# Run-Length Encoded Nondeterministic KMP and Suffix Automata

Emanuele Giaquinta<sup>1</sup>

Department of Computer Science and Engineering, Aalto University, Finland emanuele.giaquinta@aalto.fi

Abstract. We present a novel bit-parallel representation, based on the run-length encoding, of the nondeterministic KMP and suffix automata for a string P with at least two distinct symbols. Our method is targeted to the case of long strings over small alphabets and complements the method of Cantone et al. (2012), which is effective for long strings over large alphabets. Our encoding requires  $O((\sigma+m)\lceil\rho/w\rceil)$  space and allows one to simulate the automata on a string in time  $O(\lceil\rho/w\rceil)$  per transition, where  $\sigma$  is the alphabet size, m is the length of P,  $\rho$  is the length of the run-length encoding of P and w is the machine word size in bits. The input string can be given in either unencoded or run-length encoded form.

### 1 Introduction

The string matching problem consists in finding all the occurrences of a string Pof length m in a string T of length n, both over a finite alphabet  $\Sigma$  of size  $\sigma$ . The matching can be either exact or approximate, according to some metric which measures the closeness of a match. The finite automata for the language  $\Sigma^* P$ (prefix automaton) and Suff(P) (suffix automaton), where Suff(P) is the set of suffixes of P, are the main building blocks of very efficient algorithms for the exact and approximate string matching problem. Two fundamental algorithms for the exact problem, based on the deterministic version of these automata, are the KMP and BDM algorithms, which run in O(n) and O(nm) worst-case time, respectively, using O(m) space [9,4]. In the average case, the BDM algorithm achieves the optimal  $\mathcal{O}(n \log_{\sigma}(m)/m)$  time bound. The nondeterministic version of the prefix and suffix automata can be simulated using an encoding, known as bit-parallelism, based on bit-vectors and word-level parallelism [1]. The variants of the KMP algorithm based on the nondeterministic prefix automaton, known as SHIFT-OR and SHIFT-AND, run in O(n[m/w]) worst-case time and use  $O(\sigma[m/w])$  space, where w is the machine word size in bits [1,11]. Similarly, the variant of the BDM algorithm based on the nondeterministic suffix automaton, known as BNDM, runs in  $\mathcal{O}(nm[m/w])$  worst-case time and uses  $O(\sigma[m/w])$  space [10]. In the average case, the BNDM algorithm runs in  $\mathcal{O}(n \log_{\sigma}(m)/w)$  time, which is suboptimal for patterns whose length is greater than w. There also exists a variant of SHIFT-OR which achieves  $\mathcal{O}(n \log_{\sigma}(m)/w)$ time in the average case [6]. As for the approximate string matching problem,

there are also various algorithms based on the nondeterministic prefix and suffix automata [11,10,2,7,8].

In general, the bit-parallel algorithms are suboptimal if compared to the their "deterministic" counterparts in the case m > w, because of the  $\lfloor m/w \rfloor$  additional term in the time complexity. A way to overcome this problem is to use a filtering method, namely, searching for the prefix of P of length w and verifying each occurrence with the naive algorithm. Assuming uniformly random strings, the average time complexity of SHIFT-AND and BNDM with this method is O(n)and  $O(n \log_{\sigma} w/w)$ , respectively. Recently, a few approaches were proposed to improve the case of long patterns. In 2010 Durian et al. presented three variants of BNDM tuned for the case of long patterns, two of which are optimal in the average case [5]. In the same year, Cantone et al. presented a different encoding of the prefix and suffix automata, based on word-level parallelism and on a particular factorization on strings [3]. The general approach is to devise, given a factorization f on strings, a bit-parallel encoding of the automata based on f such that one transition can be performed in O([|f(P)|/w]) time instead of  $O(\lceil m/w \rceil)$ , at the price of more space. The gain is two-fold: i) if |f(P)| < 1|P|, then the overhead of the simulation is reduced. In particular, there is no overhead if  $|f(P)| \leq w$ , which is preferable if |f(P)| < |P|; ii) if we use the filtering method, we can search for the longest substring P' of P such that  $|f(P')| \leq w$ . This yields  $O(n \log_{\sigma} |P'|/|P'|)$  average time for BNDM, which is preferable if |P'| > w. The method of Cantone et al. is effective for long patterns over large alphabets. Their factorization is such that  $\lceil m/\sigma \rceil \leq |f(P)| \leq m$  and their encoding requires  $O(\sigma^2 \lceil |f(P)|/w \rceil)$  space.

In this paper we present a novel encoding of the prefix and suffix automata, based on this approach, where f(P) is the run-length encoding of P, provided that P has at least two distinct symbols. The run-length encoding of a string is a simple encoding where each maximal consecutive sequence of the same symbol is encoded as a pair consisting of the symbol plus the length of the sequence. Our encoding requires  $O((\sigma + m)\lceil \rho/w \rceil)$  space and allows one to simulate the automata in  $O(\lceil \rho/w \rceil)$  time per transition, where  $\rho$  is the length of the runlength encoding of P. While the present algorithm uses the run-length encoded form. The run-length encoding is suitable for strings over small alphabets. Therefore, our method complements the one of Cantone et al., which is effective for large alphabets.

## 2 Notions and Basic Definitions

Let  $\Sigma$  be a finite alphabet of symbols and let  $\Sigma^*$  be the set of strings over  $\Sigma$ . The empty string  $\varepsilon$  is a string of length 0. Given a string S, we denote with |S| the length of S and with S[i] the *i*-th symbol of S, for  $0 \leq i < |S|$ . The concatenation of two strings S and  $\overline{S}$  is denoted by  $S\overline{S}$ . Given two strings S and  $\overline{S}$ , S is a substring of  $\overline{S}$  if there are indices  $0 \leq i, j < |S|$  such that  $\overline{S} = S[i]...S[j]$ . If i = 0 (j = |S| - 1) then  $\overline{S}$  is a prefix (suffix) of S. The set Suff(S) is the set of all suffixes of S. We denote by S[i..j] the substring S[i]..S[j] of S. For i > j $S[i..i] = \varepsilon$ . We denote by  $S^k$  the concatenation of k S's, for  $S \in \Sigma^*$  and  $k \ge 1$ . The string  $S^r$  is the reverse of the string S, i.e.,  $S^r = S[|S| - 1]S[|S| - 2] \dots S[0]$ .

Given a string  $P \in \Sigma^*$  of length m, we denote by  $\mathcal{A}(P) = (Q, \Sigma, \delta, q_0, F)$  the nondeterministic finite automaton (NFA) for the language  $\Sigma^* P$  of all strings in  $\Sigma^*$  whose suffix of length m is P, where:

- $\begin{array}{ll} & Q = \{q_0, q_1, \dots, q_m\} & (q_0 \text{ is the initial state}) \\ & \text{the transition function } \delta : Q \times \Sigma \longrightarrow \mathscr{P}(Q) \text{ is defined by:} \end{array}$

$$\delta(q_i, c) =_{\text{Def}} \begin{cases} \{q_0, q_1\} & \text{if } i = 0 \text{ and } c = P[0] \\ \{q_0\} & \text{if } i = 0 \text{ and } c \neq P[0] \\ \{q_{i+1}\} & \text{if } 1 \le i < m \text{ and } c = P[i] \\ \emptyset & \text{otherwise} \end{cases}$$

 $-F = \{q_m\}$ (F is the set of final states).

Similarly, we denote by  $\mathcal{S}(P) = (Q, \Sigma, \delta, I, F)$  the nondeterministic suffix automaton with  $\varepsilon$ -transitions for the language Suff (P) of the suffixes of P, where:

- $\begin{array}{ll} & Q = \{I, q_0, q_1, \dots, q_m\} & (I \text{ is the initial state}) \\ & \text{the transition function } \delta : Q \times (\Sigma \cup \{\varepsilon\}) \longrightarrow \mathscr{P}(Q) \text{ is defined by:} \end{array}$

$$\delta(q,c) =_{\text{Def}} \begin{cases} \{q_{i+1}\} & \text{if } q = q_i \text{ and } c = P[i] \quad (0 \le i < m) \\ Q & \text{if } q = I \text{ and } c = \varepsilon \\ \emptyset & \text{otherwise} \end{cases}$$

 $-F = \{q_m\}$ (F is the set of final states).

We use the notation  $q_I$  to indicate the initial state of the automaton, i.e.,  $q_I$ is  $q_0$  for  $\mathcal{A}(P)$  and I for  $\mathcal{S}(P)$ . The valid configurations  $\delta^*(q_I, S)$  which are reachable by the automata  $\mathcal{A}(P)$  and  $\mathcal{S}(P)$  on input  $S \in \Sigma^*$  are defined recursively as follows:

$$\delta^*(q_I, S) =_{Def} \begin{cases} E(q_I) & \text{if } S = \varepsilon, \\ \bigcup_{q' \in \delta^*(q_I, S')} \delta(q', c) & \text{if } S = S'c, \text{ for some } c \in \Sigma \text{ and } S' \in \Sigma^*. \end{cases}$$

where  $E(q_I)$  denotes the  $\varepsilon$ -closure of  $q_I$ . Given a string P, a run of P is a maximal substring of P containing exactly one distinct symbol. The run-length encoding (RLE) of a string P, denoted by RLE(P), is a sequence of pairs (runs)  $\langle (c_0, l_0), (c_2, l_2), \dots, (c_{\rho-1}, l_{\rho-1}) \rangle$  such that  $c_i \in \Sigma$ ,  $l_i \ge 1$ ,  $c_i \ne c_{i+1}$  for  $0 \le i < \rho$ , and  $P = c_0^{l_0} c_1^{l_1} \dots c_{\rho-1}^{l_{\rho-1}}$ . The starting (ending) position in P and length of the run  $(c_i, l_i)$  are  $\alpha_P(i) = \sum_{j=0}^{i-1} l_j$   $(\beta_P(i) = \sum_{j=0}^{i} l_j - 1)$  and  $\ell_P(i) = l_i$ , for  $i = 0, \dots, \rho - 1$ . We also put  $\alpha_P(\rho) = |P|$ .

Finally, we recall the notation of some bitwise infix operators on computer words, namely the bitwise and "&", the bitwise or "|", the left shift " $\ll$ " operator (which shifts to the left its first argument by a number of bits equal to its second argument), and the unary bitwise **not** operator " $\sim$ ".



**Fig. 1.** (a) The automata  $\mathcal{A}(P)$  and  $\mathcal{S}(P)$  for the pattern P = cttcct. The state labels corresponding to the starting positions of the runs of RLE(P) are underlined.

## 3 The Shift-And and BNDM algorithms

In this section we briefly describe the SHIFT-AND and BNDM algorithms. Given a pattern P of length m and a text T of length n, the SHIFT-AND and BNDM algorithms find all the occurrences of P in T. The SHIFT-AND algorithm works by simulating the  $\mathcal{A}(P)$  automaton on T and reporting all the positions j in T such that the final state of  $\mathcal{A}(P)$  is active in the corresponding configuration  $\delta^*(q_I, T[0 \dots j])$ . Instead, the BNDM algorithm works by sliding a window of length m along T. For a given window ending at position j, the algorithm simulates the automaton  $\mathcal{S}(P^r)$  on  $(T[j-m+1 \dots j])^r$ . Based on the simulation, the algorithm computes the length k and k' of the longest suffix of  $T[j-m+1 \dots j]$ which is a prefix and a proper prefix, respectively, of P (i.e., a suffix of  $P^r$ ). If k = m then  $T[j - m + 1 \dots j] = P$  and the algorithm reports an occurrence of P at position j. The window is then shifted by m - k' positions to the right, so as to align it with the longest proper prefix of P found. The automata are simulated using an encoding based on bit-vectors and word-level parallelism. The algorithms run in O(n[m/w]) and O(nm[m/w]) time, respectively, using  $O(\sigma[m/w])$  space, where w is the word size in bits.

# 4 RLE-based encoding of the Nondeterministic KMP and suffix automata

Given a string P of length m defined over an alphabet  $\Sigma$  of size  $\sigma$ , let  $\operatorname{RLE}(P) = \langle (c_0, l_0), (c_1, l_1, ), \dots, (c_{\rho-1}, l_{\rho-1}) \rangle$  be the run-length encoding of P. In the following, we describe how to simulate the  $\mathcal{A}(P)$  and  $\mathcal{S}(P)$  automata, using word-level parallelism, on a string S of length n in  $O(\lceil \rho/w \rceil)$  time per transition and  $O((m+\sigma)\lceil \rho/w \rceil)$  space. We recall that the simulation of the automaton  $\mathcal{A}(P)$  on a string S detects all the prefixes of S whose suffix of length m is P. Similarly, the simulation of the automaton  $\mathcal{S}(P)$  detects all the prefixes of S which are suffixes of P.

Let  $\mathcal{I}(S) = \{\alpha_S(i) \mid 0 \le i \le |\text{RLE}(S)|\}$  be the set of starting positions of the runs of S, for a given string S. Note that  $0 \in \mathcal{I}(S)$ . Given a string S, we denote with  $D_j = \delta^*(q_I, S[0 \dots j-1])$  the configuration of  $\mathcal{A}(P)$  or  $\mathcal{S}(P)$  after reading  $S[0 \dots j-1]$ , for any  $0 \le j \le |S|$ . The main idea of our algorithm is to compute the configurations  $D_j$  corresponding to positions  $j \in \mathcal{I}(S)$  only. We show that,

for any two consecutive positions  $j, j' \in \mathcal{I}(T)$ , 1) we need only the states of the automaton corresponding to positions in  $\mathcal{I}(P)$  to compute  $D_{j'}$  from  $D_j$ ; 2) if the pattern includes at least distinct two symbols (i.e., if  $\rho \geq 2$ ), there can be at most one occurrence of P ending in a position between j and j' (or equivalently spanning a proper prefix of a given run of S). Similarly, there can be at most one prefix of S ending in a position between j and j' which corresponds to a suffix of P with at least two distinct symbols. We start with the following Lemma:

**Lemma 1.** Let  $j \in \mathcal{I}(S)$ . Then, for any  $q_i \in D_j$  such that  $i \notin \mathcal{I}(P)$ , we have  $\delta(q_i, S[j]) = \emptyset$ .

*Proof.* Let  $q_i \in D_j$  with  $i \notin \mathcal{I}(P)$ . By definition of  $q_i$  it follows that S[j-1] = P[i-1] and P[i-1] = P[i], respectively. Moreover, by  $j \in \mathcal{I}(S)$ , we have  $S[j] \neq S[j-1]$ . Suppose that  $\delta(q_i, S[j]) \neq \emptyset$ , which implies S[j] = P[i]. Then we have S[j] = P[i] = P[i-1] = S[j-1], which yields a contradiction.

This Lemma states that, for any  $j \in \mathcal{I}(S)$ , any state  $q_i \in D_j$  with  $i \notin \mathcal{I}(P)$  is dead, as no transition is possible from it on S[j]. Figure 1 shows the automata  $\mathcal{A}(P)$  and  $\mathcal{S}(P)$  for P = cttcct; the state labels corresponding to indexes in  $\mathcal{I}(P)$  are underlined.

We assume that P has at least two distinct symbols. The following Lemma shows that, under this assumption, there can be at most one configuration Dcontaining the final state  $q_m$  with index between  $\alpha_S(i) + 1$  and  $\beta_S(i) + 1$  in S, for any  $1 \le i \le |\text{RLE}(S)|$  (note that  $i \ge 1$  implies that the corresponding prefix of S in the language has at least two distinct symbols).

**Lemma 2.** Let  $i \in \{1, \ldots, |\text{RLE}(S)| - 1\}$ . If  $\rho \ge 2$ , there exists at most one j' in the interval  $[\alpha_S(i), \beta_S(i)]$  such that  $q_m \in D_{j'+1}$ .

Proof. Let  $j' \in [j_1, j_2]$  such that  $q_m \in D_{j'+1}$ , where  $j_1 = \alpha_S(i)$  and  $j_2 = \beta_S(i)$ . This corresponds to i)  $S[j' - |P| + 1 \dots j'] = P$  for  $\mathcal{A}(P)$  and to ii) and  $S[0 \dots j'] \in Suff(P)$  for  $\mathcal{S}(P)$ . Since  $\rho \geq 2$  and  $i \geq 1$ , in both cases, for this to hold we must have  $S[j'-k] = c_{\rho-1}$ , for  $k = 0, \dots, l_{\rho-1}-1$ , and  $S[j'-l_{\rho-1}] = c_{\rho-2}$ , where  $c_{\rho-2} \neq c_{\rho-1}$ . By definition of  $j_1, S[j_1] = S[j_1 + 1] = \dots = S[j']$  and  $S[j_1 - 1] \neq S[j_1]$ . Hence, the only possibility is  $j' = j_1 + l_{\rho-1} - 1$ .

Specifically, the only configuration containing  $q_m$ , if any, corresponds to index  $\alpha_S(i) + l_{\rho-1}$ . By definition of  $D_j$  and by Lemma 1, we have

$$D_{\alpha_{S}(j+1)} = \delta^{*}(q_{I}, S[0 \dots \beta_{S}(j)])$$
  
=  $\bigcup_{q \in D_{\alpha_{S}(j)}} \delta^{*}(q, S[\alpha_{S}(j) \dots \beta_{S}(j)])$   
=  $\bigcup_{q \in D_{\alpha_{S}(j)} \cap \{q_{i} \mid i \in \mathcal{I}(P)\}} \delta^{*}(q, S[\alpha_{S}(j)]^{\ell_{S}(j)})$  (1)

for any position  $\alpha_S(j+1)$ . Moreover, by Lemma 2, there can be at most one configuration D with index between  $\alpha_S(j) + 1$  and  $\beta_S(j) + 1$  containing the final state  $q_m$ , for  $j \ge 1$ . For j = 0 it is easy to see that: i) in the case of the prefix automaton, since  $\rho \ge 2$ ,  $q_m \notin D_{j'+1}$  for  $j' \in [\alpha_S(0), \beta_S(0)]$ ; ii) in the case of the suffix automaton, if S[0] = P[m-1] then  $q_m \in D_{j'+1}$  for  $j' \in [\alpha_S(0), \beta_S(0)]$ 

 $[0,\min(\ell_S(0), l_{\rho-1}) - 1]$ , and  $q_m \notin D_{j'+1}$  otherwise. Hence, in the case of the suffix automaton, it is enough to test if S[0] = P[m-1] to know all the matching prefixes in the interval of the first run of S. The idea is then to compute the configurations  $D_j$ , restricted to the states with index in  $\mathcal{I}(P)$ , corresponding to positions  $j \in \mathcal{I}(S)$  only by reading S run-wise. Observe that it is not possible to detect the single prefix of S in the language, if any, ending at a position between  $\alpha_S(j-1)$  and  $\beta_S(j-1)$  using  $D_{\alpha_S(j)}$ , because  $q_m \notin D_{\alpha_S(j)}$  if the prefix does not end at position  $\beta_S(j-1)$ , or equivalently if  $\ell_S(j-1) > l_{\rho-1}$ . To overcome this problem we modify the automata by adding a self-loop on  $q_m$  labeled by P[m-1]. In this way, if  $q_m$  belongs to a configuration D with index between  $\alpha_S(j-1) + 1$  and  $\beta_S(j-1) + 1$  then  $q_m$  will be in  $D_{\alpha_S(j)}$ . Observe that, if  $q_m \in D_{\alpha_S(j)}$ , then  $q_m \notin D_{\alpha_S(j+1)}$ , since  $S[\alpha_S(j)] \neq S[\alpha_S(j+1)]$ .

Let  $\overline{D}_j = \{1 \leq i \leq \rho \mid q_{\alpha_P(i)} \in D_{\alpha_S(j)}\}$ , be the encoding of the configuration of  $\mathcal{A}(P)$  after reading  $S[0 \dots \beta_S(j-1)]$ , for  $0 \leq j \leq |\text{RLE}(S)|$ . For example:

j	P = cttcct			
2	S = cttccttcct			
j	$\bar{D}_1 = \{1\}$	$\bar{D}_2$	$_{2} = \{2\}$	
$\bar{D}_3 = \{1,3\}$ $\bar{D}_4 = \{2,4\}$				
$\bar{D}_5 = \{1,3\}$ $\bar{D}_6 = \{4\}$				
_		-		•
	$0\ 1\ 2\ 3\ 4$	i	$0\ 1\ 2\ 3\ 4$	5.6
P(i)	$0\ 1\ 3\ 5\ 6$	$\alpha_{S}(i)$	$0\ 1\ 3\ 5\ 7$	9 10

Note that  $q_0$  is not represented and that  $\overline{D}_0$  is equal to  $\emptyset$  and  $\{1, \ldots, \rho\}$  for  $\mathcal{A}(P)$ and  $\mathcal{S}(P)$ , respectively. We now describe how to compute the configurations  $\overline{D}_j$ , starting with the automaton  $\mathcal{A}(P)$ . The following property easily follows from Equation 1:

**Lemma 3.** For any  $0 \le j < |\text{RLE}(S)|$  and  $1 \le i \le \rho$ ,  $q_{\alpha_P(i)} \in D_{\alpha_S(j+1)}$  if and only if either

a) 
$$q_{\alpha_P(i-1)} \in D_{\alpha_S(j)} \land S[\alpha_s(j)] = c_{i-1} \land \ell_S(j) = l_{i-1}, \text{ for } i = 2, \dots, \rho - 1;$$
  
b)  $q_{\alpha_P(i-1)} \in D_{\alpha_S(j)} \land S[\alpha_s(j)] = c_{i-1} \land \ell_S(j) \ge l_{i-1}, \text{ for } i \in \{1, \rho\}.$ 

Based on the above Lemma and the fact that there is a self-loop on  $q_0$  labeled by  $\Sigma$ , we have that

$$\bar{D}_{j+1} = \{i+1 \mid i \in \bar{D}_j \cup \{0\}\} 
\cap \{1 \le i \le \rho \mid S[\alpha_S(j)] = c_{i-1}\} 
\cap (\{2 \le i \le \rho - 1 \mid \ell_S(j) = l_{i-1}\} \cup \{i \in \{1, \rho\} \mid \ell_S(j) \ge l_{i-1}\})$$
(2)

for  $j \ge 0$ . We now show to implement Equation 2 efficiently using word-level parallelism. Let

$$B_1(c) = \{ 1 \le i \le \rho \mid c = c_{i-1} \}, \\ B_2(l) = \{ 1 \le i \le \rho \mid l = l_{i-1} \} \\ \cup \{ i \in \{1, \rho\} \mid l \ge l_{i-1} \}$$

for any  $c \in \Sigma$  and  $1 \leq l \leq |P| + 1$ . The set  $B_1(c)$  includes the indices of all the runs whose symbol is equal to c. Similarly, The set  $B_2(l)$  includes the indices of all the runs whose length is equal to l and, in addition, the index of the first and/or last run if the corresponding length is smaller than or equal to l. Note that  $B_2(l) = \{1, \rho\}$ , for any l > |P|; thus, we can define  $B_2$  up to |P| + 1 and map any integer greater than |P| onto |P| + 1. For example, for P = cttcct we have:

$$B_1(c) = \{0, 2\} B_2(1) = \{0, 3\} B_1(t) = \{1, 3\} B_2(2) = \{0, 1, 2, 3\}$$

and  $B_2(l) = \{0,3\}$ , for  $3 \le l \le 7$ . We represent the configurations D and the sets B as bit-vectors of  $\rho$  bits, denoted with D and B, respectively. Based on these two definitions, Equation 2 can be rewritten as

$$\bar{D}_{j+1} = \{i+1 \mid i \in \bar{D}_j \cup \{0\}\} \cap B_1(S[\alpha_S(j)]) \cap B_2(\min(\ell_S(j), |P|+1)),$$

which corresponds to the following bitwise operations

$$\mathsf{D}_{j+1} = ((\mathsf{D}_j \ll 1) \mid 0^{\rho-1}1) \& \mathsf{B}_1(S[\alpha_S(j)]) \& \mathsf{B}_2(\min(\ell_S(j), |P|+1)).$$

We now describe the computation of  $D_j$  for the automaton S(P). The simulation of S(P) is equivalent to the one of  $\mathcal{A}(P)$ , up to state  $q_0$  and the first transition. For j > 1, Lemma 3 holds if index 1 is handled by case a instead of b, which accounts for the fact that there is no self-loop on  $q_0$ . We can thus change Equation 2 and the definition of  $B_2$  accordingly. Instead, for j = 1 (i.e., the transition on the first run of S), we have

$$q_{\alpha_P(i)} \in D_{\alpha_S(1)} \iff S[0] = c_{i-1} \land (i = \rho \lor \ell_S(0) \le l_{i-1}), \text{ for } i = 1, \dots, \rho,$$

and

$$\bar{D}_1 = \{ 1 \le i \le \rho \mid S[0] = c_{i-1} \} \cap \left( \{ 1 \le i < \rho \mid \ell_S(0) \le l_{i-1} \} \cup \{\rho\} \right),$$

since, before the first transition, all the states are active because of the  $\varepsilon$ -transitions and so state  $q_{\alpha_P(i)}$  can be activated by any state with index between  $\alpha_P(i-1)$  and  $\alpha_P(i) - 1$  by reading a run with symbol equal to  $c_{i-1}$  and length no larger than  $l_{i-1}$ , with the exception of  $q_m$  which can be activated with a run of any length because of the self-loop. To account for this case we define the set  $B_3(l) = \{1 \leq i \leq \rho \mid l \leq l_{i-1}\} \cup \{\rho\}$ , for  $1 \leq l \leq |P| + 1$ . The set  $B_3(l)$  includes the indices of all the runs whose length is greater than or equal to l and, in addition, the index of the last run. Note that  $B_3(l) = \{\rho\}$  for  $l \geq |P| + 1$ , so we can define  $B_3$  up to |P| + 1 and map any integer greater than |P| onto |P| + 1, as done for  $B_2$ . Then, we have:

$$\overline{D}_1 = B_1(S[0]) \cap B_3(\min(\ell_S(0), |P|+1)),$$

which corresponds to the following bitwise operations

$$\mathsf{D}_1 = \mathsf{B}_1(S[0]) \& \mathsf{B}_3(\min(\ell_S(0), |P|+1)).$$



Fig. 2. The variants of SHIFT-AND and BNDM based on the run-length encoding.

The computation of a single configuration  $\overline{D}_j$  requires  $O(\lceil \rho/w \rceil)$  time. The total time complexity of the simulation is thus  $O(|S|\lceil \rho/w \rceil)$ , as the total number of configurations is  $|\text{RLE}(S)| \leq |S|$ . The bit-vectors B can be preprocessed in  $O(m + (\sigma + m)\lceil \rho/w \rceil)$  time and require  $O((\sigma + m)\lceil \rho/w \rceil)$  space. The string P or S can be given in either unencoded or run-length encoded form. In the former case its run-length encoding does not need to be stored. It can be computed on the fly in O(m) or O(|S|) time, using constant space, during the preprocessing or searching phase.

## 5 The variants of Shift-And and BNDM

The variants of the SHIFT-AND and BNDM algorithms based on the encoding described in the previous section run in  $O(n\lceil \rho/w\rceil)$  and  $O(nm\lceil \rho/w\rceil)$  time, respectively, using  $O((\sigma + m)\lceil \rho/w\rceil)$  space. The encoding of the suffix automaton is however not ideal in practice, due to the different first transition. We now describe a variant of BNDM, based on a modified suffix automaton, where the first transition of the automaton is equal to the subsequent ones. Let j be the ending position of a window of length m in T and let  $\overline{j}$  be the minimum position



**Fig. 3.** (a) The automaton  $S_R(P)$  for the pattern P = cttcct.

such that  $\overline{j} \geq j$  and  $T[\overline{j}] \neq T[\overline{j}+1]$ . In other words,  $\overline{j}$  is the ending position of the run of RLE(T) spanning T[j]. Consider the window of length  $m + \overline{j} - j$ ending at  $\overline{j}$ . By Lemma 2, if  $\rho \geq 2$ , there can be at most one occurrence of Pending at a position between j and  $\overline{j}$ . Our idea is to process this larger window with  $S(P^r)$ , finding the single occurrence of P in it, if any, and the length k' of the longest proper prefix of P ending at position  $\overline{j}$ . We then shift the window by  $m + \overline{j} - j - k'$  positions to the right. If  $\overline{j} - j < m$ , the time needed to process this window is O(m). Otherwise, if  $\overline{j} - j \geq m$ , observe that, since k' < m, the symbols of T in the interval  $[j + 1, \overline{j} - m]$  are covered by this window only and therefore the time needed to process this window is O(m), for reading the symbols of Tin the intervals [j - m + 1, j] and  $[\overline{j} - m + 1, \overline{j}]$ , plus a term which, summed over all such windows, is O(n). Hence, the time complexity of the algorithm remains  $O(nm[\rho/w])$  in the worst-case.

Suppose that we simulate the automaton  $\mathcal{S}(P^r)$  on  $(T[j-m+1...\bar{j}])^r$  and let k' be the length of the longest suffix of  $T[i - m + 1 \dots \overline{i}]$  which is a proper prefix of P. Observe that, if k' does not correspond to the ending position of a run of RLE(P), i.e., if P[k'-1] = P[k'], then the window corresponding to shift k' does not contain an occurrence of P, because  $P[k'-1] = T[\overline{j}]$  and  $T[\overline{j}] \neq T[\overline{j}+1]$ (shifting the window by  $m + \bar{j} - j - k'$  corresponds to aligning position k' with  $\overline{j}$  + 1). Indeed, it is easy to see that the smallest useful shift corresponds to the length k' of the longest suffix of  $T[j - m + 1 \dots \overline{j}]$  which is a proper prefix of P and such that  $P[k'-1] \neq P[k']$ . Hence, we can modify the automaton  $\mathcal{S}(P)$  so that it recognizes the subset of Suff(P) { $P[i \dots m-1] \mid i = 0 \lor P[i] \neq 0$ P[i-1]. To accomplish this, it is enough to remove all the  $\varepsilon$ -transitions entering a state  $q_i$  with P[i] = P[i-1]. Observe that the automaton can recognize the occurrence of P, if any, ending at a position between j and  $\overline{j}$  only if it ends at position  $\overline{j}$ . To account for this problem, we add a self-loop on  $q_0$  labeled by P[0]. We denote the resulting automaton with  $\mathcal{S}_R(P)$ . The transition function of the  $\mathcal{S}_R(P)$  automaton is defined as follows:

$$\delta(q,c) =_{\text{Def}} \begin{cases} \{q_0,q_1\} & \text{if } q = q_0 \text{ and } c = P[0] \\ \{q_{i+1}\} & \text{if } q = q_i \text{ and } c = P[i] \\ \{q_i \mid i = 0 \lor P[i] \neq P[i-1]\} & \text{if } q = I \text{ and } c = \varepsilon \\ \emptyset & \text{otherwise} \end{cases}$$

Figure 3 shows the automaton  $S_R(P)$  for P = cttcct. It is not hard to verify that Lemma 3 also holds for  $S_R(P)$  and therefore the simulation of this automaton is analogous to the one of the  $\mathcal{A}(P)$  automaton.

The pseudocode of the variants of the SHIFT-AND and BNDM algorithms based on the run-length encoding is shown in Figure 4.

### 6 Acknowledgments

We thank Jorma Tarhio for helpful comments.

## References

- R. A. Baeza-Yates and G. H. Gonnet. A new approach to text searching. Commun. ACM, 35(10):74–82, 1992.
- R. A. Baeza-Yates and G. Navarro. Faster approximate string matching. Algorithmica, 23(2):127–158, 1999.
- D. Cantone, S. Faro, and E. Giaquinta. A compact representation of nondeterministic (suffix) automata for the bit-parallel approach. *Inf. Comput.*, 213:3–12, 2012.
- 4. M. Crochemore and W. Rytter. Text Algorithms. Oxford University Press, 1994.
- B. Durian, H. Peltola, L. Salmela, and J. Tarhio. Bit-parallel search algorithms for long patterns. In *Experimental Algorithms, 9th International Symposium, SEA* 2010, Ischia Island, Naples, Italy, May 20-22, 2010. Proceedings, pages 129–140, 2010.
- K. Fredriksson and S. Grabowski. Average-optimal string matching. J. Discrete Algorithms, 7(4):579–594, 2009.
- H. Hyyrö. Improving the bit-parallel NFA of baeza-yates and navarro for approximate string matching. *Inf. Process. Lett.*, 108(5):313–319, 2008.
- H. Hyyrö and G. Navarro. Bit-parallel witnesses and their applications to approximate string matching. Algorithmica, 41(3):203–231, 2005.
- D. E. Knuth, J. H. M. Jr., and V. R. Pratt. Fast pattern matching in strings. SIAM J. Comput., 6(2):323–350, 1977.
- G. Navarro and M. Raffinot. Fast and flexible string matching by combining bitparallelism and suffix automata. ACM Journal of Experimental Algorithmics, 5:4, 2000.
- S. Wu and U. Manber. Fast text searching allowing errors. Commun. ACM, 35(10):83–91, 1992.