

VU Research Portal

Internal shortest absent word queries in constant time and linear space

Badkobeh, Golnaz; Charalampopoulos, Panagiotis; Kosolobov, Dmitry; Pissis, Solon P.

published in

Theoretical Computer Science
2022

DOI (link to publisher)

[10.1016/j.tcs.2022.04.029](https://doi.org/10.1016/j.tcs.2022.04.029)

document version

Publisher's PDF, also known as Version of record

document license

Article 25fa Dutch Copyright Act

[Link to publication in VU Research Portal](#)

citation for published version (APA)

Badkobeh, G., Charalampopoulos, P., Kosolobov, D., & Pissis, S. P. (2022). Internal shortest absent word queries in constant time and linear space. *Theoretical Computer Science*, 922, 271-282.
<https://doi.org/10.1016/j.tcs.2022.04.029>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

E-mail address:

vuresearchportal.ub@vu.nl



Internal shortest absent word queries in constant time and linear space [☆]



Golnaz Badkobeh ^a, Panagiotis Charalampopoulos ^{b,1}, Dmitry Kosolobov ^{c,*,2}, Solon P. Pissis ^{d,e,3}

^a Department of Computing, Goldsmiths University of London, UK

^b Efi Arazi School of Computer Science, Reichman University, Herzliya, Israel

^c Ural Federal University, Ekaterinburg, Russia

^d CWI, Amsterdam, the Netherlands

^e Vrije Universiteit, Amsterdam, the Netherlands

ARTICLE INFO

Article history:

Received 3 June 2021

Received in revised form 5 April 2022

Accepted 21 April 2022

Available online 26 April 2022

Communicated by M. Sciortino

Keywords:

String algorithms

Internal queries

Shortest absent word

Bit parallelism

ABSTRACT

Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size σ , we are to preprocess T so that given a range $[i, j]$, we can return a representation of a shortest string over Σ that is absent in the fragment $T[i] \cdots T[j]$ of T . We present an $\mathcal{O}(n)$ -space data structure that answers such queries in constant time and can be constructed in $\mathcal{O}(n \log_{\sigma} n)$ time.

© 2022 Elsevier B.V. All rights reserved.

1. Introduction

Range queries are a classic data structure topic [66,12,11]. In one dimension, a range query $q = f(A, i, j)$ on an array of n elements over some set U , denoted by $A[1..n]$, takes two indices $1 \leq i \leq j \leq n$, a function f defined over arrays of elements of U , and outputs $f(A[i..j]) = f(A[i], \dots, A[j])$. Range query data structures in one dimension can thus be viewed as data structures answering queries on a string in the internal setting, where U is the considered alphabet.

Internal queries on a string have received much attention in recent years. In the internal setting, we are asked to preprocess a string T of length n over an alphabet Σ of size σ , so that queries about substrings of T can be answered efficiently. Note that an arbitrary substring of T can be encoded in $\mathcal{O}(1)$ words of space by the indices i, j of its occurrence as a fragment $T[i] \cdots T[j] = T[i..j]$ of T . Data structures for answering internal queries are interesting in their own right, but

[☆] The present paper is an extended and improved version of an earlier text that appeared in the 32nd Annual Symposium on Combinatorial Pattern Matching, CPM 2021 [8].

* Corresponding author.

E-mail addresses: g.badkobeh@gold.ac.uk (G. Badkobeh), panagiotis.charalampopoulos@post.idc.ac.il (P. Charalampopoulos), dkosolobov@mail.ru (D. Kosolobov), solon.pissis@cwi.nl (S.P. Pissis).

¹ Supported by the Israel Science Foundation grants 592/17 and 810/21.

² Supported by the Ministry of Science and Higher Education of the Russian Federation (Ural Mathematical Center project No. 075-02-2022-877).

³ Supported by the ALPACA and PANGAIA projects that have received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreements no. 956229 and no. 872539.

also have numerous applications in the design of algorithms and (more sophisticated) data structures. Because of these numerous applications, we usually place particular emphasis on the construction time—other than on the tradeoff between space and query time, which is the main focus in the classic data structure literature.

In data structures on strings it is typically assumed that the input alphabet is integer and polynomially bounded, i.e., it is a subset of $\{1, 2, \dots, n^{\mathcal{O}(1)}\}$ where n is the length of the input string T . One of the most widely-used internal queries is that of asking for the *longest common prefix* of two suffixes $T[i..n]$ and $T[j..n]$ of T . The classic data structure for this problem [49] consists of the suffix tree of T [28] and a lowest common ancestor data structure [40] over the suffix tree. It occupies $\mathcal{O}(n)$ space, it can be constructed in $\mathcal{O}(n)$ time, and it answers queries in $\mathcal{O}(1)$ time. In the word RAM model of computation with word size $\Theta(\log n)$ bits the construction time is not necessarily optimal when the input alphabet is $\{1, 2, \dots, \sigma\}$ and the string is packed into $\mathcal{O}(n/\log_\sigma n)$ machine words. A sequence of results [14,13,63,48,64,57] has culminated in the recent optimal data structure of Kempa and Kociumaka [43]: it occupies $\mathcal{O}(n/\log_\sigma n)$ space, it can be constructed in $\mathcal{O}(n/\log_\sigma n)$ time, and it answers queries in $\mathcal{O}(1)$ time (see also [15] and [27]).

Another fundamental problem in this setting is the *internal pattern matching* (IPM) problem. It consists in preprocessing T so that we can efficiently compute the occurrences of a substring U of T in another substring V of T . For the decision version of the IPM problem, Keller et al. [42] presented a data structure of nearly-linear size supporting sublogarithmic-time queries. Kociumaka et al. [47] presented a data structure of linear size supporting constant-time queries when the ratio between the lengths of V and U is bounded by a constant. The $\mathcal{O}(n)$ -time construction algorithm of the latter data structure was derandomized in [45]. In fact, Kociumaka et al. [47], using their efficient IPM queries as a subroutine, managed to show efficient solutions for other internal problems, such as for computing the periods of a substring (*period queries*, introduced in [46]), and for checking whether two substrings are rotations of one another (*cyclic equivalence queries*). Other problems that have been studied in the internal setting include string alignment [65,60,20,61], approximate pattern matching [23], dictionary matching [22,21], longest common substring [4], counting palindromes [59], range longest common prefix [3, 1,50,37], the computation of the lexicographically minimal or maximal suffix, and minimal rotation [6,44], as well as of the lexicographically k th suffix [7]. We refer the interested reader to the Ph.D dissertation of Kociumaka [45], for a nice exposition.

In this work, we extend this line of research by investigating the following basic internal query, which, to the best of our knowledge, has not been studied previously. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$, preprocess T so that given a range $[i, j]$, we can return a shortest string over Σ that does not occur in $T[i..j]$. The latter shortest string is also known as a shortest absent word in the literature. We work on the standard unit-cost word RAM model with machine word-size $w = \Theta(\log n)$ bits. We measure the space used by our algorithms and data structures in machine words, unless stated otherwise. We assume that we have random access to T and so our algorithms return a constant-space representation of a shortest string (a witness) consisting of a substring of T and a letter. A naïve solution for this problem precomputes a table of size $\mathcal{O}(n^2)$ that stores the answer for every possible query $[i, j]$. Our main result is the following theorem.

Theorem 1. *Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size σ , we can construct in $\mathcal{O}(n \log_\sigma n)$ time a data structure of size $\mathcal{O}(n)$ that, for any given query $[a, b]$, can compute in $\mathcal{O}(1)$ time a shortest string over Σ that does not occur in $T[a..b]$.*

In an earlier conference version of the present paper [8], we have obtained a weaker result: a data structure of size $\mathcal{O}((n/k) \cdot \log \log_\sigma n)$ that can answer queries in $\mathcal{O}(\log \log_\sigma k)$ time, where k is a user-defined parameter from $[1, \log \log_\sigma n]$. The improved data structure presented in this manuscript combines ideas from the conference version and the utilization of succinct fusion trees introduced by Grossi et al. [39].

In the related *range shortest unique substring* problem, defined by Abedin et al. [2], the task is to construct a data structure over T to be able to answer the following type of online queries efficiently. Given a range $[i, j]$, return a shortest string with exactly one occurrence (starting position) in $[i, j]$. Abedin et al. presented a data structure of size $\mathcal{O}(n \log n)$ supporting $\mathcal{O}(\log_w n)$ -time queries, where $w = \Theta(\log n)$ is the word size. Additionally, Abedin et al. [2] presented a data structure of size $\mathcal{O}(n)$ supporting $\mathcal{O}(\sqrt{n} \log^\epsilon n)$ -time queries, where ϵ is an arbitrarily small positive constant.

Our techniques For clarity of exposition, in this overview, we skip the time-efficient construction algorithms of our data structures and only describe how to compute the *length* of a shortest absent word (without a witness) in $T[a..b]$; note that this length is at most $\log_\sigma n$. Let us also recall that the length of a shortest absent word of T can be computed in $\mathcal{O}(n)$ time using the suffix tree of T [28]. It suffices to traverse the suffix tree of T recording the shortest string-depth ℓ , where an implicit or explicit node has less than σ outgoing edges.

First approach: We precompute, for each position $b \in [1, n]$ and for each length $j \in [1, \log_\sigma n]$, the starting position p_j of the shortest suffix of $T[1..b]$ that contains an occurrence of each of the σ^j distinct words of length j . In other words, $T[p_j..b]$ is the shortest suffix of $T[1..b]$ containing all distinct words of length j from Σ . Then, a query for the length of a shortest absent word of $T[a..b]$ reduces to computing the predecessor of a among the starting positions $p_1, p_2, \dots, p_{\lfloor \log_\sigma n \rfloor}$ we have precomputed for position b . By maintaining these $\mathcal{O}(\log_\sigma n)$ starting positions in a fusion tree [35] for every position $b \in [1, n]$, we obtain a data structure of size $\mathcal{O}(n \log_\sigma n)$ supporting queries in $\mathcal{O}(\log_w \log n) = \mathcal{O}(1)$ time.

Second approach: We precompute, for each length $j \in [1, \log_\sigma n]$, all minimal fragments of T that contain an occurrence of each of the distinct σ^j words of length j . As these fragments are inclusion-free, we can encode them using two n -bit arrays storing their starting and ending positions in T , respectively. We thus require $\mathcal{O}(n)$ words of space in total over all j s. Observe that $T[a..b]$ does not have an absent word of length j if and only if it contains a minimal fragment for length j ; we can check this condition in $\mathcal{O}(1)$ time after augmenting the computed bit arrays with succinct rank and select data structures [41]. Finally, due to monotonicity (if $T[a..b]$ contains all strings of length $j + 1$ then T contains all strings of length j), we can binary search for the answer in $\mathcal{O}(\log \log_\sigma n)$ time.

Third approach: We optimize the first approach by utilizing succinct fusion trees to store the sets of size $\mathcal{O}(\log_\sigma n)$ associated with positions of T , thus reducing the space on top of the sets to $\mathcal{O}(n \log_\sigma \log n)$. Instead of storing the $\mathcal{O}(\log_\sigma n)$ -size sets explicitly, we compute their elements on demand using $\mathcal{O}(\log_\sigma n)$ select data structures, each occupying $\mathcal{O}(n)$ bits. This leads to an $\mathcal{O}(n \log_\sigma \log n)$ -space solution. In order to optimize it further, we rely on the following combinatorial observation: if the length of a shortest absent word of a string X over Σ is λ , we need to append $\Omega(\sigma^{d-1} \cdot \lambda)$ letters to X in order to obtain a string with a shortest absent word of length $\lambda + d$. (For intuition, think of $|X|$ as a constant; then, we essentially need to append the de Bruijn sequence of order d over Σ to X in order to achieve the desired result.) This observation allows us to lower the memory consumption by truncating all succinct fusion trees at positions that are not multiples of $\log \log n$, by building them only for their first $\mathcal{O}(\log n / \log \log n)$ entries. The total space thus reduces to $\mathcal{O}(n)$ words. A query for the length of a shortest absent word of $T[a..b]$ is performed by first checking whether the answer is at most $\log n / \log \log n$, which is done using the (truncated) fusion tree stored at b , and, if not, a query on $T[a..b']$ is performed, where b' is the closest multiple of $\log \log n$ after b . It can be shown using the combinatorial observation that the answer for $T[a..b]$ is within an $\mathcal{O}(1)$ -length range of the answer for $T[a..b']$, and it is computed by the data structure from the second approach.

Other related work Let us recall that a string S that does not occur in T is called *absent* from T , and if all its proper substrings appear in T it is called a *minimal absent word* of T . It should be clear that every shortest absent word is also a minimal absent word. Minimal absent words (MAWs) are used in many applications [62,56,31,38,16,54,26] and their theory is well developed [52,30,32], also from an algorithmic and data structure point of view [51,24,9,19,18,5,36,10,25]. For example, it is well known that, given two strings X and Y , one has $X = Y$ if and only if X and Y have the same set of MAWs [52].

Paper organization Section 2 provides some preliminaries. The first approach is detailed in Section 3 and the second one in Section 4. Section 5 provides the combinatorial foundations for the third approach, which is detailed in Section 6. Sections 3–5 have essentially already appeared in the conference version [8] of our paper; the main difference and novelty lie in Section 6. We conclude with open problems in Section 7.

2. Preliminaries

An *alphabet* Σ is a finite nonempty set whose elements are called *letters*. A *string* (or *word*) $S = S[1..n]$ is a sequence of length $|S| = n$ over Σ . The *empty* string ε is the string of length 0. The *concatenation* of two strings S and T is the string composed of the letters of S followed by the letters of T ; it is denoted by $S \cdot T$ or simply by ST . The set of all strings (including ε) over Σ is denoted by Σ^* . The set of all strings of length $k > 0$ over Σ is denoted by Σ^k . For $1 \leq i \leq j \leq n$, $S[i]$ denotes the i th letter of S , and the fragment $S[i..j]$ denotes an *occurrence* of the underlying *substring* $P = S[i] \cdots S[j]$. We say that P *occurs* at (starting) *position* i in S . A string P is called *absent* from S if it does not occur in S . A substring $S[i..j]$ is a *suffix* of S if $j = n$ and it is a *prefix* of S if $i = 1$.

The following proposition is straightforward (as explained in Section 1).

Proposition 1. *Let T be a string of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$. A shortest absent word of T can be computed in $\mathcal{O}(n)$ time.*

Given an array A of n items taken from a totally ordered set, the *range minimum query* $\text{RMQ}_A(\ell, r) = \arg \min A[k]$ (with $1 \leq \ell \leq k \leq r \leq n$) returns the position of the minimal element in $A[\ell..r]$. The following result is known.

Theorem 2 ([11,34]). *Let A be an array of n integers. A data structure of size $2n + o(n)$ bits that supports RMQs on A in $\mathcal{O}(1)$ time without the need to store and access A itself can be constructed in $\mathcal{O}(n)$ time.*

We make use of *rank and select* data structures constructed over bit vectors. For a bit vector H we define $\text{rank}_q(i, H) = |\{k \in [1, i] : H[k] = q\}|$ and $\text{select}_q(i, H) = \min\{k \in [1, n] : \text{rank}_q(k, H) = i\}$, for $q \in \{0, 1\}$. The following result is known.

Theorem 3 ([41,53]). *Let H be a bit vector of n bits. A data structure of $o(n)$ additional bits that supports rank and select queries on H in $\mathcal{O}(1)$ time can be constructed in $\mathcal{O}(n)$ time.*

The *static predecessor* problem consists in preprocessing a set Y of integers, over an ordered universe U , so that, for any integer $x \in U$ one can efficiently return the predecessor $\text{pred}(x) := \max\{y \in Y : y \leq x\}$ of x in Y . The successor problem is defined analogously: upon a queried integer $x \in U$, the successor $\min\{y \in Y : y \geq x\}$ of x in Y is to be returned. Willard and Fredman designed the *fusion tree* data structure for this problem [35]. In the dynamic variant of the problem, updates to Y are interleaved with predecessor and successor queries. Pătraşcu and Thorup [55] presented a dynamic version of fusion trees, which, in particular, yields an efficient construction of this data structure.

Theorem 4 ([35,55]). *Let Y be a set of at most n w -bit integers. A data structure of size $\mathcal{O}(n)$ can be constructed in $\mathcal{O}(n \log_w n)$ time supporting insertions, deletions, and predecessor queries on Y in $\mathcal{O}(\log_w n)$ time.*

We also use a succinct version of the (static) fusion tree that utilizes only $\mathcal{O}(n \log w)$ bits on top of a read-only array Y of length n (in contrast, the fusion tree from Theorem 4 uses $\mathcal{O}(nw)$ bits). In this data structure there is no need to store the array Y explicitly. Instead, Y can be “emulated” by computing its elements on demand in $\mathcal{O}(1)$ time. Albeit it is not explicitly stated in [39,17], it follows from their construction that the succinct version can be constructed from a (usual) fusion tree in linear time.

Theorem 5 ([39,17]). *Let Y be a read-only array of at most n w -bit integers and $n \leq w^{\mathcal{O}(1)}$. A data structure of size $\mathcal{O}(n \log w)$ bits can be constructed in $\mathcal{O}(n \log_w n)$ time supporting predecessor queries on the elements of Y in $\mathcal{O}(\log_w n)$ time, provided that a table computable in $o(2^w)$ time and independent of the array has been precomputed.*

Note that if we build multiple predecessor queries for sets of w -bit integers using the above theorem, they can all share a unique table computable in $o(2^w)$ time.

If $|U| = \mathcal{O}(n)$, then, after an $\mathcal{O}(n)$ -time preprocessing, we can answer predecessor queries over the integer universe U in $\mathcal{O}(1)$ time as follows. For each $y \in Y$, we set the y th bit of an initially all-zeros $|U|$ -size bit vector. We then preprocess this bit vector as in Theorem 3. Then, a predecessor query for any integer x can be answered in $\mathcal{O}(1)$ time due to the following readily verifiable formula: $\text{pred}(x) = \text{select}_1(\text{rank}_1(x))$.

The main problem considered in this paper is formally defined as follows.

INTERNAL SHORTEST ABSENT WORD (ISAW)

Input: A string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size $\sigma > 1$.

Output: Given integers a and b , with $1 \leq a \leq b \leq n$, output a shortest string in Σ^* with no occurrence in $T[a..b]$.

If $a = b$ then the answer is trivial. So, in what follows we assume that $a < b$. Let us also remark that the output (shortest absent word) can be represented in $\mathcal{O}(1)$ space using: either a range $[i, j] \subseteq [1, n]$ and a letter α of Σ , such that the shortest string in Σ^* with no occurrence in $T[a..b]$ is $T[i..j]\alpha$; or simply a range $[i, j] \subseteq [1, n]$ such that the shortest string in Σ^* with no occurrence in $T[a..b]$ is $T[i..j]$.

Example 1. Given the string $T = \text{abaabaa}\underline{\text{abbabbb}}\text{aaab}$ and the range $[a, b] = [8, 14]$ (shown in red and underlined), the only shortest absent word of $T[8..14]$ is $T[i..j] = T[7..8] = \text{aa}$.

3. $\mathcal{O}(n \log_\sigma n)$ space and $\mathcal{O}(1)$ query time

Let T be a string of length n . We define $S_T(j)$ as the function counting the cardinality of the set of length- j substrings of T . This is known as the *substring complexity* function [29,58]. Note that $S_T(j) \leq n$, for all j . We have the following simple fact.

Fact 6. *The length ℓ of a shortest absent word of a string T of length n over an alphabet of size σ is equal to the smallest j for which $S_T(j) < \sigma^j$ and hence $\ell \in [1, \lfloor \log_\sigma n \rfloor]$.*

We denote the set of shortest absent words of T by SAW_T . Recall that, by Proposition 1, a shortest absent word of T can be computed in $\mathcal{O}(n)$ time. We denote the length of the shortest absent words of T by ℓ . By Fact 6, $\ell \leq \lfloor \log_\sigma n \rfloor$. Since ℓ is an upper bound on the length of the answer for any ISAW query on T , in what follows, we consider only lengths in $[1, \ell - 1]$. Let one such length be denoted by j . By constructing and traversing the suffix tree of T , we can assign to each $T[i..i + j - 1]$ its lexicographic rank in Σ^j . The time required for each length j is $\mathcal{O}(n)$, since the suffix tree of T can be constructed within this time [28]. Thus, the total time for all lengths $j \in [1, \ell - 1]$ is $\mathcal{O}(n \log_\sigma n)$ by Fact 6.

We design the following warm-up solution to the ISAW problem. For all $j \in [1, \ell - 1]$ we store an array RNK_j of n integers such that $\text{RNK}_j[i]$ is equal to the lexicographic rank of $T[i..i + j - 1]$ in Σ^j . Then, given a range $[a, b]$, in order to check if there is an absent word of length j in $T[a..b]$ we only need to compute the number of distinct elements in $\text{RNK}_j[a..b - j + 1]$. It is folklore that using a persistent segment tree, we can preprocess an array A of n integers in

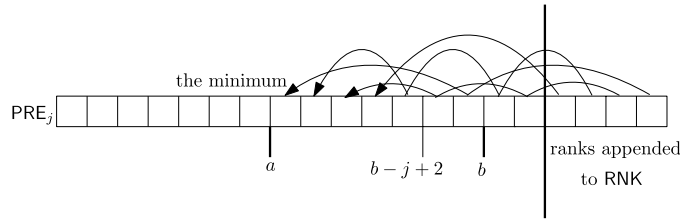


Fig. 1. Illustration of the setting in Fact 7.

$\mathcal{O}(n \log n)$ time so that upon a range query $[a, b]$ we can return the number of distinct elements in $A[a..b]$ in $\mathcal{O}(\log n)$ time.⁴ Thus, we could use this tool as a black box for every array RNK_j resulting, however, in $\Omega(\log n)$ -time queries. We improve upon this solution as follows.

We employ a range minimum query (RMQ) data structure [11] over a slight modification of RNK_j . For each j , we have an auxiliary procedure checking whether all strings from Σ^j occur in $T[a..b]$ or not (i.e., it suffices to check whether any lexicographic rank is absent from the corresponding range). Similar to the previous solution, we rank the elements of Σ^j by their lexicographic order. We append RNK_j with all integers in $[1, \sigma^j]$. Let this array be APP_j . By Fact 6, we have that $|APP_j| \leq 2n$. Then, we construct an array PRE_j of size $|APP_j|$: $PRE_j[i]$ stores the position of the rightmost occurrence of $APP_j[i]$ in $APP_j[1..i-1]$ (or 0 if such an occurrence does not exist). This can be done in $\mathcal{O}(n)$ time per j by sorting the list of pairs $(T[i..i+j-1], i)$, for all i , using the suffix tree of T to assign ranks for $T[i..i+j-1]$ and then radix sort to sort the list of pairs.

We now rely on the following fact.

Fact 7. $S_{T[a..b]}(j) = \sigma^j$ if and only if $\min\{PRE_j[i] : i \in [b-j+2, |PRE_j|]\} \geq a$.

Proof. If the smallest element in $PRE_j[b-j+2..|PRE_j|]$, say $PRE_j[k]$, is such that $PRE_j[k] \geq a$, then all ranks of elements in Σ^j occur in $APP_j[a..b-j+1]$. This is because all elements (ranks) in Σ^j occur at least once after $b-j+2$ (due to appending all integers in $[1, \sigma^j]$ to RNK_j), and thus all must have a representative occurrence after $b-j+2$. Inspect Fig. 1 for an illustration. (The opposite direction is analogous.) \square

Examples 2 and 3 illustrate the construction of arrays RNK_j , APP_j , and PRE_j as well as Fact 7.

Example 2 (Construction). Let $T = \text{abaabaaabbabbbaaab}$ and $\Sigma = \{a, b\}$. The set SAW_T of shortest absent words of T over Σ , each of length $\ell = 4$, is $\{\text{aaaa}, \text{abab}, \text{baba}, \text{bbbb}\}$. Arrays RNK_j , APP_j , and PRE_j , for all $j \in [1, \ell - 1]$, are as depicted in Table 1. For instance, $RNK_2[15] = APP_2[15] = 1$ denotes that the lexicographic rank of aa in Σ^2 is 1; and $PRE_2[15] = 7$ denotes that the previous rightmost occurrence of aa is at position 7.

Example 3 (Fact 7). Let $[a, b] = [7, 11]$ and $j = 2$ (see Example 2). The smallest element in $\{PRE_2[11], \dots, PRE_2[21]\}$ is $PRE_2[15] = 7 \geq a = 7$, which corresponds to rank $APP_2[15] = 1$. Indeed all other ranks 2, 3, 4 have at least one occurrence within $APP_2[7..11] = 1, 2, 4, 3, 2$.

To apply Fact 7, we construct, in $\mathcal{O}(n)$ time, an $\mathcal{O}(n)$ -space, $\mathcal{O}(1)$ -query-time RMQ data structure over PRE_j ; see Theorem 2. This results in $\mathcal{O}(n\ell) = \mathcal{O}(n \log_\sigma n)$ preprocessing time and space over all j .

For querying, let us observe that $\sigma^j - S_{T[a..b]}(j)$, for any T, a, b and increasing j , is non-decreasing. We can thus apply binary search on j to find the smallest length j such that $S_{T[a..b]}(j) < \sigma^j$. This results in $\mathcal{O}(\log \ell) = \mathcal{O}(\log \log_\sigma n)$ query time. We obtain the following proposition (retrieving a witness shortest absent word is detailed later).

⁴ Here, we provide details of this folklore data structure. Consider an array B of size n , all of whose entries are initially set to a dummy value, and a binary tree of height $\mathcal{O}(\log n)$ with n leaves, such that, for each integer $i \in [1, n]$, the i -th leaf of the tree corresponds to the i -th entry of B ; each node of the tree then naturally corresponds to the segment of B spanned by its leaf-descendants. Let us call an index i active when $B[i] \neq B[k]$ for all $k \in (i, n]$. For each node v of the tree, we will maintain the number $\text{active}(v)$ of its active leaf-descendants as we transform B to A by updating B 's entries in the left-to-right order. Note that each active index contributes to the value $\text{active}(v)$ if and only if v lies on the path from the root to the i -th leaf. Each update of the form " $B[i] := A[i]$ " can activate/deactivate at most one index and thus affects $\mathcal{O}(\log n)$ nodes. Further, we can efficiently precompute when all the (de)activations happen by sorting the entries of A . Thus, overall, $\mathcal{O}(n \log n)$ updates of the stored values $\text{active}(v)$ are required and they can be performed in $\mathcal{O}(n \log n)$ time in total. This data structure can be made persistent (i.e., it can allow access at its state after each update " $B[i] := A[i]$ ") as follows: the tree is implemented using pointers and, for each i , after setting $B[i]$ to be equal to $A[i]$, instead of making the $\mathcal{O}(\log n)$ required updates of values $\text{active}(v)$, we create $\mathcal{O}(\log n)$ new nodes, including a new root r_i , and update the pointers as necessary. It remains to explain how to answer a query for the number of distinct elements in $A[i..j]$: we access the state of the data structure just after update " $B[j] := A[j]$ ", partition $[i, j]$ into $\mathcal{O}(\log n)$ intervals that correspond to a set V of nodes of the tree, and return $\sum_{v \in V} \text{active}(v)$.

Table 1
Arrays RNK_j , APP_j , and PRE_j in Example 2.

| | | | | | | | | | | | | | | | | | | | | | | | | |
|---------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 |
| T | a | b | a | a | b | a | a | a | b | b | a | b | b | b | a | a | a | b | | | | | | |
| RNK_1 | 1 | 2 | 1 | 1 | 2 | 1 | 1 | 1 | 2 | 2 | 1 | 2 | 2 | 2 | 1 | 1 | 1 | 2 | | | | | | |
| APP_1 | 1 | 2 | 1 | 1 | 2 | 1 | 1 | 1 | 2 | 2 | 1 | 2 | 2 | 2 | 1 | 1 | 1 | 2 | 1 | 2 | | | | |
| PRE_1 | 0 | 0 | 1 | 3 | 2 | 4 | 6 | 7 | 5 | 9 | 8 | 10 | 12 | 13 | 11 | 15 | 16 | 14 | 17 | 18 | | | | |
| RNK_2 | 2 | 3 | 1 | 2 | 3 | 1 | 1 | 2 | 4 | 3 | 2 | 4 | 4 | 3 | 1 | 1 | 2 | | | | | | | |
| APP_2 | 2 | 3 | 1 | 2 | 3 | 1 | 1 | 2 | 4 | 3 | 2 | 4 | 4 | 3 | 1 | 1 | 2 | 1 | 2 | 3 | 4 | | | |
| PRE_2 | 0 | 0 | 0 | 1 | 2 | 3 | 6 | 4 | 0 | 5 | 8 | 9 | 12 | 10 | 7 | 15 | 11 | 16 | 17 | 14 | 13 | | | |
| RNK_3 | 3 | 5 | 2 | 3 | 5 | 1 | 2 | 4 | 7 | 6 | 4 | 8 | 7 | 5 | 1 | 2 | | | | | | | | |
| APP_3 | 3 | 5 | 2 | 3 | 5 | 1 | 2 | 4 | 7 | 6 | 4 | 8 | 7 | 5 | 1 | 2 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| PRE_3 | 0 | 0 | 0 | 1 | 2 | 0 | 3 | 0 | 0 | 0 | 8 | 0 | 9 | 5 | 6 | 7 | 15 | 16 | 4 | 11 | 14 | 10 | 13 | 12 |

Proposition 2. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{O(1)}\}$ of size σ , we can construct a data structure of size $O(n \log_\sigma n)$ in $O(n \log_\sigma n)$ time, so that if query $[a, b]$ is given, we can compute a shortest string over Σ that does not occur in $T[a..b]$ in $O(\log \log_\sigma n)$ time.

We further improve the query time via employing fusion trees as follows. We create a 2d array $FTR[1..l-1][1..n]$ of integers, where

$$FTR[j][i] = \min\{PRE_j[i - j + 2], \dots, PRE_j[|PRE_j|]\},$$

for all $j \in [1, l - 1]$ and $i \in [1, n]$. Intuitively, $FTR[j][i]$ is the rightmost index of T such that $T[FTR[j][i]..i]$ contains all strings of length j over Σ if such an index exists and 0 otherwise.

Array FTR can be constructed in $O(nl) = O(n \log_\sigma n)$ time by scanning each array PRE_j from right to left maintaining the minimum. Within the same complexities we also maintain satellite information specifying the index $k \in [i - j + 2, |PRE_j|]$ where the range minimum $FTR[j][i]$ came from in the sub-array $PRE_j[i - j + 2..|PRE_j|]$. We then construct n fusion trees, one for every collection of $l - 1$ integers in $FTR[1..l-1][i]$. This takes total preprocessing time and space $O(nl) = O(n \log_\sigma n)$ by Theorem 4. Given the range query $[a, b]$, we need to find the smallest $j \in [1, l - 1]$ such that $FTR[j][b] < a$. By Theorem 4, we find where the predecessor of a lies in $FTR[1..l-1][b]$ in $O(\log_w l)$ time, where w is the word size; this time cost is $O(1)$ since $w = \Theta(\log n)$.

We finally retrieve a witness shortest absent word as follows. If there is no $j < l$ such that $FTR[j][b] < a$, then we output any shortest absent word of length l of T arbitrarily. If such a $j < l$ exists, by the definition of $FTR[j][b]$, we output $T[FTR[j][b]..FTR[j][b] + j - 1]$ if $FTR[j][b] > 0$ or $T[k..k + j - 1]$ if $FTR[j][b] = 0$, where k is the index of PRE_j , where the minimum came from. Inspect the following illustrative example.

Example 4 (Querying). We construct array FTR for T from Example 2. For a given $[a, b]$ we look up column b , and find the topmost entry whose value is less than a . If all entries have values greater than or equal to a , we output any element from SAW_T arbitrarily.

| | | | | | | | | | | | | | | | | | | |
|----------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|
| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| T | a | b | a | a | b | a | a | a | b | b | a | b | b | b | a | a | a | b |
| $FTR[1]$ | 0 | 1 | 2 | 2 | 4 | 5 | 5 | 5 | 8 | 8 | 10 | 11 | 11 | 11 | 14 | 14 | 14 | 17 |
| $FTR[2]$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 7 | 7 | 7 | 7 | 7 | 7 | 11 | 11 | 13 |
| $FTR[3]$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 4 | 4 | 4 | 4 |

If $[a, b] = [3, 14]$ then no entry in column $b = 14$ is less than $a = 3$, which means the length of the shortest absent word is 4; we output one from $\{aaaa, abab, baba, bbbb\}$ arbitrarily. If $[a, b] = [5, 14]$ then $FTR[3][14] = 4 < 5$ so the length of a shortest absent word of $T[5..14]$ is 3; a shortest absent word is $T[FTR[3][14]..FTR[3][14] + 3 - 1] = T[4..6] = aba$.

If $[a, b] = [7, 9]$, $FTR[2][9] = 0 < 7$ so the length of a shortest absent word is 2; a shortest absent word is $T[k..k + j - 1] = T[9..10] = bb$ because $FTR[2][9] = \min\{PRE_2[9], \dots, PRE_2[|PRE_2|]\} = PRE_2[9] = 0$ tells us that the minimum in this range came from index $k = 9$.

We obtain the following proposition.

Proposition 3. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{O(1)}\}$ of size σ , we can construct a data structure of size $O(n \log_\sigma n)$ in $O(n \log_\sigma n)$ time, so that if query $[a, b]$ is given, we can compute a shortest string over Σ that does not occur in $T[a..b]$ in $O(1)$ time.

4. $\mathcal{O}(n)$ space and $\mathcal{O}(\log \log_\sigma n)$ query time

Definition 1 (*Order- j Fragment*). Given a string T over an alphabet of size σ and an integer j , V is called an *order- j fragment* of T if and only if V is a fragment of T and $S_V(j) = \sigma^j$. V is further called a *minimal order- j fragment* of T if $S_U(j) < \sigma^j$ and $S_Z(j) < \sigma^j$ for $U = V[1..|V| - 1]$ and $Z = V[2..|V|]$.

In particular, minimal order- j fragments are pairwise not included in each other. The following fact follows directly.

Fact 8. Given a string T of length n over an alphabet of size σ and an integer j we have $\mathcal{O}(n)$ minimal order- j fragments. Moreover, an arbitrary fragment F of T has $S_F[j] = \sigma^j$ if and only if it contains at least one of these minimal fragments.

For each $j \in [1, \log_\sigma n]$, we consider all minimal order- j fragments T , separately. We encode the minimal order- j fragments of T using two bit vectors SP_j and EP_j , standing for starting positions and ending positions. Inspect the following example.

Example 5. We consider T from Example 2 and $j = 2$.

| | | | | | | | | | | | | | | | | | | | | | |
|---------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|
| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| T | a | b | a | a | b | a | a | a | b | b | a | b | b | b | a | a | a | b | | | |
| APP_2 | 2 | 3 | 1 | 2 | 3 | 1 | 1 | 2 | 4 | 3 | 2 | 4 | 4 | 3 | 1 | 1 | 2 | 1 | 2 | 3 | 4 |
| PRE_2 | 0 | 0 | 0 | 1 | 2 | 3 | 6 | 4 | 0 | 5 | 8 | 9 | 12 | 10 | 7 | 15 | 11 | 16 | 17 | 14 | 13 |
| SP_2 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | | | |
| EP_2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | | | |

For instance, $SP_2[13] = 1$ and $EP_2[18] = 1$ denote the minimal order-2 fragment $V = T[13..18] = bbaaab$.

We construct a rank and select data structure on SP_j and EP_j , for all $j \in [1, \ell - 1]$ supporting $\mathcal{O}(1)$ -time queries. The overall space is $\mathcal{O}(n)$ by Theorem 3 and Fact 6.

Let us now explain how this data structure enables fast computation of absent words of length j . Given a range $[a, b]$, by Fact 8, we only need to find whether $T[a..b]$ contains a minimal order- j fragment. We can do this in $\mathcal{O}(1)$ time using one rank and one select query: $t = \text{rank}_1(a - 1, SP_j) + 1$ and $\text{select}_1(t, EP_j)$. The select query returns the ending position of the leftmost minimal order- j fragment that starts after the position $a - 1$; it remains to check whether this minimal order- j fragment is inside $[a, b]$.

Example 6. We consider T , SP_2 and EP_2 from Example 5. Let $[a, b] = [5, 14]$. We have $t = \text{rank}_1(a - 1, SP_2) + 1 = \text{rank}_1(4, SP_2) + 1 = 1$, $\text{select}_1(t, SP_2) = \text{select}_1(1, SP_2) = 5 < b = 14$ and $\text{select}_1(t, EP_2) = \text{select}_1(1, EP_2) = 10 < b = 14$, which means $T[5, 14]$ contains a minimal order-2 fragment.

Let us now describe a time-efficient construction of SP_j and EP_j . We use arrays PRE_j and APP_j of T , which are constructible in $\mathcal{O}(n)$ time (see Section 3). Recall that $PRE_j[i]$ stores the position of the rightmost occurrence of rank $APP_j[i]$ in $APP_j[1..i - 1]$ (or 0 if such an occurrence does not exist). We apply Fact 7 as follows. We start with all bits of SP_j and EP_j unset. Then, for each $b \in [1, n]$ for which $PRE_j[b - j + 1] < \min\{PRE_j[i] : i \in [b - j + 2, |PRE_j|]\} = a$, we set the b th bit of EP_j and the a th bit of SP_j . This can be done online in a right-to-left scan of PRE_j in $\mathcal{O}(n)$ time.

Example 7. We consider T , SP_2 and EP_2 from Example 5. We start by setting $b = n = 18$ and scan PRE_2 from right to left: we have $a = 13$ because $\min\{PRE_2[i] : i \in [18, 21]\} = 13$. This gives fragment $T[13..18]$, which is minimal since $PRE_2[b - 1] = PRE_2[17] < 13$. Then we set $b = n - 1 = 17$ and have $a = 11$ because $\min\{PRE_2[i] : i \in [17, 21]\} = 11$. This gives fragment $T[11..17]$, which is not minimal since $PRE_2[b - 1] = PRE_2[16] \geq 11$. Then we set $b = n - 2 = 16$ and have $a = 11$ because $\min\{PRE_2[i] : i \in [16, 21]\} = 11$. This gives fragment $T[11..16]$, which is minimal since $PRE_2[b - 1] = PRE_2[15] < 11$.

Lemma 1. SP_j and EP_j can be constructed in $\mathcal{O}(n)$ time.

For all j , the construction time is $\mathcal{O}(n\ell) = \mathcal{O}(n \log_\sigma n)$ by Theorem 3, Lemma 1, and Fact 6. All the arrays SP_j and EP_j in total occupy $\mathcal{O}(n\ell) = \mathcal{O}(n \log_\sigma n)$ bits of space, which is $\mathcal{O}(n)$ space when measured in $\Theta(\log n)$ -bit machine words. We obtain the following lemma.

Lemma 2. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size σ , we can construct a data structure of size $\mathcal{O}(n)$ in $\mathcal{O}(n \log_\sigma n)$ time, so that if query $(j, [a, b])$ is given, we can check in $\mathcal{O}(1)$ time whether there is any string in Σ^j that does not occur in $T[a..b]$, and if so return such a string.

We can now perform binary search on j using Lemma 2 to find the smallest j for which $S_{T[a..b]}(j) < \sigma^j$. This results in $\mathcal{O}(\log \ell) = \mathcal{O}(\log \log_{\sigma} n)$ query time by Fact 6. It should now be clear that when we find the j corresponding to the length of a shortest absent word, we can output the length- j suffix of the leftmost minimal order- j fragment starting after a . Note that outputting this suffix is correct by the definition of minimal order- j fragments.

Example 8. We consider T , SP_2 and EP_2 from Example 5. Let $[a, b] = [2, 7]$. The length of a shortest absent word of $T[2..7]$ is 2. We output bb , which is the length-2 suffix of the leftmost minimal order-2 fragment $T[5..10] = ba aabb$ starting after $a = 2$.

We obtain the following result.

Proposition 4. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size σ , we can construct a data structure of size $\mathcal{O}(n)$ in $\mathcal{O}(n \log_{\sigma} n)$ time, so that if query $[a, b]$ is given, we can compute a shortest string over Σ that does not occur in $T[a..b]$ in $\mathcal{O}(\log \log_{\sigma} n)$ time.

5. Combinatorial insights

A positive integer p is a *period* of a string S if $S[i] = S[i + p]$ for all $i \in [1, |S| - p]$. We refer to the smallest period as the *period* of the string. Let us state the periodicity lemma, one of the most elegant combinatorial results on strings.

Lemma 3 (Periodicity Lemma (weak version) [33]). If a string S has periods p and q such that $p + q \leq |S|$, then $\gcd(p, q)$ is also a period of S .

Lemma 4. If all strings in $\{UW : U \in \Sigma^k\}$ for $W \neq \varepsilon$ occur in some string S , then $|S| \geq |W| \cdot \sigma^k/4$.

Proof. Let p be the period of W , and let $a \in \Sigma$ be such that the period of aW is also p . All strings ZbW for a letter $b \neq a$ and $Z \in \Sigma^{k-1}$ must occur in S . Let $A = \{UW : U \in \Sigma^k\} \setminus \{ZaW : Z \in \Sigma^{k-1}\}$, and note that it is of size $\sigma^k - \sigma^{k-1} \geq \sigma^k/2$. The following claim immediately implies the statement of the lemma.

Claim. Let i and j be starting positions of occurrences of different strings $UW, VW \in A$ in S , respectively. Then, we have $|j - i| \geq |W|/2$.

Proof. Let us assume, without loss of generality, that $j > i$. Further, let us assume towards a contradiction that $j - i < |W|/2$. Then, $j - i$ is a period of W and $p + j - i \leq |W|$ since $p \leq j - i$. Therefore, due to the periodicity lemma (Lemma 3), $j - i$ must be divisible by the period p of W . Hence, V ends with the letter a and $VW \notin A$, a contradiction. \square

This concludes the proof of this lemma. \square

Lemma 5. If a shortest absent word of a string X is of length λ , then the length of a shortest absent word of XY is in $[\lambda, \lambda + \max\{10, 4 + \log_{\sigma}(|Y|/\lambda)\}]$.

Proof. Let W and W' be shortest absent words of X and XY , respectively. Further, let $d = |W'| - |W|$. In order to have $d > 0$, all strings UW for $U \in \Sigma^{d-1}$ must occur in XY , and hence in $X[|X| - |UW| + 2..|X|] \cdot Y$, since none of them occurs in X . Lemma 4 implies that $|Y| + \lambda + d > \lambda \cdot \sigma^{d-1}/4$. Then, since $\lambda + d \leq 2\lambda d$ for any positive integers λ, d , we have $|Y| > \lambda \cdot (\sigma^{d-1}/4 - 2d)$. Assuming that $d \geq 10$, and since $\sigma \geq 2$, we conclude that $|Y| > \lambda \cdot \sigma^{d-1}/8$. Consequently, $\log_{\sigma}(8|Y|/\lambda) + 1 > d$. Since $\log_{\sigma} 8 \leq 3$ we get the claimed bound. \square

Lemma 6. If a shortest absent word of XY is of length m , a shortest absent word of X is of length λ , and $|Y| \leq m \cdot \tau$, for a positive integer $\tau \geq 16$, then $m - \lambda \leq 10 + 2 \log_{\sigma} \tau$.

Proof. From Lemma 5 we have $\lambda \in [m - \max\{10, 4 + \log_{\sigma}(|Y|/\lambda)\}, m]$. If $\max\{10, 4 + \log_{\sigma}(|Y|/\lambda)\} = 10$, then $m - \lambda \leq 10$ and we are done.

In the complementary case, since $|Y| \leq m \cdot \tau$, we get the following:

$$\lambda \geq m - \log_{\sigma}(m \cdot \tau/\lambda) - 4 \iff \lambda \geq m + \log_{\sigma} \lambda - \log_{\sigma} m - \log_{\sigma} \tau - 4.$$

In particular, $\lambda \geq m - \log_{\sigma} m - \log_{\sigma} \tau - 4$.

From the above, if $m \leq \tau$, then $m - \lambda \leq 4 + 2 \log_{\sigma} \tau$.

In what follows we assume that $m > \tau \geq 16$. Rearranging the original equation, and since $\log_{\sigma}(\cdot)$ is an increasing function and $\lambda \geq m - \log_{\sigma} m - \log_{\sigma} \tau - 4$, we have

$$\begin{aligned}
 m - \lambda \leq 4 + \log_\sigma(m \cdot \tau / \lambda) &\leq 4 + \log_\sigma\left(\frac{m}{m - \log_\sigma m - \log_\sigma \tau - 4}\right) + \log_\sigma \tau \\
 &\leq 4 + \log_\sigma\left(\frac{m}{m - 2 \log_\sigma m - 4}\right) + \log_\sigma \tau.
 \end{aligned}$$

Then, we have $m - 2 \log_\sigma m - 4 \geq m/5$ since, for any $\sigma \geq 2$, $4x/5 - 2 \log_\sigma x - 4$ is an increasing function on $[16, \infty)$ and positive for $x = 16$. Hence, $m - \lambda \leq 4 + \log_\sigma 5 + \log_\sigma \tau \leq 7 + \log_\sigma \tau$.

By combining the bounds on $m - \lambda$ we get the claimed bound. \square

6. $\mathcal{O}(n)$ space and $\mathcal{O}(1)$ query time

Our linear-space solution of the ISAW problem with constant query time is an optimization of the $\mathcal{O}(n \log_\sigma n)$ -space solution from Section 3 with some “boundary” cases processed using the data structure of Section 4. Let us first describe a simpler $\mathcal{O}(n \log_\sigma \log n)$ -space data structure, which will be then optimized using the combinatorial insights from Section 5.

Recall that we denote by ℓ the length of a shortest absent word of T . The issue with the solution of Section 3 is that the 2d array $\text{FTR}[1.. \ell - 1][1.. n]$, equipped with fusion trees, occupies $\mathcal{O}(n \log_\sigma n)$ space. In order to reduce the memory consumption, we store the array FTR implicitly, computing its entries on demand, and utilize succinct fusion trees from Theorem 5 instead of usual fusion trees.

Recall that $\text{FTR}[j][i]$ is the rightmost index of T such that $T[\text{FTR}[j][i].. i]$ contains as substrings all strings of length j over Σ and it is equal to $\min\{\text{PRE}_j[i - j + 2], \dots, \text{PRE}_j[\text{PRE}_j[i]]\}$. Therefore, the content of the 2d array FTR can be “emulated” without storing it explicitly if one can compute in $\mathcal{O}(1)$ time the minima $\min\{\text{PRE}_j[a], \dots, \text{PRE}_j[\text{PRE}_j[i]]\}$, for any $a \in [1, n]$. For $j \in [1, \ell - 1]$ and $a \in [1, n]$, denote $M_{j,a} = \min\{\text{PRE}_j[a], \dots, \text{PRE}_j[\text{PRE}_j[i]]\}$. Let us fix some j . Since the sequence $M_{j,1}, M_{j,2}, \dots, M_{j,n}$ is non-decreasing, we can encode it in a $2n$ -bit array B_j using the select data structure from Theorem 3 as follows: we construct B_j (initially empty) by considering $a = 1, 2, \dots, n$ in increasing order and, for each a , we append to the end of B_j exactly $M_{j,a} - M_{j,a-1}$ zeroes followed by 1, setting $M_{j,0} = 0$ (i.e., we append the number $M_{j,a} - M_{j,a-1}$ written in unary); then, we have $M_{j,a} = \text{select}_1(a, B_j) - a$.

Example 9. We consider T from Example 2 and $j = 2$.

| | | | | | | | | | | | | | | | | | | | | | |
|----------------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|
| i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| T | a | b | a | a | b | a | a | a | b | b | a | b | b | b | a | a | a | b | | | |
| APP_2 | 2 | 3 | 1 | 2 | 3 | 1 | 1 | 2 | 4 | 3 | 2 | 4 | 4 | 3 | 1 | 1 | 2 | 1 | 2 | 3 | 4 |
| PRE_2 | 0 | 0 | 0 | 1 | 2 | 3 | 6 | 4 | 0 | 5 | 8 | 9 | 12 | 10 | 7 | 15 | 11 | 16 | 17 | 14 | 13 |
| $M_{2,i}$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 7 | 7 | 7 | 7 | 7 | 11 | 11 | 13 | | | |

In this case, we have $B_2 = 11111111100000100111111000011001$.

Besides access to the 2d array FTR , the algorithm of Section 3 also required access to the values $\text{argmin}\{\text{PRE}_j[a], \dots, \text{PRE}_j[\text{PRE}_j[i]]\}$ in order to retrieve a witness shortest absent word. To this end, we build the $2n$ -bit RMQ data structure from Theorem 2 on each array PRE_j ; the data structure does not need to store the array PRE_j itself to compute argmin . The arrays B_j , for $j \in [1, \ell - 1]$, equipped with select data structures, and the RMQ data structures on arrays PRE_j , for $j \in [1, \ell - 1]$, can be constructed in total $\mathcal{O}(n\ell) = \mathcal{O}(n \log_\sigma n)$ time and they altogether occupy $\mathcal{O}(n \log_\sigma n)$ bits of space, which is $\mathcal{O}(n)$ space when measured in machine words.

To answer a query $[a, b]$, it suffices to find the smallest j such that $\text{FTR}[j][b] < a$. We do this by finding where the predecessor of a lies in $\text{FTR}[1.. \ell - 1][b]$. To this end, we constructed n fusion trees: one per $\text{FTR}[1.. \ell - 1][i]$, resulting in a data structure of size $\Theta(n\ell) = \mathcal{O}(n \log_\sigma n)$ with $\mathcal{O}(1)$ query time. But now we do not store the arrays $\text{FTR}[1.. \ell - 1][i]$ explicitly, while still having $\mathcal{O}(1)$ -time “oracle” access to their entries on demand. Hence, we can construct a succinct fusion tree of Theorem 5, for each array $\text{FTR}[1.. \ell - 1][i]$, which takes $\mathcal{O}(\ell \log \log n)$ bits of space since the size of machine words is $w = \Theta(\log n)$ bits (a shared table mentioned in Theorem 5 is also precomputed for all the trees in $o(2^w) = o(n)$ time).

Thus, all the succinct fusion trees can be constructed in $\mathcal{O}(n \log_\sigma n)$ time and occupy $\mathcal{O}(n \log_\sigma n \log \log n)$ bits, which is $\mathcal{O}(n \log_\sigma \log n)$ space when measured in $\Theta(\log n)$ -bit machine words. The ISAW queries are answered in $\mathcal{O}(1)$ time by the same algorithm as in Section 3.

Now we are to further reduce the memory usage of the data structure. We truncate all the arrays $\text{FTR}[0.. \ell - 1][i]$ except those where i is a multiple of $\lfloor \log \log n \rfloor$ or $i = n$: namely, if i is a multiple of $\lfloor \log \log n \rfloor$ or $i = n$, then the succinct fusion tree for the whole array $\text{FTR}[1.. \ell - 1][i]$ is stored, occupying $\mathcal{O}(\ell \log \log n)$ bits, by Theorem 5; otherwise ($i \neq n$ is not a multiple of $\lfloor \log \log n \rfloor$), we store the succinct fusion tree only for the subarray $\text{FTR}[1.. \lfloor \log n / \log \log n \rfloor][i]$, thus taking $\mathcal{O}(\log n)$ bits, by Theorem 5. In total, the space used is $\mathcal{O}(\frac{n}{\log \log n} \ell \log \log n + n \log n) = \mathcal{O}(n \log n)$ in bits or $\mathcal{O}(n)$ in words.

In order to answer an ISAW query for $T[a.. b]$, we first check whether the length λ of a shortest absent word in $T[a.. b]$ is smaller than $\log n / \log \log n$ by querying the fusion tree of $\text{FTR}[1.. \lfloor \log n / \log \log n \rfloor][b]$. If it is the case, then we have computed the length λ and we find the absent word itself using RMQs exactly as in the $\mathcal{O}(n \log_\sigma \log n)$ -space solution described above.

Suppose that $\lambda \geq \log n / \log \log n$. We compute b' , the successor of b among the positions i for which we have not truncated $\text{FTR}[1.. \ell - 1][i]$: $b' = \min\{n, \lceil b / \lfloor \log \log n \rfloor \cdot \lfloor \log \log n \rfloor\}$. Observe that $[a, b] \subseteq [a, b']$. Then, using the fusion tree of $\text{FTR}[1.. \ell - 1][b']$, we compute the smallest m such that $\text{FTR}[m][b'] < a$. Then, m is the length of a shortest absent word in $T[a..b']$. Denote $X = T[a..b]$ and $T[a..b'] = XY$ where Y is a suffix of $T[a..b']$ of length $b' - b$. We obviously have $\lambda \leq m$. Since $|Y| = b' - b < \log n / \log \log n$ and $m \geq \log n / \log \log n$, we have $|Y| < m$. It follows from Lemma 6 that the answer λ is within a range of length 18 from m . Therefore, λ belongs to the range $[m - 18, m]$ and we can find it in $\mathcal{O}(1)$ time using $\mathcal{O}(1)$ queries of the $\mathcal{O}(n)$ -space data structure encapsulated by Lemma 2. We thus arrive at the main result of the paper.

Theorem 1. *Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$ of size σ , we can construct in $\mathcal{O}(n \log_\sigma n)$ time a data structure of size $\mathcal{O}(n)$ that, for any given query $[a, b]$, can compute in $\mathcal{O}(1)$ time a shortest string over Σ that does not occur in $T[a..b]$.*

7. Open problems

It remains open whether a data structure for the ISAW problem with the same query time and space complexities as the ones encapsulated in Theorem 1 can be constructed in linear time. Also, it is natural to pose the following related open problem, which may require the development of fundamentally different techniques. Given a string T of length n over an alphabet $\Sigma \subset \{1, 2, \dots, n^{\mathcal{O}(1)}\}$, preprocess T so that given a range $[i, j]$, we can return a representation of a shortest string over $\Sigma_{[i,j]}$ that is absent in the fragment $T[i] \cdots T[j]$ of T , where $\Sigma_{[i,j]}$ is the set of letters from Σ occurring in the fragment $T[i] \cdots T[j]$.

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] P. Abedin, A. Ganguly, W. Hon, K. Matsuda, Y. Nekrich, K. Sadakane, R. Shah, S.V. Thankachan, A linear-space data structure for range-lcp queries in poly-logarithmic time, *Theor. Comput. Sci.* 822 (2020) 15–22, <https://doi.org/10.1016/j.tcs.2020.04.009>.
- [2] P. Abedin, A. Ganguly, S.P. Pissis, S.V. Thankachan, Efficient data structures for range shortest unique substring queries, *Algorithms* 13 (2020) 276, <https://doi.org/10.3390/a13110276>.
- [3] A. Amir, A. Apostolico, G.M. Landau, A. Levy, M. Lewenstein, E. Porat, Range LCP, *J. Comput. Syst. Sci.* 80 (2014) 1245–1253, <https://doi.org/10.1016/j.jcss.2014.02.010>.
- [4] A. Amir, P. Charalampopoulos, S.P. Pissis, J. Radoszewski, Dynamic and internal longest common substring, *Algorithmica* 82 (2020) 3707–3743, <https://doi.org/10.1007/s00453-020-00744-0>.
- [5] L.A.K. Ayad, G. Badkobeh, G. Fici, A. Héliou, S.P. Pissis, Constructing antidictionaries of long texts in output-sensitive space, *Theory Comput. Syst.* 65 (2021) 777–797, <https://doi.org/10.1007/s00224-020-10018-5>.
- [6] M.A. Babenko, P. Gawrychowski, T. Kociumaka, I.I. Kolesnichenko, T. Starikovskaya, Computing minimal and maximal suffixes of a substring, *Theor. Comput. Sci.* 638 (2016) 112–121, <https://doi.org/10.1016/j.tcs.2015.08.023>.
- [7] M.A. Babenko, P. Gawrychowski, T. Kociumaka, T. Starikovskaya, Wavelet trees meet suffix trees, in: *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, SIAM, 2015*, pp. 572–591.
- [8] G. Badkobeh, P. Charalampopoulos, S.P. Pissis, Internal shortest absent word queries, in: *32nd Annual Symposium on Combinatorial Pattern Matching, CPM 2021, 2021*, pp. 24:1–24:18.
- [9] C. Barton, A. Héliou, L. Mouchard, S.P. Pissis, Linear-time computation of minimal absent words using suffix array, *BMC Bioinform.* 15 (2014) 388, <https://doi.org/10.1186/s12859-014-0388-9>.
- [10] C. Barton, A. Héliou, L. Mouchard, S.P. Pissis, Parallelising the computation of minimal absent words, in: *Parallel Processing and Applied Mathematics - 11th International Conference, PPAM 2015. Revised Selected Papers, Part II*, Springer, 2015, pp. 243–253.
- [11] M.A. Bender, M. Farach-Colton, The LCA problem revisited, in: *LATIN 2000: Theoretical Informatics, 4th Latin American Symposium, Proceedings*, Springer, 2000, pp. 88–94.
- [12] O. Berkman, U. Vishkin, Recursive star-tree parallel data structure, *SIAM J. Comput.* 22 (1993) 221–242, <https://doi.org/10.1137/0222017>.
- [13] P. Bille, I.L. Gørtz, M.B.T. Knudsen, M. Lewenstein, H.W. Vildhøj, Longest common extensions in sublinear space, in: *Combinatorial Pattern Matching - 26th Annual Symposium, CPM 2015, 2015*, pp. 65–76.
- [14] P. Bille, I.L. Gørtz, B. Sach, H.W. Vildhøj, Time-space trade-offs for longest common extensions, *J. Discret. Algorithms* 25 (2014) 42–50, <https://doi.org/10.1016/j.jda.2013.06.003>.
- [15] O. Birenzweige, S. Golan, E. Porat, Locally consistent parsing for text indexing in small space, in: *Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020, 2020*, pp. 607–626.
- [16] S. Chairungsee, M. Crochemore, Using minimal absent words to build phylogeny, *Theor. Comput. Sci.* 450 (2012) 109–116, <https://doi.org/10.1016/j.tcs.2012.04.031>.
- [17] T.M. Chan, K.G. Larsen, M. Pătraşcu, Orthogonal range searching on the ram, revisited, in: *Proceedings of the 27th ACM Symposium on Computational Geometry, SCG 2011, ACM, 2011*, pp. 1–10.
- [18] P. Charalampopoulos, M. Crochemore, G. Fici, R. Mercaş, S.P. Pissis, Alignment-free sequence comparison using absent words, *Inf. Comput.* 262 (2018) 57–68, <https://doi.org/10.1016/j.ic.2018.06.002>.
- [19] P. Charalampopoulos, M. Crochemore, S.P. Pissis, On extended special factors of a word, in: *String Processing and Information Retrieval - 25th International Symposium, SPIRE 2018, Springer, 2018*, pp. 131–138.
- [20] P. Charalampopoulos, P. Gawrychowski, S. Mozes, O. Weimann, An almost optimal edit distance oracle, in: *48th International Colloquium on Automata, Languages, and Programming, ICALP 2021, 2021*, pp. 48:1–48:20.

- [21] P. Charalampopoulos, T. Kociumaka, M. Mohamed, J. Radoszewski, W. Rytter, J. Straszyński, T. Waleń, W. Zuba, Counting distinct patterns in internal dictionary matching, in: 31st Annual Symposium on Combinatorial Pattern Matching, CPM 2020, 2020, pp. 8:1–8:15.
- [22] P. Charalampopoulos, T. Kociumaka, M. Mohamed, J. Radoszewski, W. Rytter, T. Waleń, Internal dictionary matching, *Algorithmica* 83 (2021) 2142–2169, <https://doi.org/10.1007/s00453-021-00821-y>.
- [23] P. Charalampopoulos, T. Kociumaka, P. Wellnitz, Faster approximate pattern matching: a unified approach, in: 61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, IEEE, 2020, pp. 978–989.
- [24] M. Crochemore, A. Héliou, G. Kucherov, L. Mouchard, S.P. Pissis, Y. Ramusat, Absent words in a sliding window with applications, *Inf. Comput.* 270 (2020), <https://doi.org/10.1016/j.ic.2019.104461>.
- [25] M. Crochemore, F. Mignosi, A. Restivo, Automata and forbidden words, *Inf. Process. Lett.* 67 (1998) 111–117, [https://doi.org/10.1016/S0020-0190\(98\)00104-5](https://doi.org/10.1016/S0020-0190(98)00104-5).
- [26] M. Crochemore, F. Mignosi, A. Restivo, S. Salemi, Data compression using antidictionaries, *Proc. IEEE* 88 (2000) 1756–1768, <https://doi.org/10.1109/5.892711>.
- [27] P. Dinklage, J. Fischer, A. Herlez, T. Kociumaka, F. Kurpicz, Practical performance of space efficient data structures for longest common extensions, in: 28th Annual European Symposium on Algorithms, ESA 2020, 2020, pp. 39:1–39:20.
- [28] M. Farach, Optimal suffix tree construction with large alphabets, in: 38th Annual Symposium on Foundations of Computer Science, FOCS 1997, IEEE Computer Society, 1997, pp. 137–143.
- [29] S. Ferenczi, Complexity of sequences and dynamical systems, *Discrete Math.* 206 (1999) 145–154, [https://doi.org/10.1016/S0012-365X\(98\)00400-2](https://doi.org/10.1016/S0012-365X(98)00400-2).
- [30] G. Fici, P. Gawrychowski, Minimal absent words in rooted and unrooted trees, in: *String Processing and Information Retrieval - 26th International Symposium, SPIRE 2019*, Springer, 2019, pp. 152–161.
- [31] G. Fici, F. Mignosi, A. Restivo, M. Sciortino, Word assembly through minimal forbidden words, *Theor. Comput. Sci.* 359 (2006) 214–230, <https://doi.org/10.1016/j.tcs.2006.03.006>.
- [32] G. Fici, A. Restivo, L. Rizzo, Minimal forbidden factors of circular words, *Theor. Comput. Sci.* 792 (2019) 144–153, <https://doi.org/10.1016/j.tcs.2018.05.037>.
- [33] N.J. Fine, H.S. Wilf, Uniqueness theorems for periodic functions, *Proc. Am. Math. Soc.* 16 (1965) 109–114, <http://www.jstor.org/stable/2034009>.
- [34] J. Fischer, V. Heun, Space-efficient preprocessing schemes for range minimum queries on static arrays, *SIAM J. Comput.* 40 (2011) 465–492, <https://doi.org/10.1137/090779759>.
- [35] M.L. Fredman, D.E. Willard, Surpassing the information theoretic bound with fusion trees, *J. Comput. Syst. Sci.* 47 (1993) 424–436, [https://doi.org/10.1016/0022-0000\(93\)90040-4](https://doi.org/10.1016/0022-0000(93)90040-4).
- [36] Y. Fujishige, Y. Tsujimaru, S. Inenaga, H. Bannai, M. Takeda, Computing DAWGs and minimal absent words in linear time for integer alphabets, in: 41st International Symposium on Mathematical Foundations of Computer Science, MFCS 2016, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2016, pp. 38:1–38:14.
- [37] A. Ganguly, M. Patil, R. Shah, S.V. Thankachan, A linear space data structure for range LCP queries, *Fundam. Inform.* 163 (2018) 245–251, <https://doi.org/10.3233/FI-2018-1741>.
- [38] S.P. Garcia, A.J. Pinho, J.M.O.S. Rodrigues, C.A.C. Bastos, P.J.S.G. Ferreira, Minimal absent words in prokaryotic and eukaryotic genomes, *PLoS ONE* 6 (2011), <https://doi.org/10.1371/journal.pone.0016065>.
- [39] R. Grossi, A. Orlandi, R. Raman, S.S. Rao, More haste, less waste: lowering the redundancy in fully indexable dictionaries, in: 26th International Symposium on Theoretical Aspects of Computer Science, STACS 2009, 2009, pp. 517–528.
- [40] D. Harel, R.E. Tarjan, Fast algorithms for finding nearest common ancestors, *SIAM J. Comput.* 13 (1984) 338–355, <https://doi.org/10.1137/0213024>.
- [41] G. Jacobson, Space-efficient static trees and graphs, in: 30th Annual Symposium on Foundations of Computer Science, FOCS 1989, IEEE Computer Society, 1989, pp. 549–554.
- [42] O. Keller, T. Kopelowitz, S.L. Feibish, M. Lewenstein, Generalized substring compression, *Theor. Comput. Sci.* 525 (2014) 42–54, <https://doi.org/10.1016/j.tcs.2013.10.010>.
- [43] D. Kempa, T. Kociumaka, String synchronizing sets: sublinear-time BWT construction and optimal LCE data structure, in: *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019*, ACM, 2019, pp. 756–767.
- [44] T. Kociumaka, Minimal suffix and rotation of a substring in optimal time, in: 27th Annual Symposium on Combinatorial Pattern Matching, CPM 2016, 2016, pp. 28:1–28:12.
- [45] T. Kociumaka, Efficient Data Structures for Internal Queries in Texts, Ph.D. thesis, University of Warsaw, 2018, <https://mimuw.edu.pl/~kociumaka/files/phd.pdf>.
- [46] T. Kociumaka, J. Radoszewski, W. Rytter, T. Waleń, Efficient data structures for the factor periodicity problem, in: *String Processing and Information Retrieval - 19th International Symposium, SPIRE 2012*, 2012, pp. 284–294.
- [47] T. Kociumaka, J. Radoszewski, W. Rytter, T. Waleń, Internal pattern matching queries in a text and applications, in: *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015*, SIAM, 2015, pp. 532–551.
- [48] D. Kosolobov, Tight lower bounds for the longest common extension problem, *Inf. Process. Lett.* 125 (2017) 26–29, <https://doi.org/10.1016/j.ipl.2017.05.003>.
- [49] G.M. Landau, U. Vishkin, Fast string matching with k differences, *J. Comput. Syst. Sci.* 37 (1988) 63–78, [https://doi.org/10.1016/0022-0000\(88\)90045-1](https://doi.org/10.1016/0022-0000(88)90045-1).
- [50] K. Matsuda, K. Sadakane, T. Starikovskaya, M. Tateshita, Compressed orthogonal search on suffix arrays with applications to range LCP, in: 31st Annual Symposium on Combinatorial Pattern Matching, CPM 2020, 2020, pp. 23:1–23:13.
- [51] T. Mieno, Y. Kuhara, T. Akagi, Y. Fujishige, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, Minimal unique substrings and minimal absent words in a sliding window, in: 46th SOFSEM, Springer, 2020, pp. 148–160.
- [52] F. Mignosi, A. Restivo, M. Sciortino, Words and forbidden factors, *Theor. Comput. Sci.* 273 (2002) 99–117, [https://doi.org/10.1016/S0304-3975\(00\)00436-9](https://doi.org/10.1016/S0304-3975(00)00436-9).
- [53] G. Navarro, *Compact Data Structures - a Practical Approach*, Cambridge University Press, 2016.
- [54] T. Ota, H. Morita, On the adaptive antidictionary code using minimal forbidden words with constant lengths, in: *Proceedings of the International Symposium on Information Theory and Its Applications, ISITA 2010*, IEEE, 2010, pp. 72–77.
- [55] M. Pătraşcu, M. Thorup, Dynamic integer sets with optimal rank, select, and predecessor search, in: 55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, IEEE Computer Society, 2014, pp. 166–175.
- [56] D. Pratas, J.M. Silva, Persistent minimal sequences of SARS-CoV-2, *Bioinformatics* (2020), <https://doi.org/10.1093/bioinformatics/btaa686>.
- [57] N. Prezza, In-place sparse suffix sorting, in: *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018*, 2018, pp. 1496–1508.
- [58] S. Raskhodnikova, D. Ron, R. Rubinfeld, A.D. Smith, Sublinear algorithms for approximating string compressibility, *Algorithmica* 65 (2013) 685–709, <https://doi.org/10.1007/s00453-012-9618-6>.
- [59] M. Rubinchik, A.M. Shur, Counting palindromes in substrings, in: *String Processing and Information Retrieval - 24th International Symposium, SPIRE 2017*, Springer, 2017, pp. 290–303.
- [60] Y. Sakai, A substring-substring LCS data structure, *Theor. Comput. Sci.* 753 (2019) 16–34, <https://doi.org/10.1016/j.tcs.2018.06.034>.

- [61] Y. Sakai, A data structure for substring–substring lcs length queries, *Theor. Comput. Sci.* 911 (2022) 41–54.
- [62] R.M. Silva, D. Pratas, L. Castro, A.J. Pinho, P.J.S.G. Ferreira, Three minimal sequences found in Ebola virus genomes and absent from human DNA, *Bioinformatics* 31 (2015) 2421–2425, <https://doi.org/10.1093/bioinformatics/btv189>.
- [63] Y. Tanimura, T. I. H. Bannai, S. Inenaga, S.J. Puglisi, M. Takeda, Deterministic sub-linear space LCE data structures with efficient construction, in: 27th Annual Symposium on Combinatorial Pattern Matching, CPM 2016, 2016, pp. 1:1–1:10.
- [64] Y. Tanimura, T. Nishimoto, H. Bannai, S. Inenaga, M. Takeda, Small-space LCE data structure with constant-time queries, in: 42nd International Symposium on Mathematical Foundations of Computer Science, MFCS 2017, 2017, pp. 10:1–10:15.
- [65] A. Tiskin, Semi-local string comparison: algorithmic techniques and applications, *Math. Comput. Sci.* 1 (2008) 571–603, <https://doi.org/10.1007/s11786-007-0033-3>.
- [66] A.C. Yao, Space-time tradeoff for answering range queries (extended abstract), in: Proceedings of the Fourteenth Annual ACM Symposium on Theory of Computing, STOC 1982, ACM, 1982, pp. 128–136.