

A Linear Time Algorithm for Seeds Computation

Tomasz Kociumaka^{1,2}, Marcin Kubica¹, Jakub Radoszewski¹, Wojciech Rytter¹, and Tomasz Walen¹

¹Institute of Informatics, University of Warsaw, Poland,
[kociumaka,kubica,jrad,rytter,walen]@mimuw.edu.pl

²Department of Computer Science, Bar-Ilan University, Ramat Gan, Israel

Abstract

A seed in a word is a relaxed version of a period in which the occurrences of the repeating subword may overlap. We show a linear-time algorithm computing a linear-size representation of all the seeds of a word (the number of seeds might be quadratic). In particular, one can easily derive the shortest seed and the number of seeds from our representation. Thus, we solve an open problem stated in the survey by Smyth (2000) and improve upon a previous $\mathcal{O}(n \log n)$ algorithm by Iliopoulos, Moore, and Park (1996). Our approach is based on combinatorial relations between seeds and subword complexity (used here for the first time in context of seeds). In the previous papers, the compact representation of seeds consisted of two independent parts operating on the suffix tree of the word and the suffix tree of the reverse of the word, respectively. Our second contribution is a simpler representation of all seeds which avoids dealing with the reversed word.

A preliminary version of this work, with a much more complex algorithm constructing the earlier representation of seeds, was presented at the 23rd Annual ACM-SIAM Symposium of Discrete Algorithms (SODA 2012).

1 Introduction

The notion of periodicity in words is widely used in many fields, such as combinatorics on words, pattern matching, data compression, automata theory, formal language theory, molecular biology etc. (see [35]). The concept of quasiperiodicity, introduced by Apostolico and Ehrenfeucht in [5], is a generalization of the notion of periodicity: A quasiperiodic word is entirely covered by occurrences of another (shorter) word, called the *quasiperiod* or the *cover*. The occurrences of the quasiperiod may overlap, while in a periodic repetition the occurrences of the period do not overlap. Hence, quasiperiodicity enables detecting repetitive structure of words which cannot be found using the classic characterizations in terms of periods.

An extension of the notion of a cover is the notion of a *seed*: a cover which is not necessarily aligned with the ends of the word being covered, but is allowed to overflow on either side; see Fig. 1. More formally, a word v is a *seed* of w if v is a subword of w and w is a subword of some word u covered by v .

Seeds were first introduced and studied by Iliopoulos, Moore, and Park [26]. The original motivation for covers and seeds comes from DNA sequence analysis (the search for regularities and common features in

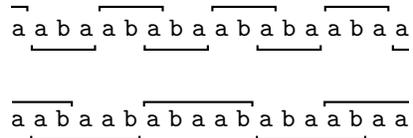


Figure 1: The word aba (above) is the shortest seed of the word $w = aabaababaababaabaa$. Another seed of w is $abaab$ (below). Two of its “overhanging” occurrences correspond to boundary subwords aab and $abaa$. In total, the word w has 35 distinct seeds, but does not have a non-trivial cover.

DNA sequences). Due to natural applications in molecular biology (a hybridization approach to analysis of a DNA sequence), both covers and seeds have also been extended in the sense that a number of factors are considered instead of a single word [28]. This way, the notions of k -covers [12, 25], λ -covers [24], and λ -seeds [22] were introduced. In applications such as molecular biology and computer-assisted music analysis, finding exact repetitions is not always sufficient; the same problem occurs for quasiperiodic repetitions. This led to the introduction of the notions of approximate covers and seeds [2, 3, 4, 40, 10], partial covers and seeds [18, 32, 31], and approximate λ -covers [23].

Previous results. Iliopoulos, Moore, and Park [26] gave an $\mathcal{O}(n \log n)$ -time algorithm computing a linear representation of all the seeds of a given word. For the next 15 years, no $o(n \log n)$ -time algorithm was known for this problem. Smyth formulated computing all the seeds of a word in linear time as an open problem in his survey [41]. Berkman et al. [7] gave a parallel algorithm computing all the seeds in $\mathcal{O}(\log n)$ time and $\mathcal{O}(n^{1+\varepsilon})$ space (for any positive ε) using n processors in the CRCW PRAM model. Much later, Christou et al. [11] proposed an alternative sequential $\mathcal{O}(n \log n)$ -time algorithm for computing the shortest seed.

In contrast, a linear-time algorithm finding the shortest cover of a word was given by Apostolico et al. [6] and later on improved into an on-line algorithm by Breslauer [8]. Moore and Smyth [37, 38] proposed a linear-time algorithm computing all the covers of a word, whereas Li and Smyth [34] afterwards developed an on-line algorithm for the problem of representing all covers of all prefixes of a word.

Another line of research is finding maximal quasiperiodic subwords of a word. This notion resembles maximal repetitions (runs) in a word [33], which is another widely studied notion of combinatorics on words. Two $\mathcal{O}(n \log n)$ -time algorithms for reporting all maximal quasiperiodic subwords of a word of length n have been proposed by Brodal and Pedersen [9] and Iliopoulos and Mouchard [27]; these results improved upon the initial $\mathcal{O}(n \log^2 n)$ -time algorithm by Apostolico and Ehrenfeucht [5].

Our result. We present a linear-time algorithm computing the set $\text{Seeds}(w)$ of all seeds of a given word w . As illustrated in Example 1.1, the number of seeds can be quadratic in the length $|w|$ (contrary to the number of covers, which is always linear). Consequently, our algorithm returns a linear-size *package representation* of the set $\text{Seeds}(w)$, which allows finding a shortest seed and the number of all seeds in a very simple way.

Example 1.1. The following word of length $4m + 3$ contains $\Theta(m^2)$ different seeds:

$$w = \mathbf{a}^m \mathbf{b} \mathbf{a}^m \mathbf{b} \mathbf{a}^m \mathbf{b} \mathbf{a}^m.$$

Those seeds are $\mathbf{a}^i \mathbf{b} \mathbf{a}^j$ with $i + j \geq m$ and $0 \leq i, j \leq m$.

Our procedure assumes that the alphabet Σ of the input word w consists of integers that are polynomial in terms of the length n of w .

Package Representation of Seeds. For a word w , by $w[i..j]$ we denote the subword $w[i] \cdots w[j]$. We introduce packages which are collections of consecutive prefixes of a subword of w . More formally, for positive integers $i \leq j_1 \leq j_2$, we define a *package*:

$$\text{pack}(i, j_1, j_2) = \{w[i..j] : j_1 \leq j \leq j_2\}.$$

If \mathcal{L} is a set of triples of integers, then we denote:

$$\text{PACK}(\mathcal{L}) = \bigcup_{(i, j_1, j_2) \in \mathcal{L}} \text{pack}(i, j_1, j_2).$$

The output of our algorithm, called the *package representation* of the set $\text{Seeds}(w)$, consists of a set \mathcal{L} of integer triples such that $\text{Seeds}(w) = \text{PACK}(\mathcal{L})$ and all the packages in the representation are pairwise disjoint. (In other words, each seed belongs to exactly one package).

complexity, which is the key combinatorial contribution behind our main recursive algorithm described in Section 6. The implementation of an auxiliary operation on package representations (applied to merge the results of recursive calls) is deferred to Section 7.

Our Techniques. Our linear-time solution relies on several combinatorial and algorithmic tools.

Compact representation of seeds: Despite its quadratic size and the failure of the naive argument (Remark 1.3), the set $\text{Seeds}(w)$ always admits a package representation of linear size (Section 3). While this fact is not essential for our algorithm (see [30]), it makes our results much simpler and cleaner.

Combinatorial properties of seeds: The connection between seeds and subword complexity gives an efficient reduction to a set of recursive calls of total size decreased by a constant factor (see Sections 5 and 6).

Interpretation of packages as paths on the suffix trie: Packages naturally correspond to paths in the (uncompacted) suffix trie, which can be easily stored using the (compacted) suffix tree. We use this interpretation both to derive the combinatorial upper bound on the package representation size (Section 3) and in the algorithmic construction of the package representation of long seeds (Section 4). The new linear-time offline algorithm answering Weighted Ancestor Queries (Section 7.1) lets us efficiently map packages on the suffix tree.

Efficient manipulation of package representations: Package representations provide a convenient way of interpreting the results of recursive calls (seeds of certain subwords) as families of subwords of the whole word w . This allows for a simple linear-time procedure aggregating the results of the recursive calls (Section 7).

2 Preliminaries

We consider *words* over a polynomially bounded integer alphabet Σ . For a word w , by $|w|$ we denote its length and by $w[i]$, for $1 \leq i \leq |w|$, we denote its i th letter. By $\text{Alph}(w)$ we denote the set of letters occurring in w . For $1 \leq i \leq j \leq |w|$, a word $u = w[i] \cdots w[j]$ is called a *subword* of w . We also say that w is a *superstring* of u . In this case, by $w[i..j]$ we denote the occurrence of u at position i , called a *fragment* of w . The set of positions where u occurs in w is denoted by $\text{Occ}(u, w)$ (or $\text{Occ}(u)$ if w is clear from the context).

A fragment of w other than the whole word w is called a *proper* fragment of w . A fragment starting at position 1 is called a *prefix* of w and a fragment ending at position $|w|$ is called a *suffix* of w ; we also say that the corresponding subword is a prefix or suffix of w , respectively.

A *border* of w is a subword of w which occurs both as a prefix and as a suffix of w . An integer p , $1 \leq p \leq n$, is a *period* of a word w if $w[i] = w[i + p]$ for $1 \leq i \leq |w| - p$. It is well known that p is a period of w if and only if w has a border of length $|w| - p$; see [13, 14]. Moreover, Fine and Wilf's Periodicity Lemma [17] asserts that if a word w has periods p and q such that $p + q - \gcd(p, q) \leq |w|$, then w also has a period $\gcd(p, q)$.

Throughout the paper, by w we denote the word which seeds are to be computed and by n we denote its length.

2.1 Tries, Suffix Trees, and Package Representations

A *trie* is a rooted tree whose nodes correspond to prefixes of words in a given (finite) family W . If ν is a node, then the corresponding prefix v is called the *value* of the node. The node with value v is called the *locus* of v . The parent-child relation in the trie is defined so that the root is the locus of the empty word, while the parent μ of a node ν is the locus of the value of ν with the last character removed. This character is the *label* of the edge from μ and ν . In general, if μ is an ancestor of ν , then the label of the path from μ to ν is the concatenation of edge labels on the path. A trie of the set W containing all the suffixes of a word *babaad* is shown in Fig. 3.

A node is *branching* if it has at least two children and *terminal* if its value belongs to W . A *compacted trie* is obtained from the underlying trie by dissolving all nodes except for the root, branching nodes, and terminal nodes. In other words, we compress paths of non-terminal vertices with single children, and thus the number of remaining nodes becomes bounded by $2|W|$. We refer to all preserved nodes of the trie as *explicit* (since they are stored explicitly) and to the dissolved ones as *implicit*. If ν is the locus of a word v in an uncompactified trie, then the locus of v in the corresponding compacted trie is defined as (μ, d) , where μ is the lowest explicit descendant of ν , and d is the distance (in the uncompactified trie) from ν to μ . Note that $\mu = \nu$ and $d = 0$ if ν is explicit. Edges of a compacted trie correspond to paths in the underlying trie and thus their labels are non-empty words, typically stored as references to fragments of the words in W .

The *suffix trie* of word w is the trie of all suffixes of w (see Fig. 3), with the locus $w[i..n]$ labelled by the position i . Consequently, there is a natural bijection between subwords of w and nodes of the suffix trie; we often use it to identify subwords of w with their loci in the suffix trie.

The *suffix tree* of w [42] is the compacted suffix trie of w . For a word of length n , it takes $\mathcal{O}(n)$ space and can be constructed in $\mathcal{O}(n)$ time either directly [15] or from the suffix array of w ; see [36, 29, 14, 13].

We say that two subwords u and v of w are *equivalent* if $Occ(u, w) = Occ(v, w)$. The equivalence classes of this relation correspond to edges of the suffix tree of w , as shown in the following observation and Fig. 3; see [13].

Observation 2.1. *Each equivalence class is of the form $E = \text{pack}(i, j_1, j_2)$ for some positions $i \leq j_1 \leq j_2$ and corresponds to the set of all nodes on an edge of the suffix tree of w (excluding the topmost explicit node). Hence, there are at most $2|w|$ equivalence classes.*

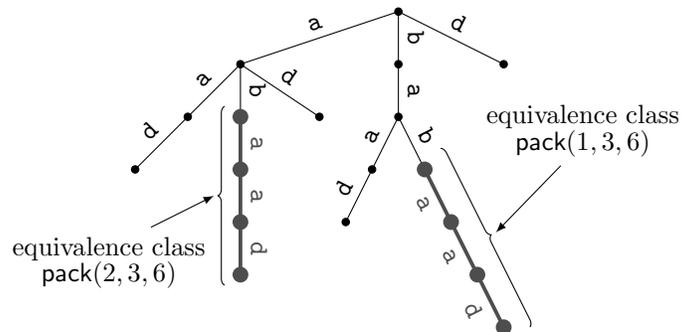


Figure 3: The suffix trie of the word **babaad**. The equivalence classes correspond to compacted edges (excluding their topmost nodes). Two of them are marked in the figure: $\{\text{ab}, \text{aba}, \text{abaa}, \text{abaaad}\} = \text{pack}(2, 3, 6)$ and $\{\text{bab}, \text{baba}, \text{babaa}, \text{babaad}\} = \text{pack}(1, 3, 6)$.

In Section 7, we heavily use the connections between suffix trees and package representation to develop the following auxiliary procedure:

Combine (R_1, \dots, R_k) : Given package representations of sets R_1, \dots, R_k of subwords of a word w , compute a smallest package representation of $\bigcap_{i=1}^k R_i$.

Note that the Combine operation is much more complicated than a simple intersection of sets of intervals. The following lemma is proved in Section 7.

Lemma 2.2. *For a word w of length n and sets R_1, \dots, R_k of subwords of w , given in package representations of total size N , $\text{Combine}(R_1, \dots, R_k)$ can be implemented in $\mathcal{O}(n+N)$ time. The size of the resulting package representation is at most N .*

2.2 Seeds: Formal Definition and Equivalent Characterizations

We say that a fragment $w[i..j]$ covers a position k (or that the position k lies within $w[i..j]$) if $i \leq k \leq j$. A word v is a *cover* of word w if the occurrences of v cover all positions in w . A word v is a *seed* of word w if it is a subword of w and a cover of a superstring of w . This definition immediately implies the following observation.

Observation 2.3. *Let v be a seed of word w . If v occurs in a subword w' of w , then v is a seed of w' .*

Denote by $\text{Seeds}_I(w)$ the set of all seeds of w with lengths in interval I . We also denote $\text{Seeds}(w) = \text{Seeds}_{[1..n]}(w)$, $\text{Seeds}_{\leq k}(w) = \text{Seeds}_{[1..k]}(w)$, $\text{Seeds}_{\geq k}(w) = \text{Seeds}_{[k..n]}(w)$, and $\text{Seeds}_k(w) = \text{Seeds}_{[k..k]}(w)$.

A *left-overhanging* occurrence of v in w is a prefix of w that matches a proper suffix of v . Symmetrically, a *right-overhanging* occurrence of v in w is a suffix of w that matches a proper prefix of v . The length of the occurrence is the length of the corresponding prefix/suffix of w . In this context, usual occurrences are sometimes called *full*, while a *generalized* occurrence is a full or an overhanging one. A more operational definition of seeds can be formulated in terms of generalized occurrences as follows; see Fig. 1.

Fact 2.4 (Alternative definition of seeds). *A subword v of word w is a seed of w if and only if each position in w is covered by a full, left-overhanging, or right-overhanging occurrence of v in w .*

Proof. Suppose that v is a seed of w , i.e., a cover of a superstring xwy of w . For each position in w , the corresponding position in xwy is covered by a full occurrence of v in xwy , which becomes a generalized occurrence of v when restricted to w . Hence, generalized occurrences of v in w cover all positions of w .

For the converse implication, we construct a superstring xwy of w which has v as a cover. For this, we extend w so that the longest left-overhanging occurrence of v and the longest right-overhanging occurrence of v become full occurrences. Shorter overhanging occurrences may be destroyed, but they do not cover any extra positions in w compared to the longest ones. Consequently, the fragment w in xwy is covered by full occurrences of v . What is more, the prefix x is covered by the occurrence of v as a prefix of xwy and the suffix y is covered by the occurrence of v as a suffix of xwy . Thus, v is a cover of xwy and a seed of w . \square

We denote by $\text{SUB}_k(w)$ the set all subwords of length k of word w . Fact 2.4 lets us show that the family $\text{SUB}_{2\ell-1}(w)$ of subwords of length $2\ell-1$ of a word w uniquely determines the length- ℓ seeds of w .

Lemma 2.5. *Let v be a non-empty word such that $2|v| - 1 \leq n$. The following conditions are equivalent:*

- (1) v is a seed of w ;
- (2) v is a seed of every subword of w of length $2|v| - 1$.

Proof. (1) \Rightarrow (2). Each $s \in \text{SUB}_{2|v|-1}(w)$ occurs as $w[i - |v| + 1 .. i + |v| - 1]$ for some position i , $|v| \leq i \leq n - |v| + 1$. The position i can only be covered by a full occurrence of v contained within this fragment, so v is a subword of each $s \in \text{SUB}_{2|v|-1}(w)$. Due to Observation 2.3, v is also a seed of every $s \in \text{SUB}_{2|v|-1}(w)$.

(2) \Rightarrow (1). Positions i satisfying $|v| \leq i \leq n - |v| + 1$ are covered by full occurrences of v due to the fact that v occurs in each subword of w of length $2|v| - 1$, including $w[i - |v| + 1 .. i + |v| - 1]$. If some position $i < |v|$ or $i > n - |v| + 1$ in w is not covered by a full occurrence of v , then it can be covered by a left-overhanging occurrence of v in $w[1 .. 2|v| - 1]$ or a right-overhanging occurrence of v in $w[n - 2|v| + 2 .. n]$, respectively. These overhanging occurrences are also present in w . \square

Corollary 2.6. *For each word w and integer k , $1 \leq k \leq \frac{n+1}{2}$, we have*

$$\text{Seeds}_{\leq k}(w) = \bigcap_{u \in \text{SUB}_{2k-1}(w)} \text{Seeds}_{\leq k}(u).$$

Proof. Consider a length ℓ not exceeding k . Note that $\text{SUB}_{2\ell-1}(w) = \bigcup_{s \in \text{SUB}_{2k-1}(w)} \text{SUB}_{2\ell-1}(s)$, so Lemma 2.5 yields

$$\text{Seeds}_{\ell}(w) = \bigcap_{u \in \text{SUB}_{2\ell-1}(w)} \text{Seeds}_{\ell}(u) = \bigcap_{s \in \text{SUB}_{2k-1}(w)} \bigcap_{u \in \text{SUB}_{2\ell-1}(s)} \text{Seeds}_{\ell}(u) = \bigcap_{s \in \text{SUB}_{2k-1}(w)} \text{Seeds}_{\ell}(s).$$

The equality holds for each length $\ell \leq k$, which concludes the proof. \square

3 Representation Theorem

Let v be a subword of a word w . Let us introduce a decomposition

$$w = v^- \cdot \hat{v} \cdot v^+$$

such that \hat{v} is the longest subword of w having v as a border. In other words,

$$\hat{v} = w[\min \text{Occ}(v) \dots \max \text{Occ}(v) + |v| - 1],$$

i.e., \hat{v} can be seen as the shortest fragment of w containing all full occurrences of v in w ; see Fig. 4.

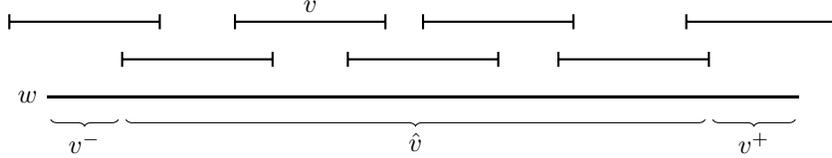


Figure 4: Generalized occurrences of a seed v of w and the decomposition of w into v^- , \hat{v} , and v^+ .

Definition 3.1 (see Fig. 5). *We say that v is a quasiseed of w if it is a cover of \hat{v} . If v is a seed of v^-v , then v is called a left candidate of w , and in case v is a seed of vv^+ , then v is called a right candidate.*

By $\text{QSeeds}(w)$, $\text{LCands}(w)$, and $\text{RCands}(w)$ we denote the sets of quasiseeds, left candidates, and right candidates of w , respectively.

Lemma 3.2. $\text{Seeds}(w) = \text{LCands}(w) \cap \text{QSeeds}(w) \cap \text{RCands}(w)$.

Proof. We apply the characterization of Fact 2.4. Observe that there are natural bijections between full occurrences of v in w and in \hat{v} , between left-overhanging occurrences of v in w and in v^-v , and between right-overhanging occurrences of v in w and in vv^+ .

Next, note that any position of w within \hat{v} covered an overhanging occurrence of v must be located within the leading or trailing $|v|$ characters of \hat{v} , so it is also covered by a full occurrence of v since v is a border of \hat{v} . Thus, v is a quasiseed of w if and only if the positions of w within \hat{v} are covered by generalized occurrences of v .

Moreover, observe that v occurs in v^-v only as a suffix, so the positions within the leading v^- of v^-v and of w can be covered by left-overhanging occurrences only. Consequently, v is a left candidate if and only if the positions within the leading v^- of w are covered by generalized occurrences of v . Symmetrically, v is a right candidate if and only if the positions within the trailing v^+ are covered by generalized occurrences of v .

Combining these three facts, we conclude that v is a seed of w if and only if it is simultaneously a quasiseed, a left candidate, and a right candidate. \square

Next, we shall characterize the quasiseeds and candidates in a computationally feasible way and bound the sizes of their package representations.

3.1 Quasiseeds

For a set $X = \{x_1, \dots, x_k\}$ of integers, $x_1 < \dots < x_k$, let us define the value $\text{maxgap}(X)$ as:

$$\text{maxgap}(X) = \begin{cases} 0 & \text{if } k \leq 1, \\ \max\{x_{i+1} - x_i : 1 \leq i < k\} & \text{if } k \geq 2. \end{cases}$$

For example, $\text{maxgap}(\{1, 3, 8, 13, 17\}) = 5$. The following easy observation relates this function to covers; see, e.g., [7].

The figure shows two rows of the word $w = \text{aabaabababababaa}$. In the first row, the subword ababaa is highlighted with a grey background and a black border, and is labeled as a quasiseed. In the second row, the subword baab is highlighted with a grey background and a black border, and is labeled as both a left and a right candidate. The word w is written in a light grey font with black outlines for each letter.

Figure 5: The word ababaa is a quasiseed of $w = \text{aabaabababababaa}$, whereas the word baab is both a left and a right candidate of w . Neither of them is a seed of w , though.

Observation 3.3. *A subword v of word w is a quasiseed of w if and only if $\text{maxgap}(\text{Occ}(v, w)) \leq |v|$.*

Corollary 3.4. *For each equivalence class E , the quasiseeds contained in E form a single package. In other words, $\text{QSeeds}(w) \cap E$ has a package representation of size 1.*

Proof. Consider the shortest subword $v \in E \cap \text{QSeeds}(w)$ and let $u \in E$ satisfy $|u| \geq |v|$. Due to Observation 3.3, $\text{maxgap}(\text{Occ}(u)) = \text{maxgap}(\text{Occ}(v)) \leq |v| \leq |u|$, so u is also a quasiseed. This completes the proof. \square

3.2 Left Candidates

The *border table* $B[0..n]$ stores at $B[i]$ the length of the longest proper border of $w[1..i]$ (we assume $B[0] = -1$). The following fact was already shown implicitly in [11]; we give its proof for completeness.

Lemma 3.5. *A subword v of word w is a left candidate of w if and only if $B[|v^-v|] \geq |v^-|$.*

Proof. Recall that positions within the prefix v^- of v^-v can only be covered by left-overhanging occurrences of v .

Consequently, if v is a seed of v^-v , then there is a prefix of v^-v of length at least $|v^-|$ equal to a suffix of v . In other words, there is a proper border of v^-v of length at least $|v^-|$.

For a proof of the converse implication, suppose that v^-v has a proper border of length at least $|v^-|$. Since v has only one occurrence in v^-v , the border must be a proper suffix of v . Consequently, v has a left-overhanging occurrence in v^-v of length at least $|v^-|$. This fragment covers all positions in v^- ; the remaining positions in v^-v are trivially covered by the occurrence of v as a suffix of v^-v . \square

A classic property of the border table is that $0 \leq B[i] \leq B[i-1] + 1$ holds for $1 \leq i \leq n$. Fine and Wilf's Periodicity Lemma [17] lets us further characterize positions where the right inequality is strict.

Lemma 3.6. *If $B[i] \leq B[i-1]$ holds for a position i of the word w , then $B[i] + B[i-1] < i - 1$.*

Proof. We assume $B[i] > 0$; otherwise, the claim holds trivially: $B[i] + B[i-1] = B[i-1] < i - 1$. Since $w[1..i]$ does not have a border of length $B[i-1] + 1$, we must have $w[B[i-1] + 1] \neq w[i] = w[B[i]]$, so $B[i-1] + 1 - B[i]$ is not a period of $w[1..i-1]$ despite the fact that both $i - B[i]$ and $i - 1 - B[i-1]$ are periods of this prefix. By the Periodicity Lemma, this yields $(i - B[i]) + (i - 1 - B[i-1]) - 1 > i - 1$, so $B[i] + B[i-1] < i - 1$ holds as claimed. \square

Let v be a subword of w and let $v = v'a$ where $a \in \Sigma$. We say that v is a *critical left candidate* if it is a left candidate, but v' is not a left candidate (in particular, v' can be the empty word) or $\min \text{Occ}(v') < \min \text{Occ}(v)$. Let

$$\text{Active}(w) = \{ j \in [1..n] : B[j-1] < B[j] \leq \frac{j-1}{2} \}.$$

The following lemma is the main technical contribution in the characterization of left candidates.

Lemma 3.7. *A fragment $w[i..j]$ is*

(a) *the leftmost occurrence of a left candidate if and only if $i - 1 \leq B[j] \leq j - i$;*

(b) *the leftmost occurrence of a critical left candidate if and only if $i = B[j] + 1$ and $j \in \text{Active}(w)$.*

Proof.

(a) If $w[i..j]$ is the leftmost occurrence of a left candidate v , then $w[1..j] = v^-v$ and $B[j] \geq |v^-| = i - 1$ by Lemma 3.5. Moreover, $B[j] < |v| = j - i + 1$, because $w[i - j + B[j]..B[j]]$ would otherwise be an earlier occurrence of v . These two conditions yield $i - 1 \leq B[j] \leq j - i$.

Conversely, assume that

$$i - 1 \leq B[j] \leq j - i. \quad (1)$$

Suppose for a proof by contradiction that $v = w[i..j]$ has an earlier occurrence at position i' , i.e., $w[i'..i' + j - i] = w[i..j]$ for some position $i' < i$. Let $w' = w[i'..j]$. The word v is a border of w' , so w' has a period $i - i'$, which we denote by p_1 . The shortest period of the whole prefix $w[1..j]$, hence a period of w' , is $j - B[j]$, which we denote by p_2 . By the first inequality of (1),

$$p_1 + p_2 = i - i' + j - B[j] \leq j - i' + 1 = |w'|.$$

Hence, p_1 and p_2 satisfy the assumption of the Periodicity Lemma and w' has period $\gcd(p_1, p_2)$, which we denote by p . Moreover, by the second inequality of (1),

$$p_1 = i - i' < i \leq j - B[j] = p_2,$$

so $p < p_2$ and p divides p_2 . Consequently, $w[1..j]$ has period p , which contradicts the fact that p_2 is its shortest period.

Therefore, $w[i..j]$ indeed is the leftmost occurrence of v and $v^-v = w[1..j]$. Due to $B[j] \geq i - 1 = |v^-|$, we conclude that v is a left candidate by Lemma 3.5.

(b) First, assume that $B[j - 1] < B[j] \leq \frac{j-1}{2}$ and $i = B[j] + 1$.

Note that

$$i - 1 = B[j] = 2B[j] - B[j] \leq j - 1 - B[j] = j - i.$$

By part (a), $w[i..j]$ is therefore the leftmost occurrence of a left candidate. On the other hand, $B[j - 1] < B[j] = i - 1$, so $w[i..j - 1]$ is not the leftmost occurrence of a left candidate. Thus, $w[i..j]$ is the leftmost occurrence of a critical left candidate.

Next, assume that $w[i..j]$ is the leftmost occurrence of a critical left candidate. Part (a) yields $i - 1 \leq B[j] \leq j - i$ (since $w[i..j]$ is the leftmost occurrence of a left candidate) and that $B[j - 1] < i - 1$ or $B[j - 1] \geq j - i$ (since $w[i..j - 1]$ is not).

In the latter case, we have $B[j - 1] + B[j] \geq j - i + i - 1 = j - 1$, so $B[j] = B[j - 1] + 1$ holds by Lemma 3.6. This yields a contradiction:

$$j - i \geq B[j] = B[j - 1] + 1 \geq j - i + 1.$$

Thus, the only possibility is that $B[j - 1] < i - 1$, i.e., $B[j - 1] \leq i - 2$. We then have

$$B[j] \leq B[j - 1] + 1 \leq i - 1 \leq B[j],$$

so $B[j - 1] < B[j] = i - 1$. Furthermore,

$$B[j] = \frac{1}{2}(B[j] + B[j]) \leq \frac{1}{2}(j - i + i - 1) = \frac{j-1}{2}$$

holds as claimed. \square

For a table $A[0..n + 1]$, assuming $A[n + 1] = -\infty$, define the *nearest smallest value* table \mathcal{N}_A such that for $0 \leq i \leq n$ we have: $\mathcal{N}_A[i] = \min\{j > i : A[j] < A[i]\}$.

Lemma 3.8. *Packages $\text{pack}(B[j] + 1, j, \mathcal{N}_B[j] - 1)$ for $j \in \text{Active}(w)$ form a package representation of the family of left candidates of w . This representation (of size at most n) can be computed in $\mathcal{O}(n)$ time.*

Proof. First, we shall prove that for each $j \in \text{Active}(w)$, the fragments $w[B[j] + 1..k]$ for $j \leq k \leq \mathcal{N}_B[j] - 1$ are leftmost occurrences of left candidates. To apply the characterization of Lemma 3.7(a), we need to show that $B[j] \leq B[k] \leq k - B[j] - 1$. The first inequality follows directly from the definition of the \mathcal{N}_B table, while for the second one, we observe that $B[j] + B[k] \leq 2B[j] + k - j \leq j - 1 + k - j = k - 1$ holds due to $B[j] \leq \frac{i-1}{2}$ and the classic property of the border table.

Consequently, the packages are disjoint and consist of left candidates only.

It remains to prove that we do not leave any left candidate behind. Suppose that $w[i..k]$ is the leftmost occurrence of a left candidate. Repeatedly trimming the trailing character, we can reach a critical left candidate whose leftmost occurrence is $w[i..j]$ (for $j \leq k$). By Lemma 3.7(b), we have that $j \in \text{Active}(w)$ and $i = B[j] + 1$. However, since $w[i..k']$ is the leftmost occurrence of a left candidate for all $j \leq k' \leq k$, Lemma 3.7(a) yields $B[k'] \geq i - 1 = B[j]$ for $j \leq k' \leq k$. Consequently, $\mathcal{N}_B[j] > k$. Thus, $w[i..k] \in \text{pack}(B[j] + 1, j, \mathcal{N}_B[j] - 1)$ indeed belongs to one of the packages we created.

The size of the package representation is $|\text{Active}(w)| \leq n$. As for the $\mathcal{O}(n)$ -time construction algorithm, we first build the border table B [39, 13, 14] and its nearest smallest value table \mathcal{N}_B (using a Cartesian tree construction algorithm [19]). Then, for each position j , we test in constant time whether $j \in \text{Active}(w)$ and, if so, we retrieve the corresponding package. \square

3.3 Right Candidates

For word w , let us define a *reverse border array* $B^R[1..n]$ such that $B^R[i]$ is the length of the longest proper border of $w[i..n]$.

Lemma 3.9. *A subword v of a word w is a right candidate of w if and only if $B^R[n - |vv^+| + 1] \geq |v^+|$.*

Proof. Follows from Lemma 3.5 by the symmetry between B and B^R as well as left and right candidates. \square

Even though Lemma 3.9 for right candidates and Lemma 3.5 for left candidates are symmetric, the former allows us to represent right candidates on each edge of the suffix tree of w as a single package. Thus a representation for right candidates is much simpler to be computed than the representation for left candidates.

Lemma 3.10. *The intersection of $\text{RCands}(w)$ with a single equivalence class forms at most one package. Moreover, a package representation (of size at most $2n$) of the set $\text{RCands}(w)$ can be computed in $\mathcal{O}(n)$ time.*

Proof. Let us consider an equivalence class E of subwords of w , with $P = \text{Occ}(v)$ for $v \in E$, and denote

$$\ell(E) = n - \max P + 1 - B^R[\max P].$$

We will show that $v \in E$ is a right candidate if and only if $|v| \geq \ell(E)$. By Lemma 3.9, v is a right candidate if and only if $B^R[n - |vv^+| + 1] \geq |v^+|$. However, $|vv^+| = n - \max P + 1$ and $|v^+| = n - \max P + 1 - |v|$, so the two conditions are equivalent.

Let us show how to compute $\text{RCands}(w) \cap E$ for each equivalence class $E = \text{pack}(i, j_1, j_2)$. We construct the array B^R [39, 13, 14] and the suffix tree of w . Moreover, for each equivalence class E , we determine the common value $\max(\text{Occ}(v))$ for $v \in E$. This lets us compute $\ell(E)$. We have proved that the right candidates in E are precisely words in E of length $\ell(E)$ or more. If $\ell(E) > |w[i..j_2]|$, there are no right candidates in E . Otherwise, we output a package

$$\text{pack}(i, \max(j_1, i + \ell(E) - 1), j_2). \quad \square$$

3.4 Representation Theorem for Seeds

Theorem 3.11. *For a word w of length n , the set $\text{Seeds}(w)$ has a package representation of size at most $3n$.*

Proof. By Lemma 3.2, we have $\text{Seeds}(w) = \text{LCands}(w) \cap \text{QSeeds}(w) \cap \text{RCands}(w)$.

First, we shall prove the package representation of $\text{QSeeds}(w) \cap \text{RCands}(w)$ consists of at most $2n$ packages. Indeed, for each equivalence class E , both $\text{QSeeds}(w) \cap E$ and $\text{RCands}(w) \cap E$ form at most one package. The intersection of two packages forms at most one package, so $\text{QSeeds}(w) \cap \text{RCands}(w)$ has a package representation with at most one package per equivalence class, i.e., of total size at most $2n$.

By Lemma 3.8, $\text{LCands}(w)$ has a package representation of size at most n .

Finally, Lemma 2.2 implies that $\text{Seeds}(w)$ has a package representation of size at most $3n$. \square

4 Computing Long Seeds

In this section, we provide a linear-time algorithm computing seeds of length $\Theta(n)$. Formally, we consider the following operation for $\ell = \Omega(n)$.

LONG-SEEDS(ℓ, w): Computes a package representation of $\text{Seeds}_{\geq \ell}(w)$.

Let us denote by $\text{QSeeds}_{\geq \ell}(w)$ the set of all quasisseeds of w with lengths at least ℓ . Our implementation is based on Theorem 3.11, and it uses the suffix tree of w to determine $\text{QSeeds}_{\geq \ell}(w)$.

Let us partition the positions of w into a family \mathcal{F} of $\mathcal{O}(n/\ell)$ disjoint *blocks* of length at most ℓ each. For a set of positions X , we define its refined version:

$$\text{refine}(X) = \bigcup_{Y \in \mathcal{F}: Y \cap X \neq \emptyset} \{\min(Y \cap X), \max(Y \cap X)\};$$

see Fig. 6. Note that $|\text{refine}(X)| \leq 2|\mathcal{F}| = \mathcal{O}(n/\ell) = \mathcal{O}(1)$.

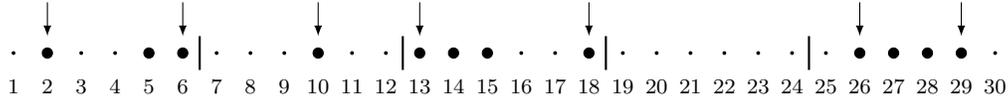


Figure 6: For $n = 30$ and $\ell = 6$, $\text{refine}(\{2, 5, 6, 10, 13, 14, 15, 18, 26, 27, 28, 29\}) = \{2, 6, 10, 13, 18, 26, 29\}$.

Lemma 4.1. *A subword v of length at least ℓ is a quasisseed if and only if $\text{maxgap}(\text{refine}(\text{Occ}(v))) \leq |v|$.*

Proof. Clearly, $\text{refine}(\text{Occ}(v)) \subseteq \text{Occ}(v)$, so $\text{maxgap}(\text{refine}(\text{Occ}(v))) \geq \text{maxgap}(\text{Occ}(v))$. Due to Observation 3.3, it remains to prove that $\text{maxgap}(\text{refine}(\text{Occ}(v))) > |v|$ yields $\text{maxgap}(\text{Occ}(v)) > |v|$. Let $i < i'$ be consecutive elements of $\text{refine}(\text{Occ}(v))$ such that $i' - i > |v|$. Since $|v| \geq \ell$, positions i and i' belong to different blocks of \mathcal{F} . Moreover, these positions are consecutive elements of $\text{refine}(\text{Occ}(v))$, so i must be the largest in its block, i' must be the smallest in its block, and all blocks in between cannot contain any element of $\text{Occ}(v)$. Consequently, $i < i'$ are consecutive elements of $\text{Occ}(v)$, i.e., $\text{maxgap}(\text{Occ}(v)) > |v|$. \square

We traverse the suffix tree of w bottom-up, computing $\text{refine}(\text{Occ}(v))$ in constant time for each subword v whose locus is an explicit node.

Fact 4.2. *The sets $\text{refine}(\text{Occ}(v))$ for all explicit nodes v of the suffix tree can be computed in $\mathcal{O}(n)$ time.*

Proof. We compute $\text{refine}(\text{Occ}(v))$ for each explicit node v and store it as a sorted list. For this, we start with an empty set and consider all explicit children u of v . For each child, we merge the current list with $\text{refine}(\text{Occ}(u))$, removing elements which are not extremes in their blocks. This takes $\mathcal{O}(|\mathcal{F}|) = \mathcal{O}(1)$ time for each edge of the suffix tree, which gives $\mathcal{O}(n)$ time in total. \square

Lemma 4.3 (LONG-SEEDS Implementation). For a threshold $\ell = \Theta(n)$ and a word w of length n , $\text{LONG-SEEDS}(\ell, w)$ can be implemented in $\mathcal{O}(n)$ time.

Proof. First, we compute a package representation of the family $\text{QSeeds}_{\geq \ell}(w)$ of long quasiseeds. Consider an equivalence class $E = \text{pack}(i, j_1, j_2)$. Let P be the common occurrence set $\text{Occ}(v)$ for $v \in E$. The longest subword in each package is represented by an explicit node, so the procedure of Fact 4.2 computes $\text{refine}(P)$. Let us define $\ell' := \max(\ell, \max\text{gap}(\text{refine}(P)))$. By Lemma 4.1, a subword $v \in E$ is a long quasiseed if and only if $|v| \geq \ell'$. If $j_2 - i + 1 < \ell'$, there are no such quasiseeds. Otherwise, we report a package

$$E \cap \text{QSeeds}_{\geq \ell}(w) = \text{pack}(i, \max(j_1, i + \ell' - 1), j_2).$$

Finally, we apply Lemma 3.2 and compute $\text{Seeds}_{\geq \ell}(w) = \text{QSeeds}_{\geq \ell}(w) \cap \text{LCands}(w) \cap \text{RCands}(w)$ using Lemma 2.2 to implement the intersection. Linear-size package representations of left candidates and right candidates are constructed using Lemmas 3.8 and 3.10, respectively. The total running time is $\mathcal{O}(n)$. \square

We also use the following operation that can be computed using LONG-SEEDS.

INT-SEEDS(I, w): Computes a package representation of $\text{Seeds}_I(w)$.

Our implementation of $\text{INT-SEEDS}(I, w)$ uses a reduction to Lemmas 2.2 and 4.3, and its running time is linear provided that I is *balanced*, i.e., if the ratio of its endpoints is bounded by a constant.

Lemma 4.4 (INT-SEEDS Implementation). For an interval $I = [\ell..r]$ and a word w of length n , $\text{INT-SEEDS}(I, w)$ can be implemented in $\mathcal{O}(n)$ time if $r = \mathcal{O}(\ell)$.

Proof. We construct a family R of fragments of length $4r$ covering w with overlaps of size $2(r - 1)$ (the last fragment might be shorter). Note that the total length of these fragments is at most $2n$. Furthermore, observe that $\text{SUB}_{2r-1}(w) = \bigcup_{s \in R} \text{SUB}_{2r-1}(s)$, so Corollary 2.6 yields

$$\text{Seeds}_{\leq r}(w) = \bigcap_{u \in \text{SUB}_{2r-1}(w)} \text{Seeds}_{\leq r}(u) = \bigcap_{s \in R} \bigcap_{u \in \text{SUB}_{2r-1}(s)} \text{Seeds}_{\leq r}(u) = \bigcap_{s \in R} \text{Seeds}_{\leq r}(s).$$

In particular, $\text{Seeds}_I(w) = \bigcap_{s \in R} \text{Seeds}_I(s)$. For each $s \in R$, we apply Lemma 4.3 to determine a package representation of $\text{Seeds}_{\geq \ell}(s)$, and we filter out seeds of length greater than r . This takes $\mathcal{O}(|s|)$ time for each $s \in R$, which is $\mathcal{O}(n)$ in total. Finally, we combine the package representations in $\mathcal{O}(n)$ time using Lemma 2.2. \square

5 Relation between Seeds, Compression and Subword Complexity

The *subword complexity* of a word w is a function which gives the number of subwords of a given length k , i.e., $|\text{SUB}_k(w)|$. Since $|\text{SUB}_k(w)|$ is not monotone in general, we define a non-decreasing sequence $(\beta_k)_{k=1}^n$ with

$$\beta_k(w) = |\text{SUB}_k(w)| + k - 1;$$

see Fig. 7.

k	1	2	3	4	5	6	7	8
$w[k]$	a	b	a	b	a	a	b	b
$ \text{SUB}_k(w) $	2	4	5	5	4	3	2	1
$\beta_k(w)$	2	5	7	8	8	8	8	8

Figure 7: Subword complexity of w and $\beta_k(w)$ for $w = \text{ababaabb}$.

A connection between the subword complexity and the existence of seeds of certain lengths is crucial in the paper. More precisely, each seed provides an upper bound on the values $\beta_k(w)$.

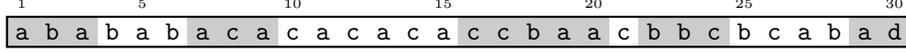


Figure 8: $\text{COMPR}_2(w) = \{w[1..3], w[7..9], w[16..20], w[22..24], w[29..30]\} = \{\text{aba}, \text{aca}, \text{ccbaa}, \text{bbc}, \text{ad}\}$. We have $\text{SUB}_2(\text{aba}) \cup \text{SUB}_2(\text{aca}) \cup \text{SUB}_2(\text{ccbaa}) \cup \text{SUB}_2(\text{bbc}) \cup \text{SUB}_2(\text{ad}) = \{\text{ab}, \text{ba}\} \cup \{\text{ac}, \text{ca}\} \cup \{\text{cc}, \text{cb}, \text{ba}, \text{aa}\} \cup \{\text{bb}, \text{bc}\} \cup \{\text{ad}\} = \text{SUB}_2(w)$.

Lemma 5.1 (Gap Lemma). *If $\beta_k(w) > \frac{2}{3}n$, then w has no seed v whose length satisfies $2k - 2 \leq |v| \leq \frac{1}{6}n$.*

Proof. We shall prove that if v is a seed of w , then $\beta_k(w) \leq |v| + (k - 1)\frac{n}{|v|}$.

Note that we may assume $|v| \geq k$; otherwise, the right-hand side is at least n . Consider length- k fragments starting at positions $i, \dots, i + \ell$, where $\ell < |v|$. We claim that at most $k - 1$ of these fragments are not covered by single occurrences of v . Let $w[i'..i' + k - 1]$, for $i \leq i' \leq i + \ell$, be the first such fragment that is not covered by any single occurrence of v . If i' does not exist or $i' + k - 1 > i + \ell$, we are done. Otherwise, the occurrence of v covering position $i' + k - 1$ must start at some position j , $i' < j < i' + k$. Consequently, the length- k fragments starting at positions i'' , $j \leq i'' \leq j + |v| - k$, are all covered by this occurrence of v . We are left with at most $k - 1$ remaining length- k fragments (starting at positions i'' such that $i' \leq i'' < j$ or $j + |v| - k < i'' \leq i + \ell$).

Thus, at most $(k - 1)\lceil \frac{n - |v| + 1}{|v|} \rceil$ length- k fragments of w are not covered by single occurrences of v . As a result, $|\text{SUB}_k(w) \setminus \text{SUB}_k(v)| \leq (k - 1)\lceil \frac{n - |v| + 1}{|v|} \rceil$, and we obtain the claimed upper bound on $\beta_k(w)$:

$$\beta_k(w) = k - 1 + |\text{SUB}_k(w)| = k - 1 + |\text{SUB}_k(v)| + |\text{SUB}_k(w) \setminus \text{SUB}_k(v)| \leq |v| + (k - 1)\lceil \frac{n - |v| + 1}{|v|} \rceil \leq |v| + (k - 1)\frac{n}{|v|}.$$

If $2k - 2 \leq |v| \leq \frac{1}{6}n$, then we conclude that

$$\beta_k(w) \leq |v| + (k - 1)\frac{n}{|v|} \leq \frac{1}{6}n + \frac{1}{2}n = \frac{2}{3}n,$$

which contradicts our assumption. □

Lemma 5.1 yields a gap in the feasible lengths of seeds of w . Seeds of length $\frac{1}{6}n$ or more can be determined using the LONG-SEEDS procedure, so we may focus on computing relatively short seeds. Due to the characterization of Lemma 2.5, we may ignore some regions of w as long as we do not miss any subword of certain length. To formalize this intuition, we define the following operation $\text{COMPR}_k(w)$, which results in a set of subwords S of w such that $\text{SUB}_k(w) = \bigcup_{s \in S} \text{SUB}_k(s)$.

Let us take the set of fragments which are leftmost occurrences of subwords in $\text{SUB}_k(w)$. Any two overlapping or consecutive fragments are joined together, and the subwords indicated by the resulting set of fragments give the family $\text{COMPR}_k(w)$; see Fig. 8 for an example.

Recall that our motivation is to reduce computing short seeds in w to the analogous operation in each $s \in \text{COMPR}_k(w)$. This is illustrated by the following result.

Lemma 5.2 (Reduction Lemma).

Consider a word w of length n . For every integer k , $1 \leq k \leq \frac{n+1}{2}$, we have

$$\text{Seeds}_{\leq k}(w) = \bigcap_{s \in \text{COMPR}_{2k-1}(w)} \text{Seeds}_{\leq k}(s).$$

Proof. Note that $\text{SUB}_{2k-1}(w) = \bigcup_{s \in \text{COMPR}_{2k-1}(w)} \text{SUB}_{2k-1}(s)$. Hence, Corollary 2.6 yields

$$\text{Seeds}_{\leq k}(w) = \bigcap_{u \in \text{SUB}_{2k-1}(w)} \text{Seeds}_{\leq k}(u) = \bigcap_{s \in \text{COMPR}_{2k-1}(w)} \bigcap_{u \in \text{SUB}_{2k-1}(s)} \text{Seeds}_{\leq k}(u) = \bigcap_{s \in \text{COMPR}_{2k-1}(w)} \text{Seeds}_{\leq k}(s).$$

This concludes the proof. □

The values $\beta_k(w)$ can be used to bound the total length of words $s \in \text{COMPR}_k(w)$, denoted $\|\text{COMPR}_k(w)\|$.

Lemma 5.3 (Compression Lemma). *For each word w and integer k , $1 \leq k \leq \frac{n+1}{2}$, we have*

$$\|\text{COMPR}_k(w)\| \leq \beta_{2k-1}(w).$$

Proof. Let P be the set of positions covered by the leftmost occurrences of subwords from $\text{COMPR}_k(w)$. By construction, $\|\text{COMPR}_k(w)\| = |P|$ and $i \in P$ if and only if i is included in the leftmost occurrence of some subword $s \in \text{SUB}_k(w)$. If $k \leq i \leq n - k + 1$, then i is the midpoint of a length- $(2k - 1)$ fragment $w[i - k + 1 .. i + k - 1]$, which covers all length- k fragments of w containing position i including the leftmost occurrence of $s \in \text{SUB}_k(w)$. Thus, $w[i - k + 1 .. i + k - 1]$ is also the leftmost occurrence of the corresponding subword of length $2k - 1$. This yields an injective mapping from the set of positions $i \in P \cap [k .. n - k + 1]$ to the family $\text{SUB}_{2k-1}(w)$. The remaining positions ($i < k$ and $i > n - k + 1$) account for the extra term $2k - 2 = \beta_{2k-1}(w) - |\text{SUB}_{2k-1}(w)|$ in the upper bound. \square

From the last two lemmas we obtain the following corollary that conveys the intuition behind the main structural theorem in the following section.

Corollary 5.4. *Let $1 \leq k \leq \frac{n+3}{4}$ and $S = \text{COMPR}_{2k-1}(w)$. Then*

$$\text{Seeds}_{\leq k}(w) = \bigcap_{s \in S} \text{Seeds}_{\leq k}(s) \quad \text{and} \quad \|S\| \leq \beta_{4k-3}(w).$$

6 Main Algorithm

The main algorithm is a recursive procedure based on a structural theorem which combines the results of the previous section.

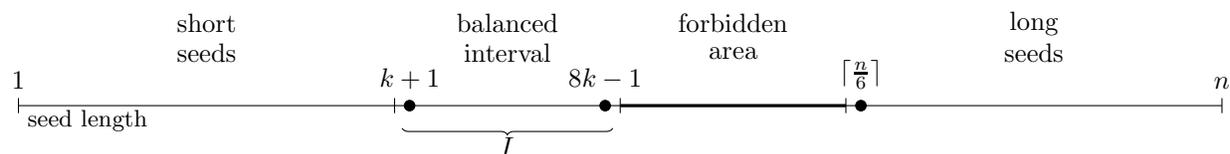


Figure 9: Illustration of Theorem 6.1(A) in case that $8k \leq \lceil \frac{n}{6} \rceil$. There is no seed of length in the forbidden area. The computation of all seeds is split into recursive computation of short seeds (of length at most k) in subwords $s \in S$, of seeds with lengths in a balanced interval I , and of long seeds.

Theorem 6.1 (Decomposition Theorem).

Consider a word w of length n . If $n \leq 6$ or $|\text{Alph}(w)| > \frac{2}{3}n$, then $\text{Seeds}(w) = \text{Seeds}_{\geq \frac{1}{6}n}(w)$. Otherwise, we can compute in linear time a balanced interval I and a family S of subwords of w such that

$$\text{(A)} \quad \text{Seeds}(w) = \text{Seeds}_{\geq \frac{1}{6}n}(w) \cup \text{Seeds}_I(w) \cup \bigcap_{s \in S} \text{Seeds}_{\leq k}(s) \quad \text{and} \quad \text{(B)} \quad \|S\| \leq \frac{2}{3}n.$$

Proof. We take

$$\begin{aligned} k &= \max \{ \ell : 1 \leq \ell < \lceil \frac{n}{6} \rceil \text{ and } \beta_{4\ell-3}(w) \leq \frac{2}{3}n \}, \\ S &= \text{COMPR}_{2k-1}(w), \\ I &= [k + 1 .. \min(8k, \lceil \frac{n}{6} \rceil) - 1]. \end{aligned}$$

As for the first part of the statement, note that $\text{Seeds}(w) = \text{Seeds}_{\geq \frac{1}{6}n}(w)$ holds trivially if $n \leq 6$. On the other hand, if $|\text{Alph}(w)| = \beta_1(w) > \frac{2}{3}n$, then the equality follows directly from the Gap Lemma (Lemma 5.1).

Hence we can assume further that $n > 6$ and $\beta_1(w) \leq \frac{2}{3}n$.

Correctness of (B). Now we can guarantee that k is well-defined and that it satisfies $\beta_{4k-3}(w) \leq \frac{2}{3}n$. Then Compression Lemma (Lemma 5.3) implies that S satisfies the condition **(B)**.

Correctness of (A). To prove condition **(A)**, we observe that the Reduction Lemma (Lemma 5.2) yields $\text{Seeds}_{\leq k}(w) = \bigcap_{s \in S} \text{Seeds}_{\leq k}(s)$. Thus, it is enough to prove that $\text{Seeds}_{\geq k+1}(w) = \text{Seeds}_{\geq \frac{1}{6}n}(w) \cup \text{Seeds}_I(w)$, i.e., that there is no seed v with $8k \leq |v| < \lceil \frac{n}{6} \rceil$. This claim holds trivially for $k = \lceil \frac{n}{6} \rceil - 1$. Otherwise, $\beta_{4k+1} > \frac{2}{3}n$, and the claim follows directly from the Gap Lemma (Lemma 5.1).

Computing k and S . Let us start with the following claim.

Claim.

The sequence $|\text{SUB}_m(w)|$ (consequently, also the sequence $\beta_m(w)$) can be computed in linear time.

Proof. Each equivalence class $E = \text{pack}(i, j_1, j_2)$ contributes a single subword in $\text{SUB}_m(w)$ for each $m \in [j_1 - i + 1 .. j_2 - i + 1]$. We obtain at most $2n$ such intervals; $|\text{SUB}_m(w)|$ is then the number of intervals that contain the element m . This quantity can be computed in $\mathcal{O}(n)$ time using an auxiliary array A of size n (initially set to zeroes). For an interval $[a .. b]$, $A[a]$ is incremented and $A[b+1]$ is decremented. Then $|\text{SUB}_m(w)| = A[1] + \dots + A[m]$. \square

To retrieve $S = \text{COMPR}_{2k-1}(w)$, we iterate again over the equivalence classes $E = \text{pack}(i, j_1, j_2)$. If $2k-1 \in [j_1 - i + 1 .. j_2 - i + 1]$, we mark the position $\min P$, where $P = \text{Occ}(v)$ for $v \in E$. Next, we scan the text marking the first $2k-2$ positions in each gap between subsequent already marked positions and the first $2k-2$ positions after the final already marked position. After these two phases, we build a word $s \in S$ out of each maximal region of marked positions. \square

The algorithm computing seeds relies on Theorem 6.1, with the LONG-SEEDS procedure applied to compute $\text{Seeds}_{\geq \frac{1}{6}n}(n)$, and recursive calls made to determine $\text{Seeds}(s)$ for each $s \in S$. Finally, a package representation of $\text{Seeds}_I(w)$ is computed using INT-SEEDS.

Algorithm 1: Recursive procedure SEEDS(w)

Input: A word w of length n .

Output: An $\mathcal{O}(n)$ -size package representation of $\text{Seeds}(w)$.

if $n \leq 6$ **or** $|\text{Alph}(w)| > \frac{2}{3}n$ **then return** LONG-SEEDS($\lceil \frac{n}{6} \rceil, w$)

Let I, S be as in the proof of the Decomposition Theorem

foreach $v \in S$ **do**

$R_v := \text{SEEDS}(v)$

$R := \text{Combine}(\{R_v : v \in S\})$

Remove from R seeds of length at least $\min I$

return $R \cup \text{LONG-SEEDS}(\lceil \frac{n}{6} \rceil, w) \cup \text{INT-SEEDS}(I, w)$

(The three package representations contain distinct seeds)

Theorem 6.2 (Main Result). *An $\mathcal{O}(n)$ -size package representation of the set $\text{Seeds}(w)$ of all seeds of a given length- n word w can be found in $\mathcal{O}(n)$ time. In particular, a shortest seed and the total number of seeds can be computed within the same time complexity.*

Proof. Correctness of the algorithm SEEDS follows immediately from Theorem 6.1. To bound the running time, let us denote by $T(n)$ the maximum number of operations performed by the SEEDS function executed

for a word of length n . Due to Theorem 6.1 and Lemmas 2.2, 4.3 and 4.4,

$$T(n) = \mathcal{O}(n) + \sum_i T(n_i), \quad \text{where } \sum_i n_i \leq \frac{2}{3}n.$$

This recurrence yields $T(n) = \mathcal{O}(n)$. □

7 Implementation of the Operation Combine

We describe here the missing part of the algorithm, which is based on some rather technical computations on weighted trees. First, our algorithm requires an efficient solution to the so-called weighted ancestor queries problem. Thus, in Section 7.1 we present an optimal offline procedure answering weighted ancestor queries in an arbitrary weighted tree with polynomially bounded weights.

7.1 Offline Weighted Ancestor Queries

In the weighted ancestor problem, introduced by Farach and Muthukrishnan [16] (see also [21]), we consider a rooted tree T with an integer weight function `weight` defined on the nodes. The weight of the root must be zero and the weight of any other node must be strictly larger than the weight of its parent. A classic example is any compacted trie with the weight of a node defined as the length of its value.

The weighed ancestor queries, given a node ν and an integer value $\ell \leq \text{weight}(\nu)$, ask for the highest ancestor μ of ν such that $\text{weight}(\mu) \geq \ell$, i.e., such an ancestor μ that $\text{weight}(\mu) \geq \ell$ and $\text{weight}(\mu)$ is smallest possible. We denote the node μ as `ancestor`(ν, ℓ).

Weighted ancestor queries in the online setting can be answered in $\mathcal{O}(\log \log n)$ time after $\mathcal{O}(n)$ -time preprocessing [1]. In the special case of the weighted tree being a suffix tree of a word, they can be answered in $\mathcal{O}(1)$ time with a data structure of $\mathcal{O}(n)$ space [21]. Nevertheless, no efficient construction of this data structure is known. Below we show that if all the queries are given offline and the weights are polynomially bounded, then q queries can be answered in $\mathcal{O}(n + q)$ time.

Let us first recall the classic union-find data structure. It maintains a partition \mathcal{S} of $[1..n]$. Each set $S \in \mathcal{S}$ has an identifier $\text{id}(S) \in S$. Initially, \mathcal{S} is a partition into singletons and $\text{id}(\{i\}) = i$ for $1 \leq i \leq n$.

The union-find data structure supports `find`(i) queries which, given an integer $i \in [1..n]$, return the identifier $\text{id}(S)$ of the set $S \in \mathcal{S}$ containing i . Moreover, a `union`(i_1, i_2) operation, given integers $i_1, i_2 \in [1..n]$, replaces the sets $S_1, S_2 \in \mathcal{S}$ such that $i_1 \in S_1$ and $i_2 \in S_2$ with their union $S_1 \cup S_2$. Note that `union`(i_1, i_2) is void if $S_1 = S_2$, i.e., if `find`(i_1) = `find`(i_2).

We will only encounter *linear* union-find instances, in which sets $S \in \mathcal{S}$ are formed by consecutive integers and $\text{id}(S) = \min(S)$. In other words, `union`(i_1, i_2) is allowed for $i_1 < i_2$ only if `find`(i_1) = `find`(`find`(i_2) - 1). For such instances, the union-find operations can be implemented in amortized $\mathcal{O}(1)$ time.

Lemma 7.1 (Gabow and Tarjan [20]). *A sequence of m linear union-find operations on a partition of $[1..n]$ can be implemented in $\mathcal{O}(n + m)$ time.*

We are now ready to describe an efficient offline procedure answering weighted ancestor queries.

Lemma 7.2. *Given a collection Q of weighted ancestor queries on a weighted tree T on n nodes with integer weights up to $(n + |Q|)^{\mathcal{O}(1)}$, all the queries from Q can be answered in $\mathcal{O}(n + |Q|)$ time.*

Proof. We process the tree and the queries according to non-increasing weights. We maintain a union-find data structure which stores a partition of the set $V(T)$ of nodes of T . After the nodes with weight ℓ have been processed, each partition class is either a singleton of a node μ such that $\text{weight}(\mu) < \ell$, or consists of the nodes of a subtree rooted at a node μ such that $\text{weight}(\mu) \geq \ell$. In either case, μ is the identifier of the set.

Note that, in order to update the partition for the next smaller value of ℓ , for each node μ with weight ℓ it suffices to union the singleton $\{\mu\}$ with the node sets of subtrees rooted at the children of μ . Moreover,

observe that after processing nodes at level ℓ , for each node ν , its ancestors μ with $\text{weight}(\mu) < \ell$ form singletons, while ancestors μ with $\text{weight}(\mu) \geq \ell$ belong to the same class as ν . Hence, the identifier of this class is the highest ancestor of μ with $\text{weight}(\mu) \geq \ell$, i.e., $\text{ancestor}(\nu, \ell) = \text{find}(\nu)$. Consequently, the weighted ancestor queries can be answered using Algorithm 2.

Algorithm 2: Offline weighted ancestors for a weighted tree T and a set of queries Q

```

 $W_T = \{\text{weight}(\mu) : \mu \in V(T)\};$ 
 $W_Q = \{\ell : (\nu, \ell) \in Q\};$ 
foreach  $\ell \in W_T \cup W_Q$  in the decreasing order do
  foreach  $\mu \in V(T) : \text{weight}(\mu) = \ell$  do
    foreach  $\nu : \text{child of } \mu \text{ in the left-to-right order}$  do
       $\text{union}(\mu, \nu);$ 
    foreach  $\nu : (\nu, \ell) \in Q$  do
      Report  $\text{ancestor}(\nu, \ell) := \text{find}(\nu);$ 

```

Next, we shall prove that Algorithm 2 can be implemented in $\mathcal{O}(n + |Q|)$ time. Since node weights and query weights are integers bounded by $(n + |Q|)^{\mathcal{O}(1)}$, they can be sorted using radix sort in $\mathcal{O}(n + |Q|)$ time. For union-find operations, we need to identify nodes with integers $[1..n]$. We use the pre-order identifiers as they guarantee that each partition class (which can be either a singleton or the node set of a subtree rooted at a given node) consists of consecutive integers. With $n - 1$ **union** operations and $|Q|$ **find** operations, Lemma 7.1 guarantees $\mathcal{O}(n + |Q|)$ overall running time of the union-find data structure. \square

We apply offline weighted ancestor queries to the suffix tree to obtain the following algorithmic tool.

Corollary 7.3. *Given a collection of subwords s_1, \dots, s_k of a word w of length n , each represented by an occurrence in w , in $\mathcal{O}(n + k)$ total time we can compute the locus of each subword s_i in the suffix tree of w . Moreover, these loci can be made explicit in $\mathcal{O}(n + k)$ extra time.*

Proof. Let T be the suffix tree of w . Assume that $s_i = w[a..b]$ and consider the following nodes: the node μ representing $w[a..n]$ (the terminal node of T annotated with a) and $\nu = \text{ancestor}(\mu, b - a + 1)$. If we denote by d the depth of ν , then the locus of s_i is $(\nu, d - (b - a + 1))$. By Lemma 7.2, the loci of s_i can be computed in $\mathcal{O}(n + k)$ time.

In order to make the corresponding implicit nodes explicit, we need to group them according to the nearest explicit descendant and sort them by distances from that node. This can be implemented in $\mathcal{O}(n + k)$ time via radix sort. Then we simply create the explicit nodes in the obtained order. \square

7.2 Intersection of Families of Paths in a Tree

Let us introduce one more abstract operation on a rooted tree T . A *path family* is a family of pairwise disjoint paths in T , each leading downwards. A path family is called *minimal* if there is no other path family covering the same set of nodes in T and consisting of a smaller number of paths.

Let P_1, \dots, P_m be path families in T . Then by $\text{PATHS}_T(P_1, \dots, P_m)$ we denote a minimal path family representing the nodes covered by all the families P_1, \dots, P_m ; see Fig. 10.

Lemma 7.4. *The family $\text{PATHS}_T(P_1, \dots, P_m)$ can be computed in linear time with respect to the size of the tree T , the number m , and total number of paths in all families P_i .*

Proof. For each node μ in T , we would like to compute a value, denoted $S[\mu]$, that is equal to the number of paths in all families P_i that contain node μ . Observe that a node μ is covered by all the families P_i if and only if $S[\mu] = m$.

For each node ν in T , we will store a counter $C[\nu]$ so that, for a node μ , $S[\mu]$ is equal to the sum of values $C[\nu]$ across the nodes ν in the subtree of μ . The counters C are initially set to zeroes. For each path

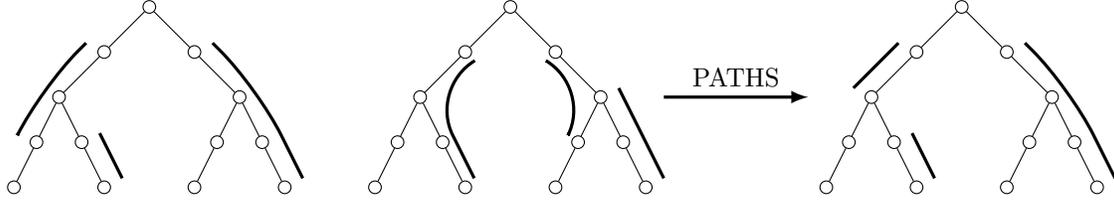


Figure 10: The $\text{PATHS}_T(P_1, P_2)$ operation: three copies of the same rooted tree T , the first two show P_1 and P_2 , and the third one shows $\text{PATHS}_T(P_1, P_2)$.

leading from μ down to ν in P_i , we increment the counter C at ν and decrement the counter at the parent of μ (unless μ is the root). Next, for each node μ , we compute $S[\mu]$ as the sum of values $C[\nu]$ across the nodes ν in the subtree of μ . This is done in a bottom-up fashion using the $S[\cdot]$ values of the children of μ .

This way, we compute all the nodes represented by $\text{PATHS}_T(P_1, \dots, P_m)$. Finally, in order to find the minimal path family covering all these nodes, we repeat the following process traversing the tree in a bottom-up order: If the value $S[\mu]$ is m , we create a single-node path $\{\mu\}$. If also $S[\nu] = m$ for a child ν of μ , we merge the paths containing these two nodes. (We choose the child arbitrarily if $S[\nu] = m$ for several children.) \square

7.3 Proof of Lemma 2.2

We are now ready to provide an efficient implementation of the Combine operation.

Lemma 2.2. *For a word w of length n and sets R_1, \dots, R_k of subwords of w , given in package representations of total size N , $\text{Combine}(R_1, \dots, R_k)$ can be implemented in $\mathcal{O}(n + N)$ time. The size of the resulting package representation is at most N .*

Proof. We reduce the problem to computing $\text{PATHS}_T(P_1, \dots, P_m)$ for some path families P_1, \dots, P_m .

We first apply Corollary 7.3 for subwords $w[i..j_1]$ and $w[i..j_2]$ across all packages $\text{pack}(i, j_1, j_2)$ in R_1, \dots, R_m , extending the set of explicit nodes of the suffix tree T of w by the obtained loci. For each package, we create a path connecting the corresponding two loci. The packages in a package representation are disjoint, so for each R_i this process results in a path family. We apply Lemma 7.4 to compute a minimal path family $P = \text{PATHS}_T(P_1, \dots, P_m)$.

Finally, for each path in P , we create a package in T . We assume that each node stores the label ℓ of any terminal node in its subtree. If a path in P connects two nodes at depths $d_1 \leq d_2$, and the second one stores the label ℓ , then we create a package $\text{pack}(\ell, \ell + d_1 - 1, \ell + d_2 - 1)$. \square

Acknowledgements The authors thank Patryk Czajka for suggesting a simplification of the algorithm in Section 3.2 (Lemma 3.7).

Jakub Radoszewski was supported by the ‘‘Algorithms for text processing with errors and uncertainties’’ project carried out within the HOMING program of the Foundation for Polish Science co-financed by the European Union under the European Regional Development Fund.

References

- [1] Amihood Amir, Gad M. Landau, Moshe Lewenstein, and Dina Sokol. Dynamic text and static pattern matching. *ACM Trans. Algorithms*, 3(2):19, 2007. doi:10.1145/1240233.1240242.
- [2] Amihood Amir, Avivit Levy, Moshe Lewenstein, Ronit Lubin, and Benny Porat. Can we recover the cover? In Juha Kärkkäinen, Jakub Radoszewski, and Wojciech Rytter, editors, *28th Annual*

Symposium on Combinatorial Pattern Matching, CPM 2017, volume 78 of *LIPICs*, pages 25:1–25:15. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2017. doi:10.4230/LIPICs.CPM.2017.25.

- [3] Amihod Amir, Avivit Levy, Ronit Lubin, and Ely Porat. Approximate cover of strings. In Juha Kärkkäinen, Jakub Radoszewski, and Wojciech Rytter, editors, *28th Annual Symposium on Combinatorial Pattern Matching, CPM 2017*, volume 78 of *LIPICs*, pages 26:1–26:14. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2017. doi:10.4230/LIPICs.CPM.2017.26.
- [4] Amihod Amir, Avivit Levy, and Ely Porat. Quasi-periodicity under mismatch errors. In Gonzalo Navarro, David Sankoff, and Binhai Zhu, editors, *Annual Symposium on Combinatorial Pattern Matching, CPM 2018*, volume 105 of *LIPICs*, pages 4:1–4:15. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2018. doi:10.4230/LIPICs.CPM.2018.4.
- [5] Alberto Apostolico and Andrzej Ehrenfeucht. Efficient detection of quasiperiodicities in strings. *Theoretical Computer Science*, 119(2):247–265, 1993. doi:10.1016/0304-3975(93)90159-Q.
- [6] Alberto Apostolico, Martin Farach, and Costas S. Iliopoulos. Optimal superprimitivity testing for strings. *Information Processing Letters*, 39(1):17–20, 1991. doi:10.1016/0020-0190(91)90056-N.
- [7] Omer Berkman, Costas S. Iliopoulos, and Kunsoo Park. The subtree max gap problem with application to parallel string covering. *Information and Computation*, 123(1):127–137, 1995. doi:10.1006/inco.1995.1162.
- [8] Dany Breslauer. An on-line string superprimitivity test. *Information Processing Letters*, 44(6):345–347, 1992. doi:10.1016/0020-0190(92)90111-8.
- [9] Gerth Stølting Brodal and Christian N. S. Pedersen. Finding maximal quasiperiodicities in strings. In Raffaele Giancarlo and David Sankoff, editors, *Combinatorial Pattern Matching, CPM 2000*, volume 1848 of *Lecture Notes in Computer Science*, pages 397–411. Springer, 2000. doi:10.1007/3-540-45123-4_33.
- [10] Manolis Christodoulakis, Costas S. Iliopoulos, Kunsoo Park, and Jeong Seop Sim. Approximate seeds of strings. *Journal of Automata, Languages and Combinatorics*, 10(5/6):609–626, 2005.
- [11] Michalis Christou, Maxime Crochemore, Costas S. Iliopoulos, Marcin Kubica, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Bartosz Szreder, and Tomasz Waleń. Efficient seed computation revisited. *Theoretical Computer Science*, 483:171–181, 2013. doi:10.1016/j.tcs.2011.12.078.
- [12] Richard Cole, Costas S. Iliopoulos, Manal Mohamed, William F. Smyth, and Lu Yang. The complexity of the minimum k-cover problem. *Journal of Automata, Languages and Combinatorics*, 10(5/6):641–653, 2005.
- [13] Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on strings*. Cambridge University Press, 2007. doi:10.1017/cbo9780511546853.
- [14] Maxime Crochemore and Wojciech Rytter. *Jewels of Stringology*. World Scientific, 2003. doi:10.1142/4838.
- [15] Martin Farach-Colton, Paolo Ferragina, and S. Muthukrishnan. On the sorting-complexity of suffix tree construction. *Journal of the ACM*, 47(6):987–1011, 2000. doi:10.1145/355541.355547.
- [16] Martin Farach-Colton and S. Muthukrishnan. Perfect hashing for strings: Formalization and algorithms. In Daniel S. Hirschberg and Eugene W. Myers, editors, *Combinatorial Pattern Matching, CPM 1996*, volume 1075 of *Lecture Notes in Computer Science*, pages 130–140. Springer, 1996. doi:10.1007/3-540-61258-0_11.

- [17] Nathan J. Fine and Herbert S. Wilf. Uniqueness theorems for periodic functions. *Proceedings of the American Mathematical Society*, 16(1):109–114, 1965. doi:10.2307/2034009.
- [18] Tomás Flouri, Costas S. Iliopoulos, Tomasz Kociumaka, Solon P. Pissis, Simon J. Puglisi, William F. Smyth, and Wojciech Tyczyński. Enhanced string covering. *Theoretical Computer Science*, 506:102–114, 2013. doi:10.1016/j.tcs.2013.08.013.
- [19] Harold N. Gabow, Jon Louis Bentley, and Robert Endre Tarjan. Scaling and related techniques for geometry problems. In Richard A. DeMillo, editor, *16th Annual ACM Symposium on Theory of Computing, STOC 1984*, pages 135–143. ACM, 1984. doi:10.1145/800057.808675.
- [20] Harold N. Gabow and Robert E. Tarjan. A linear-time algorithm for a special case of disjoint set union. *Journal of Computer and System Sciences*, 30(2):209–221, 1985. doi:10.1016/0022-0000(85)90014-5.
- [21] Paweł Gawrychowski, Moshe Lewenstein, and Patrick K. Nicholson. Weighted ancestors in suffix trees. In Andreas S. Schulz and Dorothea Wagner, editors, *Algorithms, ESA 2014*, volume 8737 of *Lecture Notes in Computer Science*, pages 455–466. Springer, 2014. doi:10.1007/978-3-662-44777-2_38.
- [22] Qing Guo, Hui Zhang, and Costas S. Iliopoulos. Computing the λ -seeds of a string. In Siu-Wing Cheng and Chung Keung Poon, editors, *Algorithmic Aspects in Information and Management, AAIM 2006*, volume 4041 of *Lecture Notes in Computer Science*, pages 303–313. Springer, 2006. doi:10.1007/11775096_28.
- [23] Qing Guo, Hui Zhang, and Costas S. Iliopoulos. Computing the minimum approximate λ -cover of a string. In Fabio Crestani, Paolo Ferragina, and Mark Sanderson, editors, *String Processing and Information Retrieval, 13th International Conference, SPIRE 2006*, volume 4209 of *Lecture Notes in Computer Science*, pages 49–60. Springer, 2006. doi:10.1007/11880561_5.
- [24] Qing Guo, Hui Zhang, and Costas S. Iliopoulos. Computing the λ -covers of a string. *Information Sciences*, 177(19):3957–3967, 2007. doi:10.1016/j.ins.2007.02.020.
- [25] Costas S. Iliopoulos, Manal Mohamed, and William F. Smyth. New complexity results for the k-covers problem. *Information Sciences*, 181(12):2571–2575, 2011. doi:10.1016/j.ins.2011.02.009.
- [26] Costas S. Iliopoulos, Dennis W. G. Moore, and Kunsoo Park. Covering a string. *Algorithmica*, 16(3):288–297, 1996. doi:10.1007/BF01955677.
- [27] Costas S. Iliopoulos and Laurent Mouchard. Quasiperiodicity: From detection to normal forms. *Journal of Automata, Languages and Combinatorics*, 4(3):213–228, 1999.
- [28] Costas S. Iliopoulos and William F. Smyth. An on-line algorithm of computing a minimum set of k-covers of a string. In *Australasian Workshop on Combinatorial Algorithms, AWOCA 1998*, pages 97–106, 1998.
- [29] Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. Linear work suffix array construction. *Journal of the ACM*, 53(6):918–936, 2006. doi:10.1145/1217856.1217858.
- [30] Tomasz Kociumaka, Marcin Kubica, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. A linear time algorithm for seeds computation. In Yuval Rabani, editor, *23rd Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012*, pages 1095–1112. SIAM, 2012. doi:10.1137/1.9781611973099.
- [31] Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. Efficient algorithms for shortest partial seeds in words. *Theoretical Computer Science*, 710:139–147, 2018. doi:10.1016/j.tcs.2016.11.035.
- [32] Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. Fast algorithm for partial covers in words. *Algorithmica*, 73(1):217–233, 2015. doi:10.1007/s00453-014-9915-3.

- [33] Roman M. Kolpakov and Gregory Kucherov. Finding maximal repetitions in a word in linear time. In *40th Annual Symposium on Foundations of Computer Science, FOCS 1999*, pages 596–604. IEEE Computer Society, 1999. doi:10.1109/SFFCS.1999.814634.
- [34] Yin Li and William F. Smyth. Computing the cover array in linear time. *Algorithmica*, 32(1):95–106, 2002. doi:10.1007/s00453-001-0062-2.
- [35] M. Lothaire. *Applied Combinatorics on Words*. Cambridge University Press, 2005. doi:10.1017/cbo9781107341005.
- [36] Udi Manber and Eugene W. Myers. Suffix arrays: A new method for on-line string searches. *SIAM Journal on Computing*, 22(5):935–948, 1993. doi:10.1137/0222058.
- [37] Dennis W. G. Moore and William F. Smyth. An optimal algorithm to compute all the covers of a string. *Information Processing Letters*, 50(5):239–246, 1994. doi:10.1016/0020-0190(94)00045-X.
- [38] Dennis W. G. Moore and William F. Smyth. A correction to "An optimal algorithm to compute all the covers of a string". *Information Processing Letters*, 54(2):101–103, 1995. doi:10.1016/0020-0190(94)00235-Q.
- [39] James H. Morris, Jr. and Vaughan R. Pratt. A linear pattern-matching algorithm. Technical Report 40, Department of Computer Science, University of California, Berkeley, 1970.
- [40] Jeong Seop Sim, Kunsoo Park, Sung-Ryul Kim, and Jee-Soo Lee. Finding approximate covers of strings. *Journal of Korea Information Science Society*, 29(1):16–21, 2002. URL: http://www.koreascience.or.kr/article/ArticleFullRecord.jsp?cn=JBGHG6_2002_v29n1_16.
- [41] William F. Smyth. Repetitive perhaps, but certainly not boring. *Theoretical Computer Science*, 249(2):343–355, 2000. doi:10.1016/S0304-3975(00)00067-0.
- [42] Peter Weiner. Linear pattern matching algorithms. In *14th Annual Symposium on Switching and Automata Theory, SWAT 1973*, pages 1–11, Washington, DC, USA, 1973. IEEE Computer Society. doi:10.1109/SWAT.1973.13.