



# Approximate cover of strings <sup>☆</sup>

A. Amir <sup>a,b</sup>, A. Levy <sup>c,\*</sup>, R. Lubin <sup>a</sup>, E. Porat <sup>a</sup>

<sup>a</sup> Department of Computer Science, Bar-Ilan University, Israel

<sup>b</sup> Johns Hopkins University, United States of America

<sup>c</sup> Department of Software Engineering, Shenkar College, 12 Anna Frank, Ramat-Gan, Israel



## ARTICLE INFO

### Article history:

Received 1 November 2017

Received in revised form 8 May 2019

Accepted 8 May 2019

Available online 28 May 2019

Communicated by P. Spirakis

### Keywords:

Periodicity

Quasi-periodicity

Cover

Approximate cover

## ABSTRACT

Regularities in strings arise in various areas of science, including coding and automata theory, formal language theory, combinatorics, molecular biology and many others. A common notion to describe regularity in a string  $T$  is a *cover*, which is a string  $C$  for which every letter of  $T$  lies within some occurrence of  $C$ . The alignment of the cover repetitions in the given text is called a *tiling*. In many applications finding exact repetitions is not sufficient, due to the presence of errors. In this paper, we use a new approach for handling errors in coverable phenomena and define the *approximate cover problem (ACP)*, in which we are given a text that is a sequence of some cover repetitions with possible mismatch errors, and we seek a string that covers the text with the minimum number of errors. We first show that the ACP is  $\mathcal{NP}$ -hard, by studying the *cover-length relaxation of the ACP*, in which the requested length of the approximate cover is also given with the input string. We show that this relaxation is already  $\mathcal{NP}$ -hard. We also study another two relaxations of the ACP, which we call the *partial-tiling relaxation of the ACP* and the *full-tiling relaxation of the ACP*, in which a tiling of the requested cover is also given with the input string. A given full tiling retains all the occurrences of the cover before the errors, while in a partial tiling there can be additional occurrences of the cover that are not marked by the tiling. We show that the partial-tiling relaxation has a polynomial time complexity and give experimental evidence that the full-tiling also has polynomial time complexity. The study of these relaxations, besides shedding another light on the complexity of the ACP, also involves a deep understanding of the properties of covers, yielding some key lemmas and observations that may be helpful for a future study of regularities in the presence of errors.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

Regularities in strings arise in various areas of science, including coding and automata theory, formal language theory, combinatorics, molecular biology and many others. A typical form of regularity is *periodicity*, meaning that a “long” string  $T$  can be represented as a concatenation of copies of a “short” string  $P$ , possibly ending in a prefix of  $P$ . Periodicity has been extensively studied in Computer Science over the years (see [31]).

<sup>☆</sup> A partial version of this paper appeared in the proceedings of CPM 2017.

\* Corresponding author.

E-mail addresses: amir@cs.biu.ac.il (A. Amir), avivitlevy@shenkar.ac.il (A. Levy), porately@cs.biu.ac.il (E. Porat).

For many phenomena, it is desirable to broaden the definition of periodicity and study wider classes of repetitive patterns in strings. One common such notion is that of a *cover*, defined as follows.

**Definition 1. [Cover]** A length  $m$  substring  $C$  of a string  $T$  of length  $n$ , is said to be a *cover* of  $T$ , if  $n > m$  and every letter of  $T$  lies within some occurrence of  $C$ .

Note that by the definition of cover, the string  $C$  is both a prefix and a suffix of the string  $T$ . For example, consider the string  $T = abaabababab$ . Clearly,  $T$  is “almost” periodic with period  $aba$ , however, as it is not completely periodic, the algorithms that exploit repetitions cannot be applied to it. On the other hand, the string  $C = aba$  is a cover of  $T$ , which allows applying to  $T$  cover-based algorithms. In this paper we study coverable phenomena in the presence of errors.

There are related regularity types and several approaches to handle errors in regularities. Quasi-periodicity was introduced by Ehrenfeucht in 1990 (according to [6]). The earliest paper in which it was studied is by Apostolico, Farach and Iliopoulos [8], which defined the *quasi-period* of a string to be the length of its shortest cover and presented an algorithm for computing the quasi-period of a given string in  $O(n)$  time and space. The new notion attracted immediately several groups of researchers (e.g. [10], [32,33], [30], [11]). An overview on the first decade of the research on covers can be found in the surveys [6,21,35].

While covers are a significant generalization of the notion of periods as formalizing regularities in strings, they are still restrictive, in the sense that it remains unlikely that an arbitrary string has a cover shorter than the word itself. Due to this reason, different variants of quasi-periodicity have been introduced. These include *seeds* [22], *maximal quasi-periodic substring* [7], the notion of  $k$ -covers [23],  $\lambda$ -cover [36], *enhanced covers* [18], *partial cover* [24] and *partial seed* [26]. Other recently explored directions include the inverse problem for cover arrays [16], extensions to strings in which not all letters are uniquely defined, such as *indeterminate strings* [5,15] or *weighted sequences* [37,9]. Some of the related problems are  $\mathcal{NP}$ -hard (see e.g., [5,12,24,25]).

In applications such as molecular biology and computer-assisted music analysis, finding exact repetitions and covers is not always sufficient. A more appropriate notion is that of *approximate repetitions*, where errors are allowed (see, e.g., [14, 17]). This notion was first studied in 1993 by Landau and Schmidt [28,29] who concentrated on approximate tandem repeats. Note that, the natural definition of an approximate repetition is not clear. One possible definition is that the distance between any two adjacent repeats is small. Another possibility is that all repeats lie at a small distance from a single “original”. Such a definition of *approximate seeds* is studied in [12,20,19]. Indeed, all these definitions along with other ones were proposed and studied (see [1,2,27,34]). Yet another possibility is that all repeats must be equal, but we allow a fixed total number of mismatches. The possibility presented in [1] is a global one, assuming that an original unknown string is a sequence of repeats without errors, but the process of sequence creation or transmission incurs errors to the sequence of repeats, and, thus, the examined input string is not a sequence of repeats. Therefore, a (smallest) repeat generating a string with the minimum total number of mismatches with the input string is sought. Extension of this definition approach to approximate covers is the topic of this paper.

### 1.1. Our results

In this paper we extend the approach of [1] to the notion of covers and define the *approximate cover problem (ACP)*, in which we are given a text that is a sequence of some cover repetitions with possible mismatch errors, and we seek a string that covers the text with the minimum number of errors. The alignment of the cover repetitions in the given text is called a *tiling*. We prove that the ACP is  $\mathcal{NP}$ -hard by studying a relaxation of this problem, which we call *the cover-length relaxation of the ACP*. In this relaxation the requested length of the approximate cover is also given with the input string. We prove that this relaxation is already  $\mathcal{NP}$ -hard, while proving the  $\mathcal{NP}$ -hardness of ACP.

We also study another two relaxations of the problem, which we call *the partial-tiling relaxation of the ACP* and *the full-tiling relaxation of the ACP*. In these relaxations a tiling of the requested cover is also given, and we seek a string such that when using the given tiling to align it with the given text, the number of mismatches is minimized. The full tiling retains all the occurrences of the cover before the errors, while in the partial tiling there can be additional occurrences of the cover that are not marked by the tiling. We examine these relaxations and show that the partial-tiling has polynomial time complexity and give experimental evidence that the full-tiling also has polynomial time complexity. The study of these relaxations, besides shedding another light on the complexity of the ACP, also involves a deep understanding of the properties of covers and seeds, yielding some key lemmas and observations (such as [3]) that may be helpful for a future study of regularities in the presence of errors. Table 1 summarizes the current knowledge of the complexity of ACP and its relaxations.

**Paper Contributions.** The main contributions of this paper are:

- Proving that the ACP is  $\mathcal{NP}$ -hard.
- Formalizing the partial-tiling relaxation of the ACP and proving it is polynomial time computable.
- Formalizing the full-tiling relaxation of the ACP and suggesting a polynomial time algorithm for its computation, while giving an experimental evidence for the correctness of this algorithm.

**Table 1**

A summary the current knowledge of the complexity of ACP and its relaxations.

The problem	Complexity	Reference
ACP	$\mathcal{NP}$ -hard	This paper
cover-length relaxation	$\mathcal{NP}$ -hard	This paper
full-tiling relaxation	$\mathcal{P}$ (conjectured)	This paper
partial-tiling relaxation	$\mathcal{P}$	This paper
candidate relaxation	$\mathcal{P}$	[4]

The paper is organized as follows. In Section 2, we give formal definitions. In Section 3, we study the cover-length relaxation of the ACP and prove the  $\mathcal{NP}$ -hardness of the ACP. In Section 4, we study the partial-tiling relaxation of the ACP and show it is polynomial-time computable. In Section 5, we study the full-tiling relaxation of the ACP, suggest a polynomial-time algorithm for this problem and experimentally test its correctness. We conclude with some open problems in Section 6.

## 2. Preliminaries

In this section we give the needed formal definitions.

**Definition 2. [Tiling]** Let  $T$  be a string over alphabet  $\Sigma$  such that the string  $C$  over alphabet  $\Sigma$  is a cover of  $T$ . Then, the sorted list of indices representing the start positions of occurrences of the cover  $C$  in the text  $T$  is called the *tiling* of  $C$  in  $T$ .

In this paper we have a text  $T$  which may have been introduced to errors and, therefore, is not coverable. However, we would like to refer to a retained tiling of an unknown string  $C$  in  $T$  although  $C$  does not cover  $T$  because of mismatch positions. The following definition makes a distinction between a list of indices that may be assumed to be a tiling of the text before mismatch errors occurred and a list of indices that cannot be such a tiling.

**Definition 3. [A Valid Tiling]** Let  $T$  be an  $n$ -length string over alphabet  $\Sigma$  and let  $L$  be a sorted list of indices  $L \subset \{1, \dots, n\}$ . Let  $m = n + 1 - L_{last}$ , where  $L_{last}$  is the last index in  $L$ . Then,  $L$  is called a *valid tiling* of  $T$ , if  $i_1 = 1$  and for every  $i_k, i_{k+1} \in L$ , it holds that  $i_{k+1} - i_k \leq m$ .

**Notation 1.** Let  $C$  be an  $m$  length string over alphabet  $\Sigma$ . Denote by  $S(C)$  a string of length  $n$ ,  $n > m$ , such that  $C$  is a cover of  $S(C)$ .

Note that  $S(C)$  is not uniquely defined even for a fixed  $n > m$ , since every different valid tiling of the  $m$ -length string  $C$  generates a different  $n$ -length string  $S(C)$ . A unique version can be obtained if a valid tiling  $L$  is also given.

**Notation 2.** Let  $T$  be an  $n$ -length string over alphabet  $\Sigma$  and let  $L$  be a valid tiling of  $T$ . Let  $m = n + 1 - L_{last}$ , where  $L_{last}$  is the last index in the tiling  $L$ . For any  $m$ -length string  $C'$ , let  $S_L(C')$  be the  $n$ -length string obtained using  $C'$  as a cover and  $L$  as the tiling as follows:  $S_L(C')$  begins with a copy of  $C'$  and for each index  $i$  in  $L$  a new copy of  $C'$  is concatenated starting from index  $i$  of  $S_L(C')$  (maybe running over a suffix of the last copy of  $C'$ ).

**Definition 4.** Let  $T$  be a string of length  $n$  over alphabet  $\Sigma$ . Let  $H$  be the Hamming distance. The *distance of  $T$  from being covered* is:

$$dist = \min_{C \in \Sigma^*, |C| < n, S(C) \in \Sigma^n} H(S(C), T).$$

We will also refer to  $dist$  as *the number of errors in  $T$* .

**Definition 5.** Let  $T$  be an  $n$ -long string over alphabet  $\Sigma$ . An  $m$ -long string  $C$  over  $\Sigma$ ,  $m \in \mathbb{N}$ ,  $m < n$ , is called an  *$m$ -length approximate cover of  $T$* , if for every string  $C'$  of length  $m$  over  $\Sigma$ ,  $\min_{S(C') \in \Sigma^n} H(S(C'), T) \geq \min_{S(C) \in \Sigma^n} H(S(C), T)$ , where  $H$  is the hamming distance of the given strings.

We refer to  $\min_{S(C) \in \Sigma^n} H(S(C), T)$  as the *number of errors of an  $m$ -length approximate cover of  $T$* .

**Definition 6. [Approximate Cover]** Let  $T$  be a string of length  $n$  over alphabet  $\Sigma$ . A string  $C$  over alphabet  $\Sigma$  is called an *approximate cover of  $T$*  if:

1.  $C$  is an  $m$ -length approximate cover of  $T$  for some  $m \in \mathbb{N}$ ,  $m < n$ , for which

$$\min_{S(C) \in \Sigma^n} H(S(C), T) = \text{dist}.$$

2. for every  $m'$ -length approximate cover of  $T$ ,  $C'$ , s.t.  $\min_{S(C') \in \Sigma^n} H(S(C'), T) = \text{dist}$ , it holds that:  $m' \geq m$ .

**Primitivity.** By definition, an approximate cover  $C$  should be *primitive*, i.e., it cannot be covered by a string other than itself (otherwise,  $T$  has a cover with a smaller length). Note that a periodic string can be covered by a smaller string (not necessarily the period), and therefore, is not primitive.

**Definition 7.** The *Approximate Cover Problem* (ACP) is the following:

*INPUT:* String  $T$  of length  $n$  over alphabet  $\Sigma$ .

*OUTPUT:* An approximate cover of  $T$ ,  $C$ , and the number of errors in  $T$ .

### 3. $\mathcal{NP}$ -hardness of the ACP

In this section we prove the  $\mathcal{NP}$ -hardness of the ACP. To this end, we study a variant of the problem where  $m$ , the length of a requested approximate cover, is also given together with the input string  $T$ , and we are requested to find a string  $C$  of length  $m$  that is an  $m$ -length approximate cover of  $T$ , i.e.,  $C$  covers  $T$  with the minimum number of errors over all strings of length  $m$ . We call this problem *the cover-length relaxation of the ACP*. We prove that the cover-length relaxation of the ACP is already  $\mathcal{NP}$ -hard, and then deduce the hardness of ACP.

Our hardness proof uses a reduction from the 3-SAT problem, in which the input is a logical formula  $\varphi$  on  $N$  variables in 3-CNF (each clause has exactly three literals), and we need to decide whether  $\varphi$  is satisfiable or not. The  $\mathcal{NP}$ -hardness of 3-SAT is well-known (see e.g. [13]).

#### 3.1. The reduction from 3-SAT

Given a 3-CNF formula  $\varphi$  on  $N$  variables,  $x_1, \dots, x_N$ , with  $\ell$  clauses. Assume without loss of generality that the literals in each clause are sorted by the index of their variables. We need to define a text  $T$  of length  $n$  over an alphabet  $\Sigma$  and to specify the length  $m$  of the requested approximate cover. We will then show that  $\varphi$  is satisfiable if and only if  $T$  has an  $m$ -approximate cover with at most some specified number of errors to be defined.

We begin by defining the alphabet  $\Sigma$  to include all the variables and their negation together with four additional dummy variables:  $x_0, x_{-1}, x_{N+1}, x_{N+2}$  and also a special padding character  $p$ . Formally,

$$\Sigma = \{x_i, \bar{x}_i | i \in [1..N]\} \cup \{x_{-1}, x_0, x_{N+1}, x_{N+2}, p\}.$$

The definition of the text  $T$  has two parts: a header and a body, where the body of  $T$  is defined according to the clauses of the given logical formula  $\varphi$ , and the header preceding this body imposes a structure on an  $m$ -approximate cover for  $T$ .

The definition of the body of  $T$  follows directly from the formula  $\varphi$ . For each clause  $C_j = L_1^j \vee L_2^j \vee L_3^j$  of  $\varphi$ ,  $1 \leq j \leq \ell$ , we add to the body of  $T$  the substring  $L_1^j L_2^j L_3^j$ , preceded and followed by a padding of  $2N + 14$  occurrences of the character  $p$ . The role of this padding is to avoid overlaps between occurrences of an approximate cover covering substrings originating from different clauses. The header is composed of  $\ell(N + 3)$  copies of the following string:  $B = p \dots p x_{N+2} x_{N+1} \bar{x}_N \dots \bar{x}_1 x_0 x_{-1} p \dots p x_{N+2} x_{N+1} x_N \dots x_1 x_0 x_{-1} p \dots p$ , where each padding contains  $N + 7$  occurrences of  $p$ .

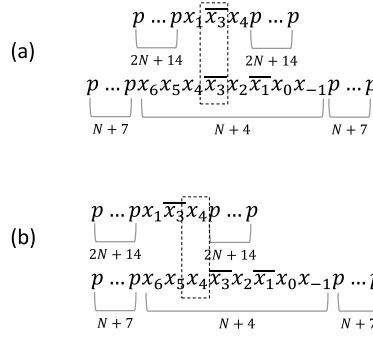
We define the length of the requested approximate cover  $m$  to be  $3N + 18$ . Note that the length of  $T$  and  $m$  as well as the complexity of their construction are polynomial in  $N$  and  $\ell$ . Lemma 1 assures the correctness of the reduction.

**Lemma 1.**  $\varphi$  is satisfiable if and only if  $T$  has an  $m$ -approximate cover with at most  $\ell(N + 3)(N + 1)$  errors.

**Proof.** Assume that  $\varphi$  is satisfiable, and let  $\mathcal{A}$  be an assignment to  $x_1, \dots, x_N$  that satisfies  $\varphi$ . Consider the following  $m$ -length string  $C$ :

1.  $C$  has a prefix of length  $N + 7$  occurrences of the padding character  $p$ .
2.  $C$  has a body of  $N + 4$ -length string that starts with  $x_{N+2} x_{N+1}$  and ends with  $x_0 x_{-1}$  and between these, for each  $i \in [N..1]$ , includes  $x_i$  if  $x_i$  gets the value true in  $\mathcal{A}$ , otherwise, it includes  $\bar{x}_i$ .
3.  $C$  has a suffix of length  $N + 7$  occurrences of the padding character  $p$ .

We show that  $C$  covers  $T$  with at most  $\ell(N + 3)(N + 1)$  errors. First,  $C$  can cover each of the  $\ell(N + 3)$  copies of  $B$  in the header of  $T$  by two consecutive copies of  $C$  (where only their  $N + 7$ -length padding suffix and prefix overlap in the middle) with exactly  $N$  errors for each copy of  $B$ . This gives a total of  $\ell N(N + 3)$  errors for covering the header of  $T$  with  $C$ . Now, since the body of  $T$  is constructed according to the  $\ell$  clauses of  $\varphi$  and each of them can be satisfied by the assignment  $\mathcal{A}$



**Fig. 1.** (a) and (b) show two possible alignments, each having exactly one matching character, of the block  $p \dots p x_1 \bar{x}_3 x_4 p \dots p$  in  $T$  with the  $m$ -length  $C$ -structure approximate cover  $p \dots p x_6 x_5 x_4 x_3 x_2 \bar{x}_1 x_0 x_{-1} p \dots p$  assuming the satisfying assignment  $\mathcal{A}(x_1, x_2, x_3, x_4) = (0, 1, 0, 1)$ .

we are assured that  $C$  can be aligned to match a character originates from a satisfied literal in each of the  $\ell$  blocks of  $T$ 's body. See Fig. 1 for an example. Note, that though there can be several satisfied literals in a clause,  $C$  can be aligned to have only one match due to the reversed order of the variable characters in  $C$ . See Fig. 1 for an example. The other text variable characters are aligned with dummy variables characters of  $C$  and the rest of the variable characters of  $C$  are aligned with occurrences of the padding character  $p$  in  $T$ , and thus add one error each, where the padding character occurrences of  $C$  can always be matched to padding character occurrences in  $T$  and add no errors. Therefore, covering the body of  $T$  with  $C$  adds  $N+3$  errors for each of its  $\ell$  blocks, and the total number of errors to cover  $T$  with  $C$  is  $\ell(N+1)(N+3)$ , as requested.

Assume that  $T$  has an  $m$ -approximate cover with at most  $\ell(N+3)(N+1)$  errors. First, note that any  $m$ -length string over  $\Sigma$  that does not have the prefix-body-suffix structure of  $C$ , has at least  $N+1$  errors for covering each of the  $\ell(N+3)$  copies of  $B$  in the header of  $T$ . This already gives a total of  $\ell(N+3)(N+1)$  errors even if it can cover the body of  $T$  with no errors. Moreover, since the clauses of  $\varphi$  are different, the text body substrings resulting from them are also different. Thus, no  $m$ -length string over alphabet  $\Sigma$  can cover the body of  $T$  with no errors. Therefore, an  $m$ -approximate cover of  $T$  must have the structure of  $C$ .

An  $m$ -length string over  $\Sigma$  with the structure of  $C$  has a total of  $\ell N(N+3)$  errors for covering the header of  $T$  and, therefore, can only have  $\ell(N+3)$  more errors for covering the body of  $T$ . Due to the reverse order of the variable characters in  $C$ , such a string cannot have more than one match for covering any of the text body substrings originating from a clause of  $\varphi$ . Therefore, such a  $C$ -structure string has exactly one match for covering any of the text body substrings originating from a clause of  $\varphi$ . Each variable character appears exactly once in any  $C$ -structure string, so we can interpret the body of this  $C$ -structure string as defining a valid assignment  $\mathcal{A}$  on  $x_1, \dots, x_n$  in the following way: if  $C$  includes  $x_i$ , then  $\mathcal{A}(x_i) = \text{true}$ , otherwise,  $C$  must include  $\bar{x}_i$ , and  $\mathcal{A}(x_i) = \text{false}$ . Since there is a match for covering each of the  $\ell$  text body substrings originating from the  $\ell$  clauses of  $\varphi$ , we have that  $\mathcal{A}$  satisfies  $\varphi$ .  $\square$

We have, therefore, proven Theorem 1.

**Theorem 1.** The cover-length relaxation of the ACP is  $\mathcal{NP}$ -hard.

We can now deduce the  $\mathcal{NP}$ -hardness of the ACP.

**Theorem 2.** The ACP is  $\mathcal{NP}$ -hard.

**Proof.** In fact, for proving the  $\mathcal{NP}$ -hardness of ACP we use the same reduction from 3-SAT described above, and observe that Lemma 1 can be strengthened to claim that  $\varphi$  is satisfiable if and only if  $T$  has an approximate cover with at most  $\ell(N+3)(N+1)$  errors. The first direction follows immediately from the proof that if  $\varphi$  is satisfiable then  $T$  has an  $m$ -approximate cover with at most  $\ell(N+3)(N+1)$  errors, because then the approximate cover, by definition, may have shorter length but cannot have more errors. The reverse direction follows by observing that the structure of  $T$  forces any string with length shorter or greater than  $m$  to have more errors than the  $m$ -approximate cover of  $T$ . Therefore, if  $T$  has an approximate cover with at most  $\ell(N+3)(N+1)$  errors, it must have an  $m$ -approximate cover with at most  $\ell(N+3)(N+1)$  errors. The proof of Lemma 1, then assures that  $\varphi$  is satisfiable in this case.  $\square$

#### 4. The partial-tiling relaxation of the ACP

In this section we study another relaxation of the approximate cover problem: the partial-tiling relaxation, in which we are given a retained tiling of the cover before the errors have occurred together with the input string itself. In order to formally define the relaxation we need Definitions 8 and 9.

**Definition 8.** Let  $T$  be an  $n$ -length string over alphabet  $\Sigma$  and let  $L$  be a valid tiling of  $T$ . Let  $m = n + 1 - L_{last}$ , where  $L_{last}$  is the last index in the tiling  $L$ . Then, an  $L$ -approximate cover of  $T$  is a primitive string  $C$  of length  $m$  such that for every string  $C'$  of length  $m$  over  $\Sigma$ ,  $H(S_L(C'), T) \geq H(S_L(C), T)$ , where  $H$  is the hamming distance of the given strings.

$\min_{C \in \Sigma^m} H(S_L(C), T)$  is the number of errors of an  $L$ -approximate cover of  $T$ .

**Definition 9.** Let  $T$  be an  $n$ -length string over alphabet  $\Sigma$ . Let  $L$  be a valid tiling of  $T$  and let  $L'$  be a valid tiling of  $T$  such that  $L \subseteq L'$ . Let  $m' = n + 1 - L'_{last}$ , where  $L'_{last}$  is the last index in the tiling  $L'$ . Then, a *partial  $L$ -approximate cover* of  $T$  is a primitive string  $C$  of length  $m'$  such that for every string  $C'$  of length  $m'$  over  $\Sigma$ ,  $H(S_{L'}(C'), T) \geq H(S_{L'}(C), T)$ , where  $H$  is the hamming distance of the given strings.

$\min_{C \in \Sigma^{m'}} H(S_{L'}(C), T)$  is the number of errors of a partial  $L$ -approximate cover of  $T$ .

**Definition 10. [The Partial-Tiling Relaxation of the ACP]**

*INPUT:* String  $T$  of length  $n$  over alphabet  $\Sigma$ , and a valid tiling  $L$  of  $T$ .

*OUTPUT:* A partial  $L$ -approximate cover  $C$  of  $T$ .

We describe an algorithm for the partial-tiling relaxation of the approximate cover problem in two parts. We first describe the mandatory part of the algorithm, which we call the Histogram Greedy Algorithm. This algorithm does the main work in finding an approximate cover subject to the tiling  $L$ . It returns a candidate for the final  $L$  approximate cover to be output. This candidate is legal if it is primitive and illegal, otherwise. We then describe the second part, which we call the Partial-Tiling Primitivity Coercion. In this part, the legality of the candidate is checked, and if needed, the candidate is corrected in order to coerce the primitivity requirement.

**4.1. The histogram greedy algorithm**

This part of the algorithm performs the following steps given the text  $T$  and the valid tiling  $L$ :

1. Find  $m$ , the length of an approximate cover subject to the tiling  $L$ , by computing the difference between  $n + 1$ , and the last index in the tiling  $L$ ,  $L_{last}$ , which indicates the last occurrence of the cover in  $T$ .
2. Compute the  $m$ -length mask  $M$  of an approximate cover, by initializing  $M$  to zeroes, setting  $M[1] = 1$ , then reading the tiling  $L$  from beginning to end and for each  $i_k, i_{k+1} \in L$  setting  $M[i_{k+1} - i_k] = 1$ .
3. Compute the  $m$ -long string  $V_C$  of variables from an auxiliary alphabet

$$\Sigma_V = \{v_1, v_2, \dots, v_m\}.$$

First, we initialize the  $m$ -long string  $V_C$  to  $v_1 v_2 \dots v_m$ . Then, we read the mask  $M$  from end to beginning, and for every  $j$  such that  $M[j] = 1$ , we update the string  $V_C$  by equalizing the substrings  $V_C[1..m - j + 1]$  and  $V_C[j..m]$ . In the equalization process, when we obtain an equation  $v_k = v_\ell$  for  $k < \ell$ , we replace both letters by  $v_k$ . The resulting string  $V_C$  represents  $C$  in the following sense: for any pair of indices  $1 \leq i < j \leq m$ , if  $V_C[i] = V_C[j]$  then  $C[i] = C[j]$ . However, it can be that  $V_C[i] \neq V_C[j]$ , while  $C[i] = C[j]$ . In other words,  $V_C$  carries the information on equalities imposed by the mask  $M$  between indices of  $C$ .

4. Compute the  $n$ -long string  $V_T$  of variables from the auxiliary alphabet  $\Sigma_V$ , which is a string covered by  $V_C$  according to the tiling  $L$  of  $T$ .  $V_C$  is computed using the tiling  $L$  and  $V_C$  as follows: it begins with a copy of  $V_C$  and for each index  $i$  in  $L$  a new copy of  $V_C$  is concatenated starting from index  $i$  of  $V_T$  (maybe running over a suffix of the last copy of  $V_C$ ).
5. Compute the histogram  $Hist_{V_C, \Sigma}$  using the alignment of  $T$  with  $V_T$  and counting for each variable  $V \in V_C$  and each  $\sigma \in \Sigma$ , the number of indices  $i$  in  $T$ ,  $V_T[i] = V$  and  $T[i] = \sigma$ .
6. Compute an  $L$ -approximate cover candidate  $C$  greedily according to the histogram  $Hist_{V_C, \Sigma}$ , as follows: for every index  $1 \leq i \leq m$ , set  $C[i] = \sigma_0$ , where  $Hist_{V_C, \Sigma}[V_C[i], \sigma_0] = \max_{\sigma \in \Sigma} Hist_{V_C, \Sigma}[V_C[i], \sigma]$ , i.e., for each index in  $C$  we choose the alphabet symbol that minimizes the number of mismatch errors between  $S_L(C)$  and  $T$  in the relevant indices according to the tiling  $L$ .

The algorithm outputs the  $m$ -length string  $C$  from its last step and the histogram table  $Hist_{V_C, \Sigma}$ .

Lemma 2 describes a property of the output  $C$  returned by the Histogram Greedy algorithm, and immediately follows from the greedy criterion used in step 6 of the algorithm. Lemma 3 describes time complexity of the algorithm.

**Lemma 2.** Let  $C$  be the output of the Histogram Greedy algorithm. Then,

$$H(T, S_L(C)) = \min_{C' \in \Sigma^m} H(T, S_L(C')).$$

**Lemma 3.** The time complexity of the Histogram Greedy algorithm is:  $O(|\Sigma| \cdot m + n)$ .



Despite Lemma 2, the output  $C$  of the Histogram Greedy algorithm might not be an  $L$ -approximate cover of  $T$ , because it might not be primitive, as the following example shows.

**Example.** Assume that  $V_C = XYZWXY$  and  $\Sigma = \{a, b\}$  and that the histogram  $Hist_{V_C, \Sigma}$  computed by the algorithm is the following:

$V_C \setminus \Sigma$	a	b
X	4	1
Y	2	3
Z	2	1
W	0	3

Then, the Histogram Greedy algorithm chooses:  $X = a$ ,  $Y = b$ ,  $Z = a$ ,  $W = b$ , and outputs  $C = ababab$ , which cannot be considered a legal cover since it is not primitive, i.e.,  $C$  itself can be covered by the shorter string  $ab$ . However, the partial  $L$ -approximate cover can have a tiling  $L'$ , such that  $L \subseteq L'$ , which exactly is the case with  $ab$ . Therefore,  $ab$  should be returned as the partial  $L$ -approximate cover of  $T$ . The Partial-Tiling Primitivity Coercion algorithm described in Subsection 4.2 is responsible for checking the legality of the output string received from the Histogram Greedy algorithm and returning a partial  $L$ -approximate cover.

Note, that the input tiling  $L$  requires an  $m$ -length string as an output. Therefore, the (primitive) 2-length approximate cover  $ab$  is precluded as an  $L$ -approximate cover. Assuming that the input tiling  $L$  is the retained tiling of the cover of the original text before the errors occurred, such a case means that, though  $ab$  is a string covering  $T$  subject to a partial tiling  $L$  with the least number of errors, it does not cover  $T$  with  $L$  as a full tiling. In this sense,  $L$  is an evidence that the original cover is of larger length than  $ab$  and that more errors actually happened. Section 5 is devoted to finding an  $L$ -approximate cover.

#### 4.2. The partial-tiling primitivity coercion algorithm

This part of the algorithm gets as input the string  $C$  returned by the Histogram Greedy algorithm and performs the following steps:

1. Check the primitivity of  $C$  (using the linear-time algorithm of [8]). If  $C$  is primitive, return  $C$ .
2. Else, return the primitive cover  $C'$  of  $C$  (found using the linear-time algorithm of [8] in the first step).

The time complexity of the Partial-Tiling Primitivity Coercion algorithm is immediate from the linear-time complexity of the algorithm in [8]. Thus, we get:

**Lemma 4.** *The time complexity of the Partial-Tiling Primitivity Coercion algorithm is  $O(m)$ .*

Theorem 3 follows.

**Theorem 3.** *Given a text  $T$  of length  $n$  over alphabet  $\Sigma$  and a valid tiling  $L$ . Let  $L_{last}$  be the last index in  $L$ . Then, the partial-tiling relaxation of the approximate cover problem of  $T$  can be solved in  $O(|\Sigma| \cdot m + n)$  time, where  $m = n + 1 - L_{last}$ .*

**Proof.** First, note that if the output is returned in the first step of the Partial-Tiling Coercion algorithm, then the returned string is an  $L$ -approximate cover and, therefore, also a trivial partial  $L$ -approximate cover, and the correctness in this case is assured by Lemma 2. We now prove the correctness in case the answer is returned in the second step of the Partial-Tiling Coercion algorithm. Let  $m'$  be the length of the returned string  $C'$ , in this case. Note that  $m' < m$ . Clearly, it holds that  $L' \subset L$ , where  $L'$  is the tiling of  $C'$  in  $T$ , because  $C'$  is a cover of the string  $C$  returned from the Histogram Greedy algorithm which obeys the given tiling  $L$ . It remains to prove that  $C'$  has the least number of errors over all strings of length  $m'$ . The important observation here is that, given  $T$  and  $L'$  instead of  $L$ , the Histogram Greedy algorithm would return  $C'$  with the same number of errors as it returned for  $C$  with the tiling  $L$ . Thus, Lemma 2 assures that

$$H(T, S_{L'}(C')) = \min_{C'' \in \Sigma^{m'}} H(T, S_{L'}(C'')).$$

The time complexity is composed of the time complexities of the Histogram Greedy algorithm and the Partial-Tiling Coercion algorithm and is given by Lemmas 3 and 4.  $\square$

## 5. The full-tiling relaxation of the ACP

In this section we study another relaxation of the approximate cover problem: the full-tiling relaxation, in which we are given a retained tiling of the cover before the errors have occurred together with the input string itself. Unlike the situation in the problem of the previous section, this tiling is assumed to be exact. Therefore, the algorithm cannot return as cover a string that in order to cover  $T$  must have repetitions that are not marked in the tiling  $L$ . The formal definition of the problem is as follows.

### Definition 11. [The Full-Tiling Relaxation of the ACP]

**INPUT:** String  $T$  of length  $n$  over alphabet  $\Sigma$ , and a valid tiling  $L$  of  $T$ .

**OUTPUT:** An  $L$ -approximate cover  $C$  of  $T$ .

In order to impose the requirement of the definition of an  $L$ -approximate cover of  $T$  to be a primitive string such that all its repetitions to cover  $T$  (with minimum number of errors) are marked in the tiling  $L$ , we need a different primitivity coercion algorithm than the one described in the previous section. This algorithm is described in Subsection 5.1. Unfortunately, proving the correctness of this algorithm requires a deep understanding of the properties of coverability in the presence of mismatch errors. Although we are making progress in proving this needed background (see, for example [3]), a lack in the complete understanding of the phenomenon prevents us from proving the correctness formally. Hence, in Subsection 5.2, we resort to experimental evidence of the correctness.

#### 5.1. The full-tiling primitivity coercion algorithm

This part of the algorithm gets as input the string  $C$  returned by the Histogram Greedy algorithm (Subsection 4.1) and performs the following steps:

1. Check the primitivity of  $C$  (using the linear-time algorithm of [8]). If  $C$  is primitive, return  $C$ .
2. Else, find  $V_k \in V_C$  such that if the assignment of  $V_k$  is changed from the symbol with the largest value in the row of  $V_k$  in  $Hist_{V_C, \Sigma}$  to the symbol with the second largest value in this row, thus obtaining a new  $m$ -length candidate string  $C'$ , we have that:
  - (a)  $C'$  is primitive, and
  - (b)  $H(S_L(C'), T) - H(S_L(C), T)$  is minimized over all such primitive strings  $C'$ .

Lemma 5 below describes the time complexity of the Full-Tiling Primitivity Coercion algorithm and immediately follows from the linear-time complexity of the algorithm [8] we use in the first and second step and the description of the second step.

**Lemma 5.** *The time complexity of the Full-Tiling Primitivity Coercion algorithm is  $O((|\Sigma| + m) \cdot m)$ .*

**Remark.** Note that we can use a different algorithm that instead of checking the change of single variables to the second best assignment and choosing the one that gives primitivity with the least number of errors (as our algorithm does), checks the changing to the second best assignment of all subsets of variables and chooses the set that gives primitivity with the least number of errors. This algorithm is obviously correct, i.e., assures primitivity with the least number of errors, however, it has an exponential-time complexity. On the other hand, our algorithm is assured to have polynomial-time complexity, so a proof of its correctness will assure the polynomial-time complexity of the full-tiling relaxation of the ACP.

#### 5.2. Experimental tests of the full-tiling relaxation algorithm

Experiments were designed to test the full-tiling relaxation algorithm, which is composed of the algorithms of Subsections 4.1 and 5.1. In particular, we also wanted to experimentally test how many times the full-tiling primitivity coercion is necessary. Note that, due to the result of [4], this algorithm is only of interest to test under a rather high error rate, in which there is an error in every occurrence of the approximate cover of the text, otherwise, the dynamic programming algorithm solving the candidate-relaxation of the ACP is applicable, where trying every substring of  $T$  as a candidate cover [4]. In order to comprehensively test the algorithm, the inputs for the tests were classified according to the following criteria:

**cover length:** A cover  $C$  of length  $m$  is constructed, where  $m$  is small (less than 10), medium (10-100) or large (100-400). Covers of length more than 400 were not created due to space limitations.

**alphabet size:** The alphabet size was chosen to be either small (at most  $\sqrt{m}$ ) or large (more than  $\sqrt{m}$ ).

**tiling style:** Given a cover  $C$  and its mask  $M$ , a tiling  $L$  for the text  $S_L(C)$  is constructed where the decision of the next index in  $L$  is made according to the following styles: random – an equal priority is given to every set bit in  $M$ ,



**Table 2**

Experimental Tests of the Full-Tiling Relaxation Algorithm for Small Alphabets.

Cover Length	Tiling Style	Error Rate	Error Style	Percent of Inputs	Identical	Primitive	Non-Primitive
small	left	medium	random	1.57	80.89	13.04	6.07
small	left	medium	priority	1.57	80.55	31.51	5.95
small	random	medium	random	1.57	78.80	15.28	5.93
small	random	medium	priority	1.57	78.30	15.46	6.24
small	right	medium	random	1.57	76.20	17.47	6.32
small	right	medium	priority	1.57	76.13	17.82	6.05
small	left	high	random	1.57	5.55	77.78	16.67
small	left	high	priority	1.57	1.87	81.39	16.74
small	random	high	random	1.57	5.24	78.57	16.19
small	random	high	priority	1.57	1.68	82.18	16.13
small	right	high	random	1.57	4.74	80.41	14.85
small	right	high	priority	1.57	1.27	84.28	14.45
medium	left	medium	random	3.22	100	0	0
medium	left	medium	priority	3.22	100	0	0
medium	random	medium	random	3.22	100	0	0
medium	random	medium	priority	3.22	100	0	0
medium	right	medium	random	3.22	100	0	0
medium	right	medium	priority	3.22	100	0	0
medium	left	high	random	3.22	89.80	10.17	0.03
medium	left	high	priority	3.22	87.87	12.09	0.03
medium	random	high	random	3.22	89.39	10.59	0.02
medium	random	high	priority	3.22	87.42	12.54	0.05
medium	right	high	random	3.22	89.07	10.90	0.33
medium	right	high	priority	3.22	86.63	13.35	0.03
large	left	medium	random	0.81	100	0	0
large	left	medium	priority	0.81	100	0	0
large	random	medium	random	0.81	100	0	0
large	random	medium	priority	0.81	100	0	0
large	right	medium	random	0.81	100	0	0
large	right	medium	priority	0.81	100	0	0
large	left	high	random	0.81	100	0	0
large	left	high	priority	0.81	100	0	0
large	random	high	random	0.81	100	0	0
large	random	high	priority	0.81	100	0	0
large	right	high	random	0.81	100	0	0
large	right	high	priority	0.81	100	0	0

left priority – a decreasing priority is given to the set bits in  $M$ , right priority – an increasing priority is given to the set bits in  $M$ .

**error rate:** The input string  $T$  is constructed from  $S_L(C)$  by inserting mismatch errors according to the following error rates: medium (in every  $m$  characters at least one error), high (in every  $m$  characters at least  $\sqrt{m}$  errors).

**error style:** The mismatching character is determined according to the following style: random (replacing by a uniformly at random choice of another character from the alphabet) or priority (replacing by another character with priority to the first character in the alphabet, and if the first character is to be replaced then by a uniformly at random chosen different character).

These criteria guarantee that the inputs created for testing the algorithm all have a coverable original string, that its valid tiling is retained. This original string is then introduced with a sufficiently high error rate to produce the current input string together with the valid tiling as inputs for the tiling relaxation algorithm. Therefore, all the tested inputs have an  $L$ -approximate cover and our tiling relaxation algorithm is indeed applicable for them. Moreover, the above criteria for input generation also aim at neutralizing the effect of the cover length, the alphabet size, the tiling style, the error rate or the error style on the validity of the algorithm, by exhaustively using all reasonable alternatives.

A total of 372000 texts  $T$  were constructed as described above and served as inputs (together with the tiling  $L$ ) to the full-tiling relaxation algorithm. The results are given in Tables 2 and 3. The column “Percent of Inputs” describes how many of the input texts had each row’s characteristics. Numbers are rounded to two digits after decimal point. The column “Identical” describes in how many of the input texts the Histogram Greedy algorithm of Subsection 4.1 returned the original cover  $C$  of the text  $S_L(C)$  built prior to the error insertion process. The column “Primitive” describes in how many of the input texts the Histogram Greedy algorithm of Subsection 4.1 returned a primitive cover and there was no need to proceed with the second phase of the Full-Tiling Primitivity Coercion algorithm of Subsection 5.1. The column “Non-Primitive” describes in how many of the input texts the Histogram Greedy algorithm of Subsection 4.1 returned a non-primitive string and, therefore, the second phase of the Full-Tiling Primitivity Coercion algorithm of Subsection 5.1 was performed. This latter case happened in 8912 texts, which are about 2% of the texts.

**Table 3**  
Experimental Tests of the Full-Tiling Relaxation Algorithm for Large Alphabets.

Cover Length	Tiling Style	Error Rate	Error Style	Percent of Inputs	Identical	Primitive	Non-Primitive
small	left	medium	random	0.59	89.45	9.59	0.96
small	left	medium	priority	0.59	76.51	21.93	1.56
small	random	medium	random	0.59	87.57	11.24	1.19
small	random	medium	priority	0.59	76.15	22.39	1.47
small	right	medium	random	0.59	86.56	12.89	0.55
small	right	medium	priority	0.59	75.51	23.40	1.10
small	left	high	random	0.59	25.55	70.78	3.67
small	left	high	priority	0.59	1.84	84.45	13.72
small	random	high	random	0.59	24.82	70.96	4.22
small	random	high	priority	0.59	2.06	86.15	11.79
small	right	high	random	0.59	23.76	72.75	3.49
small	right	high	priority	0.59	1.88	85.32	12.80
medium	left	medium	random	1.62	100	0	0
medium	left	medium	priority	1.62	100	0	0
medium	random	medium	random	1.62	100	0	0
medium	random	medium	priority	1.62	100	0	0
medium	right	medium	random	1.62	100	0	0
medium	right	medium	priority	1.62	100	0	0
medium	left	high	random	1.62	99.77	0.23	0
medium	left	high	priority	1.62	85.11	14.89	0
medium	random	high	random	1.62	99.75	0.25	0
medium	random	high	priority	1.62	84.19	15.81	0
medium	right	high	random	1.62	99.90	0.10	0
medium	right	high	priority	1.62	84.11	15.90	0
large	left	medium	random	0.54	100	0	0
large	left	medium	priority	0.54	100	0	0
large	random	medium	random	0.54	100	0	0
large	random	medium	priority	0.54	100	0	0
large	right	medium	random	0.54	100	0	0
large	right	medium	priority	0.54	100	0	0
large	left	high	random	0.54	100	0	0
large	left	high	priority	0.54	100	0	0
large	random	high	random	0.54	100	0	0
large	random	high	priority	0.54	100	0	0
large	right	high	random	0.54	100	0	0
large	right	high	priority	0.54	100	0	0

**Experiments Conclusion:** Primitivity coercion was necessary in 2% of the total tested inputs. In a 100% of the tests the returned string after the Full-Tiling Primitivity Coercion algorithm was indeed an  $L$ -approximate cover of the input string.

## 6. Open problems

In this paper we initiated the study of the approximate cover problem using a new approach. We proved that the ACP is  $\mathcal{NP}$ -hard. We also broadened the study to consider several relaxations of the ACP, proving that some of them (the cover-length relaxation) are already  $\mathcal{NP}$ -hard, while others (the partial-tiling relaxation) are polynomial-time computable. Some interesting questions and open problems are:

- Our  $\mathcal{NP}$ -hardness proof uses unbounded-size alphabet. Is the ACP still  $\mathcal{NP}$ -hard for constant-size alphabet?
- It is interesting to define other relaxations of the ACP and to study their complexity in order to have a deeper understanding of the ACP.
- In this paper we only experimentally checked the correctness of our full-tiling relaxation algorithm. We would like to have a formal proof of its correctness.
- In this paper we considered the Hamming distance as a metric in the definition of approximate cover. Other string metrics can be considered as well. It is interesting to see if and how the complexity of the problem changes with the use of other string metrics.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

- [1] A. Amir, E. Eisenberg, A. Levy, Approximate periodicity, in: Proc. ISAAC 2010, in: LNCS, vol. 6506, Springer, 2010, pp. 25–36.
- [2] A. Amir, E. Eisenberg, A. Levy, E. Porat, N. Shapira, Cycle detection and correction, *ACM Trans. Algorithms* 9 (1) (2012) 13.
- [3] A. Amir, C.S. Iliopoulos, J. Radoszewski, Two strings at Hamming distance 1 cannot be both quasiperiodic, manuscript, 2016.
- [4] A. Amir, A. Levy, M. Lewenstein, R. Lubin, B. Porat, Can we recover the cover?, *CPM*, 2017.
- [5] P. Antoniou, M. Crochemore, C.S. Iliopoulos, I. Jayasekera, G.M. Landau, Conservative string covering of indeterminate strings, in: *Proc. Stringology*, 2008, pp. 108–115.
- [6] A. Apostolico, D. Breslauer, Of periods, quasiperiods, repetitions and covers, in: *Proc. Structures in Logic and Computer Science*, in: LNCS, vol. 1261, Springer, 1997, pp. 236–248.
- [7] A. Apostolico, A. Ehrenfeucht, Efficient detection of quasiperiodicities in strings, *Theor. Comput. Sci.* 119 (1993) 247–265.
- [8] A. Apostolico, M. Farach, C.S. Iliopoulos, Optimal superprimitivity testing for strings, *Inf. Process. Lett.* 39 (1991) 17–20.
- [9] C. Barton, T. Kociumaka, S.P. Pissis, J. Radoszewski, Efficient index for weighted sequences, in: *Proc. CPM*, 2016, pp. 4:1–4:13.
- [10] D. Breslauer, An on-line string superprimitivity test, *Inf. Process. Lett.* 44 (1992) 345–347.
- [11] D. Breslauer, Testing string superprimitivity in parallel, *Inf. Process. Lett.* 49 (5) (1994) 235–241.
- [12] M. Christodoulakis, C.S. Iliopoulos, K. Park, J.S. Sim, Approximate seeds of strings, *J. Autom. Lang. Comb.* 10 (2005) 609–626.
- [13] T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein, NP-Completeness. Introduction to Algorithms, second edition, MIT Press, 2001, pp. 966–1021.
- [14] T. Crawford, C.S. Iliopoulos, R. Raman, String matching techniques for musical similarity and melodic recognition, *Comput. Musicol.* 11 (1998) 73–100.
- [15] M. Crochemore, C.S. Iliopoulos, T. Kociumaka, J. Radoszewski, W. Rytter, T. Walen, Covering problems for partial words and for indeterminate strings, *Theor. Comput. Sci.* 698 (2017) 25–39.
- [16] M. Crochemore, C.S. Iliopoulos, S.P. Pissis, G. Tischler, Cover array string reconstruction, in: *Proc. CPM*, 2010, pp. 251–259.
- [17] M. Crochemore, C.S. Iliopoulos, H. Yu, Algorithms for computing evolutionary chains in molecular and musical sequences, in: *Proc. 9th Austral. Workshop on Combinatorial Algorithms*, 1998, pp. 172–185.
- [18] T. Flouri, C.S. Iliopoulos, T. Kociumaka, S.P. Pissis, S.J. Puglisi, W.F. Smyth, W. Tyczynski, Enhanced string covering, *Theor. Comput. Sci.* 506 (2013) 102–114.
- [19] O. Guth, B. Melichar, Using finite automata approach for searching approximate seeds of strings, in: *Intelligent Automation and Computer Engineering*, Springer Netherlands, ISBN 978-90-481-3517-2, 2010, pp. 347–360.
- [20] O. Guth, B. Melichar, M. Balik, Searching all approximate covers and their distance using finite automata, in: *Information Technologies – Applications and Theory*, ISBN 978-80-969184-9-2, 2009, pp. 21–26.
- [21] C.S. Iliopoulos, L. Mouchard, Quasiperiodicity and string covering, *Theor. Comput. Sci.* 218 (1) (1999) 205–216.
- [22] C.S. Iliopoulos, D.W.G. Moore, K. Park, Covering a string, *Algorithmica* 16 (3) (1996) 288–297.
- [23] C.S. Iliopoulos, W.F. Smyth, An on-line algorithm of computing a minimum set of  $k$ -covers of a string, in: *Proc. 9th Australasian Workshop on Combinatorial Algorithms*, AWOCA, 1998, pp. 97–106.
- [24] T. Kociumaka, S.P. Pissis, J. Radoszewski, W. Rytter, T. Walen, Fast algorithm for partial covers in words, in: *Proc. CPM*, 2013, pp. 177–188.
- [25] T. Kociumaka, S.P. Pissis, J. Radoszewski, W. Rytter, T. Walen, Fast algorithm for partial covers in words, *Algorithmica* 73 (2015) 217–233.
- [26] T. Kociumaka, S.P. Pissis, J. Radoszewski, W. Rytter, T. Walen, Efficient algorithms for shortest partial seeds in words, *Theor. Comput. Sci.* 710 (2018) 139–147.
- [27] R.M. Kolpakov, G. Kucherov, Finding approximate repetitions under hamming distance, *Theor. Comput. Sci.* 303 (2003) 135–156.
- [28] G.M. Landau, J.P. Schmidt, An algorithm for approximate tandem repeats, in: *Proc. 4th Symp. Combinatorial Pattern Matching*, in: LNCS, vol. 648, Springer, Berlin, 1993, pp. 120–133.
- [29] G.M. Landau, J.P. Schmidt, D. Sokol, An algorithm for approximate tandem repeats, *J. Comput. Biol.* 8 (1) (2001) 1–18.
- [30] Y. Li, W.F. Smyth, Computing the cover array in linear time, *Algorithmica* 32 (1) (2002) 95–106.
- [31] M. Lothaire, Combinatorics on Words, Addison-Wesley, Reading, Mass, 1983.
- [32] D. Moore, W.F. Smyth, An optimal algorithm to compute all the covers of a string, *Inf. Process. Lett.* 50 (5) (1994) 239–246.
- [33] D. Moore, W.F. Smyth, A correction to: an optimal algorithm to compute all the covers of a string, *Inf. Process. Lett.* 54 (1995) 101–103.
- [34] J.S. Sim, C.S. Iliopoulos, K. Park, W.F. Smyth, Approximate periods of strings, *Theor. Comput. Sci.* 262 (2001) 557–568.
- [35] W.F. Smyth, Repetitive perhaps, but certainly not boring, *Theor. Comput. Sci.* 249 (2) (2000) 343–355.
- [36] H. Zhang, Q. Guo, C.S. Iliopoulos, Algorithms for computing the lambda-regularities in Strings, *Fundam. Inform.* 84 (1) (2008) 33–49.
- [37] H. Zhang, Q. Guo, C.S. Iliopoulos, Varieties of regularities in weighted sequences, in: *Proc. AAIM 2010*, in: LNCS, vol. 6124, Springer, 2010, pp. 271–280.