A Simple Streaming Bit-parallel Algorithm for Swap Pattern Matching^{*}

Václav Blažej^{**}, Ondřej Suchý^{* * *}, and Tomáš Valla[†]

Faculty of Information Technology, Czech Technical University in Prague, Prague, Czech Republic

Abstract. The pattern matching problem with swaps is to find all occurrences of a pattern in a text while allowing the pattern to swap adjacent symbols. The goal is to design fast matching algorithm that takes advantage of the bit parallelism of bitwise machine instructions and has only streaming access to the input. We introduce a new approach to solve this problem based on the graph theoretic model and compare its performance to previously known algorithms. We also show that an approach using deterministic finite automata cannot achieve similarly efficient algorithms. Furthermore, we describe a fatal flaw in some of the previously published algorithms based on the same model. Finally, we provide experimental evaluation of our algorithm on real-world data.

1 Introduction

In the Pattern Matching problem with Swaps (Swap Matching, for short), the goal is to find all occurrences of any swapped version of a pattern P in a text T, where P and T are strings of length p and t over an alphabet Σ , respectively. By the swapped version of a pattern P we mean a string of symbols created from P by swapping adjacent symbols while ensuring that each symbol is swapped at most once (see Section 2 for formal definitions). The solution of Swap Matching is a set of indices which represent where occurrences swapped version of P in T begin. Swap Matching is intensively studied due to its use in practical applications such as text and music retrieval, data mining, network security and biological computing [7].

The swap of two consecutive symbols is one of the most typical typing errors. It also represent a simpler version of swaps that appear in nature. In particular,

** Supported by the OP VVV MEYS funded project CZ.02.1.01/0.0/0.0/16_019/0000765 "Research Center for Informatics" and by the SGS CTU project SGS17/209/OHK3/3T/18.

^{*} An extended abstract of this work appeared in the Proceedings of the 7th International Conference on Mathematical Aspects of Computer and Information Sciences, MACIS 2017 [8].

 $^{^{\}star\,\star\,\star}$ Supported by grant 17-20065S of the Czech Science Foundation.

[†] Supported by the Centre of Excellence – Inst. for Theor. Comp. Sci. 79 (project P202/12/G061 of the Czech Science Foundation.)

the phenomenon of swaps occurs in gene mutations and duplications such as in the region of human chromosome 5 that is implicated in the disease called spinal muscular Atrophy, a common recessive form of muscular dystrophy [18]. While the biological swaps occur at a gene level and have several additional constraints and characteristics, which make the problem much more difficult, they do serve as a convincing pointer to the theoretical study of swaps as a natural edit operation for the approximation metric [2]. Indeed Lowrance and Wagner [21] suggested to add the swap operation when considering the edit distance of two strings.

Swap Matching was introduced in 1995 as an open problem in non-standard string matching [20]. The first result was reported by Amir et al. [2] in 1997, who provided an $O(tp^{\frac{1}{3}} \log p)$ -time solution for alphabets of size 2, while also showing that alphabets of size exceeding 2 can be reduced to size 2 with a little overhead. Amir et al. [5] came up with solution with $O(t\log^2 p)$ time complexity for some very restrictive cases. Several years later Amir et al. [3] showed that Swap Matching can be solved by an algorithm for the overlap matching achieving the running time of $O(t\log p \log |\Sigma|)$. This algorithm as well as all the previous ones is based on fast Fourier transformation (FFT).

In 2008 Iliopoulos and Rahman [17] introduced a new graph theoretic approach to model the Swap Matching problem and came up with the first efficient solution to Swap Matching without using FFT (we show it to be incorrect). Their algorithm based on bit parallelism runs in $O((t + p) \log p)$ time if the pattern length is similar to the word-size of the target machine. One year later Cantone and Faro [10] presented a dynamic programming algorithm named Cross Sampling solving Swap Matching in O(t) time and $O(|\Sigma|)$ space, assuming that the pattern length is similar to the word-size in the target machine. In the same year Campanelli et al. [9] enhanced the Cross Sampling algorithm using notions from Backward directed acyclic word graph matching algorithm and named the new algorithm Backward Cross Sampling. This algorithm also assumes short pattern length. Although Backward Cross Sampling has $O(|\Sigma|)$ space and O(tp) time complexity, which is worse than that of Cross Sampling, it improves the real-world performance.

In 2013 Faro [14] presented a new model to solve Swap Matching using reactive automata and also presented a new algorithm with O(t) time complexity assuming short patterns. The same year Chedid [12] reformulated the dynamic programming solution by Cantone and Faro [10] which results in more intuitive algorithms. In 2014 a minor improvement by Fredriksson and Giaquinta [15] appeared, yielding slightly (at most factor $|\Sigma|$) better asymptotic time complexity (and also slightly worse space complexity) for special cases of patterns. The same year Ahmed et al. [1] took ideas of the algorithm by Iliopoulos and Rahman [17] and devised two algorithms named SMALGO-I and SMALGO-II which both run in O(t) for short patterns, but bear the same error as the original algorithm.

Another remarkable effort related to Swap Matching is to actually count the number of swaps needed to match the pattern at the location [6]. This is more often studied with an extra operation of character change allowed [4,13,19].

Our Contribution. We design a simple algorithm which solves the Swap Matching problem. The goal is to design a streaming algorithm, which is given one symbol per each execution step until the end-of-input arrives, and thus does not need access to the whole input. This algorithm has $O(\lceil \frac{p}{w} \rceil (|\Sigma| + t) + p)$ time and $O(\lceil \frac{p}{w} \rceil |\Sigma|)$ space complexity where w is the word-size of the machine. We would like to stress that our solution, as based on the graph theoretic approach, does not use FFT. Therefore, it yields a much simpler non-recursive algorithm allowing bit parallelism and is not suffering from the disadvantages of the convolution-based methods. While our algorithm matches the best asymptotic complexity bounds of the previous results [10,15] (up to a $|\Sigma|$ factor), we believe that its strength lies in the applications where the alphabet is small and the pattern length is at most the word-size, as it can be implemented using only $7 + |\Sigma|$ CPU registers and few machine instructions. This makes it practical for tasks like DNA sequences scanning. Also, as far as we know, our algorithm is currently the only known streaming algorithm for the swap matching problem.

We continue by proving that any deterministic finite automaton that solves Swap Matching has number of states exponential in the length of the pattern.

We also describe the SMALGO (swap matching algorithm) by Iliopoulos and Rahman [17] in detail. Unfortunately, we have discovered that SMALGO and derived algorithms contain a flaw which cause false positives to appear. We have prepared implementations of SMALGO-I, Cross Sampling, Backward Cross Sampling and our own algorithm, measured the running times and the rate of false positives for the SMALGO-I algorithm. All of the sources are available for download.¹

This paper is organized as follows. First we introduce all the basic definitions, and also recall the graph theoretic model introduced in [17] and its use for matching in Section 2. In Section 3 we show our algorithm for Swap Matching problem and follow it in Section 4 with the proof that Swap Matching cannot be solved efficiently by deterministic finite automata. Then we describe the SMALGO algorithms in detail in Section 5 and finish with the experimental evaluation of the algorithms in Section 6.

2 Basic Definitions and the Graph Theoretic Model

In this section we state the basic definitions, present the graph theoretic model and show a basic algorithm that solves Swap Matching using the model.

2.1 Notations and Basic Definitions

We use the word-RAM as our computational model. That means we have access to memory cells of fixed capacity w (e.g., 64 bits). A standard set of arithmetic and bitwise instructions include And (&), Or (|), Left bitwise-shift (LShift or $\ll 1$) and Right bitwise-shift (RShift or $\gg 1$). Each of the standard operations on

¹ http://users.fit.cvut.cz/blazeva1/gsm.html



Fig. 1. *P*-graph \mathcal{P}_P for the pattern P = abcbbac

words takes single unit of time. In order to compare to other existing algorithms, which are not streaming, we define the access to the input in a less restrictive way – the input is read from a read-only part of memory and the output is written to a write-only part of memory. However, it will be easy to observe that our algorithm accesses the input sequentially. We do not include the input and the output into the space complexity analysis.

A string S over an alphabet Σ is a finite sequence of symbols from Σ and |S|is its length. By S_i we mean the *i*-th symbol of S and we define a substring $S_{[i,j]} = S_i S_{i+1} \dots S_j$ for $1 \leq i \leq j \leq |S|$, and prefix $S_{[1,i]}$ for $1 \leq i \leq |S|$. String P prefix matches string T k symbols on position i if $P_{[1,k]} = T_{[i,i+k-1]}$.

Next we formally introduce a swapped version of a string.

Definition 1 (Campanelli et al. [9]). A swap permutation for S is a permutation $\pi : \{1, \ldots, n\} \rightarrow \{1, \ldots, n\}$, where n = |S|, such that:

- (i) if $\pi(i) = j$ then $\pi(j) = i$ (symbols at positions i and j are swapped),
- (ii) for all $i, \pi(i) \in \{i 1, i, i + 1\}$ (only adjacent symbols are swapped),
- (iii) if $\pi(i) \neq i$ then $S_{\pi(i)} \neq S_i$ (identical symbols are not swapped).

For a string S a swapped version $\pi(S)$ is a string $\pi(S) = S_{\pi(1)}S_{\pi(2)}\dots S_{\pi(n)}$ where π is a swap permutation for S.

Now we formalize the version of matching we are interested in.

Definition 2. Given a text $T = T_1T_2...T_t$ and a pattern $P = P_1P_2...P_p$, the pattern P is said to swap match T at location i if there exists a swapped version $\pi(P)$ of P that matches T at location i, i.e., $\pi(P) = T_{[i,i+p-1]}$.

2.2 A Graph Theoretic Model

The algorithms in this paper are based on a model introduced by Iliopoulos and Rahman [17]. In this section we briefly describe this model.

For a pattern P of length p we construct a labeled graph $\mathcal{P}_P = (V, E, \sigma)$ with vertices V, edges E, and a vertex labeling function $\sigma : V \to \Sigma$ (see Fig. 1 for an example). Let $V = V' \setminus \{m_{-1,1}, m_{1,p}\}$ where $V' = \{m_{r,c} \mid r \in \{-1, 0, 1\}, c \in \}$

Algorithm 1 The basic matching algorithm (BMA)

Input: Labeled directed acyclic graph $G = (V, E, \sigma)$, set $Q_0 \subseteq V$ of starting vertices, set $F \subseteq V$ of accepting vertices, text T, and position k.

- 1: Let $D'_1 := Q_0$.
- 2: for $i = 1, 2, 3, \dots, p$ do

3: Let $D_i := \{x \mid x \in D'_i, \sigma(x) = T_{k+i-1}\}.$

- 4: **if** $D_i = \emptyset$ **then** finish.
- 5: **if** $D_i \cap F \neq \emptyset$ **then** we have found a match and finish.
- 6: Define the next iteration set D'_{i+1} as vertices which are successors of D_i , i.e., $D'_{i+1} := \{ d \in V(\mathcal{P}_P) \mid (v, d) \in E(\mathcal{P}_P) \text{ for some } v \in D_i \}.$

 $\{1, 2, \ldots, p\}\}$. For $m_{r,c} \in V$ we set $\sigma(m_{r,c}) = P_{r+c}$. Each vertex $m_{r,c}$ is identified with an element of a $3 \times p$ grid. We set $E' := E'_1 \cup E'_2 \cup \cdots \cup E'_{p-1}$, where $E'_j := \{(m_{k,j}, m_{i,j+1}) \mid k \in \{-1, 0\}, i \in \{0, 1\}\} \cup \{(m_{1,j}, m_{-1,j+1})\}$, and let $E = E' \cap V \times V$. We call \mathcal{P}_P the *P*-graph. Note that \mathcal{P}_P is directed acyclic graph, $|V(\mathcal{P}_P)| = 3p - 2$, and $|E(\mathcal{P}_P)| = 5(p-1) - 4$.

The idea behind the construction of \mathcal{P}_P is as follows. We create vertices V' and edges E' which represent every swap pattern without unnecessary restrictions (equal symbols can be swapped). We remove vertices $m_{-1,1}$ and $m_{1,p}$ which represent symbols from invalid indices 0 and p + 1.

The *P*-graph now represents all possible swap permutations of the pattern *P* in the following sense. Vertices $m_{0,j}$ represent ends of prefixes of swapped version of the pattern which end by a non-swapped symbol. Possible swap of symbols P_j and P_{j+1} is represented by vertices $m_{1,j}$ and $m_{-1,j+1}$. Edges represent symbols which can be consecutive. Each path from column 1 to column *p* represents a swap pattern and each swap pattern is represented this way.

Definition 3. For a given Σ -labeled directed acyclic graph $G = (V, E, \sigma)$ vertices $s, e \in V$ and a directed path $f = v_1, v_2, \ldots, v_k$ from $v_1 = s$ to $v_k = e$, we call $S = \sigma(f) = \sigma(v_1)\sigma(v_2)\ldots\sigma(v_k) \in \Sigma^*$ a path string of f.

2.3 Using Graph Theoretic Model for Matching

In this section we describe an algorithm called *Basic Matching Algorithm* (BMA) which can determine whether there is a match of pattern P in text T on a position k using any graph model which satisfies the following conditions.

- It is a directed acyclic graph,
- $-V = V_1 \uplus V_2 \uplus \cdots \uplus V_p$ (we can divide vertices to columns),
- $E \subseteq \{(u, w) \mid u \in V_i, w \in V_{i+1}, 1 \le i < p\} \text{ (edges lead to next column)}.$

Let $Q_0 = V_1$ be the starting vertices and $F = V_p$ be the accepting vertices. BMA is designed to run on any graph which satisfies these conditions. Since P-graph satisfies these assumptions we can use BMA for \mathcal{P}_P .

The algorithm runs as follows (see also Algorithm 1). We initialize the algorithm by setting $D'_1 := Q_0$ (Step 1). D'_1 now holds information about vertices

6

Algorithm 2 BMA in terms of prefix match signals

1: Let $I^0(v) := 1$ for each $v \in Q_0$ and $I^0(v) := 0$ for each $v \notin Q_0$. 2: **for** i = 0, 1, 2, 3, ..., p - 1 **do** 3: Filter signals by a symbol T_{k+i} . 4: **if** $I^i(v) = 0$ for every $v \in \mathcal{P}_P$ **then** finish. 5: **if** $I^i(v) = 1$ for any $v \in F$ **then** we have found a match and finish. 6: Propagate signals along the edges.

which are the end of some path f starting in Q_0 for which $\sigma(f)$ possibly prefix matches 1 symbol of $T_{[k,k+p-1]}$. To make sure that the path f represents a prefix match we need to check whether the label of the last vertex of the path fmatches the symbol T_k (Step 3). If no prefix match is left we did not find a match (Step 4). If some prefix match is left we need to check whether we already have a complete match (Step 5). If the algorithm did not stop it means that we have some prefix match but it is not a complete match yet. Therefore we can try to extend this prefix match by one symbol (Step 6) and check whether it is a valid prefix match (Step 3). Since we extend the matched prefix in each step, we repeat these steps until the prefix match is as long as the pattern (Step 2).

Having vertices in sets is not handy for computing so we present another way to describe this algorithm. We use their characteristic vectors instead.

Definition 4. A Boolean labeling function $I : V \to \{0, 1\}$ of vertices of \mathcal{P}_P is called a prefix match signal.

The algorithm can be easily divided into *iterations* according to the value of i in Step 2. We denote the value of the prefix match signal in j-th iteration as I^{j} and we define the following operations:

- propagate signal along the edges, is an operation which sets $I^{j}(v) := 1$ if and only if there exists an edge $(u, v) \in E$ with $I^{j-1}(u) = 1$,
- filter signal by a symbol $x \in \Sigma$, is an operation which sets I(v) := 0 for each v where $\sigma(v) \neq x$,
- match check, is an operation which checks whether there exists $v \in F$ such that I(v) = 1 and if so reports a match.

With these definitions in hand we can describe BMA in terms of prefix match signals as Algorithm 2. See Fig. 2 for an example of use of BMA to figure out whether P = acbab swap matches T = babcabc at a position 2.

2.4 Shift-And Algorithm

The following description is based on [11, Chapter 5] describing the Shift-Or algorithm.

For a pattern P and a text T of length p and t, respectively, let R be a bit array of size p and R^j its value after text symbol T_j has been processed. It contains information about all matches of prefixes of P that end at the position j

7



Fig. 2. BMA of $T_{[2,6]} = abcab$ on a *P*-graph of the pattern P = acbab. The prefix match signal propagates along the dashed edges. Index *j* above a vertex *v* represent that $I^{j}(v) = 1$, otherwise $I^{j}(v) = 0$.

in the text. For $1 \leq i \leq p$, $R_i^j = 1$ if $P_{[1,i]} = T_{[j-i+1,j]}$ and 0 otherwise. The vector R^{j+1} can be computed from R^j as follows. For each positive i we have $R_{i+1}^{j+1} = 1$ if $R_i^j = 1$ and $P_{i+1} = T_{j+1}$, and $R_{i+1}^{j+1} = 0$ otherwise. Furthermore, $R_1^{j+1} = 1$ if $P_1 = T_{j+1}$ and 0 otherwise. If $R_p^{j+1} = 1$ then a complete match can be reported.

The transition from R^j to R^{j+1} can be computed very fast as follows. For each $x \in \Sigma$ let D^x be a bit array of size p such that for $1 \leq i \leq p, D_i^x = 1$ if and only if $P_i = x$. The array D^x denotes the positions of the symbol x in the pattern P. Each D^x can be preprocessed before the search. The computation of R^{j+1} is then reduced to three bitwise operations, namely $R^{j+1} = (\text{LShift}(R^j) \mid 1) \& D^{T_{j+1}}$. When $R_p^j = 1$, the algorithm reports a match on a position j - p + 1.

3 Our Algorithm

In this section we will show an algorithm which solves Swap Matching. We call the algorithm GSM (Graph Swap Matching). GSM uses the graph theoretic model presented in Section 2.2 and is based on the Shift-And algorithm from Section 2.4.

The basic idea of the GSM algorithm is to represent prefix match signals (see Definition 4) from the basic matching algorithm (Section 2.3) over \mathcal{P}_P in bit vectors. The GSM algorithm represents all signals I in the bitmaps RX formed by three vectors, one for each row. Each time GSM processes a symbol of T, it first propagates the signal along the edges, then filters the signal and finally checks for matches. All these operations can be done very quickly thanks to bitwise parallelism.

First, we make the concept of GSM more familiar by presenting a way to interpret the Shift-And algorithm by means of the basic matching algorithm (BMA) from Section 2.3 to solve the (ordinary) Pattern Matching problem. Then we expand this idea to Swap Matching by using the graph theoretic model.

3.1 Graph Theoretic View of the Shift-And Algorithm

Let T and P be a text and a pattern of lengths t and p, respectively. We create the T-graph $\mathcal{T}_P = (V, E, \sigma)$ of the pattern P.

8

Definition 5. Let S be a string. The T-graph of S is a graph $\mathcal{T}_S = (V, E, \sigma)$ where $V = \{v_i \mid 1 \leq i \leq |S|\}, E = \{(v_i, v_{i+1}) \mid 1 \leq i \leq |S-1|\}$ and $\sigma : V \to \Sigma$ such that $\sigma(v_i) = S_i$.

Note that the *T*-graph is directed acyclic graph which can be divided into columns $V_i, 1 \leq i \leq p$ (each of them containing one vertex v_i) such that the edges lead from V_j to V_{j+1} . This means that the *T*-graph satisfies all assumptions of BMA. We apply BMA to \mathcal{T}_P to figure out whether *P* matches *T* at a position *j*. We get a correct result because for each $i \in \{1, \ldots, p\}$ we check whether $T_{j+i-1} = \sigma(v_i) = P_i$.

To find every occurrence of P in T we would have to run BMA for each position separately. This is basically the naive approach to solve the pattern matching. We can improve the algorithm significantly when we parallelize the computations of p runs of BMA in the following way.

The algorithm processes one symbol at a time starting from T_1 . We say that the algorithm is in the *j*-th step when a symbol T_j has been processed. BMA represents a prefix match as a prefix match signal $I : V \to \{0, 1\}$. Its value in the *j*-th step is denoted I^j . Since one run of the BMA uses only one column of the *T*-graph at any time we can use other vertices to represent different runs of the BMA. We represent all prefix match indicators in one vector so that we can manipulate them easily. To do that we prepare a bit vector *R*. Its value in *j*-th step is denoted R^j and defined as $R_i^j = I^j(v_i)$.

First operation which is used in BMA (propagate signal along the edges) can be done easily by setting the signal of v_i to value of the signal of its predecessor v_{i-1} in the previous step. I.e., for $i \in \{1, \ldots, p\}$ we set $I^j(v_i) = 1$ if i = 1 and $I^j(v_i) = I^{j-1}(v_{i-1})$ otherwise. In terms of R^j this means just $R^j = \text{LSO}(R^{j-1})$, where LSO is defined as $\text{LSO}(x) = \text{LShift}(x) \mid 1$.

We also need a way to set $I(v_i) := 0$ for each v_i for which $\sigma(v_i) \neq T_{j+i}$ which is another basic BMA operation (filter signal by a symbol). We can do this using the bit vector D^x from Section 2.4 and taking $R \& D^x$. I.e., the algorithm computes R^j as $R^j = \text{LSO}(R^{j-1}) \& D^{T_{j+1}}$.

The last BMA operation we have to define is the *match detection*. We do this by checking whether $R_p^j = 1$ and if this is the case then a match starting at position j - p + 1 occurred.

3.2 Our Algorithm for Swap Matching Using the Graph Theoretic Model

Now we are ready to describe the GSM algorithm.

We again let $\mathcal{P}_P = (V, E, \sigma)$ be the *P*-graph of the pattern *P*, apply BMA to \mathcal{P}_P to figure out whether *P* matches *T* at a position *j*, and parallelize *p* runs of BMA on \mathcal{P}_P .

Again, the algorithm processes one symbol at a time and it is in the *j*-th step when symbol T_j is being processed. We again denote the value of the prefix match signal $I : V \to \{0, 1\}$ of BMA in the *j*-th step by I^j . I.e., the semantic

Algorithm 3 The graph swap matching (GSM)

Input: Pattern P of length p and text T of length t over alphabet Σ . Output: Positions of all swap matches. 1: Let $RU^0 := RM^0 := RD^0 := 0^p$. 2: Let $D^x := 0^p$, for all $x \in \Sigma$. 3: for $i = 1, 2, 3, \dots, p$ do $D_i^{P_i} := 1$ 4: 5: for $j = 1, 2, 3, \dots, t$ do $RU'^j := \mathrm{LSO}(RD^{j-1}).$ 6: $RM^{\prime j} := \mathrm{LSO}(RM^{j-1} \mid RU^{j-1}).$ 7: $RD'^{j} := LSO(RM^{j-1} \mid RU^{j-1}).$ 8: $RU^j := RU'^j \& \mathrm{LShift}(D^{T_j}).$ 9: $RM^j := RM'^j \& D^{T_j}.$ 10: $RD^j := RD'^j \& \operatorname{RShift}(D^{T_j}).$ 11: if $RU_p^j = 1$ or $RM_p^j = 1$ then 12:13:report a match on position j - p + 1.

meaning of $I^{j}(m_{r,c})$ is that $I^{j}(m_{r,c}) = 1$ if there exists a swap permutation π such that $\pi(c) = c + r$ and $\pi(P)_{[1,c]} = T_{[j-c+1,j]}$. Otherwise $I^{j}(m_{r,c})$ is 0.

We want to represent all prefix match indicators in vectors so that we can manipulate them easily. We can do this by mapping the values of I for rows $r \in \{-1, 0, 1\}$ of the *P*-graph to vectors RU, RM, and RD, respectively. We denote value of the vector $RX \in \{RU, RM, RD\}$ in *j*-th step as RX^j . We define values of the vectors as $RU_i^j = I^j(m_{-1,i}), RM_i^j = I^j(m_{0,i}), \text{ and } RD_i^j = I^j(m_{1,i}),$ where the value of $I^j(v) = 0$ for every $v \notin V$.

We define BMA propagate signal along the edges operation as setting the signal of $m_{r,c}$ to 1 if at least one of its predecessors have signal set to 1. I.e., we set $I^{j+1}(m_{-1,i}) := I^j(m_{1,i-1}), I^{j+1}(m_{0,i}) := I^j(m_{-1,i-1}) | I^j(m_{0,i-1}), I^{j+1}(m_{0,1}) := 1, I^{j+1}(m_{1,i}) := I^j(m_{-1,i-1}) | I^j(m_{0,i-1}), \text{ and } I^{j+1}(m_{1,1}) := 1$. We can perform the above operation using the LSO(R) operation. We obtain the propagate signal along the edges operation in the form $RU'^{j+1} := \text{LSO}(RD^j), RM'^{j+1} := \text{LSO}(RM^j | RU^j)$, and $RD'^{j+1} := \text{LSO}(RM^j | RU^j)$.

The operation filter signal by a symbol can be done by first constructing a bit vector D^x for each $x \in \Sigma$ as $D_i^x = 1$ if $x = P_i$ and $D_i^x = 0$ otherwise. Then we use these vectors to filter signal by a symbol x by taking $RU^j := RU'^j \& LShift(D^{T_j})$, $RM^j := RM'^j \& D^{T_j}$, and $RD^j := RD'^j \& RShift(D^{T_j})$.

The last operation we define is the match detection. We do this by checking whether $RU_p^j = 1$ or $RM_p^j = 1$ and if this is the case, then a match starting at a position j - p + 1 occurred.

The final GSM algorithm (Algorithm 3) first prepares the D-masks D^x for every $x \in \Sigma$ and initializes $RU^0 := RM^0 := RD^0 := 0$ (Steps 1–4). Then the algorithm computes the value of vectors RU^j , RM^j , and RD^j for $j \in \{1, \ldots, t\}$ by first using the above formula for signal propagation (Steps 6–8) and then the formula for signal filtering (Steps 9–11) and checks whether $RU_p^j = 1$ or $RM_p^j = 1$ and if this is the case the algorithm reports a match (Steps 12 and 13).

Observe that Algorithm 3 accesses the input sequentially and thus it is a streaming algorithm. We now prove correctness of our algorithm. To ease the notation let us define $R^{j}(m_{r,c})$ to be RU_{c}^{j} if r = -1, RM_{c}^{j} if r = 0, and RD_{c}^{j} if r = 1. We define $R^{\prime j}(m_{r,c})$ analogously. Similarly, we define $D^{x}(m_{r,c})$ as $(\text{LShift}(D^{x}))_{c} = D_{c-1}^{x}$ if r = -1, D_{c}^{x} if r = 0, and $(\text{RShift}(D^{x}))_{c} = D_{c+1}^{x}$ if r = 1. By the way the masks D^{x} are computed on lines 2–4 of Algorithm 3, we get the following observation.

Observation 1 For every $m_{r,i} \in V$ and every $j \in \{i, \ldots, t\}$ we have $D^{T_j}(m_{r,i}) = 1$ if and only if $T_j = P_{r+i}$.

The following lemma constitutes the crucial part of the correctness proof.

Lemma 1. For every $m_{r,i} \in V$ and every $j \in \{i, \ldots t\}$ we have $R^j(m_{r,i}) = 1$ if and only if there exists a swap permutation π such that $\pi(P)_{[1,i]} = T_{[j-i+1,j]}$ and $\pi(i) = i + r$.

Proof. Let us start with the "if" part. We prove the claim by induction on *i*. If i = 1 and there is a swap permutation π such that $\pi(1) = 1 + r$ and $P_{1+r} = T_j$, then the algorithm sets $R'^j(m_{r,1})$ to 1 on line 6, 7, or 8 (recall the definition of LSO). As $P_{1+r} = T_j$, we have $D^{T_j}(m_{r,1}) = 1$ by Observation 1 and, therefore, by lines 9–11, also $R^j(m_{r,1})$.

Now assume that i > 1 and that the claim is true for every smaller *i*. Assume that there exists a swap permutation π such that $\pi(P)_{[1,i]} = T_{[j-i+1,j]}$ and $\pi(i) = i + r$. By induction hypothesis we have that $R^{j-1}(m_{r',i-1}) = 1$, where $r' = i-1-\pi(i-1)$. Since *r* equals -1 if and only if r' equals +1 by Definition 1, we have $(r,r') \in \{(-1,1), (0,-1), (0,0), (1,-1), (1,0)\}$. Therefore the algorithm sets $R^{j}(m_{r,i})$ to 1 on line 6, 7, or 8. Moreover, since $P_{i+r} = T_j$, we have $D^{T_j}(m_{r,i}) = 1$ by Observation 1 and the algorithm sets $R^j(m_{r,i})$ to 1 on one of the lines 9–11.

Now we prove the "only if" part again by induction on i. If i = 1 and $R^{j}(m_{r,i}) = 1$, then we must have $D^{T_{j}}(m_{r,1}) = 1$ and, by Observation 1, also $P_{1+r} = T_{j}$. We obtain π by setting $\pi(1) = 1 + r$, $\pi(2) = 2 - r$ and $\pi(i') = i'$ for every $i' \in \{2, \ldots, p\}$. It is easy to verify that this is a swap permutation for P and has the desired properties.

Now assume that i > 1 and that the claim is true for every smaller i. Assume that $R^j(m_{r,i}) = 1$. Then, due to lines 9–11 we must have $D^{T_j}(m_{r,i}) = 1$ and, hence, by Observation 1, also $P_{i+r} = T_j$. Moreover, we must have $R'^j(m_{r,i}) = 1$ and, hence, by lines 6–8 of the algorithm also $R^{j-1}(m_{r',i-1}) = 1$ for some r' with $(r,r') \in \{(-1,1), (0,-1), (0,0), (1,-1), (1,0)\}$. By induction hypothesis there exists a swap permutation π' for P such that $\pi'(P)_{[1,i-1]} = T_{[j-i+1,j-1]}$ and $\pi'(i-1) = i - 1 + r'$. If $\pi'(i) = i + r$, then setting $\pi = \pi'$ finishes the proof. Otherwise we have either r = 0 or r = 1 and i < p. In the former case we let $\pi(i') = i'$ for every $i' \in \{i, \ldots, p\}$ and in the later case we let $\pi(i) = i + 1$, $\pi(i+1) = i$ and $\pi(i') = i'$ for every $i' \in \{i, \ldots, i-1\}$. It is again easy to verify that π is a swap permutation for P with the desired properties.

A Simple Streaming Bit-parallel Algorithm for Swap Pattern Matching

11

Theorem 2. The GSM algorithm is correct.

Proof. Our GSM algorithm reports a match on position j - p + 1 if and only if $R^{j}(m_{p,-1}) = 1$ or $R^{j}(m_{p,0}) = 1$. However, by Lemma 1, this happens if and only if there is a swap match of P on position j - p + 1 in T. Hence, the algorithm is correct.

Theorem 3. The GSM algorithm runs in $O(\lceil \frac{p}{w} \rceil(|\Sigma| + t) + p)$ time and uses $O(\lceil \frac{p}{w} \rceil |\Sigma|)$ memory cells (not counting the input and output cells), where t is the length of the input text, p length of the input pattern, w is the word-size of the machine, and $|\Sigma|$ size of the alphabet.²

Proof. The initialization of RX and D^x masks (lines 1 and 2) takes $O(\lceil \frac{p}{w} \rceil |\Sigma|)$ time. The bits in D^x masks are set according to the pattern in O(p) time (lines 3 and 4). The main cycle of the algorithm (lines 5–13) makes t iterations. Each iteration consists of computing values of RX in 13 bitwise operations, i.e., in $O(\lceil \frac{p}{w} \rceil)$ machine operations, and checking for the result in O(1) time. This gives $O(\lceil \frac{p}{w} \rceil)(|\Sigma| + t) + p)$ time in total. The algorithm saves 3 RX masks (using the same space for all j and also for RX' masks), $|\Sigma| D^x$ masks, and constant number of variables for other uses (iteration counters, temporary variable, etc.). Thus, in total the GSM algorithm needs $O(\lceil \frac{p}{w} \rceil |\Sigma|)$ memory cells.

Corollary 1. If p = cw for some constant c, then the GSM algorithm runs in $O(|\Sigma| + p + t)$ time and has $O(|\Sigma|)$ space complexity. Moreover, if $p \le w$, then the GSM algorithm can be implemented using only $7 + |\Sigma|$ memory cells.

Proof. The first part follows directly from Theorem 3. Let us show the second part. We need $|\Sigma|$ cells for all D-masks, 3 cells for R vectors (reusing the space also for R' vectors), one pointer to the text, one iteration counter, one constant for the match check and one temporary variable for the computation of the more complex parts of the algorithm. Alltogether, we need only $7 + |\Sigma|$ memory cells to run the GSM algorithm.

From the space complexity analysis we see that for some sufficiently small alphabets (e.g. DNA sequences) the GSM algorithm can be implemented in practice using solely CPU registers with the exception of text which has to be loaded from the RAM.

4 Limitations of the Finite Deterministic Automata Approach

Many of the string matching problems can be solved by finite automata. The construction of a non-deterministic finite automaton that solves Swap Matching can be done by a simple modification of the P-graph. An alternative approach

² To simplify the analysis, we assume that $\log t < w$, i.e., the iteration counter fits into one memory cell.

Table 1. An example of the construction from proof of Theorem 4 for k = 3.

to solve the Swap Matching would thus be to determinize and execute this automaton. The drawback is that the determinization process may lead to an exponential number of states. We show that in some cases it actually does, contradicting the conjecture of Holub [16], stating that the number of states of this determinized automaton is O(p).

Theorem 4. There is an infinite family F of patterns such that any deterministic finite automaton A_P accepting the language $L_S(P) = \{u\pi(P) \mid u \in \Sigma^*, \pi \text{ is a swap permutation for } P\}$ for $P \in F$ has $2^{\Omega(|P|)}$ states.

Proof. For any integer k we define the pattern $P_k := ac(abc)^k$. Note that the length of P_k is $\Theta(k)$. Suppose that the automaton A_P recognizing language L(P)has s states such that $s < 2^k$. We consider a set of strings T_0, \ldots, T_{2^k-1} where T_i is defined as follows. Let $b_{k-1}^i, b_{k-2}^i \dots b_0^i$ be the binary representation of the number *i*. Let $B_j^i = abc$ if $b_j^i = 0$ and let $B_j^i = bac$ if $b_j^i = 1$. Then, let $T_i := acB_{k-1}^i B_{k-2}^i \dots B_0^i$. See Table 1 for an example. Note that each $T_i, i \in$ $\{0,\ldots,2^k-1\}$ is a swapped version of $P = T_0$. Since $s < 2^k$, there exist $0 \leq i < j \leq 2^k - 1$ such that both T_i and T_j are accepted by the same accepting state q of the automaton A. Let m be the minimum number such that $b_{k-1-m}^{i} \neq b_{k-1-m}^{j}$. Note that $b_{m}^{i} = 0$ and $b_{m}^{j} = 1$. Now we define $T_{i}^{i} = T_{i}(abc)^{(m+1)}$ and $T_{j}^{i} = T_{j}(abc)^{(m+1)}$. Let $X = (T_{i}^{i})_{[3(m+1)+1,3(m+1+k)+2]}$ and $Y = (T'_j)_{[3(m+1)+1,3(m+1+k)+2]}$ be the suffices of the strings T'_i and T'_j both of length 3k+2. Note that X begins with $bc \ldots$ and Y begins with $ac \ldots$ and that block *abc* or *bac* repeats for k times in both. Therefore pattern P swap matches Yand does not swap match X. Since for the last symbol of both T_i and T_j the automaton is in the same state q, the computation for T'_i and T'_i must end in the same state q'. However as X should not be accepted and Y should be accepted we obtain contradiction with the correctness of the automaton A. Hence, we may define the family F as $F = \{P_1, P_2, \dots\}$, concluding the proof.

This proof shows the necessity for specially designed algorithms which solve the Swap Matching. We presented one in the previous section and now we reiterate on the existing algorithms.

A Simple Streaming Bit-parallel Algorithm for Swap Pattern Matching

5 Smalgo Algorithm

In this section we discuss how SMALGO by Iliopoulos and Rahman [17] and SMALGO-I and SMALGO-II by Ahmed et al. [1] work. Since SMALGO-I is bitwise inverse of SMALGO, we will introduce them both in terms of operations used in SMALGO-I. After that we will describe and analyze SMALGO-II.

Before we show how these algorithms work, we need one more definition.

Definition 6. A degenerate symbol w over an alphabet Σ is a nonempty set of symbols from alphabet Σ . A degenerate string S is a string built over an alphabet of degenerate symbols. We say that a degenerate string \tilde{P} matches a text T at a position j if $T_{j+i-1} \in \tilde{P}_i$ for every $1 \le i \le p$.

5.1 Smalgo-I

The SMALGO-I [1] algorithm is a modification of the Shift-And algorithm from Section 2.4 for Swap Matching. The algorithm uses the graph theoretic model introduced in Section 2.2.

First let $\tilde{P} = \{P_1, P_2\} \dots \{P_{x-1}, P_x, P_{x+1}\} \dots \{P_{p-1}, P_p\}$ be a degenerate version of pattern P. The symbol on position i in \tilde{P} represents the set of symbols of P which can swap to that position. To accommodate the Shift-And algorithm to match degenerate patterns we need to change the way the D^x masks are defined. For each $x \in \Sigma$ let \tilde{D}_i^x be the bit array of size p such that for $1 \leq i \leq p, \tilde{D}^x = 1$ if and only if $x \in \tilde{P}_i$.

While a match of the degenerate pattern \widetilde{P} is a necessary condition for a swap match of P, it is clearly not sufficient. The way the SMALGO algorithms try to fix this is by introducing P-mask $P(x_1, x_2, x_3)$ which is defined as $P(x_1, x_2, x_3)_i = 1$ if i = 1 or if there exist vertices u_1, u_2 , and u_3 and edges $(u_1, u_2), (u_2, u_3)$ in \mathcal{P}_P for which $u_2 = m_{r,i}$ for some $r \in \{-1, 0, 1\}$ and $\sigma(u_n) = x_n$ for $1 \le n \le 3$, and $P(x_1, x_2, x_3)_i = 0$ otherwise. One P-mask called P(x, x, x) is used to represent the P-masks for triples (x_1, x_2, x_3) which only contain 1 in the first column.

Now, whenever checking whether P prefix swap matches $T \ k + 1$ symbols at position j we check for a match of \tilde{P} in T and we also check whether $P(T_{j+k-1}, T_{j+k}, T_{j+k+1})_{k+1} = 1$. This ensures that the symbols are able to swap to respective positions and that those three symbols of the text T are present in some $\pi(P)$.

With the P-masks completed we initialize $R^1 = 1 \& \widetilde{D}^{T_1}$. Then for every j = 1 to t we repeat the following. We compute R^{j+1} as $R^{j+1} = \text{LSO}(R^j) \& \widetilde{D}^{T_{j+1}} \&$ RShift $(\widetilde{D}^{T_{j+2}}) \& P(T_j, T_{j+1}, T_{j+2})$. To check whether or not a swap match occurred we check whether $R_{p-1}^j = 1$. This is claimed to be sufficient because during the processing we are in fact considering not only the next symbol T_{j+1} but also the symbol T_{j+2} .

5.2 The Flaw in the Smalgo, Smalgo-I and Smalgo-II

We shall see that for a pattern P = abab and a text T = aaba all SMALGO versions give false positives.

14Václav Blažej, Ondřej Suchý, and Tomáš Valla



Fig. 3. SMALGO flaw represented in the *P*-graph for P = abab

The concept of SMALGO is based on the assumption that we can find a path in \mathcal{P}_P by searching for consecutive paths of length 3 (triplets), where each two consecutive share two columns and can partially overlap. However, this only works if the consecutive triplets *actually* share the two vertices in the common columns. If the assumption is not true then the found substring of the text might not match any swapped version of P.

The above input gives such a configuration (see Fig. 3) and therefore the assumption is false. The SMALGO-I algorithm actually reports match of pattern P = abab on a position 1 of text T = aaba. This is obviously a false positive, as the pattern has two b symbols while the text has only one.

The reason behind the false positive match is as follows. The algorithm checks whether the first triplet of symbols (a, a, b) matches. It can match the swap pattern *aabb*. Next it checks the second triplet of symbols (a, b, a), which can match baba. We know that baba is not possible since it did not appear in the previous check, but the algorithm cannot distinguish them since it only checks for triplets existence. Since each step gave us a positive match the algorithm reports a swap match of the pattern in the text.

In the Fig. 3 we see the two triplets which SMALGO assumes have two vertices in common. The SMALGO-II algorithm saves space by maintaining less information, however it simulates how SMALGO-I works and so it contains the same flaw.

The Run of Smalgo-I Resulting in the False Positive 5.3

In Tables 2 and 3 we can see the step by step execution of SMALGO-I algorithm on pattern P = abab and text T = aaba. In Table 3 we see that R^3 has 1 in the 3-rd row which means that the algorithm reports a pattern match on a position 1. This is a false positive, because it is not possible to swap match the pattern with two b symbols in the text with only one b symbol.

Description of Smalgo-II $\mathbf{5.4}$

To explain the SMALGO-II algorithm in more detail, we first introduce a notion of change. An upward change corresponds to (the BMA) going to vertex $m_{-1,i}$ for some i, a downward change corresponds to going to vertex $m_{\pm 1,i}$, and a middle-ward change corresponds to going to vertex $m_{0,i}$.

Table 2. \widetilde{D} -masks and P-masks for P = abab. A column xyz contains values $P(x, y, z)_i$.

i	\widetilde{P}_i	\widetilde{D}_i^a	\widetilde{D}_i^b	aaa	aab	aba	baa	abb	bab	bba	bbb
1	[ab]	1	1	1	1	1	1	1	1	1	1
2	[ba]	1	1	0	1	1	1	1	1	0	0
3	[ab]	1	1	0	1	1	0	1	1	1	0
4	[ba]	1	1	0	0	0	0	0	0	0	0

Table 3. SMALGO-I algorithm execution for P = abab and T = aaba. The column RD^x denotes the values of $RShift(\tilde{D}^x)$.

i	\mathbb{R}^1	$\mathrm{LSO}(\mathbb{R}^1)$	\widetilde{D}^a	RD^{b}	P(a, a, b)	\mathbb{R}^2	$\mathrm{LSO}(\mathbb{R}^2)$	\widetilde{D}^b	RD^a	P(a, b, a)	R^3
1	1	1	1	1	1	1	1	1	1	1	1
2	0	1	1	1	1	1	1	1	1	1	1
3	0	0	1	1	1	0	1	1	1	1	1
4	0	0	1	0	0	0	0	1	0	0	0

If a *downward change* has occurred, then we have to check whether an *upward* change occurs at the next position. If an *upward change* has occurred, then we have to check whether a *downward* or *middle-ward change* occurs at the next position. The main problem here is how to tell whether the changes *actually* occur.

To this end, the authors of the algorithm introduce three new types of masks, namely up-masks $up_{(x,y)}$, down-masks $down_{(x,y)}$, and middle-masks $middle_{(x,y)}$, which express whether an upward, a downward, and a middle-ward change can occur at the particular position, respectively, with the endpoints of the edge having labels x and y.

The authors of the algorithm now claim that to perform the above checks, it is enough to save the previous *down-mask* and match its value with current *upmask* and R_j , or to save the previous *up-mask* and match its value with current *down-mask*, *middle-mask*, and R_j , respectively. However, this way in both cases we only check whether the change can occur, not whether it actually occurred. This would lead not only to false positives (as shown in Section 5), but also to false negatives.

Unfortunately, no more details are available about the algorithm in the original paper. The pseudocode of SMALGO-II (which contains numerous errors) performs something different and we include its analysis in the next section for completeness. Nevertheless, the example presented in the Section 5 and in the previous section still makes the pseudocode (with the small errors corrected) report a false positive.

5.5 Analysis of the Pseudocode of Smalgo-II

In this section we analyze the pseudocode of the SMALGO-II algorithm as given by Ahmed et al. in [1], we will perform equivalent transformation on the pseu-

docode in order to understand the meaning of the checks the pseudocode actually performs.

The original pseudocode is as follows.

\mathbf{Al}	gorithm 4 Smalgo-II
	Require: Text T, up-mask up, down-mask down, middle-mask middle,
	P-mask pmask, D-mask D for given pattern p
1:	$R_0 \leftarrow 2^{\text{patternLength}-1}$
2:	$checkUp \leftarrow checkDown \leftarrow 0$
3:	$R_0 \leftarrow R_0 \& D_{T_0}$
4:	$R_1 \leftarrow R_0 \gg 1$
5:	for $j = 0$ to $(n-2)$ do
6:	$R_j \leftarrow R_j \& pmask_{(T_j, T_{j+1})} \& D_{T_{j+1}}$
7:	$\mathrm{temp} \gets \mathrm{prevCheckUp} \gg 1$
8:	$checkUp \leftarrow checkUp \mid up_{(T_i, T_{i+1})}$
9:	$\operatorname{check}\operatorname{Up} \leftarrow \operatorname{check}\operatorname{Up} \& \sim \operatorname{down}_{(T_j,T_{j+1})} \& \sim \operatorname{middle}_{(T_j,T_{j+1})}$
10:	$prevCheckUp \leftarrow checkUp$
11:	$R_j \leftarrow \sim (\text{temp \& checkUp}) \& R_j$
12:	$\mathrm{temp} \leftarrow \mathrm{prevCheckDown} \gg 1$
13:	$\operatorname{checkDown} \leftarrow \operatorname{checkDown} \mid \operatorname{down}_{(T_j,T_{j+1})}$
14:	$\operatorname{checkDown} \leftarrow \operatorname{checkDown} \& \sim \operatorname{up}_{(T_j, T_{j+1})}$
15:	$\operatorname{prevCheckDown} \leftarrow \operatorname{checkDown}$
16:	$R_j \leftarrow \sim (\text{temp \& checkDown}) \& R_j$
17:	$\mathbf{if} \ (R_j \ \& \ 1) = 1 \ \mathbf{then}$
18:	Match found ending at position $(j-1)$
19:	$R_{j+1} \leftarrow R_j \gg 1$
20:	$checkUp \leftarrow checkUp \gg 1$
21:	$checkDown \leftarrow checkDown \gg 1$

The pseudocode has several problems. First, in the first iteration of the cycle, the algorithm uses the value of the variable prevCheckUp which was never initialized. Second, the algorithm never adds new ones to the variable R and, hence, can never report a match after position patternLength of the text. Third, if the text is of the same length as the pattern, the algorithm only applies the shift patternLength -2 times to the original value of $2^{\text{patternLength-1}}$ (note that in the first iteration it uses R_0 and overwrites the value of R_1) before the last match check. Therefore, at the last check, the value could only drop to $2^{\text{patternLength-1-patternLength+2} = 2^1 = 2$ and the match check cannot be successful. Also the reported position of the match does not make much sense.

Let us first correct all these easy problems.

Algorithm 5 SMALGO-II

1: $R_0 \leftarrow 2^{\text{patternLength}-1}$ 2: prevCheckUp \leftarrow prevCheckDown \leftarrow checkUp \leftarrow checkDown \leftarrow 0 3: $R_0 \leftarrow R_0 \& D_{T_0}$ 4: $R_1 \leftarrow R_0 \gg 1$ 5: for j = 0 to (n - 2) do $R_{j+1} \leftarrow R_{j+1} \& pmask_{(T_j,T_{j+1})} \& D_{T_{j+1}}$ 6: 7: $temp \leftarrow prevCheckUp \gg 1$ $checkUp \leftarrow checkUp \mid up_{(T_i, T_{i+1})}$ 8: checkUp \leftarrow checkUp & $\sim \operatorname{down}_{(T_j, T_{j+1})}$ & $\sim \operatorname{middle}_{(T_j, T_{j+1})}$ 9: 10: $prevCheckUp \gets checkUp$ $R_{j+1} \leftarrow \sim (\text{temp \& checkUp}) \& R_{j+1}$ 11: $\mathrm{temp} \gets \mathrm{prevCheckDown} \gg 1$ 12: $checkDown \leftarrow checkDown \mid down_{(T_i, T_{i+1})}$ 13:checkDown \leftarrow checkDown & $\sim up_{(T_j, T_{j+1})}$ 14: $prevCheckDown \leftarrow checkDown$ 15:16: $R_{j+1} \leftarrow \sim (\text{temp \& checkDown}) \& R_{j+1}$ 17:if $(R_{j+1} \& 1) = 1$ then Match found ending at position (j + 1)18: $R_{j+2} \leftarrow (R_{j+1} \gg 1) \mid 2^{\text{patternLength}-1}$ 19: $checkUp \leftarrow checkUp \gg 1$ 20: $checkDown \leftarrow checkDown \gg 1$ 21:

If we now move the line setting prevCheckUp to checkUp after the line where the check with the temp variable is performed and similarly with prevCheckDown, we do not need the temp variable anymore. We also move the shifts of checkUp and checkDown closer to where this variables are used. We only show the important part of the algorithm.

Algorithm 6 SMALGO-II

5: for j = 0 to (n - 2) do $R_{j+1} \leftarrow R_{j+1} \And pmask_{(T_j,T_{j+1})} \And D_{T_{j+1}}$ 6: $checkUp \leftarrow checkUp \mid up_{(T_i, T_{j+1})}$ 7: checkUp \leftarrow checkUp & $\sim \operatorname{down}_{(T_j, T_{j+1})}$ & $\sim \operatorname{middle}_{(T_j, T_{j+1})}$ 8: 9: $R_{j+1} \leftarrow \sim (\text{prevCheckUp} \gg 1 \& \text{checkUp}) \& R_{j+1}$ 10: $prevCheckUp \leftarrow checkUp$ $checkUp \leftarrow checkUp \gg 1$ 11: checkDown \leftarrow checkDown | down_(T_i,T_{i+1}) 12:checkDown \leftarrow checkDown & $\sim up_{(T_j,T_{j+1})}$ 13:14: $R_{i+1} \leftarrow \sim (\text{prevCheckDown} \gg 1 \& \text{checkDown}) \& R_{i+1}$ $prevCheckDown \leftarrow checkDown$ 15: $checkDown \leftarrow checkDown \gg 1$ 16:17:if $(R_{i+1} \& 1) = 1$ then Match found ending at position (j + 1)18: $R_{j+2} \leftarrow (R_{j+1} \gg 1) \mid 2^{\text{patternLength}-1}$ 19:

Now we swap the order of setting prevCheckUp to checkUp and the shift of checkUp. As this makes prevCheckUp shifted by one, we remove the additional shift in the check. Similarly for checkDown.

Algorithm 7 SMALGO-II

7: checkUp \leftarrow checkUp | up_(T_j,T_{j+1}) 8: checkUp \leftarrow checkUp & $\sim \text{down}_{(T_j,T_{j+1})}$ & $\sim \text{middle}_{(T_j,T_{j+1})}$ 9: $R_{j+1} \leftarrow \sim (\text{prevCheckUp} \& \text{checkUp}) \& R_{j+1}$ 10: checkUp \leftarrow checkUp $\gg 1$ 11: prevCheckUp \leftarrow checkUp 12: checkDown \leftarrow checkDown | down_(T_j,T_{j+1}) 13: checkDown \leftarrow checkDown & $\sim \text{up}_{(T_j,T_{j+1})}$ 14: $R_{j+1} \leftarrow \sim (\text{prevCheckDown} \& \text{checkDown}) \& R_{j+1}$ 15: checkDown \leftarrow checkDown $\gg 1$ 16: prevCheckDown \leftarrow checkDown ...

Now we institute checkUp into the check and move its computation after the check.

Algorithm 8 SMALGO-II

- 6: $R_{j+1} \leftarrow R_{j+1} \& pmask_{(T_j,T_{j+1})} \& D_{T_{j+1}}$
- 7: $R_{j+1} \leftarrow \sim (\operatorname{prevCheckUp} \& (\operatorname{checkUp} | \operatorname{up}_{(T_j, T_{j+1})}) \& \sim \operatorname{down}_{(T_j, T_{j+1})} \& \sim \operatorname{middle}_{(T_j, T_{j+1})}) \& R_{j+1}$
- 8: checkUp \leftarrow (checkUp | up_(T_j,T_{j+1})) & \sim down_(T_j,T_{j+1}) & \sim middle_(T_j,T_{j+1})
- 9: checkUp \leftarrow checkUp $\gg 1$
- 10: prevCheckUp \leftarrow checkUp
- 11: $R_{j+1} \leftarrow \sim (\text{prevCheckDown & (checkDown | down_{(T_j,T_{j+1})}) & \sim up_{(T_j,T_{j+1})}) & R_{j+1}$
- 12: checkDown \leftarrow (checkDown | down_(T_j,T_{j+1})) & \sim up_(T_j,T_{j+1})
- 13: checkDown \leftarrow checkDown $\gg 1$
- 14: prevCheckDown \leftarrow checkDown

```
...
```

Now note that during the check, the content of prevCheckUp is exactly the same as the content of checkUp, so we can remove prevCheckUp completely.

Algorithm 9 SMALGO-II

1: $R_0 \leftarrow 2^{\text{patternLength}-1}$ 2: checkUp \leftarrow checkDown $\leftarrow 0$ 3: $R_0 \leftarrow R_0 \& D_{T_0}$ 4: $R_1 \leftarrow R_0 \gg 1$ 5: for j = 0 to (n - 2) do $R_{j+1} \leftarrow R_{j+1} \& pmask_{(T_j,T_{j+1})} \& D_{T_{j+1}}$ 6: $R_{j+1} \leftarrow \sim (\operatorname{checkUp} \& (\operatorname{checkUp} \mid \operatorname{up}_{(T_i, T_{i+1})}) \& \sim \operatorname{down}_{(T_j, T_{j+1})} \& \sim$ 7: $\operatorname{middle}_{(T_i,T_{i+1})}) \& R_{j+1}$ $\operatorname{checkUp} \leftarrow (\operatorname{checkUp} \mid \operatorname{up}_{(T_i, T_{i+1})}) \& \sim \operatorname{down}_{(T_j, T_{j+1})} \& \sim \operatorname{middle}_{(T_j, T_{j+1})}$ 8: $checkUp \gets checkUp \gg 1$ 9: $R_{j+1} \leftarrow \sim (\text{checkDown \& (checkDown | down_{(T_i, T_{i+1})}) \& \sim up_{(T_i, T_{i+1})}) \&$ 10: R_{j+1} checkDown \leftarrow (checkDown | down_(T_i,T_{j+1})) & ~up_(T_i,T_{j+1}) 11: 12: $\mathrm{checkDown} \gets \mathrm{checkDown} \gg 1$ if $(R_{j+1} \& 1) = 1$ then 13:14:Match found ending at position (j + 1) $R_{j+2} \leftarrow (R_{j+1} \gg 1) \mid 2^{\text{patternLength}-1}$ 15:

Now we modify the expressions by laws of logic to arrive at the following formulation.

Algorithm 10 SMALGO-II

 $\begin{array}{l} 7: \ R_{j+1} \leftarrow R_{j+1} \& \left(\ \sim \operatorname{checkUp} \mid \operatorname{down}_{(T_j,T_{j+1})} \mid \operatorname{middle}_{(T_j,T_{j+1})} \right) \\ 8: \ \operatorname{checkUp} \leftarrow \left(\operatorname{checkUp} \& \ \sim \operatorname{down}_{(T_j,T_{j+1})} \& \ \sim \operatorname{middle}_{(T_j,T_{j+1})} \right) \\ & \sim \operatorname{down}_{(T_j,T_{j+1})} \& \ \sim \operatorname{middle}_{(T_j,T_{j+1})} \right) \\ 9: \ \operatorname{checkUp} \leftarrow \operatorname{checkUp} \gg 1 \\ 10: \ R_{j+1} \leftarrow R_{j+1} \& \left(\ \sim \operatorname{checkDown} \mid \operatorname{up}_{(T_j,T_{j+1})} \right) \\ 11: \ \operatorname{checkDown} \leftarrow \left(\operatorname{checkDown} \& \ \sim \operatorname{up}_{(T_j,T_{j+1})} \right) \mid \left(\operatorname{down}_{(T_j,T_{j+1})} \& \ \sim \operatorname{up}_{(T_j,T_{j+1})} \right) \\ 12: \ \operatorname{checkDown} \leftarrow \operatorname{checkDown} \gg 1 \\ \dots \end{array}$

Now, if the first subexpression in the logical OR setting the new value of checkUp is true, then the appropriate bit of R_{j+1} was just set to 0 on the previous line and filtrating this bit again in future is useless. Hence, we can omit this part of the expression. We arrive at the following resulting pseudocode.

Algorithm 11 SMALGO-II

1: $R_0 \leftarrow 2^{\text{patternLength}-1}$ 2: checkUp \leftarrow checkDown $\leftarrow 0$ 3: $R_0 \leftarrow R_0 \& D_{T_0}$ 4: $R_1 \leftarrow R_0 \gg 1$ 5: for j = 0 to (n - 2) do $R_{j+1} \leftarrow R_{j+1} \& pmask_{(T_j,T_{j+1})} \& D_{T_{j+1}}$ 6: $R_{j+1} \leftarrow R_{j+1} \& (\sim \operatorname{checkUp} | \operatorname{down}_{(T_i, T_{j+1})} | \operatorname{middle}_{(T_j, T_{j+1})})$ 7: $\operatorname{checkUp} \leftarrow \operatorname{up}_{(T_j, T_{j+1})} \& \ \sim \operatorname{down}_{(T_j, T_{j+1})} \& \ \sim \operatorname{middle}_{(T_j, T_{j+1})}$ 8: 9: $checkUp \leftarrow checkUp \gg 1$ $R_{j+1} \leftarrow R_{j+1} \& (\sim \text{checkDown} \mid \text{up}_{(T_j, T_{j+1})})$ 10: checkDown $\leftarrow \operatorname{down}_{(T_i, T_{i+1})} \& \sim \operatorname{up}_{(T_i, T_{i+1})}$ 11: $\mathrm{checkDown} \gets \mathrm{checkDown} \gg 1$ 12:13:if $(R_{j+1} \& 1) = 1$ then Match found ending at position (j + 1)14: $R_{j+2} \leftarrow (R_{j+1} \gg 1) \mid 2^{\text{patternLength}-1}$ 15:

Now it is easy to see, that checkUp stores the information on whether an *upward-change* must have occurred in the previous step (provided that there was a prefix match) and this is compared with the information whether *downward-change* or *middle-change* can occur. Similarly for the *downward-change*. This is not sufficient to avoid false positives since sometimes both *upward-change* and *downward-change* can occur (e.g., as in our counterexample), in which case no filtration is performed at all.

5.6 Why the Flaw is Not Easily Repairable

Consider the following attempt to fix the SMALGO-I or SMALGO-II. After each reported match we check for the validity of the result using a single linear-time algorithm. This approach would rule out false positives but it ruins the time complexity of the algorithms, since there are texts of arbitrary length t with $\Theta(t)$ of reported occurrences.

Namely consider the text $T = aa(baa)^n$ for some positive n, pattern P = abab, and let t = |T|. Note that $n = (t-2)/3 = \Theta(t)$. Text T contains string aabaon positions $1, 4, 7, \ldots, 3(n-1) + 1$ (n occurrences in total) and string baab on positions $3, 6, \ldots, 3(n-1)$ (n-1 occurrences in total). Thus there are 2n-1occurrences which need to be checked since the n occurrences of aaba are reported by the algorithms although they are not valid matches. Even if the checking for correctness was done in linear time O(p), the algorithms will report up to $\Theta(t)$ occurrences which means we have to run the checking algorithm $\Theta(t)$ times. Therefore the time complexity of a version of the SMALGO algorithm corrected this way is O(tp) even for a pattern length similar to the word-size of the target machine.

Also, the flaw cannot be resolved by checking for subpaths of length 4 or any larger constant, due to the following. Consider a pattern $P = (ab)^n$ and a text

21

 $T = aa(ba)^{n-1}$ for any positive *n*. Obviously *P* does not swap match *T*, as they are of the same length 2*n*, but *T* contains more *a*'s than *P*. However, there is a swap permutation π for *P* such that $(\pi(P))_{[1...(2n-1)]} = T_{[1...(2n-1)]}$ and also a swap permutation π' for *P* such that $(\pi'(P))_{[2...(2n)]} = T_{[2...(2n)]}$. For example if we have P = abab and a text T = aabaabaabaa both SMALGO algorithms report swap matches on positions $\{1, 3, 4, 6, 7\}$ while the correct output would be $\{3, 6\}$.

6 Experiments

We implemented our Algorithm 3 (GSM), described in Section 3.2, the Bitwise Parallel Cross Sampling (BPCS) algorithm by Cantone and Faro [10], the Bitwise Parallel Backward Cross Sampling (BPBCS) algorithm by Campanelli et al. [9], and the faulty SMALGO algorithm by Iliopoulos and Rahman [17]. All these implementations are available online.³

We tested the implementations on three real-world datasets. The first dataset (CH) is the 7th chromosome of the human genome⁴ which consists of 159 M characters from the standard ACTG nucleobases and N as for non-determined. Second dataset (HS) is a partial genome of Homo sapiens from the Protein Corpus⁵ with 3.3 M characters representing proteins encoded in 19 different symbols. The last dataset (BIB) is the Bible text of the Cantenbury Corpus⁶ with 4.0 M characters containing 62 different symbols. For each length from 3, 4, 5, 6, 8, 9, 10, 12, 16, and 32 we randomly selected 10,000 patterns from each text and processed each of them with each implemented algorithm.

All measured algorithms were implemented in C++ and compiled with -O3 in gcc 6.3.0. Measurements were carried on an Intel Core i7-4700HQ processor with 2.4 GHz base frequency and 3.4 GHz turbo with 8 GiB of DDR3 memory at 1.6 GHz. Time was measured using std::chrono::high_resolution_clock::now() from the C++ chrono library. The resulting running times, shown in Table 4, were averaged over the 10,000 patterns of the given length.

The results show, that the GSM algorithm runs approximately 23% faster than SMALGO (ignoring the fact that SMALGO is faulty by design). Also, the performance of GSM and BPCS is almost indistinguishable and according to our experiments, it varies in the span of units of percents depending on the exact CPU, cache, RAM and compiler setting. The seemingly superior average performance of BPBCS is caused by the heuristics BPBCS uses; however, while the worst-case performance of GSM is guaranteed, the performance of BPBCS for certain patterns is worse than that of GSM. Also note that GSM is a streaming algorithm while the others are not.

Table 5 visualizes the accurateness of SMALGO-I with respect to its flaw by comparing the number of occurrences found by the respective algorithms. The

³ http://users.fit.cvut.cz/blazeva1/gsm.html

⁴ ftp://ftp.ensembl.org/pub/release-90/fasta/homo_sapiens/dna/

⁵ http://www.data-compression.info/Corpora/ProteinCorpus/

⁶ http://corpus.canterbury.ac.nz/descriptions/large/bible.html

Data	Algor	Pattern Length										
(Σ)	Algor.	3	4	5	6	8	9	10	12	16	32	
	SMALGO	426	376	355	350	347	347	344	347	345	345	
CH	BPCS	398	353	335	332	329	329	326	328	329	327	
(5)	BPBCS	824	675	555	472	366	328	297	257	199	112	
	GSM	394	354	338	333	332	331	329	333	331	333	
	SMALGO	4.80	4.73	4.72	4.74	4.70	4.71	4.71	4.71	4.72	4.70	
HS	BPCS	4.43	4.36	4.36	4.36	4.34	4.33	4.34	4.34	4.35	4.34	
(19)	BPBCS	7.16	5.80	4.79	4.05	3.03	2.70	2.44	2.06	1.62	0.95	
	GSM	4.42	4.38	4.41	4.46	4.45	4.45	4.45	4.44	4.53	4.48	
	SMALGO	8.60	8.38	8.29	8.34	8.32	8.33	8.30	8.35	8.35	8.33	
BIB	BPCS	7.53	7.36	7.28	7.29	7.26	7.27	7.26	7.28	7.29	7.25	
(62)	BPBCS	12.43	10.03	8.26	7.03	5.44	4.93	4.52	3.93	3.19	1.88	
	GSM	7.52	7.37	7.31	7.35	7.38	7.40	7.38	7.42	7.44	7.40	

Table 4. Comparison of the running times in milliseconds. Each value is the average over 10,000 patterns randomly selected from the text.

Table 5. Found occurrences across datasets: The value is simply the sum of occurrences over all the patterns.

Algorithm	Dataset								
Algorithm	CH	HS	BIB						
SMALGO	86243500784	51136419	315612770						
rest	84411799892	51034766	315606151						

ratio of false positives to true positives for the SMALGO-I was: CH 2.17%, HS 0.20% and BIB 0.002%.

References

- Ahmed, P., Iliopoulos, C.S., Islam, A.S., Rahman, M.S.: The swap matching problem revisited. Theoretical Computer Science 557, 34–49 (2014)
- Amir, A., Aumann, Y., Landau, G.M., Lewenstein, M., Lewenstein, N.: Pattern matching with swaps. Journal of Algorithms 37(2), 247–266 (2000)
- Amir, A., Cole, R., Hariharan, R., Lewenstein, M., Porat, E.: Overlap matching. Information and Computation 181(1), 57–74 (2003)
- 4. Amir, A., Eisenberg, E., Porat, E.: Swap and mismatch edit distance. Algorithmica 45(1), 109–120 (2006)
- Amir, A., Landau, G.M., Lewenstein, M., Lewenstein, N.: Efficient special cases of pattern matching with swaps. Information Processing Letters 68(3), 125–132 (1998)
- Amir, A., Lewenstein, M., Porat, E.: Approximate swapped matching. Information Processing Letters 83(1), 33–39 (2002)

A Simple Streaming Bit-parallel Algorithm for Swap Pattern Matching

- Antoniou, P., Iliopoulos, C.S., Jayasekera, I., Rahman, M.S.: Implementation of a swap matching algorithm using a graph theoretic model. In: Bioinformatics Research and Development, BIRD 2008, CCIS, vol. 13, pp. 446–455. Springer (2008)
- Blažej, V., Suchý, O., Valla, T.: A simple streaming bit-parallel algorithm for swap pattern matching. In: Mathematical Aspects of Computer and Information Sciences, MACIS 2017. LNCS, vol. 10693, pp. 333–348. Springer (2017)
- Campanelli, M., Cantone, D., Faro, S.: A new algorithm for efficient pattern matching with swaps. In: International Workshop on Combinatorial Algorithms, IWOCA 2009, LNCS, vol. 5874, pp. 230–241. Springer (2009)
- Cantone, D., Faro, S.: Pattern matching with swaps for short patterns in linear time. In: International Conference on Current Trends in Theory and Practice of Computer Science, SOFSEM 2009, LNCS, vol. 5404, pp. 255–266. Springer (2009)
- 11. Charras, C., Lecroq, T.: Handbook of Exact String Matching Algorithms. King's College Publications (2004)
- Chedid, F.: On pattern matching with swaps. In: IEEE/ACS International Conference on Computer Systems and Applications, AICCSA 2013. pp. 1–5. IEEE (2013)
- Dombb, Y., Lipsky, O., Porat, B., Porat, E., Tsur, A.: The approximate swap and mismatch edit distance. Theoretical Computer Science 411(43), 3814–3822 (2010)
- 14. Faro, S.: Swap matching in strings by simulating reactive automata. In: Proceedings of the Prague Stringology Conference 2013, pp. 7–20. CTU in Prague (2013)
- Fredriksson, K., Giaquinta, E.: On a compact encoding of the swap automaton. Information Processing Letters 114(7), 392–396 (2014)
- 16. Holub, J.: Personal communication (2015)
- Iliopoulos, C.S., Rahman, M.S.: A new model to solve the swap matching problem and efficient algorithms for short patterns. In: International Conference on Current Trends in Theory and Practice of Computer Science, SOFSEM 2008, LNCS, vol. 4910, pp. 316–327. Springer (2008)
- 18. Lewin, B.: Genes for SMA: Multum in parvo. Cell 80(1), 1-5 (1995)
- Lipsky, O., Porat, B., Porat, E., Shalom, B.R., Tzur, A.: String matching with up to k swaps and mismatches. Information and Computation 208(9), 1020–1030 (2010)
- Muthukrishnan, S.: New results and open problems related to non-standard stringology. In: Annual Symposium on Combinatorial Pattern Matching, CPM 95, LNCS, vol. 937, pp. 298–317. Springer (1995)
- Wagner, R.A., Lowrance, R.: An extension of the string-to-string correction problem. Journal of the ACM 22(2), 177–183 (1975)