



Fast detection of specific fragments against a set of sequences

Marie-Pierre Béal¹  and Maxime Crochemore¹ 

Univ. Gustave Eiffel, CNRS, LIGM, F-77454 Marne-la-Vallée, France
`{marie-pierre.beal,maxime.crochemore}@univ-eiffel.fr`

Abstract. We design alignment-free techniques for comparing a sequence or word, called a target, against a set of words, called a reference. A target-specific factor of a target T against a reference R is a factor w of a word in T which is not a factor of a word of R and such that any proper factor of w is a factor of a word of R . We first address the computation of the set of target-specific factors of a target T against a reference R , where T and R are finite sets of sequences. The result is the construction of an automaton accepting the set of all considered target-specific factors. The construction algorithm runs in linear time according to the size of $T \cup R$. The second result consists of the design of an algorithm to compute all the occurrences in a single sequence T of its target-specific factors against a reference R . The algorithm runs in real-time on the target sequence, independently of the number of occurrences of target-specific factors.

Keywords: Specific word · Minimal forbidden word · Suffix automaton.

1 Introduction

The goal of this article is to design an alignment-free technique for comparing a sequence or word, called a target, against a set of words, called a reference.

The motivation comes from the analysis of genomic sequences as done for example by Khorsand et al. in [15] in which authors introduce the notion of sample-specific strings. To avoid alignments but to extract interesting elements that differentiate the target from the reference, the chosen specific fragments are minimal forbidden factors, also called minimal absent factors. Target-specific words are factors of the target that are minimal forbidden factors of the reference. These types of factors have already been applied to compare efficiently sequences (see for example [8] and references therein), to build phylogenies of biological molecular sequences using a distance based on absent words (see [7,6],...), to discover remarkable patterns in some genomic sequences (see for example [21]) and to improve pattern matching methods (see [11],...), to quote only a few applications. In bioinformatics target-specific words act as signatures for newly sequenced biological molecules and help find their characteristics.

The notion of minimal absent factors was introduced by Mignosi et al. [18] (see also [2]) in relation to combinatorial aspects of some sequences. It has then

been extended to regular languages in [1], which obviously applies to a finite set of (finite) sequences. The first linear-time computation is described in [12] (see also [10]) and, due to the important role of the notion, the efficient computation of minimal forbidden factors has attracted quite a lot of works (see for example [20] and references therein).

In the article, we continue exploring the approach of target-specific words as done in [15] by introducing new other algorithmic techniques to detect them. See also the more general view on the usefulness of formal languages to analyze several genomes using pangenomics graphs by Bonizzoni et al. in [5].

The results. First, we address the computation of the set of target-specific factors of a target T against a reference R , where T and R are finite sets of sequences. The result is the construction of an automaton accepting the set of all considered target-specific factors. The construction algorithm runs in linear time according to the size of $T \cup R$.

The second result consists of the design of an algorithm to compute all the occurrences in a single sequence T of its target-specific factors against a reference R . The algorithm runs in real-time on the target sequence, independently of the number of occurrences of target-specific factors, after a standard processing of the reference. This improves on the result in [15], where the running time of the main algorithm depends on the number of occurrences of sought factors.

The design of both algorithms uses the notion of suffix links that are used for building efficiently indexing data structures, like suffix trees (see [14]) and DAWGs also called suffix automata (see [3,10]). The links can also be simulated with suffix arrays [17] and their implementations, for example, the FM-index [13]. The algorithm in [15] uses the FMD index by Li [16]. All these data structures can accommodate the sequences and their reverse complements.

Definitions. Let A be a finite alphabet and A^* be the set of the finite words drawn from the alphabet A , including the empty word ε . A *factor* of a word $u \in A^*$ is a word $v \in A^*$ that satisfies $u = wvt$ for some words $w, t \in A^*$. A *proper factor* of a word u is a factor distinct from the whole word. If P is a set of words, we denote by $\text{Fact}(P)$ the set of factors of words in P , and, if P is finite, $\text{size}(P)$ denotes the sum of lengths of the words in P .

A *minimal forbidden word* (also called a minimal absent word) for a given set of words $L \subseteq A^*$ with respect to a given alphabet B containing A is a word of B^* that does not belong to L but that all proper factors do.

Let R, T be two sets of finite words. A *T -specific word with respect to R* is a word u for which: u is a factor of a word of T , u is not a factor of a word in R and any proper factor of u is a factor of a word in R . The set R is called the *reference* and T the *target* of the problem.

Note that a word is a T -specific word with respect to R if and only if it is a minimal forbidden word of $\text{Fact}(R)$ with respect to the alphabet of letters occurring in $R \cup T$ and is also in $\text{Fact}(T)$. As a consequence, the set of T -specific words with respect to R is both prefix-free and suffix-free.

It follows from the definition that the set S of T -specific words with respect to R is:

$$A \text{Fact}(R) \cap \text{Fact}(R)A \cap (A^* - \text{Fact}(R)) \cap \text{Fact}(T),$$

where A is the alphabet of letters of words R and T . It is thus a regular set when R and T are regular, in particular when R and T are finite.

A *finite deterministic automaton* is denoted by $\mathcal{A} = (Q, A, i, F, \delta)$ where A is a finite alphabet, Q is a finite set of states, $i \in Q$ is the unique initial state, $F \subseteq Q$ is the set of final states and δ is the partial function from $Q \times A$ to Q representing the transitions of the automaton. The partial function δ extends to $Q \times A^*$ and a word u is accepted by \mathcal{A} if and only if $\delta(i, u)$ is defined and belongs to F .

2 Background: directed acyclic word graph

In this section, we recall the definition and the construction of the directed acyclic word graph of a finite set of words. This description already appears in [1].

Let $P = \{x_1, x_2, \dots, x_r\}$ be a finite set of words of size r . A linear-time construction of a deterministic finite state automaton recognizing $\text{Fact}(P)$ has been obtained by Blumer *et al.* in [3], [4], see also [19]. Their construction is an extension of the well-known incremental construction of the suffix automaton of a single word (see for instance [9,10]). The words are added one by one to the automaton. In the sequel, we call this algorithm the DAWG algorithm since it outputs a deterministic automaton called a *directed acyclic word graph*. Let us denote by $DAWG(P) = (Q, A, i, Q, \delta)$ this automaton. Let $\text{Suff}(v)$ denote the set of suffixes of a word v and $\text{Suff}(P)$ the union of all $\text{Suff}(v)$ for $v \in P$. The states of $DAWG(P)$ are the equivalence classes of the right invariant equivalence $\equiv_{\text{Suff}(P)}$ defined as follows. If $u, v \in \text{Fact}(P)$,

$$u \equiv_{\text{Suff}(P)} v \text{ iff } \forall i, 1 \leq i \leq r \text{ and } u^{-1}\text{Suff}(x_i) = v^{-1}\text{Suff}(x_i).$$

and there is a transition labeled by a from the class of a word u to the class of ua . The automaton $DAWG(P)$ has a unique initial state, which is the class of the empty word, and all its states are final. Note that the syntactic congruence \sim defining the minimal automaton of the language is

$$u \sim v \text{ iff } \bigcup_{i=1}^r u^{-1}\text{Suff}(x_i) = \bigcup_{i=1}^r v^{-1}\text{Suff}(x_i)$$

and is not the same as the above equivalence. In other words, $DAWG(P)$ is not always a minimal automaton.

The construction of $DAWG(P)$ is performed in time $O(\text{size}(P) \times \log |A|)$. A time complexity of $O(\text{size}(P))$ can be obtained with an implementation of automata with sparse matrices (see [10]).

Example 1. The deterministic acyclic word graph obtained with the DAWG algorithm from $P = \{abbab, abaab\}$ is displayed in Figure 1 where dashed edges represent the suffix links. Note that this deterministic automaton is not minimal since states 3 and 7, 5 and 9, and 6 and 10 can be merged pairwise.

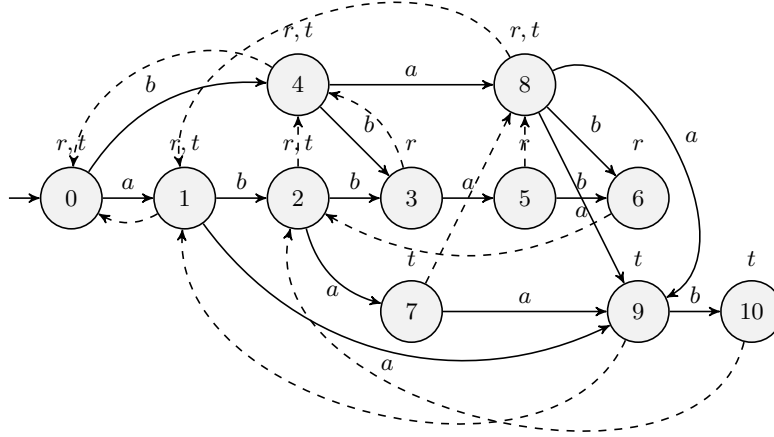


Fig. 1. Automaton $DAWG(P)$ for $P = \{abbab, abaab\}$. Marks r, t above states are defined in Section 3

We let s denote the suffix link function associated with $DAWG(P)$. We first define the function s' from $\text{Fact}(P) \setminus \{\varepsilon\}$ to $\text{Fact}(P)$ as follows: for $v \in \text{Fact}(P) \setminus \{\varepsilon\}$, $s'(v)$ is the longest word $u \in \text{Fact}(P)$ that is a suffix of v and for which $u \not\equiv_{\text{Suff}(P)} v$. Then, if $p = \delta(i, v)$, $s(p)$ is the state $\delta(i, s'(v))$.

3 Computing the set of T -specific words

In this section, we assume that the reference R and the target T are two finite sets of words and our goal is to compute the set of T -specific factors of T against R . To do so, We first compute the directed acyclic word graph $DAWG(R \cup T) = (Q, A, i, Q, \delta)$ of $R \cup T$. Further, we compute a table mark indexed by the set of states Q that satisfies: for each state p in Q , $\text{mark}[p]$ is one of the three values r , t or both r, t according to the fact that each word labeling a path from i to q is a factor of some word in R and not of a word in T , or is a factor of a word in T and not of a word in R , or is a factor of a word in R and of a word in T .

This information can be obtained during the construction of the directed acyclic word graph without increasing the time and space complexity.

The following algorithm outputs a trie (digital tree) of the set of T -specific words with respect to T and R .

```

SPECIFIC-TRIE( $(Q, A, i, Q, \delta)$  DAWG of  $(R \cup T)$ ,  $s$  its suffix link)
1  for each  $p \in Q$  with  $\text{mark}[p] = r, t$  in width-first search from  $i$ 
   and for each  $a \in A$  do
2      if  $(\delta(p, a)$  defined and  $\text{mark}[\delta(p, a)] = t$ ) and  $((p = i)$  or
         $(\delta(s(p), a)$  defined and  $\text{mark}[\delta(s(p), a)] = r$  or  $r, t$ )) then
3           $\delta'(p, a) \leftarrow$  new sink
4      else if  $(\delta(p, a) = q$  with  $\text{mark}[q] = r, t$ )
        and  $(q$  not already reached) then
5           $\delta'(p, a) \leftarrow q$ 
6  return  $\mathcal{A}$ , the automaton  $(Q, A, i, \{\text{sinks}\}, \delta')$ 

```

Example 2. The automaton $\text{DAWG}(R \cup T)$ with the input $R = \{abab\}$, $T = \{abaab\}$ is shown in Figure 1. The output of algorithm SPECIFIC-TRIE on $\text{DAWG}(R \cup T)$ is shown in Figure 2 where the squares are final or sink states of the trie. The set of T -specific words with respect to R is $\{aa, aba\}$.

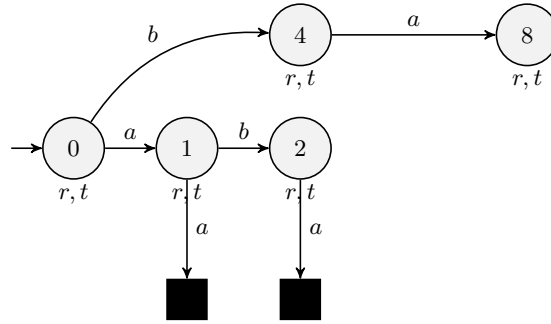


Fig. 2. The trie of T -specific words with respect to R .

Proposition 1. Let $\text{DAWG}(R \cup T)$ be the output of algorithm DAWG on the finite set of words $R \cup T$, let s be its suffix function, and let mark be the table defined as above. Algorithm SPECIFIC-TRIE builds the trie recognizing the set of T -specific words with respect to R .

Proof. Let S be the set of T -specific words with respect to R .

Consider a word ua ($a \in A$) accepted by \mathcal{A} . Note that \mathcal{A} accepts only nonempty words. Let $p = \delta'(i, u)$. Since the DAWG automaton is processed

with a width-first search, u is the shortest word for which $\delta(i, u) = p$. Therefore, if $u = bv$ with $b \in A$, we have $\delta(i, v) = s(p)$ by definition of the suffix function s . When the test “ $(\delta(p, a)$ defined and $\text{mark}[\delta(p, a)] = t$) and $(\delta(s(p), a)$ defined and $\text{mark}[\delta(s(p), a)] = r$ or r, t)” is satisfied, this implies that $va \in \text{Fact}(R)$. Thus, $bva \notin \text{Fact}(R)$, while $bv, va \in \text{Fact}(R)$ and $bva \in \text{Fact}(T)$. So, ua is a T -specific word with respect to R . If u is the empty word, then $p = i$. The transition from i to the sink labeled by a is created under the condition “ $\delta(p, a)$ defined and $\text{mark}[\delta(p, a)] = t$ ”, which means that $a \in \text{Fact}(T)$. The word a is again a T -specific word with respect to R . Thus the words accepted by \mathcal{A} are T -specific words with respect to R .

Conversely, let $ua \in S$. If u is the empty word, this means that a does not occur in $\text{Fact}(R)$ and occurs in $\text{Fact}(T)$ therefore there is a transition labeled by a from i in $\text{DAWG}(R \cup T)$ to a state marked t . Thus a transition from i to a sink state in \mathcal{A} created Line 3 and a is accepted by \mathcal{A} . Now assume that $u = bv$. The word u is in $\text{Fact}(R)$. So let $p = \delta(i, u)$. Note that u is the shortest word for which $p = \delta(i, u)$, because all such words are suffixes of each other in the DAWG automaton. The word ua is not in $\text{Fact}(R)$ and is in $\text{Fact}(T)$, so the condition “ $\delta(p, a)$ defined and $\text{mark}[p, a] = t$ ” is satisfied. Let $q = s(p)$. We have $q = \delta(i, v)$ because of the minimality of the length of u and the definition of s . Since va is in $\text{Fact}(R)$, the condition “ $\delta(s(p), a)$ defined and $\text{mark}[\delta(s(p), a)] = r$ or r, t ” at Line 2 is satisfied which yields the creation of a transition at Line 3 to make \mathcal{A} accept ua as wanted.

A main point in algorithm SPECIFIC-TRIE is that it uses the function s defined on states of the input DAWG. It is not possible to proceed similarly when considering the minimal factor automaton of $\text{Fact}(R \cup T)$ because there is no analogue function s . However, it is possible to reduce the automaton $\text{DAWG}(R \cup T)$ by merging states having the same future (right context) and the same image by s . For example, on the DAWG of Figure 1, states 6 and 10 can be merged because $s(6) = s(10) = 2$. States 3 and 7, nor states 5 and 9 cannot be merged with the same argument.

Proposition 2. *Algorithms DAWG and SPECIFIC-TRIE together run in time $O(\text{size}(R \cup T) \times |A|)$ with input two finite sets of words R, T , if the transition functions are implemented by transition matrices.*

If P is a set of words, we denote by A_P the set of letters occurring in P .

Proposition 3. *Let R, T be two finite sets of words. The number of T -specific words with respect to R is no more than $(2 \text{size}(R) - 2)(|A_R| - 1) + |A_T \setminus A_R| - |A_R| + m$, if $\text{size}(R) > 1$, where m the number of words in R . The bound becomes $|A_T \setminus A_R|$ when $\text{size}(R) \leq 1$.*

Proof. We let S denote the set of T -specific words with respect to R . Since S is included in the set of minimal forbidden words of $\text{Fact}(R)$ with respect to the alphabet $A = A_R \cup A_T$, the bound comes from [1, Corollary 4.1].

4 Computing occurrences of target-specific factors: the T -specific table

In this section, we consider that R and T are just words. The goal of the section is to design an algorithm that computes all the occurrences of T -specific words in T . To do so, we define the T -specific table associated with the pair R, T of words of the problem.

A letter of T at position k is denoted by $T[k]$ and $T[i..j]$ denotes the factor $T[i]T[i+1]\dots T[j]$ of T . Then, the T -specific table Ts is defined, for $i = 0, \dots, |T| - 1$, by

$$Ts[i] = \begin{cases} j, & \text{if } T[i..j] \text{ is } T\text{-specific, } i \leq j, \\ -1, & \text{else.} \end{cases}$$

Note 1. Since the set of T -specific factors is both prefix-free and suffix-free, for each position k on T there is at most one T -specific factor of T starting at k and for each position j on T there is at most one T -specific factor of T ending at j .

Note 2. Instead of computing the T -specific table Ts , in a straightforward way, the algorithm below can be transformed to compute the list of pairs (i, j) of positions on T for which $Ts[i] = j$ and $j \neq -1$.

To compute the table we use \mathcal{R} , the suffix automaton of R , with its transition function δ and equipped with both the suffix link s (used here as a failure link) and the length function ℓ defined on states by: $\ell[p] = \max\{z \in A^* \mid \delta(i, z) = p\}$. Functions s and ℓ transform the automaton into a search machine, see [10, Section 6.6].

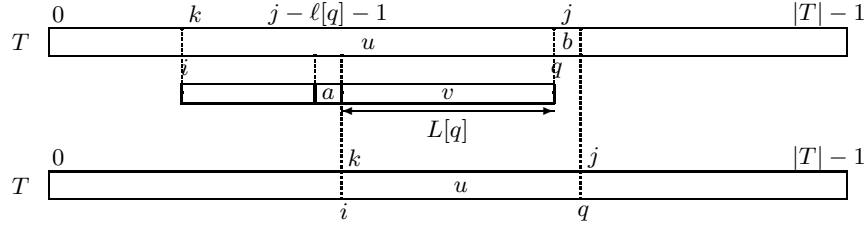


Fig. 3. A T -specific word found: when $u \in \text{Fact}(R)$ and $ub \notin \text{Fact}(R)$, either avb or b is a T -specific factor with respect to R (a, b are letters).

Figure 3 illustrates the principle of Algorithm TsTABLE. Let us assume the factor $u = T[k..j-1]$ is a factor of R but ub is not for some letter b . Then, let v be the longest suffix of u for which vb is a factor of R . If it exists, then clearly avb , with a letter preceding v , is T -specific. Indeed, $av, vb \in \text{Fact}(R)$ and $avb \notin \text{Fact}(R)$, which means that avb is a minimal forbidden word of R while

occurring in T . Therefore, setting $q = \delta(i, u)$, $Ts[j - \ell[q] - 1] = j$ since $\ell[q] = |v|$ due to a property of the DAWG of \mathcal{R} . If there is no suffix of u satisfying the condition, the letter b alone is T -specific and $Ts[j] = j$.

```

TsTABLE( $T$  target word,  $\mathcal{R}$  DAWG( $R$ ),  $i$  initial( $\mathcal{R}$ ))
1  ( $q, j$ )  $\leftarrow (i, 0)$ 
2  while  $j < |T|$  do
3       $Ts[j] \leftarrow -1$ 
4      if  $\delta(q, T[j])$  undefined then
5          while  $q \neq i$  and  $\delta(q, T[j])$  undefined do
6               $q \leftarrow s[q]$ 
7              if  $\delta(q, T[j])$  undefined then  $\triangleright q = i$ 
8                   $Ts[j] \leftarrow j$ 
9                   $j \leftarrow j + 1$ 
10             else  $Ts[j - \ell[q] - 1] \leftarrow j$ 
11             ( $q, j$ )  $\leftarrow (\delta(q, T[j]), j + 1)$ 
12         else ( $q, j$ )  $\leftarrow (\delta(q, T[j]), j + 1)$ 
13 return  $Ts$ 

```

Theorem 1. *The DAWG of the reference set R of words being preprocessed, applied to a word T , Algorithm TsTABLE computes its T -specific table with respect to R and runs in linear time, i.e. $O(|T|)$ on a fixed-size alphabet.*

Proof. The algorithm implements the ideas detailed above. A more formal proof relies on the invariant of the while loop: $q = \delta(i, u)$, where i is the initial state of the suffix automaton of R and $u = T[k..j]$ for a position $k \leq j$. Since $k = j - |u|$, it is left implicit in the algorithm. The length $|u|$ could be computed and then incremented when j is. It is made explicit only at line 10 as $L[q] + 1$ after computing the suffix v of u .

For example, when v exists, u is changed to vb and j is incremented, which maintains the equality.

As for the running time, note that instructions at lines 1 and 7-12 execute in constant time for each value of j . All the executions of the instruction at line 6 execute in time $O(|T|)$ because the link s reduces strictly the potential length of the T -specific word ending at j , that is, it virtually increments the starting position of v in the picture.

Thus the whole execution is done in time $O(|T|)$.

Algorithm TsTABLE can be improved to run in real-time on a fixed-size alphabet. This is done by optimizing the suffix link s defined on the automaton \mathcal{R} . To do so, let us define, for each state q of \mathcal{R} ,

$$Out(q) = \{a \mid \delta(q, a) \text{ defined for letter } a\}.$$

Then, the optimised suffix link G is defined by $G[\text{initial}(\mathcal{R})] = \mathbf{nil}$ and, for any other state q of \mathcal{R} , by

$$G[q] = \begin{cases} s[q], & \text{if } \text{Out}(q) \subset \text{Out}(s[q]), \\ G[s[q]], & \text{else.} \end{cases}$$

Note that, since we always have $\text{Out}(q) \subseteq \text{Out}(s[q])$, the definition of G can be reformulated as

$$G[q] = \begin{cases} s[q], & \text{if } \deg(q) < \deg(s[q]), \\ G[s[q]], & \text{else,} \end{cases}$$

where \deg is the outgoing degree of a state. Therefore, its computation can be realized in linear time with respect to the number of states of \mathcal{R} . After substituting G for s in Algorithm TSTABLE, when the alphabet is of size α the instruction at line 6 executes no more than α times for each value of q . So the time to process a given state q is constant. This is summarized in the next corollary.

Corollary 1. *When using the optimized suffix link, Algorithm TSTABLE runs in real-time on a fixed-size alphabet.*

On a more general alphabet of size α , the processing of a given state of the automaton can be done in time $\log \alpha$.

References

1. M. Béal, M. Crochemore, F. Mignosi, A. Restivo, and M. Sciortino. Computing forbidden words of regular languages. *Fundam. Informaticae*, 56(1-2):121–135, 2003.
2. M. Béal, F. Mignosi, A. Restivo, and M. Sciortino. Forbidden words in symbolic dynamics. *Adv. Appl. Math.*, 25(2):163–193, 2000.
3. A. Blumer, J. Blumer, A. Ehrenfeucht, D. Haussler, and R. M. McConnell. Building the minimal DFA for the set of all subwords of a word on-line in linear time. In J. Paredaens, editor, *Automata, Languages and Programming, 11th Colloquium, Antwerp, Belgium, July 16-20, 1984, Proceedings*, volume 172 of *Lecture Notes in Computer Science*, pages 109–118. Springer, 1984.
4. A. Blumer, J. Blumer, D. Haussler, R. McConnell, and A. Ehrenfeucht. Complete inverted files for efficient text retrieval and analysis. *Journal of the ACM*, 34(3):578–595, 1987.
5. P. Bonizzoni, C. D. Felice, Y. Pirola, R. Rizzi, R. Zaccagnino, and R. Zizza. Can formal languages help pangenomics to represent and analyze multiple genomes? In V. Diekert and M. V. Volkov, editors, *Developments in Language Theory - 26th International Conference, DLT 2022, Tampa, FL, USA, May 9-13, 2022, Proceedings*, volume 13257 of *Lecture Notes in Computer Science*, pages 3–12. Springer, 2022.
6. G. Castiglione, J. Gao, S. Mantaci, and A. Restivo. A new distance based on minimal absent words and applications to biological sequences. *CoRR*, abs/2105.14990, 2021.

7. S. Chairungsee and M. Crochemore. Using minimal absent words to build phylogeny. *Theor. Comput. Sci.*, 450:109–116, 2012.
8. P. Charalampopoulos, M. Crochemore, G. Fici, R. Mercas, and S. P. Pissis. Alignment-free sequence comparison using absent words. *Inf. Comput.*, 262:57–68, 2018.
9. M. Crochemore. Transducers and repetitions. *Theoretical Computer Science*, 45(1):63–86, 1986.
10. M. Crochemore, C. Hancart, and T. Lecroq. *Algorithms on Strings*. Cambridge University Press, 2007. 392 pages.
11. M. Crochemore, A. Héliou, G. Kucherov, L. Mouchard, S. P. Pissis, and Y. Ramusat. Absent words in a sliding window with applications. *Inf. Comput.*, 270, 2020.
12. M. Crochemore, F. Mignosi, and A. Restivo. Automata and forbidden words. *Inf. Process. Lett.*, 67(3):111–117, 1998.
13. P. Ferragina and G. Manzini. Opportunistic data structures with applications. In *41st Annual Symposium on Foundations of Computer Science, FOCS 2000, 12-14 November 2000, Redondo Beach, California, USA*, pages 390–398. IEEE Computer Society, 2000.
14. D. Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computational Biology*. Cambridge University Press, 1997.
15. P. Khorsand, L. Denti, H. G. S. V. Consortium, P. Bonizzoni, R. Chikhi, and F. Hormozdiari. Comparative genome analysis using sample-specific string detection in accurate long reads. *Bioinformatics Advances*, 1(1), 05 2021.
16. H. Li. Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly. *Bioinformatics*, 28(14):1838–1844, 05 2012.
17. U. Manber and E. W. Myers. Suffix arrays: A new method for on-line string searches. *SIAM J. Comput.*, 22(5):935–948, 1993.
18. F. Mignosi, A. Restivo, and M. Sciortino. Forbidden factors in finite and infinite words. In J. Karhumäki, H. A. Maurer, G. Paun, and G. Rozenberg, editors, *Jewels are Forever, Contributions on Theoretical Computer Science in Honor of Arto Salomaa*, pages 339–350. Springer, 1999.
19. G. Navarro and M. Raffinot. *Flexible pattern matching in strings—practical on-line search algorithms for texts and biological sequences*. Cambridge University Press, 2002. 232 pages.
20. A. J. Pinho, P. J. S. G. Ferreira, S. P. Garcia, and J. M. O. S. Rodrigues. On finding minimal absent words. *BMC Bioinform.*, 10, 2009.
21. R. M. Silva, D. Pratas, L. Castro, A. J. Pinho, and P. J. S. G. Ferreira. Three minimal sequences found in ebola virus genomes and absent from human DNA. *Bioinform.*, 31(15):2421–2425, 2015.