambda$ . The pattern  $P$  matches the text substring  $T[j \dots j+m-1]$ , iff for all  $i \in \{0 \dots m-1\}$  it holds that  $M_j(P[i]) = T[j+i]$ , where  $M_j(\cdot)$  is one-to-one mapping on  $\Sigma \cup \Lambda$  such that the mapping is identity on  $\Sigma$ , but on  $\Lambda$  can be different for each text position  $j$ . We give efficient algorithms that find all parameterized occurrences of  $P$  in  $T$ . The algorithms are based on generalizing Shift-Or and Backward DAWG Matching (BDM) algorithms. The latter can be used for searching  $r$  patterns simultaneously. The Shift-Or based algorithm runs in  $O(n \lceil m/w \rceil)$  worst case time, while the average case for fixed alphabets and under some mild and realistic assumptions is  $O(n \log_\sigma(m)/w)$ , where  $w$  is the number of bits in computer word. The BDM based algorithm runs in  $O(n \log_\sigma(rm)/m)$  average time. This is optimal within a constant factor. For general alphabets the times increase by a factor  $O(\log(m))$ .

**Keywords:** algorithms, parameterized string matching, bit-parallelism, suffix automaton

**ACM Classification:** F.2.2 [Analysis of algorithms and problem complexity]: Non-numerical algorithms and problems — *Pattern matching, Computations on discrete structures*; H.3.3 [Information storage and retrieval]: Information Search and Retrieval — *Search process*.

## 1 Introduction

In traditional string matching problem one is interested in finding the occurrences of a pattern  $P$  from a text  $T$ , where  $P$  and  $T$  are strings over some alphabet  $\Sigma$ . Many variations of this basic problem setting exist, such as searching multiple patterns simultaneously, and/or allowing some limited number of errors in the matches, and indexed searching, where  $T$  can be preprocessed to allow efficient queries of  $P$ . See e.g. [9, 11, 6] for an overview and references. Yet another variation is *parameterized matching* [4]. In this variant we have two disjoint alphabets,  $\Sigma$  for *fixed* symbols, and  $\Lambda$  for *parameter* symbols. In this setting we search *parameterized* occurrences of  $P$ , where the symbols from  $\Sigma$  must match exactly, while the symbols in  $\Lambda$  can be also renamed. This problem has important applications e.g. in software maintenance and plagiarism detection [4], where the symbols of the strings can be e.g. reserved words and identifier or parameter names of some (possibly tokenized) programming language source code. Hence one might be interested in finding code snippets that are the same up to some systematical variable renaming.

A myriad of algorithms have been developed for the classical problem, but only a few exist for parameterized matching. Linear exact on-line matching algorithms have been developed for single [3, 1] and multiple patterns [10]. Sublinear algorithm on average for single pattern was developed in [3]. Other algorithms exist for the off-line problem [4, 5]. In this paper we develop algorithms that run in sublinear time on average, are simple to implement and

---

\*Supported by the Academy of Finland, grant 202281.

perform well in practice. Our algorithms are based on generalizing the well known Shift-Or [2] and Backward DAWG Matching algorithms [7]. Our algorithms generalize for the multipattern matching as well.

## 2 Preliminaries

We use the following notation. The *pattern* is  $P[0 \dots m - 1]$  and the *text* is  $T[0 \dots n - 1]$ . The symbols of  $P$  and  $T$  are taken from two disjoint finite alphabets  $\Sigma$  of size  $\sigma$  and  $\Lambda$  of size  $\lambda$ . The pattern  $P$  matches the text substring  $T[j \dots j + m - 1]$ , iff for all  $i \in \{0 \dots m - 1\}$  it holds that  $M_j(P[i]) = T[j + i]$ , where  $M_j(\cdot)$  is one-to-one mapping on  $\Sigma \cup \Lambda$ . Moreover, the mapping must be identity on  $\Sigma$ , but on  $\Lambda$  can be different for each text position  $j$ . For example, assume that  $\Sigma = \{A, B\}$ ,  $\Lambda = \{X, Y, Z\}$  and  $P = \text{AAZYABXYZAX}$ . Then  $P$  matches the text substring  $\text{AAZYABXYZAX}$  with identity mapping, and  $\text{AAXYXABZYXAZ}$  with parameter mapping  $X \mapsto Z$ ,  $Y \mapsto Y$ , and  $Z \mapsto X$ . This mapping is simple with *prev* encoding [4]. For a string  $S$ ,  $\text{prev}(S)$  maps all parameter symbols  $s$  in  $S$  to a non-negative integer  $p$ , where  $p$  is the number of symbols since the last occurrence of symbol  $s$  in  $S$ . The first occurrence of the parameter is encoded as 0. If  $s$  belongs to  $\Sigma$ , it is mapped to itself ( $s$ ). For our example pattern,  $\text{prev}(P) = \text{AA002AB055A4}$ . This is the same as the encoding for the two example substrings, i.e.  $\text{prev}(\text{AAZYABXYZAX}) = \text{prev}(\text{AAXYXABZYXAZ})$ . Hence the problem is reduced to exact string matching, where we match  $\text{prev}(P)$  against  $\text{prev}(T[j \dots j + m - 1])$  for all  $j = 0 \dots n - m$ . The string  $\text{prev}(S)$  can be easily computed in linear time for constant size alphabets. The only remaining problem then is how to maintain  $\text{prev}(T[j \dots j + m - 1])$  (and any algorithmic parameters that depend on it) efficiently as  $j$  increases. The key is the following Lemma [4].

**Lemma 1** . *Let  $S' = \text{prev}(S)$ . Then for  $S'' = \text{prev}(S[j \dots j + m - 1])$  for all  $i$  such that  $S[i] \in \Lambda$  it holds that  $S''[i] = S'[i]$  iff  $S'[i] < m$ . Otherwise  $S''[i] = 0$ .*

We are now ready to present our algorithms. For simplicity we assume that  $\Sigma$  and  $\Lambda$  are finite constant size alphabets. For large alphabets all our time bounds hold if we multiply them by  $O(\log(m))$ . Moreover, without loss of generality, we assume the symbols to be atomic, i.e. they can be accessed in constant time.

## 3 Parameterized bit-parallel matching

In this section we present bit-parallel approach for parameterized matching, based in Shift-Or algorithm [2]. For the bit-parallel operations we adopt the following notation. A machine word has  $w$  bits, numbered from the least significant bit to the most significant bit. We use C-like notation for the bit-wise operations of words;  $\&$  is bit-wise **and**,  $|$  is **or**,  $\wedge$  is **xor**,  $\sim$  negates all bits,  $\ll$  is shift to left, and  $\gg$  shift to right, both with zero padding. For brevity, we make the assumption that  $m \leq w$ , unless explicitly stated otherwise. We first review the standard Shift-Or algorithm, and then show how it can be adapted to parameterized matching, and finally improve its average case running time.

### 3.1 Standard Shift-Or

The standard Shift-Or automaton is constructed as follows. The automaton has states  $0, 1, \dots, m$ . The state 0 is the initial state, state  $m$  is the final (accepting) state, and for  $i = 0, \dots, m - 1$  there is a transition from the state  $i$  to the state  $i + 1$  for character  $P[i]$ . In addition, there is a transition for every  $c \in \Sigma$  from and to the initial state, which makes the automaton non-deterministic. The preprocessing algorithm builds a table  $B$ , having one bit-mask entry for each  $c \in \Sigma$ . For  $0 \leq i \leq m - 1$ , the mask  $B[c]$  has  $i$ th bit set to 0, iff  $P[i] = c$ . These correspond to the transitions of the implicit automaton. That is, if the bit  $i$  in  $B[c]$  is 0, then there is a

---

**Alg. 1** Shift-Or( $T, n, P, m$ ).

---

```
1  for  $i \leftarrow 0$  to  $\sigma - 1$  do  $B[i] \leftarrow \sim 0 \gg (w - m)$ 
2  for  $i \leftarrow 0$  to  $m - 1$  do  $B[P[i]] \leftarrow B[P[i]] \& \sim(1 \ll i)$ 
3   $D \leftarrow \sim 0$ ;  $mm \leftarrow 1 \ll (m - 1)$ 
4  for  $i \leftarrow 0$  to  $n - 1$  do
5       $D \leftarrow (D \ll 1) | B[T[i]]$ 
6      if  $(D \& mm) \neq mm$  then report match
```

---

---

**Alg. 2** Encode( $P, m$ ).

---

```
1  for  $i \leftarrow 0$  to  $m - 1$  do if  $P[i] \in \Lambda$  then  $prev[P[i]] \leftarrow -1$ 
2  for  $i \leftarrow 0$  to  $m - 1$  do
3      if  $P[i] \in \Lambda$  then
4          if  $prev[P[i]] = -1$  then  $P'[i] \leftarrow \sigma$  else  $P'[i] \leftarrow i - prev[P[i]] + \sigma$ 
5           $prev[P[i]] \leftarrow i$ 
6      else
7           $P'[i] \leftarrow P[i]$ 
8  return  $P'$ 
```

---

transition from the state  $i$  to the state  $i + 1$  with character  $c$ . The bit-vector  $D$  encodes the states of the automaton. The  $i$ th bit of the state vector is set to 0, iff the state  $i$  is active, i.e. the pattern prefix  $P[0 \dots i]$  matches the current text position. Initially each bit is set to 1. For each text symbol  $c$  the vector is updated by  $D \leftarrow (D \ll 1) | B[c]$ . This simulates all the possible transitions of the nondeterministic automaton in a single step. If after the update the  $m$ th bit of  $d$  is zero, then there is an occurrence of  $P$ . Alg. 1 gives the code. If  $m \leq w$ , then the algorithm runs in time  $O(n)$ .

### 3.2 Parameterized Shift-Or

In order to generalize Shift-Or for parameterized matching, we must take care of three things: (i)  $P$  must be encoded with  $prev$ ; (ii)  $prev(T[j \dots j + m - 1])$  must be maintained in  $O(1)$  time per text position; (iii) the table  $B$  must be built so that all parameterized pattern prefixes can be searched in parallel. The items (i) and (ii) are trivial, while (iii) is a bit more tricky. To compute  $prev(P)$  we just maintain an array  $prev[c]$  that for each symbol  $c \in \Lambda$  stores the position of its last occurrence. Then  $prev(P)$  can be computed in  $O(m)$  time by a linear scan over  $P$ . To simplify indexing in the array  $B$ , we assume that  $\Sigma = \{0 \dots \sigma - 1\}$ , and map the  $prev$  encoded parameter offsets into the range  $\{\sigma \dots \sigma + m - 1\}$ . Alg. 2 gives the pseudo code. The text is encoded in the same way, but the encoding is embedded into the search code, see Alg. 3. The only difference is that we apply Lemma 1 to reset offsets that are greater than  $m - 1$  (i.e. offsets that are for parameters that are outside of the current text window) to zero. Otherwise the search algorithms is exactly the same as for normal Shift-Or.

The tricky part is the preprocessing phase. We denote the  $prev$  encoded pattern as  $P'$ . At first  $P'$  is preprocessed just as  $P$  in the normal Shift-Or algorithm. This includes the parameter offsets, which are handled as any other symbol. However, this is not enough. We illustrate the problem by an example. Let  $P = \text{xaxax}$  and  $T = \text{zzazazaz}$ . In encoded forms these are  $P' = 0a2a2$  and  $T' = 01a2a2a2$ . Clearly  $P$  has two (overlapping) parameterized matches in  $T$ . However,  $P'$  does not match in  $T'$  at all.

The problem is that as the algorithm searches all the  $m$  prefixes of the pattern in parallel, then some non-zero encoded offset  $p$  (of some text symbol) should be interpreted as zero in some cases. These prefixes have lengths from 1 to  $m$ , and in able to successfully apply Lemma 1 we should be able to apply it in parallel to all  $m$  substrings. In other words, any non-zero parameter offset  $p$  must be treated as zero for all pattern prefixes whose length  $h$  is less than  $p$ , since by Lemma 1 the parameter with offset  $p$  is dropped out of the window of length  $h$ .

This problem can be solved as follows. The bit-vector  $B[\sigma + i]$  is the match vector for offset  $i$ . If the  $j$  bit of this vector is zero, it means by definition that  $P'[j] = i$ . If any of the  $i$  least significant bits of  $B[\sigma]$  are zero, we clear the corresponding bits of  $B[\sigma + i]$  as well. More

---

**Alg. 3** P-Shift-Or( $T, n, P, m$ ).

---

```
1   $P' \leftarrow \text{Encode}(P, m)$ 
2  for  $i \leftarrow 0$  to  $\sigma + m - 1$  do  $B[i] \leftarrow \sim 0 \gg (w - m)$ 
3  for  $i \leftarrow 0$  to  $\lambda - 1$  do  $prv[\sigma + i] \leftarrow -\infty$ 
4  for  $i \leftarrow 0$  to  $m - 1$  do  $B[P'[i]] \leftarrow B[P[i]] \& \sim(1 \ll i)$ 
5  for  $i \leftarrow 1$  to  $m - 1$  do  $B[\sigma + i] \leftarrow B[\sigma + i] \& (B[\sigma] | (\sim 0 \ll i))$ 
6   $D \leftarrow \sim 0$ ;  $mm \leftarrow 1 \ll (m - 1)$ 
7  for  $i \leftarrow 0$  to  $n - 1$  do
8       $c \leftarrow T[i]$ 
9      if  $c \in \Lambda$  then
10          $c \leftarrow i - prv[T[i]] + \sigma$ 
11         if  $c > \sigma + m - 1$  then  $c \leftarrow \sigma$ 
12          $prv[T[i]] \leftarrow i$ 
13          $D \leftarrow (D \ll 1) | B[c]$ 
14         if  $(D \& mm) \neq mm$  then report match
```

---

precisely, we set

$$B[\sigma + i] \leftarrow B[\sigma + i] \& (B[\sigma] | (\sim 0 \ll i)).$$

This means that the offset  $i$  is treated as offset  $i$  for prefixes whose length is greater than  $i$ , and as zero for the shorter prefixes, satisfying the condition of Lemma 1.

Alg. 3 gives the complete code. The algorithm clearly runs in  $O(n \lceil m/w \rceil)$  worst case time. For long patterns one can search just a length  $w$  prefix of the pattern, and verify with the whole pattern whenever the prefix matches, giving  $O(n)$  average time. However, note that a long variable name (string) is just one symbol (token) in typical applications, hence  $w$  bits is usually plenty. Finally, note that for unbounded alphabets we cannot use arrays for  $prv$  and  $B$ . We can use balanced trees instead, but then the time bounds must be multiplied by  $O(\log(m))$ .

### 3.3 Skipping text symbols

Standard Shift-Or can be improved to run in optimal  $O(n \log_{\sigma}(m)/m)$  average time [8]. The algorithm takes a parameter  $q$ , and from the original pattern generates a set  $\mathcal{P}$  of  $q$  new patterns  $\mathcal{P} = \{P^0, \dots, P^{q-1}\}$ , each of length  $m' = \lfloor m/q \rfloor$ , where  $P^j[i] = P[j + iq]$  for  $i = 0 \dots \lfloor m/q \rfloor - 1$ . In other words, the algorithm generates  $q$  different alignments of the original pattern  $P$ , each alignment containing only every  $q$ th character. The total length of the patterns in  $\mathcal{P}$  is  $q \lfloor m/q \rfloor \leq m$ . For example, if  $P = \text{ABCDEF}$  and  $q = 3$ , then  $P^0 = \text{AD}$ ,  $P^1 = \text{BE}$  and  $P^2 = \text{CF}$ . Assume now that  $P$  occurs at  $T[i..i + m - 1]$ . From the definition of  $P^j$  it directly follows that  $P^j[h] = T[i + j + hq]$ , where  $j = i \bmod q$  and  $h = 0 \dots m' - 1$ . This means that we can use the set  $\mathcal{P}$  as a filter for the pattern  $P$ , and that the filter needs only to scan every  $q$ th character of  $T$ . All the patterns must be searched simultaneously. Whenever an occurrence of  $P^j$  is found in the text, we must verify if  $P$  also occurs, with the corresponding alignment.

This method clearly works for parameterized matching as well. We generate the set of patterns  $\mathcal{P}$ , and also *prev*-encode them. For example for  $P = \text{AAZYABXYZAX}$  and  $q = 3$  we process the pattern set  $prev(\{\text{AYBZ}, \text{AZXA}, \text{ZAYX}\}) = \{\text{A0B0}, \text{A00A}, \text{0A00}\}$ . In the search phase the text is also encoded on-line, encoding only every  $q$ th symbol, but assuming that they are consecutive. In other words, every parameter offset is effectively divided by  $q$  to agree with the encoding of the patterns. Finally, the verification phase checks if  $prev(P) = prev(T[v \dots v + m - 1])$ , where  $v$  is the starting position of a potential match.

The search of the pattern set can be done using the parameterized Shift-Or algorithm. This is possible by concatenating and packing the set of patterns into a single machine word [8, 2]. Another alternative is to use the parameterized version [10] of Aho-Corasic algorithm. Both lead to the same average case running time, but the latter does not require that  $m \leq w$ , as it is not based on bit-parallelism. We denote the Shift-Or based algorithm as PFSO. Alg. 4 shows the pseudo code for the filtering phase, and Alg. 5 the verification code.

The filtering time is  $O(n/q)$ . The filter searches the exact matches of  $q$  patterns, each of length  $\lfloor m/q \rfloor$ . We are not able to analyze the exact effect of the parameter alphabet to the

---

**Alg. 4** Average-Optimal-P-Shift-Or( $T, n, P, m, q$ ).

---

```
1   $P' \leftarrow \text{Encode}(P, m)$ 
2   $h \leftarrow 0; mm \leftarrow 0$ 
3  for  $j \leftarrow 0$  to  $q - 1$  do
4    for  $i \leftarrow 0$  to  $\lfloor m/q \rfloor - 1$  do
5       $\mathcal{P}[j][i] \leftarrow P[iq + j]$ 
6       $h \leftarrow h + \lfloor m/q \rfloor$ 
7       $mm \leftarrow mm \mid (1 \ll (h - 1))$ 
8       $\mathcal{P}[j] \leftarrow \text{Encode}(\mathcal{P}[j], \lfloor m/q \rfloor)$ 
9  for  $i \leftarrow 0$  to  $\sigma + \lfloor m/q \rfloor q - 1$  do  $B[i] \leftarrow \sim 0 \gg (w - \lfloor m/q \rfloor q)$ 
10 for  $j \leftarrow 0$  to  $q - 1$  do
11   for  $i \leftarrow 0$  to  $\lfloor m/q \rfloor - 1$  do
12      $B[\mathcal{P}[j][i]] \leftarrow B[\mathcal{P}[j][i]] \& \sim(1 \ll (j \lfloor m/q \rfloor + i))$ 
13 for  $j \leftarrow 0$  to  $q - 1$  do
14    $msk \leftarrow \sim 0$ 
15   for  $i \leftarrow 1$  to  $\lfloor m/q \rfloor - 1$  do
16      $msk \leftarrow msk \wedge (1 \ll (j \lfloor m/q \rfloor + i - 1))$ 
17      $B[\sigma + i] \leftarrow B[\sigma + i] \& (B[\sigma] \mid msk)$ 
18 for  $i \leftarrow 0$  to  $\lambda$  do  $prev[\sigma + i] \leftarrow -\infty$ 
19  $D \leftarrow \sim 0; i \leftarrow 0$ 
20 while  $i < n$  do
21    $c \leftarrow T[i]$ 
22   if  $c \in \Lambda$  then
23      $c \leftarrow i/q - prev[T[i]] + \sigma$ 
24     if  $c > \sigma + \lfloor m/q \rfloor - 1$  then  $c \leftarrow \sigma$ 
25      $prev[T[i]] \leftarrow i/q$ 
26      $D \leftarrow ((D \& \sim mm) \ll 1) \mid B[c]$ 
27     if  $(D \& mm) \neq mm$  then  $\text{Verify}(T, i, n, P', m, q, D, mm)$ 
28    $i \leftarrow i + q$ 
```

---

---

**Alg. 5** Verify( $T, i, n, P, m, q, D, mm$ ).

---

```
1   $D \leftarrow (D \& mm) \wedge mm$ 
2  while  $D \neq 0$  do
3     $s \leftarrow \lfloor \log_2(D) \rfloor$ 
4     $c \leftarrow -(\lfloor m/q \rfloor - 1)q - \lfloor s/\lfloor m/q \rfloor \rfloor$ 
5    if  $P = \text{Encode}(T[i + c \dots i + c + m - 1], m)$  then report match
6     $D \leftarrow D \& \sim(1 \ll s)$ 
```

---

probability that two randomly picked symbols match. However, if we assume that a constant fraction  $\varepsilon$  of the pattern positions are randomly selected to have a randomly selected symbol from  $\Sigma$ , then the probability that  $P^j$  occurs in a given text position is  $O((1/\sigma)^{\lfloor \varepsilon m/q \rfloor})$ . A brute force verification cost is in the worst case  $O(m)$  (but only  $O(1)$  on average). To keep the total time at most  $O(n/q)$  on average, we select  $q$  so that  $n/q = mn/\sigma^{\varepsilon m/q}$ , i.e.  $q = O(m/\log_\sigma(m))$ . The total average time is therefore  $O(n \log_\sigma(m)/m)$ . This is optimal [12] within a constant factor.

## 4 Parameterized backward trie matching

We now present an algorithm based on Backward DAWG Matching (BDM) [7]. BDM is optimal on average, i.e. it runs in  $O(n \log_\sigma(m)/m)$  average time. We call our parameterized version of BDM as Parameterized Backward Trie Matching, PBTM for short. In the preprocessing phase PBTM builds a trie for the encoded suffixes of the reversed pattern. Trie is a rooted tree, where each edge is labeled by a symbol. The edges of the path from the root node to some leaf node then spells out the string of symbols stored into that leaf. The pattern in reverse is denoted by  $P^r$ . The set of its suffixes is  $\{P^r[i \dots m - 1] \mid 0 \leq i < m\}$  (note that this corresponds to the prefixes of the original pattern). Each suffix is then encoded with *prev*, and the encoded strings are inserted into a trie. For example, if  $P = \text{AZBZXBXY}$ , then the set of stored strings is  $\{00b20b2a, 0b20b2a, b00b2a, 00b2a, 0b2a, b0a, 0a, a\}$ . The trie allows efficient searching of any pattern substring that occurs in  $P^r$ . A brute force algorithm for this takes  $O(m^2)$  time, but can be improved to  $O(m)$  by using efficient suffix tree construction algorithms for parameterized

---

**Alg. 6** PBTM( $T, n, P, m$ ).

---

```
1  root ← EncSTrie( $P^r$ )
2  for  $i \leftarrow 0$  to  $\lambda - 1$  do  $prev[\sigma + i] \leftarrow -\infty$ 
3   $i \leftarrow 0$ 
4  while  $i < n - m$  do
5       $j \leftarrow m$ ;  $shift \leftarrow m$ ;  $u \leftarrow root$ 
6      while  $u \neq null$  do
7           $c \leftarrow T[i + j - 1]$ 
8          if  $c \in \Lambda$  then
9               $c \leftarrow m - j - prev[T[i + j - 1]] + \sigma$ 
10             if  $c > \sigma + m - 1$  then  $c \leftarrow \sigma$ 
11              $prev[T[i + j - 1]] \leftarrow m - j$ 
12              $j \leftarrow j - 1$ 
13              $u \leftarrow child(u, c)$ 
14             if  $u \neq null$  AND  $issuffix(u)$  then
15                 if  $j > 0$  then  $shift \leftarrow j$  else report match
16         for  $k \leftarrow i + j$  to  $i + m - 1$  do if  $T[k] \in \Lambda$  then  $prev[T[k]] \leftarrow -\infty$ 
17          $i \leftarrow i + shift$ 
```

---

strings [5]. An alternative to the trie is suffix array, i.e. the trie can be replaced with sorted array of  $prev$  encoded suffixes of the reverse pattern. Following an edge in the trie can then be simulated by a binary search in the array. We call the resulting algorithm PBAM. The benefit is that the array based method is easy to implement space efficiently since only one pointer is needed for each suffix.

We now show how this can be used for efficient search. Assume that we are scanning the text window  $T[i \dots i + m - 1]$  *backwards*. The invariant is that all occurrences that start before the position  $i$  are already reported. The text window is  $prev$ -encoded (backwards as well) as we go, and the read substring of this window is matched against the trie. This is continued as long as the substring can be extended without a mismatch, or we reach the beginning of the window. If the whole window can be matched against the trie, then the pattern occurs in that window. If the pattern does not match, some of the occurrences may still overlap with the current window. However, in this case one of the suffixes stored into the trie must match, since the reverse suffixes are also the prefixes of the original pattern. The algorithm remembers the longest such suffix, that is not the whole pattern, found from the window. The window is then shifted so that its starting position will become aligned with the last symbol of that suffix. This is the position of the next possible pattern occurrence. If the length of that longest suffix was  $\ell$ , the next window to be searched is  $T[i + m - \ell \dots i + m - 1 + m - \ell]$ . This process is repeated until the whole text is scanned.

Some care must be taken to be able to do the encoding of the text window in  $O(1)$  time per read symbol. To achieve constant time per symbol we must use an auxiliary array  $prev$  (as before) to store the position of the last occurrence for each symbol. We cannot afford to initialize the whole array for each window, so before shifting the window we rescan the symbols just read in the current window, and reinitialize the array only for those symbols. This ensures  $O(1)$  total time for each symbol read. Alg. 6 gives the code.

The average case running time of this algorithm depends on how many symbols  $x$  are examined in each window. Again, if we make the simplifying assumption that a constant fraction of the pattern positions are randomly selected to have a randomly selected symbol from  $\Sigma$ , then the original analysis of BDM holds for PBTM as well, and the average case running time is  $O(n \log_\sigma(m)/m)$ . For general alphabets and for the PBAM version the time must be multiplied by  $O(\log(m))$ . Finally, this algorithm can be easily modified to search  $r$  patterns simultaneously. Basically, if all the patterns are of the same length, this generalization requires just storing all the suffixes of all the patterns into the same trie. This results in  $O(n \log_\sigma(rm)/m)$  average time. With modest additional complexity patterns of different lengths can be handled as well in the same way as with regular BDM [6].

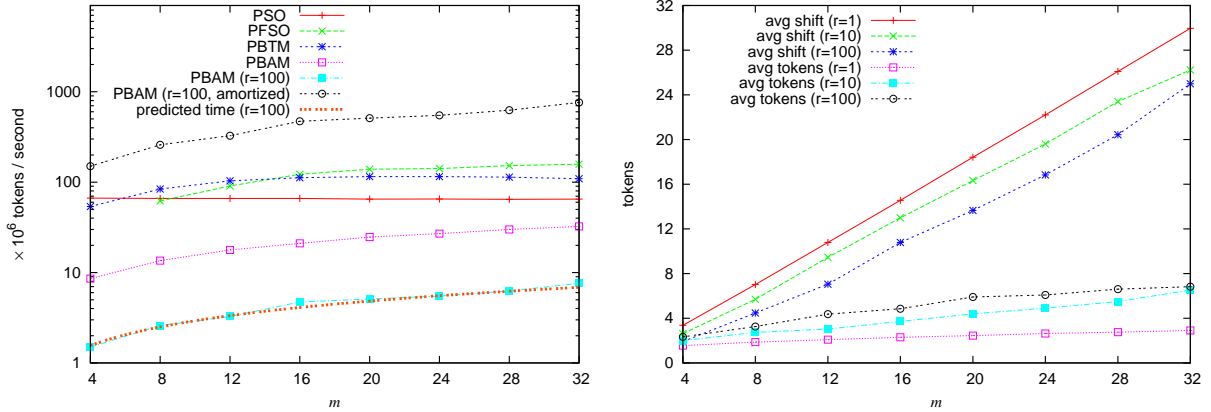


Figure 1: Left: the search speed in  $10^6$  tokens / second. Right: the average shift and average number of tokens inspected in each window of length  $m$ .

## 5 Experimental results

We have implemented the algorithms in C++, and compiled them with Borland C++Builder 6. We performed the experiments on the AMD Sempron 2600+ (1.88 GHz) machine with 768 MB RAM, running Windows XP. A tokenized string of concatenated Java source files (taken from various open source projects, such as jPOS, smppapi, and TM4J) was used as a text to be searched. The tokenization procedure (based on JavaCC<sup>1</sup> parser) converted an input file into a sequence of two-byte codes, representing single characters, reserved Java words and distinct identifiers. The initial string had a size of 5.48MB, and after encoding it consisted of 1259799 tokens, including 51 reserved Java words and 10213 unique identifiers. A set of 100 patterns for each length reported was randomly extracted from the input text. We report the average number of tokens searched per second for each algorithm.

Fig. 1 summarizes the results. PSO denotes the basic parameterized shift-or algorithm, PFSO the fast parameterized shift-or, PBTM the parameterized backward trie matching algorithm, and PBAM the suffix array version of PBTM. For short patterns plain PSO and PBTM give the best results. For longer patterns PFSO wins in case of optimal  $q$  selection. For  $m \in \{8, 12, 16, 20, 24, 28, 32\}$  we used  $q = \{2, 3, 4, 4, 4, 5, 6\}$ , respectively. For short patterns PBTM is faster than PFSO. For long patterns PBTM suffers from the large alphabet size. In our implementation we used arrays to implement the trie nodes and for long patterns the trie requires a lot of initialization time and memory, not fitting into the CPU cache. PBAM does not have this flaw, but the binary search step needed for each accessed text symbol makes it comparatively slow.

We also experimented with the multipattern version of PBAM, searching  $r = 100$  patterns simultaneously. The plot shows that while the raw speed is reduced, the amortized speed per pattern is clearly better than for any of the single pattern matching algorithms. The time also coincides nicely with the theoretical curve  $O(n \log_{\sigma}(rm) \log_2(rm)/m)$ , supporting our analysis. This is also clear given the right plot, showing the average number of tokens inspected in each text window, and the average shift for  $r = 1, 10, 100$ . These behave like in random texts supporting our assumptions in the analysis.

## 6 Conclusions

We have shown how two well-known algorithms, namely Shift-Or and BDM, can be generalized for parameterized matching. The algorithms are easy to implement, and work well in practice.

<sup>1</sup><https://javacc.dev.java.net/>



We have concentrated on obtaining fast average case times. However, the worst case times can be improved to  $O(n)$  using known results [7, 8] and standard tricks.

## References

- [1] A. Amir, M. Farach, and S. Muthukrishnan. Alphabet dependence in parameterized matching. *Inf. Process. Lett.*, 49(3):111–115, 1994.
- [2] R. A. Baeza-Yates and G. H. Gonnet. A new approach to text searching. *Commun. ACM*, 35(10):74–82, 1992.
- [3] B. S. Baker. Parameterized pattern matching by Boyer-Moore-type algorithms. In *Proceedings of ACM-SODA'95*, pages 541–550, 1995.
- [4] B. S. Baker. Parameterized duplication in strings: algorithms and an application to software maintenance. *SIAM J. Comput.*, 26(5):1343–1362, 1997.
- [5] R. Cole and R. Hariharan. Faster suffix tree construction with missing suffix links. In *Proceedings of ACM-STOC'00*, pages 407–415, Portland, Oregon, 2000.
- [6] M. Crochemore and W. Rytter. *Text algorithms*. Oxford University Press, 1994.
- [7] M. Crochemore *et al.* Speeding up two string matching algorithms. *Algorithmica*, 12(4/5):247–267, 1994.
- [8] K. Fredriksson and Sz. Grabowski. Practical and optimal string matching. In *Proceedings of SPIRE'2005*, LNCS 3772, pages 374–385. Springer-Verlag, 2005.
- [9] D. Gusfield. *Algorithms on strings, trees and sequences: computer science and computational biology*. Cambridge University Press, Cambridge, 1997.
- [10] R. M. Idury and A. A. Schäffer. Multiple matching of parameterized patterns. *Theor. Comput. Sci.*, 154(2):203–224, 1996.
- [11] G. Navarro and M. Raffinot. *Flexible Pattern Matching in Strings*. Cambridge University Press, 2002.
- [12] A. C. Yao. The complexity of pattern matching for a random string. *SIAM J. Comput.*, 8(3):368–387, 1979.