Optimal Square Detection Over General Alphabets

Jonas Ellert^{*1}, Paweł Gawrychowski^{†2}, and Garance Gourdel^{†3}

¹Technical University of Dortmund, Germany ²Institute of Computer Science, University of Wrocław, Poland

³DI/ENS, PSL Research University, IRISA Inria Rennes, France

Abstract

Squares (fragments of the form xx, for some string x) are arguably the most natural type of repetition in strings. The basic algorithmic question concerning squares is to check if a given string of length n is square-free, that is, does not contain a fragment of such form. Main and Lorentz [J. Algorithms 1984] designed an $\mathcal{O}(n \log n)$ time algorithm for this problem, and proved a matching lower bound assuming the so-called general alphabet, meaning that the algorithm is only allowed to check if two characters are equal. However, their lower bound also assumes that there are $\Omega(n)$ distinct symbols in the string. As an open question, they asked if there is a faster algorithm if one restricts the size of the alphabet. Crochemore [Theor. Comput. Sci. 1986] designed a linear-time algorithm for constant-size alphabets, and combined with more recent results his approach in fact implies such an algorithm for linearly-sortable alphabets. Very recently, Ellert and Fischer [ICALP 2021] significantly relaxed this assumption by designing a linear-time algorithm for general ordered alphabets, that is, assuming a linear order on the characters that permits constant time order comparisons. However, the open question of Main and Lorentz from 1984 remained unresolved for general (unordered) alphabets. In this paper, we show that testing square-freeness of a length-n string over general alphabet of size σ can be done with $\mathcal{O}(n\log\sigma)$ comparisons, and cannot be done with $o(n\log\sigma)$ comparisons. We complement this result with an $\mathcal{O}(n \log \sigma)$ time algorithm in the Word RAM model. Finally, we extend the algorithm to reporting all the runs (maximal repetitions) in the same complexity.

1 Introduction

The notion of repetition is a central concept in combinatorics on words and algorithms on strings. In this context, a word or a string is simply a sequence of characters from some finite alphabet Σ . In the most basic version, a repetition consists of two (or more) consecutive occurrences of the same fragment. Repetitions are interesting not only from a purely theoretical point of view, but are also very relevant in bioinformatics [51]. A repetition could be a square, defined as two consecutive occurrences of the same fragment, a higher power (for example, a cube), or a run, which is a length-wise maximal periodic substring. For example, both anan and nana are squares with two occurrences each in banananas, and they belong to the same run ananana. In this paper, we start by focusing on squares, then generalize our results for runs.

^{*}Partially supported by the German Research Association (DFG) within the Collaborative Research Center SFB 876, project A6.

[†]Partially supported by the grant ANR-20-CE48-0001 from the French National Research Agency (ANR).

The study of squares in strings goes back to the work of Thue published in 1906 [73], who considered the question of constructing an infinite word with no squares. It is easy to see that any sufficiently long binary word must contain a square, and Thue proved that there exists an infinite ternary word with no squares. His result has been rediscovered multiple times, and in 1979 Bean, Ehrenfeucht and McNulty [9] started a systematic study of the so-called avoidable repetitions, see for example the survey by Currie [27].

Combinatorics on words. The basic tool in the area of combinatorics on words is the so-called periodicity lemma. A period of a string T[1..n] is an integer d such that T[i] = T[i+d] for every $i \in [1, n-d]$, and the periodicity lemma states that if p and q are both such periods and $p+q \leq n + \gcd(p,q)$ then $\gcd(p,q)$ is also a period [32]. This was generalised in a myriad of ways, for strings [17,49,74], partial words (words with don't cares) [11–13,47,50,69,70], Abelian periods [14,20], parametrized periods [6], order-preserving periods [42,64], approximate periods [2–4]. Now, a square can be defined as a fragment of length twice its period. The string \mathbf{a}^n contains $\Omega(n^2)$ such fragments, thus from the combinatorial point of view it is natural to count only distinct squares. Fraenkel and Simpson [35] showed an upper bound of 2n and a lower bound of n – $\Theta(\sqrt{n})$ for the maximum number of distinct squares in a length-n string. After a sequence of improvements [28, 48, 72], the upper bound was very recently improved to n [16]. The last result was already generalised to higher powers [59]. Another way to avoid the trivial examples such as \mathbf{a}^n is to count only maximal periodic fragments, that is, fragments with period at most half of their length and that cannot be extended to the left or to the right without breaking the period. Such fragments are usually called runs. Kolpakov and Kucherov [52] showed an upper bound of $\mathcal{O}(n)$ on their number, and this started a long line of work on determining the exact constant [24,25,40,41,67,68], culminating in the paper of Bannai et al. [8] showing an upper bound of n, and followed by even better upper bounds for binary strings [33, 45]. This was complemented by a sequence of lower bounds [36, 62, 63, 71].

Algorithms on strings. In this paper, we are interested in the algorithmic aspects of detecting repetitions in strings. The most basic question in this direction is checking if a given length-n string contains at least one square, while the most general version asks for computing all the runs. Testing square-freeness was first considered by Main and Lorentz [61], who designed an $\mathcal{O}(n \log n)$ time algorithm based on a divide-and-conquer approach and a linear-time procedure for finding all new squares obtained when concatenating two strings. In fact, their algorithm can be used to find (a compact representation of) all squares in a given string within the same time complexity. They also proved that any algorithm based on comparisons of characters needs $\Omega(n \log n)$ such operations to test square-freeness in the worst case. Here, comparisons of characters means checking if characters at two positions of the input string are equal. However, to obtain the lower bound they had to consider instances consisting of even up to n distinct characters, that is, over alphabet of size n. This is somewhat unsatisfactory, and motivates the following open question that was explicitly asked by Main and Lorentz [61]:

Question 1.1. Is there a faster algorithm to determine if a string is square-free if we restrict the size of the alphabet?

Another $\mathcal{O}(n \log n)$ time algorithm for finding all repetitions was given by Crochemore [22], who also showed that for constant-size alphabets testing square-freeness can be done in $\mathcal{O}(n)$ time [23]. In fact, the latter algorithm works in $\mathcal{O}(n \log \sigma)$ time for alphabets of size σ with a linear order on the characters. That is, it needs to test if the character at some position is smaller than the character at another position. In the remaining part of the paper, we will refer to this model as general ordered alphabet, while the model in which we can only test equality of characters will be called general (unordered) alphabet. Later, Kosaraju [53] showed that in fact, assuming constant-size alphabet, $\mathcal{O}(n)$ time is enough to find the shortest square starting at each position of the input string. Apostolico and Preparata [7] provide another $\mathcal{O}(n \log n)$ time algorithm assuming a general ordered alphabet, based more on data structure considerations than combinatorial properties of words. Finally, a number of alternative $\mathcal{O}(n \log n)$ and $\mathcal{O}(n \log \sigma)$ time algorithms (respectively, for general unordered and general ordered alphabets) can be obtained from the work on online [46, 54, 56] and parallel [5] square detection (interestingly, this cannot be done efficiently in the related streaming model [65, 66]).

Faster algorithms for testing square-freeness of strings over general ordered alphabets were obtained as a byproduct of the more general results on finding all runs. Kolpakov and Kucherov [52] not only proved that any length-*n* string contains only $\mathcal{O}(n)$ runs, but also showed how to find them in the same time assuming linearly-sortable alphabet. Every square is contained in a run, and every run contains at least one square, thus this in particular implies a linear-time algorithm for testing square-freeness over such alphabets. For general ordered alphabets, Kosolobov [55] showed that the decision tree complexity of this problem is only $\mathcal{O}(n)$, and later complemented this with an efficient $\mathcal{O}(n(\log n)^{2/3})$ time algorithm [57] (still using only $\mathcal{O}(n)$ comparisons). The time complexity was then improved to $\mathcal{O}(n \log \log n)$ by providing a general mechanism for answering longest common extension (LCE) queries for general ordered alphabets [39], and next to $\mathcal{O}(n\alpha(n))$ by observing that the LCE queries have additional structure [26]. Finally, Ellert and Fischer provided an elegant $\mathcal{O}(n)$ time algorithm, thus fully resolving the complexity of square detection for general ordered alphabets. However, for general (unordered) alphabets the question of Main and Lorentz remains unresolved, with the best upper bound being $\mathcal{O}(n \log n)$, and only known to be asymptotically tight for alphabets of size $\Theta(n)$.

General alphabets. While in many applications one can without losing generality assume some ordering on the characters of the alphabet, no such ordering is necessary for defining what a square is. Thus, it is natural from the mathematical point of view to seek algorithms that do not require such an ordering to efficiently test square-freeness. Similar considerations have lead to multiple beautiful results concerning the pattern matching problem, such as constant-space algorithms [15, 38], or the works on the exact number of required equality comparisons [18, 19] More recent examples include the work of Duval, Lecroq, and Lefebvre [29] on computing the unbordered conjugate/rotation, and Kosolobov [58] on finding the leftmost critical point.

Main results. We consider the complexity of checking if a given string T[1..n] containing σ distinct characters is square-free. The input string can be only accessed by issuing comparisons $T[i] \stackrel{?}{=} T[j]$, and the value of σ is not assumed to be known. We start by analysing the decision tree complexity of the problem. That is, we only consider the required and necessary number of comparisons, without worrying about an efficient implementation. We show that, even if the value of σ is assumed to be known, $\Omega(n \log \sigma)$ comparisons are required.

Theorem 1.1. For any integers n and σ with $8 \leq \sigma \leq n$, there is no deterministic algorithm that performs at most $n \ln \sigma - 3.6n = \mathcal{O}(n \ln \sigma)$ comparisons in the worst case, and determines whether a length-n string with at most σ distinct symbols from a general unordered alphabet is square-free.

Next, we show that $\mathcal{O}(n \log \sigma)$ comparisons are sufficient. We stress that the value of σ is not assumed to be known. In fact, as a warm-up for the above theorem, we first prove that finding

a sublinear multiplicative approximation of this value requires $\Omega(n\sigma)$ comparisons. This does not contradict the claimed upper bound, as we are only saying that the number of comparisons used on a particular input string is at most $\mathcal{O}(n \log \sigma)$, but might actually be smaller. Thus, it is not possible to extract any meaningful approximation of the value of σ from the number of used comparisons.

Theorem 1.2. Testing square-freenes of a length-n string that contains σ distinct symbols from a general unordered alphabet can be done with $\mathcal{O}(n \log \sigma)$ comparisons.

The proof of the above result is not efficient in the sense that it only restricts the overall number of comparisons, and not the time to actually figure out which comparisons should be used. A direct implementation results in a quadratic time algorithm. We first show how to improve this to $\mathcal{O}(n \log \sigma + n \log^* n)$ time (while still keeping the asymptotically optimal $\mathcal{O}(n \log \sigma)$ number of comparisons), and finally to $\mathcal{O}(n \log \sigma)$. In this part of the paper, we assume the Word RAM model with word of length $\Omega(\log n)$. We stress that the input string is still assumed to consist of characters that can be only tested for equality, that is, one should think that we are given oracle access to a functions that, given *i* and *j*, checks whether T[i] = T[j].

Theorem 1.3. Testing square-freeness of a length-n string that contains σ distinct symbols from a general unordered alphabet can be implemented in $\mathcal{O}(n \log \sigma)$ comparisons and time.

Finally, we also generalize this result to the computation of runs.

Theorem 1.4. Computing all runs in a length-n string that contains σ distinct symbols from a general unordered alphabet can be implemented in $\mathcal{O}(n \log \sigma)$ comparisons and time.

Altogether, our results fully resolve the open question of Main and Lorentz for the case of general unordered alphabets and deterministic algorithms. We leave extending our lowerbound to randomised algorithms as an open question.

Overview of the methods. As mentioned before, Main and Lorentz [61] designed an $\mathcal{O}(n \log n)$ time algorithm for testing square-freeness of length-n strings over general alphabets. The high-level idea of their algorithm goes as follows. They first designed a procedure for checking, given two strings x and y, if their concatenation contains a square that is not fully contained in x nor y in $\mathcal{O}(|x| + |y|)$ time. Then, a divide-and-conquer approach can be used to detect a square in the whole input string in $\mathcal{O}(n \log n)$ total time. For general alphabets of unbounded size this cannot be improved, but Crochemore [23] showed that, for general ordered alphabets of size σ , a faster $\mathcal{O}(n \log \sigma)$ time algorithm exists. The gist of his approach is to first obtain the so-called f-factorisation of the input string (related to the well-known Lempel-Ziv factorisation), that in a certain sense "discovers" repetitive fragments. Then, this factorisation can be used to apply the procedure of Main and Lorentz on appropriately selected fragments of the input strings in such a way that the leftmost occurrence of every distinct square is detected, and the total length of the strings on which we apply the procedure is only $\mathcal{O}(n)$. The factorisation can be found in $\mathcal{O}(n \log \sigma)$ time for general ordered alphabets of size σ by, roughly speaking, constructing some kind of suffix structure (suffix array, suffix tree or suffix automaton).

For general (unordered) alphabets, computing the *f*-factorisation (or anything similar) seems problematic, and in fact we show (as a corollary of our lower bound on approximating the alphabet size) that computing the *f*-factorisation or Lempel-Ziv-factorisation (LZ-factorisation) of a given length-*n* string containing σ distinct characters requires $\Omega(n\sigma)$ equality tests. Thus, we need another approach. Additionally, the $\mathcal{O}(n)$ time algorithm of Ellert and Fischer [30] hinges on the notion of Lyndon words, which is simply not defined for strings over general alphabets. Thus, at first glance it might seem that $\Theta(n\sigma)$ is the right time complexity for testing square-freeness over length*n* strings over general alphabets of size σ . However, due to the $\Omega(n \log n)$ lower bound of Main and Lorentz for testing square-freeness of length-*n* string consisting of up to *n* distinct characters, one might hope for an $\mathcal{O}(n \log \sigma)$ time algorithm when there are only σ distinct characters.

We begin our paper with a lower bound of $\Theta(n \log \sigma)$ for such strings. Intuitively, we show that testing square-freeness has the direct sum property: $\frac{n}{\sigma}$ instances over length- σ strings can be combined into a single instance over length-n string. As in the proof of Main and Lorentz, we use the adversarial method. While the underlying calculation is essentially the same, we need to appropriately combine the smaller instances, which is done using the infinite square-free Prouhet-Thue-Morse sequence, and use significantly more complex rules for resolving the subsequent equality tests. As a warm-up for the adversarial method, we prove that computing any meaningful approximation of the number of distinct characters requires $\Omega(n\sigma)$ such tests, and that this implies the same lower bound on computing the f-factorisation and the Lempel-Ziv factorisation (if the size of the alphabet is unknown in advance).

We then move to designing an approach that uses $\mathcal{O}(n \log \sigma)$ equality comparisons to test squarefreeness. As discussed earlier, one way of detecting squares uses the *f*-factorisation of the string, which is similar to its LZ factorisation. However, as we prove in Corollary 1 and 2, we cannot compute either of these factorisations over a general unordered alphabet in $o(n\sigma)$ comparisons. Therefore, we will instead use a novel type of factorisation, Δ -approximate LZ factorisation, that can be seen as an approximate version of the LZ factorisation. Intuitively, its goal is to "capture" all sufficiently long squares, while the original LZ factorisation (or *f*-factorisation) captures all squares. Each phrase in a Δ -approximate LZ factorisation consists of a head of length at most Δ and a tail (possibly empty) that must occur at least once before, such that the whole phrase is at least as long as the classical LZ phrase starting at the same position. Contrary to the classical LZ factorisation, this factorisation is not unique. The advantage of our modification is that there are fewer phrases (and there is more flexibility as to what they should be), and hence one can hope to compute such factorisation more efficiently.

To design an efficient construction method for Δ -approximate LZ factorisation, we first show how to compute a sparse suffix tree while trying to use only a few symbol comparisons. This is then applied on a set of positions from a so-called difference cover with some convenient synchronizing properties. Then, a Δ -approximate LZ factorisation allows us to detect squares of length $\geq 8\Delta$.

The first warm-up algorithm fixes Δ depending on n and σ (assuming that σ is known), and uses the approximate LZ factorisation to find all squares of length at least 8Δ . It then finds all the shorter squares by dividing the string in blocks of length 8Δ , and applying the original algorithm by Main and Lorentz on each block pair. Our choice of Δ leads to $\mathcal{O}(n(\lg \sigma + \lg \lg n))$ comparisons.

The improved algorithm does not need to know σ , and instead starts with a large $\Delta = \Omega(n)$, and then progressively decreases Δ in at most $\mathcal{O}(\lg \lg n)$ phases, where later phases detect shorter squares. As soon as we notice that there are many distinct characters in the alphabet, by carefully adjusting the parameters we can afford switching to the approach of Main and Lorentz on sufficiently short fragments of the input string. Since we cannot afford $\Omega(n)$ comparisons per phase, we use a deactivation technique, where whenever we perform a large number of comparisons in a phase, we will discard a large part of the string in all following phases. More precisely, during a given phase, we avoid looking for squares in a fragment fully contained in a tail from an earlier phase. This leads to optimal $\mathcal{O}(n \lg \sigma)$ comparisons.

The above approach uses an asymptotically optimal number of equality tests in the worst case, but does not result in an efficient algorithm. The main bottleneck is constructing the sparse suffix trees. However, it is not hard to provide an efficient implementation using the general mechanism for answering LCE queries for strings over general alphabets [39]. Unfortunately, the best known approach for answering such queries incurs an additional $\mathcal{O}(n \log^* n)$ in the time complexity, even if the size of the alphabet is constant. We overcome this technical hurdle by carefully deactivating fragments of the text to account for the performed work.

Many of our techniques can easily be modified to compute all runs rather than detecting squares. We exploit that the approximate factorisation reveals long substrings with an earlier occurrence. Hence we compute runs only for the first occurrence of such substrings, while for later occurrences we simply copy the already computed runs. By carefully arranging the order of the computation, we ensure that the total time for copying is bounded by the number of runs, which is known to be $\mathcal{O}(n)$. This way, we achieve $\mathcal{O}(n \lg \sigma)$ time and comparisons to compute all runs.

2 Preliminaries

Strings. A string of length n is a sequence $T[1] \ldots T[n]$ of characters from a finite alphabet Σ of size σ . The substring T[i..j] is the string $T[i] \cdots T[j]$, whereas the fragment T[i..j] refers to the specific occurrence of T[i..j] starting at position i in T. If i > j, then T[i..j] is the empty string. A suffix of T has the form T[i..n]. We say that a fragment T[i'..j'] is properly contained in another fragment T[i..j] if $i < i' \leq j' < j$. A substring is properly contained in T[i..j], if it equals a fragment that is properly contained in T[i..j]. We write T[i..j) a shortcut for T[i..j-1]. Similarly, we write [i, j] = [i, j + 1) as a shortcut for the integer interval $\{i, \ldots, j\}$. Given two positions $i \leq j$, their longest common extension (LCE) is the length of the longest common prefix between suffixes T[i..n] and T[j..n], formally defined as $LCE(i, j) = LCE(j, i) = \max\{\ell \in \{0, \ldots, n - j + 1\} \mid T[i..i + \ell) = T[j..j + \ell)\}$.

Definition 2.1. A positive integer p is a period of a string T[1..n] if T[i] = T[i + p] for every $i \in \{1, ..., n - p\}$. The smallest such p is called the period of T[1..n], and we call a string periodic if its period p is at most $\frac{n}{2}$.

Computational model. For a general unordered alphabet Σ , the only allowed operation on the characters is comparing for equality. In particular, there is no linear order on the alphabet. Unless explicitly stated otherwise, we will only use such comparisons. A general ordered alphabet has a total order, such that comparisons of the type less-equals are possible.

In the algorithmic part of the paper, we assume the standard unit-cost Word RAM model with words of length $\Omega(\log n)$, but the algorithm is only allowed to access the input string T[1..n] by comparisons $T[i] \stackrel{?}{=} T[j]$, which are assumed to take constant time. We say that a string of length n is over a linearly-sortable alphabet, if we can sort the n symbols of the string in $\mathcal{O}(n)$ time. Note that whether or not an alphabet is linearly-sortable depends not only on the alphabet, but also on the string. For example, the alphabet $\Sigma = \{1, \ldots, m^{\mathcal{O}(1)}\}$ is linearly-sortable for strings of length $n = \Omega(m)$ (e.g., using radix sort), but it is unknown whether it is linearly-sortable for all strings of length n = o(m) [44]. Our algorithm will internally use strings over linearly-sortable alphabets. We stress that in such strings the characters are not the characters from the input string, but simply integers calculated by the algorithm. Note that every linearly-sortable alphabet is also a general ordered alphabet.

Squares and runs. A square is a length- 2ℓ fragment of period ℓ . The following theorem is a classical result by Main and Lorentz [61].

Theorem 2.1. Testing square-freenes of T[1..n] over a general alphabet can be implemented in $\mathcal{O}(n \log n)$ time and comparisons.

The proof of the above theorem is based on running a divide-and-conquer procedure using the following lemma.

Lemma 1. Given two strings x and y over a general alphabet, we can test if there is a square in xy that is not fully contained in x nor y in O(|x| + |y|) time and comparisons.

A repetition is a length- ℓ fragment of period at most $\frac{\ell}{2}$. A run is a maximal repetition. Formally, a repetition in T[1..n] is a triple $\langle s, e, p \rangle$ with $s, e \in [1, n]$ and $p \in [1, \frac{e-s+1}{2}]$ such that p is the smallest period of T[s..e]. A run is a repetition $\langle s, e, p \rangle$ that cannot be extended to the left nor to the right with the same period, in other words s = 1 or $T[s-1] \neq T[s-1+p]$ and e = n or $T[e+1] \neq T[e+1-p]$. The celebrated runs conjecture, proven by Bannai et al. [8], states that the number of runs is any length-n string is less than n. Ellert and Fischer [30] showed that all runs in a string over a general ordered alphabet can be computed in $\mathcal{O}(n)$ time. As mentioned earlier, each run contains a square, and each square is contained in a run. Thus, the string contains a square if and only if it contains a run, and it follows:

Theorem 2.2. Computing all runs (and thus testing square-freeness) of T[1..n] over a general ordered alphabet can be implemented in $\mathcal{O}(n)$ time.

Lempel-Ziv factorisation. The unique LZ phrase starting at position s of T[1..n] is a fragment T[s..e] such that T[s..(e-1)] occurs at least twice in T[1..(e-1)] and either e = n or T[s..e] occurs only once in T[1..e]. The Lempel-Ziv factorisation of T consists of z phrases f_1, \ldots, f_z such that the concatenation $f_1 \ldots f_z$ is equal to T[1..n] and each f_i is the unique LZ phrase starting at position $1 + \sum_{j=1}^{i-1} |f_j|$.

Tries. Given a collection $S = \{T_1, \ldots, T_k\}$ of strings over some alphabet Σ , its trie is a rooted tree with edge labels from Σ . For any node v, the concatenation of the edge labels from the root to the node *spells* a string. The string-depth of a node is the length of the string that it spells. No two nodes spell the same string, i.e., for any node, the labels of the edges to its children are pairwise distinct. Each leaf spells one of the T_i , and each T_i is spelled by either an internal node or a leaf.

The compacted trie of S can be obtained from its (non-compacted) trie by contracting each path between a leaf or a branching node and its closest branching ancestor into a single edge (i.e., by contraction we eliminate all non-branching internal nodes). The label of the new edge is the concatenation of the edge labels of the contracted path in root to leaf direction. Since there are at most k leaves and all internal nodes are branching, there are O(k) nodes in the compacted trie. Each edge label is some substring $T_i[s..e]$ of the string collection, and we can avoid explicitly storing the label by instead storing the reference (i, s, e). Thus O(k) words are sufficient for storing the compacted trie. Consider a string T' that is spelled by a node of the non-compacted trie. We say that T' is explicit, if and only if it is spelled by a node of the compacted trie. Otherwise T' is implicit.

The suffix tree of a string T[1..n] is the compacted trie containing exactly its suffixes, i.e., a trie over the string collection $\{T[i..n] \mid i \in \{1, ..., n\}\}$. It is one of the most fundamental data structures in string algorithmics, and is widely used, e.g., for compression and indexing [43]. The suffix tree can be stored in $\mathcal{O}(n)$ words of memory, and for linearly-sortable alphabets it can be computed in $\mathcal{O}(n)$ time [31]. The sparse suffix tree of T for some set $B \subseteq \{1, ..., n\}$ of sample positions is the compacted trie containing exactly the suffixes $\{T[i..n] \mid i \in B\}$. It can be stored in $\mathcal{O}(|B|)$ words of memory. We assume that T is terminated by some special symbol T[n] =that occurs nowhere else in T. This ensures that each suffix is spelled by a leaf, and we label the leaves with the respective starting positions of the suffixes. Note that for any two leaves $i \neq j$, their lowest common ancestor (i.e., the deepest node that is an ancestor of both i and j) spells a string of length LCE(i, j).

3 Lower Bounds

In this section, we show lower bounds on the number of symbol comparisons required to compute a meaningful approximation of the alphabet size (Section 3.1) and to test square-freeness (Section 3.2). For both bounds we use an adversarial method, which we briefly outline now.

The present model of computation may be interpreted as follows. An algorithm working on a string over a general unordered alphabet has no access to the actual string. Instead, it can only ask an oracle whether or not there are identical symbols at two positions. The number of questions asked is exactly the number of performed comparisons. In order to show a lower bound on the number of comparisons required to solve some problem, we describe an adversary that takes over the role of the oracle, forcing the algorithm to perform as many symbol comparisons as possible.

We use a conflict graph G = (V, E) with $V = \{1, \ldots, n\}$ and $E \subseteq V^2$ to keep track of the answers given by the adversary. The nodes directly correspond to the positions of the string. Initially, we have $E = \emptyset$ and all nodes are colorless, which formally means that they have color $\gamma(i) = \bot$. During the algorithm execution, the adversary may assign colors from the set $\Sigma = \{0, \ldots, n-1\}$ to the nodes, which can be seen as permanently fixing the alphabet symbol at the corresponding position (i.e., each node gets colored at most once). The rule used for coloring nodes depends on the lower bound that we want to show (we describe this in detail in the respective sections). Apart from this coloring rule, the general behaviour of the adversary is as follows. Whenever the algorithm asks whether T[i] = T[j] holds, the adversary answers "yes" if and only if $\gamma(i) = \gamma(j) \neq \bot$. Otherwise, it answers "no" and inserts an edge (i, j) into E. Whenever the adversary assigns the color of a node, it has to choose a color that is not used by any of the adjacent nodes in the conflict graph. This ensures that the coloring does not contradict the answers given in the past.

Let us define a set $\mathcal{T} \subseteq \Sigma^n$ of strings that is consistent with the answers given by the adversary. A string $T \in \Sigma^n$ is a member of \mathcal{T} if

$$\forall i \in V : \gamma(i) \in \{\bot, T[i]\} \quad \land \quad \forall i, j \in V : (T[i] = T[j]) \implies (i, j) \notin E.$$

Note that \mathcal{T} changes over time. Initially (before the algorithm starts), we have $\mathcal{T} = \Sigma^n$. With every question asked, the algorithm might eliminate some strings from \mathcal{T} . However, there is always at least on string in \mathcal{T} , which can be obtained by coloring each colorless node in a previously entirely unused color.

3.1 Approximating the Alphabet Size

Given a string T[1..n] of unknown alphabet size $\sigma \geq 2$, assume that we want to compute an approximation of σ . We show that if an algorithm takes at most $\frac{n\sigma}{8}$ comparisons in the worst-case, then it cannot distinguish strings with at most σ distinct symbols from strings with at least $\frac{n}{2}$ distinct symbols. Thus, any meaningful approximation of σ requires $\Omega(n\sigma)$ comparisons.

For the sake of the proof, consider an algorithm that performs at most $\frac{n\sigma}{8}$ comparisons when given a length-*n* string with at most $\sigma \geq 2$ distinct symbols. We use an adversary as described at the beginning of Section 3, and ensure that the set \mathcal{T} of strings consistent with the adversary's answers always contains a string with at most σ distinct symbols. Thus, the algorithm terminates after at most $\frac{n\sigma}{8}$ comparisons. At the same time, we ensure that \mathcal{T} also contains a string with at least $\frac{n}{2}$ distinct symbols, which yields the desired result. The adversary is equipped with the following coloring rule. All colors are from $\{1, \ldots, \sigma\}$. Whenever the degree of a node in the conflict graph becomes $\sigma - 1$, we assign its color. We avoid the colors of the $\sigma - 1$ adjacent nodes in the conflict graph. At any moment in time, we could hypothetically complete the coloring by assigning one of the colors $\{1, \ldots, \sigma\}$ to each colorless node, avoiding the colors of adjacent nodes. This way, each node gets assigned one of the σ colors, which means that \mathcal{T} contains a string with at most σ distinct symbols. It follows that the algorithm terminates after at most $\frac{n\sigma}{8}$ comparisons. Each comparison may increase the degree of two nodes by one. Thus, after $\frac{n\sigma}{8}$ comparisons, there are at most $\frac{n\sigma}{8} \cdot \frac{2}{\sigma-1} \leq \frac{n}{2}$ nodes with degree at least $\sigma - 1$. Therefore, at least $\frac{n}{2}$ nodes are colorless. We could hypothetically color them in $\frac{n}{2}$ distinct colors, which means that \mathcal{T} contains a string with at least $\frac{n}{2}$ distinct symbols. This leads to the following result.

Theorem 3.1. For any integers n and σ with $2 \leq \sigma < \frac{n}{2}$, there is no deterministic algorithm that performs at most $\frac{n\sigma}{8}$ equality-comparisons in the worst case, and is able to distinguish length-n strings with at most σ distinct symbols from length-n strings with at least $\frac{n}{2}$ distinct symbols.

The theorem implies lower bounds on the number of comparisons needed to compute the LZ factorisation (as defined in Section 2) and the *f*-factorisation. In the unique *f*-factorisation $T = f_1 f_2 \dots f_z$, each factor f_i is either a single symbol that does not occur in $f_1 \dots f_{i-1}$, or it is the fragment of maximal length such that f_i occurs twice in $f_1 \dots f_i$.

Corollary 1. For any integers n and σ with $2 \leq \sigma < \frac{n}{4}$, there is no deterministic algorithm that performs at most $\frac{(n-1)\sigma}{16}$ equality-comparisons in the worst case, and computes the f-factorisation of a length-n string with at most σ distinct symbols.

Proof. For some string $T = T[1]T[2] \dots T[\frac{n}{2}]$ with σ distinct symbols, consider the length-*n* string $T' = T[1]T[1]T[2]T[2] \dots T[\frac{n}{2}]T[\frac{n}{2}]$ with σ distinct symbols constructed by doubling each character of *T*. The alphabet size of *T* is exactly the number of length-one phrases in the *f*-factorisation of *T'* starting at odd positions in *T'*. Thus, by Theorem 3.1, we need $\frac{n\sigma}{16} = \frac{|T|\sigma}{8}$ comparisons to find the *f*-factorisation of *T'*. We assumed that *n* is even, and account for odd *n* by adjusting the bound to $\frac{(n-1)\sigma}{16}$.

Corollary 2. For any integers n and σ with $3 \leq \sigma < \frac{n}{6} + 1$, there is no deterministic algorithm that performs at most $\frac{(n-2)(\sigma-1)}{24}$ equality-comparisons in the worst case, and computes the Lempel-Ziv factorisation of a length-n string with at most σ distinct symbols.

Proof. For some string $T = T[1]T[2] \dots T[\frac{n}{3}]$ with $\sigma - 1$ distinct symbols, let T' be the length-n string with σ distinct symbols constructed by doubling every character of T with a separator in between, i.e., $T' = T[1]T[1]\#T[2]T[2]\# \dots \#T[\frac{n}{3}]T[\frac{n}{3}]\#$. The first occurrence of character x in T corresponds to the first occurrence of xx# in T', thus the preceding phrase (possibly of length one) ends at the first x in the first occurrence of xx#, and the subsequent phrase must be x#. Then, for the later occurrences of xx# we cannot have that x# is a phrase. Consequently, the alphabet size of T is exactly the number of length-two phrases in the Lempel-Ziv factorisation of T' starting at positions $i \equiv 2 \pmod{3}$ in T'. Thus, by Theorem 3.1, we need $\frac{n(\sigma-1)}{24} = \frac{|T|(\sigma-1)}{8}$ comparisons to find the LZ factorisation of T'. We assumed that n is divisible by 3, and account for this by adjusting the bound to $\frac{(n-2)(\sigma-1)}{24}$.



Figure 1: Example conflict graph of the adversary described in Section 3.2. The alphabet $\{0, \ldots, 15\}$ is of size $\sigma = 16$. The blocks are of length $\frac{\sigma}{4} = 4$. The gray nodes are exactly the starting positions of the blocks and contain the symbols of the ternary Thue-Morse sequence $v = 2, 1, 0, 2, 0, 1, 2, \ldots$, which is square-free. We assume that the colored nodes were colored in the following order: 2, 6, 8, 7, 15, 16, 14. At the time of coloring node 8, we had to avoid colors 0, 1, 2 (because they are reserved for the separator positions), 3 (because the adjacent node 2 already has color 3), and 4 (because node 6 is in the same block and already has color 4). The algorithm has not eliminated all squares yet. For example, nodes 10 and 11 with absent edge $(10, 11) \notin E$ are adjacent to nodes of colors $\{3, 6, 5\} \cup \{5, 3, 4\}$. Thus, any of the colors $\{0, 1, 2\} \cup \{7, \ldots, 15\}$ can be assigned to both nodes, enforcing the square T[10..11]. As visualized on the right, an edge of length ℓ eliminates at most ℓ squares.

3.2 Testing Square-Freeness

In this section, we prove that testing square-freeness requires at least $n \ln \sigma - 3.6n$ comparisons (even if σ is known). The proof combines the idea behind the original $\Omega(n \lg n)$ lower bound by Main and Lorentz [61] with the adversary described at the beginning of Section 3. This time, we ensure that \mathcal{T} always contains a square-free string with at most σ distinct symbols. At the same time, we try to ensure that \mathcal{T} also contains a string with at least one square. We will show that we can maintain this state until at least $n \ln \sigma - 3n$ comparisons have been performed.

The string (or rather family of strings) constructed by the adversary is organized in $\left\lceil \frac{4n}{\sigma} \right\rceil$ nonoverlapping blocks of length $\frac{\sigma}{4}$ (we assume $\frac{\sigma}{4} \in \mathbb{N}$ and $8 \leq \sigma \leq n$). Each block begins with a special separator symbol. More precisely, the first symbol of the k-th block is the k-th symbol of a ternary square-free word over the alphabet $\{0, 1, 2\}$ (e.g., the distance between the k^{th} and $(k + 1)^{\text{th}}$ occurrence of 0 in the Prouhet-Thue-Morse sequence, also known as the ternary Thue-Morse-Sequence, see [1, Corollary 1]). Initially, the adversary colors the nodes that correspond to the separator positions in their respective colors from $\{0, 1, 2\}$. All remaining nodes will later get colors other than $\{0, 1, 2\}$. Any fragment crossing a block boundary can be projected on the colors $\{0, 1, 2\}$, and by construction the string cannot contain a square. Thus, the separator symbols ensure that there is no square crossed by a block boundary, which implies that the string is squarefree if and only if each of its blocks is square-free.

During the algorithm execution, we use the following coloring rule. The available colors are $\{3, \ldots, \sigma - 1\}$. Whenever the degree of a node becomes $\frac{\sigma}{4}$, we assign its color. We avoid not only the at most $\frac{\sigma}{4}$ colors of already colored neighbors in the conflict graph, but also the less than $\frac{\sigma}{4}$ colors of nodes within the same block (due to $\sigma \geq 8$, there are at least $\sigma - 3 - \frac{\sigma}{2} \geq 1$ colors available). An example of the conflict graph is provided in Fig. 1. At any moment in time, we could hypothetically complete the coloring by assigning one of the colors $\{3, \ldots, \sigma - 1\}$ to each colorless node, avoiding colors of adjacent nodes and colors of nodes in the same block. Afterwards, each node holds one of the σ colors, but no two nodes within the same block have the same color. Thus, each block is square-free, and therefore \mathcal{T} always contains a square-free string with at most σ distinct symbols.

Now we consider the state of the conflict graph *after the algorithm has terminated*. We are particularly concerned with consecutive ranges of colorless nodes. The following lemma states that for each such range, the algorithm either performed many comparisons, or we can enforce a square

within the range.

Lemma 2. Let $R = \{i, \ldots, j\} \subset V$ be a consecutive range of m = j - i + 1 colorless nodes in the conflict graph. Then either $|E \cap R^2| \geq \sum_{\ell=1}^{\lfloor m/2 \rfloor} \frac{m-2\ell+1}{\ell}$, or there is a string $T \in \mathcal{T}$ with at most σ distinct symbols such that T[i..j] contains a square.

Proof. We say that an integer interval $[x, x + 2\ell - 1]$ with $i \leq x < (x + 2\ell - 1) \leq j$ has been *eliminated*, if for some y with $x \leq y < x + \ell$ there is an edge $(y, y + \ell)$ in the conflict graph. If such an edge exists, then (by the definition of \mathcal{T}) all strings $T \in \mathcal{T}$ satisfy $T[y] \neq T[y + \ell]$. Thus $T[x..x + 2\ell - 1]$ is not a square for any of them.

Now we show that if $[x, x + 2\ell - 1]$ has not been eliminated, then there exists a string $T \in \mathcal{T}$ such that $T[x..x + 2\ell - 1]$ is a square. For this purpose, consider any position y with $x \leq y < x + \ell$, i.e., a position in the first half of the potential square. Since $[x, x + 2\ell - 1]$ has not been eliminated, $(y, y + \ell)$ is not an edge in the conflict graph. It follows that we could assign the same color to y and $y + \ell$. We only have to avoid the at most $2 \cdot (\frac{\sigma}{4} - 1)$ colors of adjacent nodes of both y and $y + \ell$ in the conflict graph. Thus there are $\frac{\sigma}{2} + 2$ appropriate colors that can be assigned to both nodes. Unlike during the algorithm execution, we do not need to avoid the special separator colors or the colors in the same block; since we are trying to enforce a square, we do not have to worry about accidentally creating one. By applying this coloring scheme for all possible choices of y, we enforce that all strings $T \in \mathcal{T}$ have a square $T[x..x + 2\ell - 1]$. Note that by coloring additional nodes after the algorithm terminated, we only remove elements from \mathcal{T} . Thus, the strings with square $T[x..x + 2\ell - 1]$ were already in \mathcal{T} when the algorithm terminated. It follows that, if the algorithm actually guarantees square-freeness, then it must have eliminated all possible intervals $[x, x + 2\ell - 1]$ with $i \leq x < (x + 2\ell - 1) \leq j$.

While each interval needs at least one edge to be eliminated, a single edge eliminates multiple intervals. However, all the intervals eliminated by an edge must be of the same length. Now we give a lower bound on the number of edges needed to eliminate all intervals of length 2ℓ . Any edge $(y, y + \ell)$ eliminates ℓ intervals, namely the intervals $[x, x + 2\ell - 1]$ that satisfy $x \leq y < x + \ell$. Within R, we have to eliminate $m - 2\ell + 1$ intervals of length 2ℓ , namely the intervals $[x, x + 2\ell - 1]$ that satisfy $i \leq x \leq j - 2\ell + 1$ (see right side of Fig. 1). Thus we need at least $\frac{m-2\ell+1}{\ell}$ edges to eliminate all intervals in R. Note that the edges used for elimination have both endpoints in R, and are thus contained in $E \cap R^2$. Consequently, if $|E \cap R^2| < \sum_{\ell=1}^{\lfloor m/2 \rfloor} \frac{m-2\ell+1}{\ell}$, then not all intervals have been eliminated, and there is a string in \mathcal{T} that contains a square. \Box

Finally, we show that the algorithm either performed at least $\Omega(n \lg \sigma)$ comparisons, or there is a string $T \in \mathcal{T}$ that contains a square. Let c_1, c_2, \ldots, c_k be exactly the colored nodes. Initially (before the algorithm execution), the adversary colored $\left\lceil \frac{4n}{\sigma} \right\rceil$ nodes. Thus $k \ge \left\lceil \frac{4n}{\sigma} \right\rceil$, and there are $k - \left\lceil \frac{4n}{\sigma} \right\rceil$ nodes that have been colored after their degree reached $\frac{\sigma}{4}$. Therefore, the sum of degrees of all colored nodes is at least $\left(k - \left\lceil \frac{4n}{\sigma} \right\rceil\right) \cdot \frac{\sigma}{4} \ge \frac{\sigma k - 4n - \sigma}{4} \ge \frac{\sigma k - 5n}{4}$. Each comparison may increase the degree of two nodes by one. Thus, the colored nodes account for at least $\frac{\sigma k - 5n}{8}$ comparisons. There are k non-overlapping maximal colorless ranges of nodes, namely $\{c_i + 1, \ldots, c_{i+1} - 1\}$ for $1 \le i \le k$ with auxiliary value $c_{k+1} = n+1$. According to Lemma 2, each respective range accounts for $e_i = \sum_{\ell=1}^{\lfloor m_i/2 \rfloor} \frac{m_i - 2\ell + 1}{\ell}$ edges, where $m_i = c_{i+1} - c_i - 1$. (No edge gets counted more than once because the ranges are non-overlapping, and both endpoints of the respective edges are within the range.) Thus, in order to verify square-freeness, the algorithm must have performed at least $\sum_{i=1}^k e_i + \frac{\sigma k - 5n}{8}$ comparisons. The remainder of the proof consists of simple algebra. First, we provide a convenient lower bound for e_i (explained below):

$$\begin{split} e_i &= \sum_{\ell=1}^{\lfloor m_i/2 \rfloor} \frac{m_i - 2\ell + 1}{\ell} = \sum_{\ell=1}^{\lceil m_i/2 \rceil} \frac{m_i - 2\ell + 1}{\ell} \ge (m_i + 1) \left(\left(\sum_{\ell=1}^{\lceil m_i/2 \rceil} \frac{1}{\ell} \right) - 1 \right) \\ &> (m_i + 1) \cdot \left(\ln \frac{m_i}{2} - \frac{1}{2} \right) \\ &= (m_i + 1) \cdot \ln \frac{m_i}{2\sqrt{e}} \\ &\ge (m_i + 1) \cdot \ln \frac{m_i + 1}{2.5\sqrt{e}} \end{split}$$

We can replace $\lfloor m_i/2 \rfloor$ with $\lceil m_i/2 \rceil$ because if m_i is odd the additional summand equals zero. The first inequality uses simple arithmetic operations. The second inequality uses the classical lower bound $(\ln x + \frac{1}{2}) < H_x$ of harmonic numbers. The last inequality holds for $m_i \ge 4$. For $m_i < 4$ the result becomes negative and is thus still a correct lower bound for the number of comparisons. We obtain:

$$\underbrace{\sum_{i=1}^{k} (m_i+1) \cdot \ln \frac{m_i+1}{2.5\sqrt{e}}}_{\text{comparisons within colorless ranges}} + \underbrace{\frac{\sigma k - 5n}{8}}_{\text{comparisons for colored nodes}} \geq n \cdot \ln \frac{n}{2.5\sqrt{e}k} + \frac{\sigma k - 5n}{8}$$
$$= n \cdot \ln \frac{\sigma}{2.5\sqrt{e}x} + \frac{xn - 5n}{8}$$
$$= n \cdot \ln \sigma + n \cdot \left(\frac{x - 5}{8} - \ln 2.5\sqrt{e}x\right)$$
$$> n \cdot \ln \sigma - 3.12074n$$

The first step follows from $\sum_{i=1}^{k} (m_i + 1) = n$ and the log sum inequality (see [21, Theorem 2.7.1]). In the second step we replace k by using $x = \frac{\sigma k}{n}$. The third step uses simple arithmetic operations. The last step is reached by substituting x = 8, which minimizes the equation. Finally, we assumed that σ is divisible by 4. We account for this by adjusting the lower bound to $n \ln(\sigma - 3) - 3.12074n$, which is larger than $n \ln \sigma - 3.6n$ for $\sigma \geq 8$.

Theorem 1.1. For any integers n and σ with $8 \leq \sigma \leq n$, there is no deterministic algorithm that performs at most $n \ln \sigma - 3.6n = \mathcal{O}(n \ln \sigma)$ comparisons in the worst case, and determines whether a length-n string with at most σ distinct symbols from a general unordered alphabet is square-free.

In this section, we consider the problem of testing square-freeness of a given string. We introduce an algorithm that decides whether or not a string is square-free using only $\mathcal{O}(n \lg \sigma)$ comparisons, matching the lower bound from Section 3.2. Note that this algorithm is not yet time efficient because, apart from the performed symbol comparisons, it uses other operations that are expensive in the Word RAM model. A time efficient implementation of the algorithm will be presented in Section 4, where we first achieve $\mathcal{O}(n \lg \sigma + n \log^* n)$ time, and then improve this to $\mathcal{O}(n \lg \sigma)$ time. In Section 5, we generalize the result to compute all runs in the same time complexity.

3.3 Sparse Suffix Trees and Difference Covers

Lemma 3. The sparse suffix tree containing any b suffixes $T[i_1..n], ..., T[i_b..n]$ of T[1..n] can be constructed using $\mathcal{O}(b\sigma \log b)$ comparisons plus $\mathcal{O}(n)$ comparisons shared by all invocations of the lemma.

Proof. We maintain a union-find structure over the positions of T[1..n]. Initially, each position is in a separate component. Before issuing a query $T[x] \stackrel{?}{=} T[y]$, we check if x and y are in the same component of the union-find structure, and if so immediately return that T[x] = T[y] without performing any comparisons. Otherwise, we issue the query and if it returns that T[x] = T[y] we merge the components of x and y. Thus, the total number of issued queries with positive answer, over all invocations of the lemma, is less than n, and it remains to bound the number of issued queries with negative answer.

We insert the suffixes $T[i_i...n]$ one-by-one into an initially empty sparse suffix tree. To insert the next suffix, we descend from the root of the tree to identify the node u that corresponds to the longest common prefix between $T[i_1..n]$ and any of the already inserted suffixes. We then make u explicit unless it is explicit already, and add an edge from u to a new leaf corresponding to the whole $T[i_1..n]$. We say that the insertion procedure terminates at u. Node u can be identified with only $\mathcal{O}(\sigma \log b)$ comparisons with negative answers as follows. Let v be the current node (initially, the root of the tree), and let v_1, \ldots, v_d be its children, where $d \leq \sigma$. Here, v can be either explicit or implicit, in the latter case d = 1. We arrange the children of v so that the number of leaves in the subtree rooted at v_1 is at least as large as the number of leaves in the subtree rooted at any other child of v. Then, we compare the character on the edge leading to v_1 with the corresponding character of the current suffix. If they are equal we continue with v_1 , otherwise we compare the characters on the edges leading to v_2, \ldots, v_d with the corresponding character of the current suffix one-by-one. Then, we either continue with some v_j , $j \ge 2$, or terminate at v. To bound the number of comparisons with negative answer, observe that such comparisons only occur when we either terminate at v or continue with $v_j, j \ge 2$. Whenever we continue with $v_j, j \ge 2$, the number of leaves in the current subtree rooted at v_i decreases at least by a factor of 2 compared to subtree rooted at v (as the subtree rooted at v_1 had the largest number of leaves). Thus, during the whole descent from the root performed during an insertion this can happen only at most $1 + \log b$ times. Every time we do not continue in the subtree v_1 we might have up to $d \leq \sigma$ comparisons with negative answer, thus the total number of such comparisons is as claimed¹.

Now we describe the sample positions that we will later use to compute the approximate LZ factorisation. A set $\mathbf{S} \subseteq \mathbb{N}$ is called a *t*-cover of $\{1, \ldots, n\}$ if there is a constant-time computable function *h* such that, for any $1 \leq i, j \leq n-t+1$, we have $0 \leq h(i, j) < t$ and $i+h(i, j), j+h(i, j) \in \mathbf{S}$. A possible construction of *t*-covers is given by the lemma below, and visualized in Fig. 2.

Lemma 4. For any n and $t \leq n$, there exists a t-cover $\mathbf{D}(t)$ of $\{1, \ldots, n\}$ with size $\mathcal{O}(n/\sqrt{t})$. Furthermore, its elements can be enumerated in time proportional to their number.

Proof. We use the well-known combinatorial construction known as difference covers, see e.g. [60]. Let $r = \lfloor \sqrt{t} \rfloor$ and define $\mathbf{D}(t) = \{i \in \{1, \ldots, n\} : i \mod r = 0 \text{ or } i \mod r^2 \in \{0, \ldots, r-1\}\}$. By definition, $|\mathbf{D}(t)| \leq \lfloor n/r \rfloor + \lfloor n/r^2 \rfloor r = \mathcal{O}(n/r) = \mathcal{O}(n/\sqrt{t})$. The function h(i, j) is defined as $a + b \cdot r$, where $a = (r - i) \mod r$ and $b = (r - \lfloor (j + a)/r \rfloor) \mod r$. Note that $i + h(i, j) \leq n$ and $j + h(i, j) \leq n$. Then, $i + (a + b \cdot r) = 0 \pmod{r}$, while $\lfloor (j + (a + b \cdot r))/r \rfloor = \lfloor (j + a)/r + b \rfloor = 0 \mod r$ implies $j + h(i, j) \mod r^2 \in \{0, \ldots, r - 1\}\}$, thus $i + h(i, j), j + h(i, j) \in \mathbf{D}(t)$ as required.

3.4 Detecting Squares with a Δ -Approximate LZ Factorisation

A crucial notion in our algorithm is the following variation on the standard Lempel-Ziv factorisation:

¹In the descent, if all children are sorted according to their subtree size, the number of comparisons decreases to $\mathcal{O}(b(\sigma/\log\sigma)\log b)$, but this appears irrelevant for our final algorithm.



Figure 2: Positions in a Δ -difference cover.



Figure 3: Illustration of the definition of a LZ-phrase and a Δ -approximate phrase.

Definition 3.1 (Δ -approximate LZ factorisation). For a positive integer parameter Δ , the fragment T[s..e] is a Δ -approximate LZ phrase if it can be split into a head and a tail T[s..e] =head(T[s..e])tail(T[s..e]) such that $|head(T[s..e])| < \Delta$ and additionally

- tail(T[s..e]) is either empty or occurs at least twice in T[1..e], and
- the unique (standard) LZ phrase T[s..e'] starting at position s satisfies $e' 1 \le e$.

In a Δ -approximate LZ factorisation $T = b_1 b_2 \dots b_z$, each factor b_i is a Δ -approximate phrase T[s..e] with $s = 1 + \sum_{j=1}^{i-1} |b_j|$ and $e = \sum_{j=1}^{i} |b_j|$.

Note that a standard LZ phrase is not a Δ -approximate phrase. Also, while the LZ phrase starting at each position (and thus also the LZ factorisation) is uniquely defined, there may be multiple different Δ -approximate phrases starting at each position. This also means that a single string can have multiple different Δ -approximate factorisations. The definitions of both standard and Δ -approximate LZ phrases are illustrated in Fig. 3.

The intuition behind the above definition is that constructing the Δ -approximate LZ factorisation becomes easier for larger values of Δ . In particular, for $\Delta = n$ one phrase is enough. We formalise this in the following lemma, which is made more general for the purpose of obtaining the final result in this section.

Lemma 5. For any parameter $\Delta \in [1, m]$, a Δ -approximate LZ factorisation of any fragment T[x..y] of length m can be computed with $\mathcal{O}(m\sigma \log m/\sqrt{\Delta})$ comparisons plus $\mathcal{O}(n)$ comparisons shared by all invocations of the lemma.

Proof. By Lemma 4, there exists a Δ -cover $\mathbf{D}(\Delta)$ of $\{1, \ldots, n\}$ with size $\mathcal{O}(n/\sqrt{\Delta})$. Let $S = \mathbf{D}(\Delta) \cap \{x, x+1, \ldots, y\}$. Let $S = \{i_1, i_2, \ldots, i_b\}$. It is straightforward to verify that the construction additionally guarantees $b = \mathcal{O}(m/\sqrt{\Delta})$. We apply Lemma 3 on the suffixes $T[i_1..n], \ldots, T[i_b..n]$ to obtain their sparse suffix tree T with $\mathcal{O}(b\sigma \log b)$ comparisons plus $\mathcal{O}(n)$ comparisons shared by all invocations of the lemma. T allows us to obtain the longest common prefix of any two fragments T[i..y] and T[j..y], for $i, j \in S$, with no additional comparisons. By the properties of $\mathbf{D}(\Delta)$, for any $i, j \in \{x, x+1, \ldots, y - \Delta + 1\}$ we have $0 \leq h(i, j) < \Delta$ and $i + h(i, j), j + h(i, j) \in S$.

We compute the Δ -approximate LZ factorisation of T[x..y] phrase-by-phrase. Denoting the remaining suffix of the whole T[x..y] by T[x'..y], we need to find $x' \leq y' \leq y$ such that T[x'..y'] is a Δ -approximate phrase. This is done as follows. We iterate over every $x' \leq x'' < x' + \Delta$ such that $x'' \in S$. For every such x'', we consider every $x \leq a' < x'$ such that $a' \in S$, and compute the length ℓ of the longest common prefix of T[x''..y] and T[a'..y]. Among all such x'', a' we choose the pair that results in the largest value of $x'' - x' + \ell - 1$ and choose the next phrase to be $T[x'..(x'' + \ell - 1)]$, with the head being T[x'..(x'' - 1)] and the tail $T[x''..(x'' + \ell) - 1]$. Finally, if there is no such pair, or the value of $x'' - x' + \ell - 1$ corresponding to the found pair is less than $\Delta - 2$, we take the next phrase to be $T[x'..\min\{x' + \Delta - 1, y\}]$ (with empty tail). Selecting such a pair requires no extra comparisons, as for every $x'', a' \in S$ we can use the sparse suffix tree to compute ℓ . While it is clear that the generated Δ -approximate phrase has the required form, we need to establish that it is sufficiently long.

Let T[x'..y''] be the (unique) standard LZ phrase of T[x..y] that is prefix of T[x'..y]. If $y'' < x' + \Delta - 1$ then we only need to ensure that the generated Δ -approximate phrase is of length at least $\min\{\Delta-1, y-x'+1\}$, which is indeed the case. Therefore, it remains to consider the situation when $y'' \ge x' + \Delta - 1$. Let T[a..b] be the previous occurrence of T[x'..(y''-1)] in T[x..y] (because T[x'..y''] is a phrase this is well defined). Thus, T[a..b] = T[x'..(y''-1)] and a < x'. Because $y'' \ge x' + \Delta - 1$ and $y'' \le y$, as explained above $0 \le h(a, x') < \Delta$ and $a + h(a, x'), x' + h(a, x') \in S$. We will consider x'' = x' + h(a, x') and a' = a + h(a, x') in the above procedure. Next, T[a'..b] = T[x''..(y''-1)], so when considering this pair we will obtain $\ell \ge |T[x''..(y''-1)]|$. Thus, for the found pair we will have $x'' + \ell - 1 \ge y'' - 1$ as required in the definition of a Δ -approximate phrase.

Next, we show that even though the Δ -approximate LZ factorisation does not capture all distinct squares, as it is the case for the standard LZ factorisation, it is still helpful in detecting all sufficiently long squares. A crucial component is the following property of the Δ -approximate LZ factorisation.

Lemma 6. Let $b_1b_2...b_z$ be a Δ -approximate LZ factorisation of a string T. For every square $T[s..s+2\ell-1]$ of length $2\ell \geq 8\Delta$, there is at least one phrase b_i with $|\mathsf{tail}(b_i)| \geq \frac{\ell}{4} \geq \Delta$ such that $\mathsf{tail}(b_i)$ and the right-hand side $T[s+\ell..s+2\ell-1]$ of the square intersect.

Proof. Assume that all tails that intersect $T[s + \ell..s + 2\ell - 1]$ are of length less than $\frac{\ell}{4}$, then the respective phrases of these tails are of length at most $\frac{\ell}{4} + \Delta - 1$ (because each head is of length less than Δ). This means that $T[s + \ell..s + 2\ell - 1]$ intersects at least $\left\lceil \ell/(\frac{\ell}{4} + \Delta - 1) \right\rceil \ge \left\lceil \ell/(\frac{\ell}{2} - 1) \right\rceil = 3$ phrases. Thus there is some phrase $b_i = T[x..y]$ properly contained in $T[s + \ell..s + 2\ell - 1]$, formally $s + \ell < x \le y < s + 2\ell - 1$. However, this contradicts the definition of the Δ -approximate LZ factorisation because $T[x..s + 2\ell]$ is the prefix of a standard LZ phrase (due to $T[x..s + 2\ell - 1] = T[x - \ell..s + \ell - 1]$), and the Δ -approximate phrase $b_i = T[x..y]$ must satisfy $y \ge s + 2\ell - 1$. The contradiction implies that $T[s + \ell..s + 2\ell - 1]$ intersects a tail of length at least $\frac{\ell}{4} \ge \Delta$.

Lemma 7. Given a Δ -approximate LZ factorisation $T = b_1 b_2 \dots b_z$, we can detect a square of size $\geq 8\Delta$ in $\mathcal{O}\left(\sum_{|\mathsf{tail}(b_i)|\geq\Delta} |\mathsf{tail}(b_i)|+z\right)$ time and $\mathcal{O}\left(\sum_{|\mathsf{tail}(b_i)|\geq\Delta} |\mathsf{tail}(b_i)|\right)$ comparisons.

Proof. We consider each phrase $b_i = T[a_1..a_3]$ with $\mathsf{head}(b_i) = T[a_1..a_2 - 1]$ and $\mathsf{tail}(b_i) = T[a_2..a_3]$ separately. Let $k = |\mathsf{tail}(b_i)|$. If $k \ge \Delta$, we apply Lemma 1 to $x_1 = T[a_2 - 8k..a_2 - 1]$ and $y_1 = T[a_2..a_3 + 4k - 1]$, as well as $x_2 = T[a_2 - 8k..a_3 - 1]$ and $y_2 = T[a_3..a_3 + 4k - 1]$ trimmed to T[1..n]. This takes $\mathcal{O}(|\mathsf{tail}(b_i)|)$ time and comparisons, or $\mathcal{O}\left(\sum_{|\mathsf{tail}(b_i)|\ge\Delta}|\mathsf{tail}(b_i)|\right)$ time and comparisons for all phrases. We need additional $\mathcal{O}(z)$ time to check if $k \ge \Delta$ for each phrase.

Now we show that the described strategy detects a square of size at least 8Δ . Let $T[s..s+2\ell-1]$ be any such square. Due to Lemma 6, the right-hand side $T[s+\ell..s+2\ell-1]$ of this square intersects some tail tail $(b_i) = T[a_2..a_3]$ of length $k = |\text{tail}(b_i)| \ge \frac{\ell}{4} \ge \Delta$. Due to the intersection, we have $a_2 \le s + 2\ell - 1$ and $a_3 \ge s + \ell$. Thus, when processing b_i and applying Lemma 1, the starting position of x_1 and x_2 satisfies $a_2 - 8k \le s + 2\ell - 1 - 8\frac{\ell}{4} = s - 1$, while the end position of y_1 and y_2 satisfies $a_3 + 4k - 1 \ge s + \ell + 4\frac{\ell}{4} - 1 = s + 2\ell - 1$. Therefore, the square is entirely contained in the respective fragments corresponding to x_1y_1 and x_2y_2 . If $s < a_2 \le s + 2\ell - 1$, we find the square with our choice of x_1 and y_1 . If $s < a_3 \le s + 2\ell - 1$, we find the square with our choice of x_2 and y_2 . Otherwise, $T[s..s+2\ell-1]$ is entirely contained in tail $T[a_2..a_3]$, and we find another occurrence of the square further to the left.

3.5 Simple Algorithm for Detecting Squares

Now we have all the tools to introduce our simple method for testing square-freeness of T[1..n] using $\mathcal{O}(n(\log \sigma + \log \log n))$ comparisons, assuming that σ is known in advance. Let $\Delta = (\sigma \log n)^2$. We partition T[1..n] into blocks of length 8Δ , and denote the k^{th} block by B_k . A square of length at most 8Δ can be found by invoking Theorem 2.1 on B_1B_2 , B_2B_3 , and so on. This takes $\mathcal{O}(\Delta \log \Delta) = \mathcal{O}(\Delta(\log \sigma + \log \log n))$ comparisons for each pair of adjacent blocks, or $\mathcal{O}(n(\log \sigma + \log \log n))$ comparisons in total. It remains to test for squares of length exceeding 8Δ . This is done by first invoking Lemma 5 to compute a Δ -approximate LZ factorisation of T[1..n] with $\mathcal{O}(n\sigma \log n/\sqrt{\Delta}) = \mathcal{O}(n)$ comparisons, and then using Lemma 7, which adds another $\mathcal{O}(n)$ comparisons. The total number of comparisons is dominated by the $\mathcal{O}(n(\log \sigma + \log \log n))$ comparisons needed to apply Theorem 2.1 to the block pairs.

3.6 Improved Algorithm for Detecting Squares

We are now ready to describe the algorithm that uses only $\mathcal{O}(n \log \sigma)$ comparisons without knowing the value of σ . Intuitively, we will proceed in phases, trying to "guess" the value of σ . We first observe that Lemma 5 can be extended to obtain the following.

Lemma 8. There is an algorithm that, given any parameter $\Delta \in [1, m]$, estimate $\tilde{\sigma}$ and fragment T[x..y] of length m, uses $\mathcal{O}(m\tilde{\sigma} \log m/\sqrt{\Delta})$ comparisons plus $\mathcal{O}(n)$ comparisons shared by all invocations of the lemma, and either computes a Δ -approximate LZ factorisation of T[x..y] or determines that $\sigma > \tilde{\sigma}$.

Proof. We run the procedure described in the proof of Lemma 5 and keep track of the number of comparisons with negative answer. As soon as it exceeds $\mathcal{O}(m\tilde{\sigma}\log m/\sqrt{\Delta})$ (where the constant follows from the complexity analysis) we know that necessarily $\sigma > \tilde{\sigma}$, so we can terminate. Otherwise, the algorithm obtains a Δ -approximate LZ factorisation with $\mathcal{O}(m\tilde{\sigma}\log m/\sqrt{\Delta})$ comparisons. Comparisons with positive answer are paid for globally.

Now we describe how to find any square using $\mathcal{O}(n \lg \sigma)$ comparisons. We define the sequence $\sigma_t = 2^{2^{\lceil \log \log n \rceil - t}}$, for $t = 0, 1, \ldots, \lceil \log \log n \rceil$. We observe that $\sigma_{t-1} = (\sigma_t)^2$, and proceed in phases corresponding to the values of t. In the t^{th} phase we are guaranteed that any square of length at least $(\sigma_t)^2$ has been already detected, and we aim to detect square of length less than $(\sigma_t)^2$, and at least σ_t . We partition the whole T[1..n] into blocks of length $(\sigma_t)^2$, and denote the k^{th} block by B_k . A square of length less than $(\sigma_t)^2$ is fully contained within some two consecutive blocks $B_i B_{i+1}$, hence we consider each such pair $B_1 B_2$, $B_2 B_3$, and so on. We first apply Lemma 8 with $\Delta = \sigma_t/8$ and $\tilde{\sigma} = (\sigma_t)^{1/4} / \log(\sigma_t)$ to find an $(\sigma_t/8)$ -approximate LZ factorisation of the corresponding fragment

of T[1..n], and then use Lemma 7 to detect squares of length at least σ_t . We cannot always afford to apply Lemma 7 to all block pairs. Thus, we have to deactivate some of the blocks, which we explain when analysing the number of comparisons performed by the algorithm. If any of the calls to Lemma 8 in the current phase detects that $\sigma > \tilde{\sigma}$, we switch to applying Theorem 2.1 on every pair of blocks $B_i B_{i+1}$ of the current phase and then terminate the whole algorithm.

We now analyse the total number of comparisons, ignoring the $\mathcal{O}(n)$ comparisons shared by all invocations of Lemma 8. Throughout the t^{th} phase, we use $\mathcal{O}(n \cdot \tilde{\sigma} \log \sigma_t / \sqrt{\Delta}) = \mathcal{O}(n \cdot (\sigma_t)^{1/4} / \log(\sigma_t) \cdot \log(\sigma_t) / \sqrt{\sigma_t}) = \mathcal{O}(n/(\sigma_t)^{1/4})$ comparisons to construct the Δ -approximate factorisations (using Lemma 8) until we either process all pairs of blocks or detect that $\sigma > (\sigma_t)^{1/4} / \log(\sigma_t)$. In the latter case, we finish off the whole computation with $\mathcal{O}(n \log(\sigma_t))$ comparisons (using Theorem 2.1), and by assumption on σ this is $\mathcal{O}(n \log \sigma)$ as required. Until this happens (or until we reach phase $t = \lceil \log \log n \rceil - 3$ where $\sigma_t \leq 256$), we use $\mathcal{O}(\sum_{t=0}^{t'} n/(\sigma_t)^{1/4})$ comparisons to construct the Δ -approximate factorisations, for some $0 \leq t' \leq \lceil \log \log n \rceil$. To analyse the sum, we need the following bound (made more general for the purpose of the next section).

Lemma 9. For any $0 \le x \le y$ and $c \ge 0$ we have $\sum_{i=x}^{y} 2^{ic}/2^{2^{i}} = \mathcal{O}(2^{xc}/2^{2^{x}})$.

Proof. We observe that the sequence of exponents 2^i is strictly increasing from i = 0, hence

$$\sum_{i=x}^{y} \frac{2^{ic}}{2^{2^{i}}} \le \sum_{i=2^{x}}^{2^{y}} \frac{i^{c}}{2^{i}} \le \sum_{i=2^{x}}^{\infty} \frac{i^{c}}{2^{i}} = \sum_{i=0}^{\infty} \frac{(2^{x}+i)^{c}}{2^{(2^{x}+i)}} \le \sum_{i=0}^{\infty} \frac{2^{xc} \cdot (i+1)^{c}}{2^{(2^{x}+i)}} = \frac{2^{xc}}{2^{2^{x}}} \cdot \sum_{i=0}^{\infty} \frac{(i+1)^{c}}{2^{i}}$$

 $\sum_{i=0}^{\infty} \frac{(i+1)^c}{2^i}$ is a series of positive terms, we thus use Alembert's ratio test $\frac{(i+2)^c}{2^{i+1}} \cdot \frac{2^i}{(i+1)^c} = \frac{1}{2} \frac{(i+2)^c}{(i+1)^c}$ which tends to $\frac{1}{2}$ when *i* goes to the infinity, thus the series converges to a constant.

Corollary 3. For any $0 \le t' \le \lceil \log \log n \rceil$, it holds that $\sum_{t=0}^{t'} n \cdot \operatorname{polylog}(\sigma_t) / (\sigma_t)^{1/4} = \mathcal{O}(n)$.

Proof. We have to show that $\sum_{t=0}^{t'} n \log^c(\sigma_t) / (\sigma_t)^{1/4} = \mathcal{O}(n)$ for any constant $c \ge 0$. We achieve this by splitting the sum and applying Lemma 9.

$$\sum_{t=0}^{t'} \frac{n \log^c(\sigma_t)}{(\sigma_t)^{1/4}} \leq \sum_{t=0}^{\lfloor \log \log n \rfloor} \frac{n \cdot (2^{\lceil \log \log n \rceil - t})^c}{(2^{2^{\lceil \log \log n \rceil - t}})^{1/4}} = \sum_{t=0}^{\lfloor \log \log n \rceil} \frac{n \cdot (2^t)^c}{(2^{2^t})^{1/4}}$$
$$= n \cdot \sum_{t=0}^{\lceil \log \log n \rceil} \frac{2^{tc}}{2^{2^{t-2}}} = n \cdot \left(\frac{1}{2^{2^{-2}}} + \frac{2^c}{2^{2^{-1}}} + 4 \cdot \sum_{t=0}^{\lceil \log \log n \rceil - 2} \frac{2^{tc}}{2^{2^t}}\right) = \mathcal{O}(n)$$

Thus, all invocations of Lemma 8 cause $\mathcal{O}(\sum_{t=0}^{t'} n/(\sigma_t)^{1/4}) = \mathcal{O}(n)$ comparisons.

Deactivating Block Pairs

It remains to analyse the number of comparisons used by Lemma 7 throughout all phases. As mentioned earlier, we cannot actually afford to apply Lemma 7 to all block pairs. Thus, we introduce a mechanism that deactivates some of the pairs.

introduce a mechanism that deactivates some of the pairs. First, note that there are $\mathcal{O}(\sum_{t=0}^{t'} n/(\sigma_t)^2) \subseteq \mathcal{O}(\sum_{t=0}^{t'} n/(\sigma_t)^{1/4}) = \mathcal{O}(n)$ block pairs in all phases. For each pair, we store whether it has been deactivated or not, where being deactivated broadly means that we do not have to investigate the pair because it does not contain a leftmost distinct square. For each block pair $B_i B_{i+1}$ in the current phase t, we first check if it has been marked as deactivated. If not, we also check if it has been *implicitly* deactivated, i.e., if any of the two pairs from the previous phase that contain B_iB_{i+1} are marked as deactivated. If B_iB_{i+1} has been implicitly deactivated, then we mark it as deactivated and do not apply Lemma 8 and Lemma 7 (the implicit deactivation serves the purpose of propagating the deactivation to all later phases). Note that if some position of the string is not contained in any active block pair in some phase, then it is also not contained in any active block pair in all later phases. This is because it always holds that $\sigma_{t-1} = (\sigma_t)^2$ (with no rounding required), which guarantees that block boundaries of earlier phases do not intersect blocks of later phases.

We only apply Lemma 8 and then Lemma 7 to $B_i B_{i+1}$ if the pair has neither explicitly nor implicitly been deactivated. When applying Lemma 7, a tail $T[a..a+\ell)$ contributes $\mathcal{O}(\ell)$ comparisons if $\ell \geq \Delta = \sigma_t/8$ (and otherwise it contributes no comparisons). As the whole $T[a..a+\ell)$ occurs earlier, it cannot contain the leftmost occurrence of a square in the whole T. Thus, any block pair (of any phase) contained in $T[a..a+\ell)$ also cannot contain such an occurrence, and thus such block pairs can be deactivated.

The mechanism used for deactivation works as follows. Let $T[a..a+\ell)$ be a tail contributing $\mathcal{O}(\ell)$ comparisons with $\ell \geq \Delta = \sigma_t/8$ in phase t. We mark all block pairs of phase t+2 that are entirely contained in $T[a..a+\ell)$ as deactivated. Note that blocks in phase t+2 are of length $\sqrt{\sigma_t}$, and consider the fragment $T[a + 2\sqrt{\sigma_t}...a + \ell - 2\sqrt{\sigma_t})$. In phase t + 2, and by implicit deactivation in all later phases, this fragment overlaps (either partially or fully) only block pairs that have been deactivated. Thus, after phase t+1, we will never inspect any of the symbols in $T[a+2\sqrt{\sigma_t}..a+\ell-2\sqrt{\sigma_t}]$ again. We say that tail $T[a..a + \ell)$ deactivated the fragment of length $\ell - 4\sqrt{\sigma_t} = \Omega(\ell)$, which is positive until phase $t = \lfloor \log \log n \rfloor - 3$ because $\sigma_t > 256$. Since the number of deactivated positions is linear in the number of comparisons that the tail contributes to Lemma 7, it suffices to show that each position gets deactivated at most a constant number of times. In a single phase, any position gets deactivated at most twice. This is because the tails of each factorisation do not overlap by definition, but the tails of the two factorisations of adjacent block pairs $B_i B_{i+1}$ and $B_{i+1} B_{i+2}$ can overlap. If a position gets deactivated for the first time in phase t, then (as explained earlier) we will not consider it in any of the phases $t' \ge t+2$. Thus, it can only be that we deactivate the position again in phase t + 1, but not in any later phases. In total, each position gets deactivated at most four times. Hence Lemma 7 contributes $\mathcal{O}(n)$ comparisons in total.

We have shown:

Theorem 1.2. Testing square-freenes of a length-n string that contains σ distinct symbols from a general unordered alphabet can be done with $\mathcal{O}(n \log \sigma)$ comparisons.

4 Algorithm

In this section we show how to implement the approach described in the previous section to work in $\mathcal{O}(n \log \sigma)$ time. The main difficulty is in efficiently implementing the sparse suffix tree construction algorithm, and then computing a Δ -approximate factorisation. We first how to obtain an $\mathcal{O}(n \log \sigma + n \log^* n)$ time algorithm that still uses only $\mathcal{O}(n \log \sigma)$ comparisons, and then further improve its running time to $\mathcal{O}(n \log \sigma)$.

4.1 Constructing Suffix Tree and Δ -Approximate Factorisation

To give an efficient algorithmic construction of the sparse suffix tree from Lemma 3, we will use a restricted version of LCEs, where a query ShortLCE_x(i, j) (for any positive integer x) returns min(x, LCE(i, j))). The following result was given by Gawrychowski et al. [39]: **Lemma 10** (Lemma 14 in [39]). For a length-*n* string over a general unordered alphabet², a sequence of *q* queries ShortLCE_{4^{k_i} for $i \in \{1, ..., q\}$ can be answered online in total time $\mathcal{O}(n \log^* n + s)$ and $\mathcal{O}(n + q)$ comparisons³, where $s = \sum_{i=1}^{q} (k_i + 1)$.}

In the lemma, apart from the $\mathcal{O}(n \log^* n)$ time, each ShortLCE_{4^{k_i} query accounts for $\mathcal{O}(k_i + 1)$ time. Note that we can answer the queries online, without prior knowledge of the number and length of the queries. Also, computing an LCE in a fragment T[x..y] of length m trivially reduces to a ShortLCE_{4[log4m]} query on T. Thus, we have:</sub>}

Corollary 4. A sequence of q longest common extension queries on a fragment T[x..y] of length m over a general unordered alphabet can be answered in $\mathcal{O}(q \log m)$ time plus $\mathcal{O}(n \log^* n)$ time shared by all invocations of the lemma. The number of comparisons is $\mathcal{O}(q)$, plus $\mathcal{O}(n)$ comparisons shared by all invocations of the lemma.

While constructing the sparse suffix tree, we will maintain a heavy-light decomposition using a rebuilding scheme introduced by Gabow [37]. Let L(u) denote the number of leaves in the subtree of a node u. We use the following recursive construction of a heavy-light decomposition: Starting from a node r (initially the root of the tree), we find the deepest descendant node e such that $L(e) \geq \frac{5}{6}L(r)$ (possibly e = r). The path p from the root r(p) = r to e(p) = e is a heavy path. Any edge (u, v) on this path satisfies $L(v) \geq \frac{5}{6}L(u)$, and we call those edges heavy. As a consequence, a node u can have at most one child v such that (u, v) is heavy. For each edge (u, v) where u is on the heavy path and v is not, we recursively build a new heavy path construction starting from v.

When inserting a new suffix in our tree, we keep track of the insertion in the following way: for every root of a heavy path, we maintain I(u) the number of insertions made in the subtree of usince we built the heavy-light decomposition of this subtree. When $I(u) \geq \frac{1}{6}L(u)$ we recalculate the values of L(v) for all nodes v in the subtree of u and rebuild the heavy-light decomposition for the subtree of u.

This insures that, despite insertion, for any heavy path starting at node r and a node u on that heavy path, $L(e) \geq \frac{2}{3}L(r)$. When crossing a non-heavy edge the number of nodes in the subtree reduces by a factor at least $\frac{5}{6}$ which leads to the following property:

Observation 4.1. The path from any node to the root crosses at most $\mathcal{O}(\log m)$ heavy paths.

Additionally, rebuilding a subtree of size s takes $\mathcal{O}(s)$ time and adding a suffix $T[i_j..y]$ to the tree increases I(r) for each path p from the root r to the new leaf. Those are at most $\mathcal{O}(\log m)$ nodes, and thus maintaining the heavy path decomposition takes amortized time $\mathcal{O}(\log n)$ time per insertion.

With these building blocks now clearly defined, we are ready to describe the construction of the sparse suffix tree.

Lemma 11. The sparse suffix tree containing any b suffixes $T[i_1..y], \ldots, T[i_b..y]$ of T[x..y] with m = |T[x..y]| can be constructed using $\mathcal{O}(b\sigma \log b \log m)$ time plus $\mathcal{O}(n \log^* n)$ time shared by all invocations of the lemma.

Proof. As in the proof of Lemma 3, we consider the insertion of a suffix $T[i_j..y]$ into the sparse suffix tree with suffixes $T[i_1..y], T[i_2..y] \cdots T[i_{j-1}..y]$. At all times, we maintain the heavy path

 $^{^{2}}$ Lemma 14 in [39] does not explicitly mention that it works for general unordered alphabet. However, the proof of the lemma relies solely on equality tests.

³Lemma 14 in [39] does not explicitly mention that it requires $\mathcal{O}(n+q)$ comparisons. However, they use a unionfind approach where there can be at most $\mathcal{O}(n)$ comparisons with outcome "equal", and each LCE query performs only one comparison with outcome "not-equal", similarly to what we describe in the proof of Lemma 3.

decomposition. Additionally, we maintain for each heavy path a predecessor data structure, where given some length ℓ , we can quickly identify the deepest explicit node on the heavy path that spells a string of length at least ℓ . The data structure can, e.g., be a balanced binary search tree with insertion and search operations in $\mathcal{O}(\log b)$ time (the final sparse suffix tree and thus each heavy path contains $\mathcal{O}(b)$ nodes). When rebuilding a subtree of the heavy path decomposition, we also have to rebuild the predecessor data structure for each of its heavy paths. Thus, rebuilding a size-q subtree takes $\mathcal{O}(q \log b)$ time (each node is on exactly one heavy path and has to be inserted into one predecessor data structure), and the amortized insertion time increases from $\mathcal{O}(\log m)$ to $\mathcal{O}(\log m \cdot \log b)$. Whenever we insert a suffix, we make at most one node explicit, and thus have to perform at most one insertion into a predecessor data structure. The time for this is $\mathcal{O}(\log b)$, which is dominated be the previous term.

When inserting $T[i_j..y]$, we look for the node u corresponding to the longest common prefix between $T[i_j..y]$ and the inserted suffixes, make u explicit if necessary and add a new leaf corresponding to $T[i_j..y]$ attached to u. Let v be the current node (initialized by the root, and always an explicit node) and v_1, \dots, v_d be its (explicit) children. If there is a heavy edge (v, v_a) for $1 \le a \le d$, let p be the corresponding heavy path. For each heavy path p, we store the label of one leaf (i.e., the starting position of one suffix) that is contained in the subtree of e(p). Thus, we can use Corollary 4 to compute the longest common extension between the string spelled by e(p) and $T[i_j..y]$. Now we use the predecessor data structure on the heavy path to find the deepest (either explicit or implicit) node v' on the path that spells a prefix of $T[i_j..y]$. If v' is implicit, we make it explicit and add the leaf. If v' is explicit and $v' \ne v$, we use v' as the new current node and continue. Otherwise, we have v' = v, i.e., the suffix does not belong to the subtree rooted in v_a . In this case, we issue d LCE queries between $T[i_j..y]$ and each of the strings spelled by the nodes v_1, \ldots, v_d . This either reveals that we can continue using one of the v_a as the new current node, or that we can create a new explicit node on some (v, v_a) edge and attach the leaf to it, or that we can simply attach a new leaf to v.

Now we analyse the time spent while inserting one suffix. We spent $\mathcal{O}(b \cdot \log m \cdot \log b)$ total time for inserting $\mathcal{O}(b)$ nodes into the dynamic heavy path decomposition and the predecessor data structures. In each step of the insertion process, we either (i) move as far as possible along some heavy path or (ii) move along some non-heavy edge. For (i), we issue one LCE query and one predecessor query. For (ii) we issue $\mathcal{O}(\sigma)$ LCE queries. Due to Observation 4.1, both (i) and (ii) happen at most $\mathcal{O}(\log b)$ times per suffix. Thus, for all suffixes, we perform $\mathcal{O}(b \log b)$ predecessor queries and $\mathcal{O}(b\sigma \log b)$ LCE queries. The total time is $\mathcal{O}(b \log^2 b)$ for predecessor queries, and $\mathcal{O}(b\sigma \log b \log m)$ for LCE queries (apart from the $n \log^* n$ time shared by all invocations of Corollary 4).

Lemma 12. For any parameter $\Delta \in [1, m]$, a Δ -approximate LZ factorisation of any fragment T[x..y] of length m can be computed in $\mathcal{O}(m\sigma \log^2 m/\sqrt{\Delta})$ time plus $\mathcal{O}(n \log^* n)$ time shared by all invocations of the lemma.

Proof. Let T' = T[x..y], and let $\{i_1, i_2, \ldots, i_b\}$ be a Δ -cover of $\{1, \ldots, m\}$, which implies $b = \Theta(m/\sqrt{\Delta})$. We obtain a sparse suffix tree of the suffixes $T'[i_1..m], \ldots, T'[i_b..m]$, which takes $\mathcal{O}(b\sigma \log b \log m) \subseteq \mathcal{O}(m\sigma \log^2 m/\sqrt{\Delta})$ time according to Lemma 11, plus $\mathcal{O}(n \log^* n)$ time shared by all invocations of the lemma. Now we compute a Δ -approximate LZ factorisation of T' from the spare suffix tree in $\mathcal{O}(b)$ time.

In the following proof, we use i_1, i_2, \ldots, i_b interchangeably to denote both the difference cover positions, as well as their corresponding leaves in the sparse suffix tree. Assume that the order of difference cover positions is $i_1 < i_2 < \cdots < i_b$. First, we determine for each $i_k > i_1$, the position $\operatorname{src}(i_k) = i_h$ and the length $\operatorname{len}(i_k) = \operatorname{LCE}(i_h, i_k)$, where $i_h \in \{i_1, \ldots, i_{k-1}\}$ is a position that maximizes $\operatorname{LCE}(i_h, i_k)$. This is similar to what was done in [34] for the LZ77 factorisation. We start by assigning labels from $\{1, \ldots, b\}$ to the nodes of the sparse suffix tree. A node has label kif and only of i_k is its smallest descendant leaf. We assign the labels as follows. Initially, all nodes are unlabelled. We assign label 1 to each node on the path from i_1 to the root. Then, we process the remaining leaves i_2, \ldots, i_b in increasing order. For each i_k , we follow the path from i_k to the root. We assign label k to each unlabelled node that we encounter. As soon as we reach a node that has already been labelled, say, with label h and string-depth ℓ , we are done processing leaf i_k . It should be easy to see that i_h is also exactly the desired index that maximizes $\operatorname{LCE}(i_h, i_k)$, and we have $\operatorname{LCE}(i_h, i_k) = \ell$. Thus, we have found $\operatorname{src}(i_k) = i_h$ and $\operatorname{len}(i_k) = \ell$. The total time needed is linear in the number of sparse suffix tree nodes, which is $\mathcal{O}(b)$.

Finally, we obtain a Δ -approximate LZ factorisation using src and len. The previously computed values can be interpreted as follows: i_k could become the starting position of a length-len (i_k) tail (with previous occurrence at position $\operatorname{src}(i_k)$). For the Δ -approximate LZ factorisation, we will create the factors greedily in a left-to-right manner. Assume that we already factorised T'[1..s-1], then the next phrase starts at position s, and thus the next tail starts within $T'[s..s + \Delta)$ (as a reminder, the head is by definition shorter than Δ). Let $S = \{i_1, i_2, \ldots, i_b\} \cap \{s, \ldots, s + \Delta - 1\}$. If there is no $i_k \in S$ with $i_k + \operatorname{len}(i_k) > s + \Delta - 1$, then the next phrase is simply $T'[s..\min(|T'|, s + \Delta - 1))$ with empty tail. Otherwise, the next phrase has (possibly empty) head $T'[s..i_k)$ and tail $T'[i_k..i_k + \operatorname{len}(i_k))$ (with previous occurrence $\operatorname{src}(i_k)$), where i_k is chosen from S such that it maximizes $i_k + \operatorname{len}(i_k)$. Creating the phrase in this way clearly takes $\mathcal{O}(|S|)$ time. Since the next phrase starts at least at position $s + \Delta - 1$, none of the positions from $S \setminus \{s + \Delta - 1\}$ will ever be considered as starting positions of other tails. Thus, every i_k is considered during the creation of at most two phrases, and the total time needed to create all phrases is $\mathcal{O}(b)$.

It remains to be shown that the computed factorisation is indeed a Δ -approximate LZ factorisation, i.e., if we output a phrase T'[s..e], then the unique (non-approximate) LZ phrase T'[s..e']starting at position s satisfies $e' - 1 \leq e$. First, note that for the created approximate phrases (except possibly the last phrase of T) we have $s + \Delta - 2 \leq e$. Assume $e' < s + \Delta$, then clearly $e' - 1 \leq e$. Thus, we only have to consider $e' > s + \Delta - 1$. Since T'[s..e'] is an LZ phrase, there is some s' < s such that LCE(s', s) = e' - s. Let h be the constant-time computable function that defines the Δ -cover, and let $i_{k'} = s' + h(s', s)$ and $i_k = s + h(s', s)$. Note that $i_{k'} \in \{i_1, i_2, \ldots, i_{k-1}\}$ and $i_k \in \{i_1, i_2, \ldots, i_b\} \cap \{s, \ldots, s + \Delta - 1\}$. Therefore, we have $len(i_k) \geq LCE(i_{k'}, i_k) = LCE(s', s) - h(s', s) = (e' - s) - (i_k - s) = e' - i_k$. While computing the Δ -approximate phrase T'[s..e], we considered i_k as the starting positions of the tail, which implies $e \geq i_k + len(i_k) - 1 \geq e' - 1$.

Lemma 13. There is an algorithm that, given any parameter $\Delta \in [1, m]$, estimate $\tilde{\sigma}$ and fragment T[x..y] of length m, takes $\mathcal{O}(m\tilde{\sigma}\log^2 m/\sqrt{\Delta})$ time plus $\mathcal{O}(n\log^* n)$ time shared by all invocations of the lemma, and either computes a Δ -approximate LZ factorisation of T[x..y] or determines $\sigma > \tilde{\sigma}$.

Proof. We simply use Lemma 12 to compute the factorisation. In the first step, we have to construct the sparse suffix tree using the algorithm from Lemma 11. While this algorithm takes $\mathcal{O}(m\sigma \log^2 m/\sqrt{\Delta})$ time, it is easy to see that a more accurate time bound is $\mathcal{O}(md \log^2 m/\sqrt{\Delta})$, where d is the maximum degree of any node in the sparse suffix tree. If during construction the maximum degree of a node becomes $\tilde{\sigma} + 1$, we immediately stop and return that $\sigma > \tilde{\sigma}$. Otherwise, we finish the construction in the desired time.

Now we can describe the algorithm that detects squares in $\mathcal{O}(n \lg \sigma + n \log^* n)$ time and $\mathcal{O}(n \lg \sigma)$ comparisons. We simply use the algorithm from Section 3.2, but use Lemma 13 in-

stead of Lemma 8. Next, we analyse the time needed apart from the $\mathcal{O}(n \log^* n)$ time shared by all invocations of Lemma 13. Throughout the t^{th} phase, we use $\mathcal{O}(n \cdot \tilde{\sigma} \cdot \log^2(\sigma_t)/\sqrt{\Delta}) =$ $\mathcal{O}(n \cdot (\sigma_t)^{1/4}/\log(\sigma_t) \cdot \log^2(\sigma_t)/\sqrt{\sigma_t}) = \mathcal{O}(n\log(\sigma_t)/(\sigma_t)^{1/4})$ comparisons to construct all the Δ approximate factorisations. As before, if at any time we discover that $\tilde{\sigma} > (\sigma_t)^{1/4}/\log(\sigma_t)$, then we use Theorem 2.1 to finish the computation in $\mathcal{O}(n\lg\sigma_t) = \mathcal{O}(n\log\sigma)$ time. Until then (or until we finished all $\lceil \log \log n \rceil$ phases), we use $\mathcal{O}(\sum_{t=0}^{t'} n\log(\sigma_t)/(\sigma_t)^{1/4})$ time, and by Corollary 3 this is $\mathcal{O}(n)$. For detecting squares, we still use Lemma 7, which as explained in Section 3.2 takes $\mathcal{O}(n)$ time and comparisons in total, plus additional $\mathcal{O}(Z)$ time, where Z is the number of approximate LZ factors considered during all invocations of the lemma. We apply the lemma to each approximate LZ factorisation exactly once, and by construction each factor in phase t has size at least $\Delta = \Omega(\sigma_t)$. Also, each text position is covered by at most two tails per phase. Hence $Z = \mathcal{O}(\sum_{t=0}^{t'} n/\sigma_t)$, which is $\mathcal{O}(n)$ by Corollary 3.

The last thing that remains to be shown is how to implement the bookkeeping of blocks, i.e., in each phase we have to efficiently deactivate block pairs as described at the end of Section 3.2. We maintain the block pairs in $\lceil \log \log n \rceil$ bitvectors of total length $\mathcal{O}(n)$, where a set bit means that a block pair has been deactivated (recall that there are $\mathcal{O}(n)$ pairs in total). Bitvector t contains at position j the bit corresponding to block pair $B_j B_{j+1} = T[i..i + 2(\sigma_t)^2)$ with $i = (\sigma_t)^2 \cdot (j-1)$. Note that translating between i and j takes constant time. For each sufficiently long tail in phase t, we simply iterate over the relevant block pairs in phase t + 2 and deactivate them, i.e., we set the corresponding bit. This takes time linear in the number of deactivated blocks. Since there are $\mathcal{O}(n)$ block pairs, and each block pair gets deactivated at most a constant number of times, the total cost for this bookkeeping is $\mathcal{O}(n)$.

The number of comparisons is dominated by the $\mathcal{O}(n \log \sigma)$ comparisons used when finishing the computation with Theorem 2.1. The only other comparisons are performed by Lemma 7, which we already bounded by $\mathcal{O}(n)$, and by LCE queries via Corollary 4. Since we ask $\mathcal{O}(n)$ such queries in total, the number of comparisons is also $\mathcal{O}(n)$. We have shown:

Lemma 14. The square detection algorithm from Section 3.2 can be implemented in $\mathcal{O}(n \lg \sigma + n \log^* n)$ time and $\mathcal{O}(n \log \sigma)$ comparisons.

4.2 Final Improvement

For our final improvement we need to replace the LCE queries implemented by Corollary 4 with our own mechanism. The goal will remain the same, that is, given a parameter Δ and estimate $\tilde{\sigma}$ of the alphabet size, find a Δ -approximate LZ factorisation of any fragment T[x..y] in $\mathcal{O}(m\tilde{\sigma} \log m/\sqrt{\Delta})$ time, where m = |T[x..y]| (with $m = \Theta(\Delta^2)$), as otherwise we are not required to detect anything).

As in the previous section, the algorithm might detect that the size of the alphabet is larger than $\tilde{\sigma}$, and in such case we revert to the divide-and-conquer algorithm. Let $\tau = \lfloor \sqrt{\Delta} \rfloor$.

Initially, we only consider some fragments of T[x..y]. We say that $T[i \cdot \tau^2 ..i \cdot \tau^2 + \tau)$ is a dense fragment. We start by remapping the characters in all dense fragments that intersect T[x..y] to a linearly-sortable alphabet. This can be done in $\mathcal{O}(\tilde{\sigma})$ time for each position by maintaining a list of the already seen distinct characters. For each position in a dense fragment, we iterate over the characters in the list, and possibly append a new character to the list if it is not present. As soon as the size of the list exceeds $\tilde{\sigma}$, we terminate the procedure and revert to the divide-and-conquer algorithm. Otherwise, we replace each character by its position in the list. Overall, there are $\mathcal{O}(m/\sqrt{\Delta})$ positions in the dense fragments of T[x..y], and the remapping takes $\mathcal{O}(m\tilde{\sigma}/\sqrt{\Delta})$ time.

Next, we construct two generalised suffix trees [43], the first one of all dense fragments, and the second one of their reversals. (The generalised suffix tree of a collection of strings is the compacted



(a) Sampling dense fragments and cutting the text into chunks. Dotted lines indicate chunk boundaries, and $h_x = (j+x)\cdot\tau$ for some integer j and $x \in [0,13]$ are positions of chunk boundaries. The dense fragments are $D_1 = T[h_2..h_3)$, $D_2 = T[h_7..h_8)$, and $D_3 = T[h_{12}..h_{13})$. The primary occurrences of dense fragments are grey, while the secondary occurrences (the ones that we aim to find) are white. A purple box in the text, and the matching purple line underneath the text, correspond to some substring $T[j\cdot\tau-r_{j-1}..j\cdot\tau)$. Similarly, the orange boxes and lines correspond to substrings $T[j\cdot\tau..j\cdot\tau+\ell_j)$.



(b) The string T' used to find all the occurrences of dense fragments. Each position \hat{h}_x maps to position h_x in Fig. 4a. The substring indicated by the purple box preceding $h_x = (j+x)\tau$ and the orange box succeding h_x is exactly $T[h_x - r_{j+x-1}..h_x + \ell_{j+x})$. Each \mathfrak{s} is a distinct separator symbol that is unique within T'.

Figure 4: Supplementary drawings for Section 4.2.

trie that contains all suffixes of all strings in the collection.) Again, because we now work with a linearly-sortable alphabet this takes only $\mathcal{O}(m/\sqrt{\Delta})$ time [31]. We consider fragments of the form $T[i \cdot \tau .. (i+1) \cdot \tau)$ having non-empty intersection with T[x..y]. We call such fragments chunks. We note that there are $\mathcal{O}(m/\sqrt{\Delta})$ chunks, and their total length is $\mathcal{O}(m)$. For each chunk, we find its longest prefix $T[i \cdot \tau .. i \cdot \tau + \ell_i)$ and longest suffix $T[(i+1) \cdot \tau - r_i .. (i+1) \cdot \tau)$ that occur in one of the dense fragments. Fig. 4a visualizes the dense fragments, chunks, and longest prefixes and suffixes. This can be done efficiently by following the heavy path decomposition of the generalised suffix tree of all dense fragments and their reversals, respectively. On each current heavy path, we just naively match the characters as long as possible. In case of a mismatch, we spend $\mathcal{O}(\tilde{\sigma})$ time to descend to the appropriate subtree, which happens at most $\mathcal{O}(\log m)$ times due to the heavy path decomposition. After having found ℓ_i and r_i , we test square-freenes of $T[i \cdot \tau .. i \cdot \tau + \ell_i)$ and $T[(i+1)\cdot\tau - r_{i}..(i+1)\cdot\tau)$. Because they both occur in dense fragments, and we have remapped the alphabet of all dense fragments, we can use Theorem 2.2 to implement this in $\mathcal{O}(\ell_i + r_i)$ time. Thus, the total time per chunk is thus $\mathcal{O}(\tilde{\sigma} \log m)$ plus $\mathcal{O}(\ell_i + r_i)$. The former sums up to $\mathcal{O}(m\tilde{\sigma}\log m/\sqrt{\Delta})$, and we will later show that the latter can be amortised by deactivating blocks on the lower levels.

The situation so far is that we have remapped the alphabet of all dense fragments to linearly-

sortable, and for every chunk we know its longest prefix and suffix that occur in one of the dense fragments. We concatenate all fragments of the form $T[i \cdot \tau - r_{i-1} ... i \cdot \tau + \ell_i)$ (intersected with T[x..y]) while adding distinct separators in between to form a new string T'. We stress that, because we have remapped the alphabet of all dense fragments, and the found longest prefix and suffix of each chunk also occur in some dense fragment, T' is over linearly-sortable alphabet. Thus, we can build the suffix tree ST of T' in $\mathcal{O}(|T'|)$ time [31]. A visualization of T' is provided in Fig. 4b

Let $\mathcal{D} = \{D_1, D_2, \ldots\}$ be the set of distinct dense fragments. We would like to construct the set of all occurrences of the strings in \mathcal{D} in T[x..y]. Using the suffix tree of T' we can retrieve all occurrences of every D_j in T'. We observe that, because of how we have defined $T[i \cdot \tau ..i \cdot \tau + \ell_i)$ and $T[(i+1) \cdot \tau - r_i..(i+1) \cdot \tau)$, this will in fact give us all occurrences of every D_j in the original T[x..y]. To implement this efficiently, we proceed as follows. First, for every i we traverse STstarting from its root to find the (explicit or implicit) node corresponding to the dense fragment $T[i \cdot \tau^2 ..i \cdot \tau^2 + \tau)$. This takes only $\mathcal{O}(m\tilde{\sigma}/\sqrt{\Delta})$ time. Then, all leaves in every subtree rooted at such a node correspond to occurrences of some D_j , and can be reported by traversing the subtree in time proportional to its size, so at most $\mathcal{O}(|T'|)$ in total. Finally, remapping the occurrences back to T[x..y] can be done in constant time per occurrence by precomputing, for every position in T', its corresponding position in T[x..y], which can be done in $\mathcal{O}(|T'|)$ time when constructing T'. Thus, in $\mathcal{O}(|T'|)$ time, we obtain the set S of starting positions of all occurrences of the strings in \mathcal{D} . We summarize the properties of S below.

Proposition 1. S admits the following properties:

- 1. For every $i \in [x, y]$ such that $i = 0 \pmod{\tau^2}$, $i \in S$.
- 2. For every $i \in [x, y \tau]$, $i \in S$ if and only if $T[i..i + \tau) \in \mathcal{D}$.
- 3. $|S| \leq |T'|$.

We now define a parsing of T[x..y] based on S. Let $i_1 < i_2 < \ldots i_k$ be all the positions in S, that is, $(i_j, i_{j+1}) \cap S = \emptyset$ for every $j = 1, 2, \ldots, k - 1$. For every $j = 1, 2, \ldots, k - 1$, we create the phrase $T[i_j..i_{j+1} + \tau)$. We add the last phrase $T[i_k...y]$. We stress that consecutive phrases overlap by τ characters, and each phrase begins with a length- τ fragment starting at a position in S. This, together with property 2 of S, implies the following property.

Observation 4.2. The set of distinct phrases is prefix-free.

We would like to construct the compacted trie \mathcal{T}_{phrase} of all such phrases, so that (in particular) we identify identical phrases. We first notice that each phrase begins with a fragment $T[i_j..i_j + \tau)$ that has its corresponding occurrence in T'. We note that, given a set of positions P in T, we can find their corresponding positions in T' (if they exist) by sorting and scanning in $\mathcal{O}(|P|+|T'|)$ time.

Thus, we can assume that for each i_j we know its corresponding position i'_j in T'. Next, for each node of ST we precompute its unique ancestor at string depth τ in $\mathcal{O}(|T'|)$ time. Then, for every fragment $T[i_j..i_j + \tau)$ we can access its corresponding (implicit or explicit) node of ST. This allows us to partition all phrases according to their prefixes of length τ . In fact, this gives us the top part of $\mathcal{T}_{\text{phrase}}$ containing all such prefixes in $\mathcal{O}(m/\sqrt{\Delta})$ time, and for each phrase we can assume that we know the node of $\mathcal{T}_{\text{phrase}}$ corresponding to its length- τ prefix.

To build the remaining part of \mathcal{T}_{phrase} , we partition the phrases into short and long. $T[i_j...i_{j+1}+\tau)$ is short when $i_{j+1} \leq i_j + \tau$ (meaning that its length is at most 2τ), and long otherwise.

We begin with constructing the compacted trie $\mathcal{T}'_{\text{phrase}}$ of all short phrases. This can be done similarly to constructing the top part of $\mathcal{T}_{\text{phrase}}$, except that now the fragments have possibly different lengths. However, every short phrase $T[i_j..i_{j+1} + \tau)$ occurs in T' as $T'[i'_j..i'_{j+1} + \tau)$. We claim that the nodes of ST corresponding to every $T'[i'_j..i'_{j+1} + \tau)$ can be found in $\mathcal{O}(|T'|)$ time. This can be done by traversing ST in the depth-first order while maintaining a stack of all explicit nodes with string depth at least τ on the current path. Then, when visiting the leaf corresponding to the suffix of T' starting at position i'_j , we iterate over the current stack to find the sought node. This takes at most $\mathcal{O}(|T'[i_j + \tau..i_{j+1} + \tau]|)$ time, which sums up to $\mathcal{O}(|T'|)$. Having found the node of ST corresponding to $T[i_j..i_{j+1} + \tau)$, we extract $\mathcal{T}'_{\text{phrase}}$ from ST in $\mathcal{O}(|T'|)$ time.

With $\mathcal{T}'_{\text{phrase}}$ in hand, we construct the whole $\mathcal{T}_{\text{phrase}}$ as follows. We begin with taking the union of $\mathcal{T}'_{\text{phrase}}$ and the already obtained top part of $\mathcal{T}_{\text{phrase}}$, this can be obtained in $\mathcal{O}(|T'|)$ time. For each long phrase $T[i_j..i_{j+1} + \tau)$, we know the node corresponding to $T[i_j..i_j + \tau)$ and would like to insert the whole string $T[i_j..i_{j+1} + \tau)$ into \mathcal{T}_{phrase} . We perform the insertions in increasing order of i_i (this will be crucial for amortising the time later). This is implemented with a dynamic heavy path decomposition similarly as in Section 4.1, however with one important change. Namely, we fix a heavy path decomposition of the part of \mathcal{T}_{phrase} corresponding to the union of \mathcal{T}'_{phrase} and the top part of \mathcal{T}_{phrase} , and maintain a dynamic heavy path decomposition of every subtree hanging off from this part. Thanks to this change, the time to maintain the dynamic trie and all heavy path decompositions is $\mathcal{O}(m \log m/\sqrt{\Delta})$, as there are only $\mathcal{O}(m/\sqrt{\Delta})$ long phrases. Next, for each long phrase $T[i_{j}..i_{j+1} + \tau)$, we begin the insertion at the already known node corresponding to $T[i_j..i_j + \tau)$, and continue the insertion by following the heavy paths, first in the static heavy path decomposition in the part of $\mathcal{T}_{\text{phrase}}$ corresponding to $\mathcal{T}'_{\text{phrase}}$, second in the dynamic heavy path decomposition in the appropriate subtree. On each heavy path, we naively match the characters as long as possible. The time to insert a single phrase $T[i_j..i_{j+1} + \tau)$ is $\mathcal{O}(\log m)$ (twice) plus the length of the longest prefix of $T[i_j + \tau .. i_{j+1} + \tau)$ equal to a prefix of $T[i_{j'} + \tau .. i_{j'+1} + \tau)$, for some j' < j. The former sums up to another $\mathcal{O}(m \log m/\sqrt{\Delta})$, and we will later show that the latter can be amortised by deactivating blocks on the lower levels.

 $\mathcal{T}_{\text{phrase}}$ allows us to form metacharacters corresponding to the phrases, and transform T[x,y]into a string T_{parse} of length $\mathcal{O}(|T'|)$ consisting of these metacharacters. We build a suffix tree $\mathcal{S}_{\text{parse}}$ over this string over linearly-sortable metacharacters in $\mathcal{O}(|T'|)$ time. Next, we convert it into the sparse suffix tree $\mathcal{S}'_{\text{parse}}$ of all suffixes $T[i_j..y]$ as follows. Consider an explicit node $u \in \mathcal{S}_{\text{parse}}$ with children $v_1, v_2, \ldots, v_d, d \geq 2$. We first compute the subtree \mathcal{T}_u of $\mathcal{T}_{\text{phrase}}$ induced by the leaves corresponding to the first metacharacters on the edges (u, v_i) , for $i = 1, 2, \ldots, d$, and connect every v_i to the appropriate leaf of \mathcal{T}_u . This can be implemented in $\mathcal{O}(d)$ time, assuming constant-time lowest common ancestor queries on $\mathcal{T}_{\text{phrase}}$ [10] and processing the leaves from left to right with a stack, similarly as in the Cartesian tree construction algorithm [75]. We note that the order on the leaves is the same as the order on the metacharacters, and hence no extra sorting is necessary. Overall, this sums up to $\mathcal{O}(|T'|)$. Next, we observe that, unless u is the root of $\mathcal{S}_{\text{parse}}$, all metacharacters on the edges (u, v_i) correspond to strings starting with the same prefix of length τ . We obtain the subtree \mathcal{T}'_u by truncating this prefix (or taking \mathcal{T}_u if u is the root). Finally, we identify the root of \mathcal{T}'_u with u, and every child v_i with its corresponding leaf of \mathcal{T}'_u . Because we truncate the overlapping prefixes of length τ , after this procedure is executed on every node of $\mathcal{S}_{\text{parse}}$ we obtain a tree $\mathcal{S}'_{\text{parse}}$ with the property that each leaf corresponds to a suffix $T[i_j..y]$. Also, by Observation 4.2, the edges outgoing from every node start with different characters as required.

By following an argument from the proof of Lemma 12, S'_{parse} allows us to determine, for every suffix $T[i_j..y]$, its longest prefix equal to a prefix of some T[i'..y] with $i' < i_j$, as long as its length is at least τ . Indeed, in such case we must have $i' \in S$ by property 2, so in fact $i' = i_{j'}$ and it is enough to maximise the length of the common prefix with all earlier positions in S, which can be done using S'_{parse} . Thus, we either know that the length of this longest prefix is less than τ , or know its exact value (and the corresponding position $i' \in S$). **Lemma 15.** For any parameter $\Delta \in [1, m]$ and estimate $\tilde{\sigma}$ of the alphabet size, a $(\Delta + \tau)$ approximate LZ factorisation of any fragment T[x..y] can be computed in $\mathcal{O}(m/\sqrt{\Delta})$ time with m = |T[x..y]| (assuming the preprocessing described earlier in this section).

Proof. Let $e \in [x, y]$ and suppose we have already constructed the factorisation of T[x..e-1] and are now trying to construct the next phrase. Let e' be the next multiple of τ^2 , we have that $e' - e < \tau^2 \leq \Delta$ and $T[e'..e' + \tau)$ is a dense fragment. Thus, by property 1 we have $e' \in S$.

The first possibility is that the longest common prefix between T[e'..y] and any suffix starting at an earlier position is shorter than τ . In this case, we can simply set the head of the new phrase to be $T[e..e' + \tau)$ and the tail to be empty. Otherwise, we know the length ℓ of this longest prefix by the preprocessing described above. We set the head of the new phrase to be T[e..e') and the tail to be $T[e'..e'' + \ell)$. This takes constant time per phrase, and each phrase is of length at least τ , giving the claimed overall time complexity. It remains to argue correctness of every step.

Let T[e..s] be the longest LZ phrase starting at position e, to show that we obtain a valid $(\Delta + \tau)$ -approximate phrase it suffices to show that $s \leq e' + \max(\tau, \ell)$. Let the previous occurrence of T[e..s) be at position p < e. If $s - e' < \tau$ then there is nothing to prove. Otherwise, T[e'..s) is a string of length at least τ that also occurs starting earlier at position p + e' - e < e'. Thus, we will correctly determine that $\ell \geq \tau$, and find a previous occurrence of the string maximising the value of ℓ . In particular, we will have $\ell \geq s - e'$ as required.

To achieve the bound of Theorem 1.3, we now proceed as in Section 3.6, except that instead of Lemma 13 we use Lemma 15. For every T[x..y] with m = |T[x..y]| this takes $\mathcal{O}(m\tilde{\sigma} \log m/\sqrt{\Delta})$ time plus the time used for computing the longest prefix and suffix of each chunk (the latter also accounts for constructing the suffix tree ST and other steps that have been estimated as taking $\mathcal{O}(|T'|)$ in the above reasoning) plus the time for inserting $T[i_j + \tau ...i_{j+1} + \tau)$ into $\mathcal{T}_{\text{phrase}}$ when $i_{j+1} \geq i_j + \tau$.

We observe that we can deactivate any block pair fully contained in $T[i \cdot \tau .. i \cdot \tau + \ell_i)$ and $T[(i+1) \cdot \tau - r_i .. (i+1) \cdot \tau)$, as we have already checked that these fragments are square-free. Also, we can deactivate any block pair fully contained in the longest prefix of $T[i_j + \tau .. i_{j+1} + \tau)$ equal to $T[i_{j'} + \tau .. i_{j'+1} + \tau)$, for some j' < j, because such fragment cannot contain the leftmost occurrence of a square.

There are $\mathcal{O}(m/\sqrt{\Delta})$ chunks and long phrases. If a chunk or a long phrase contributes $x = \Omega(\sqrt[4]{\Delta})$ to the total time, then we explicitly deactivate the block pairs in phase t + 3 that are entirely contained in the corresponding fragment. Block pairs in phase t + 3 are of length $\mathcal{O}(\sqrt[4]{\Delta})$, and thus we deactivate $\Omega(x)$ positions. Therefore, the time spent on such chunks and long phrases in all phases sums to $\mathcal{O}(n)$. The remaining chunks and long phrases contribute $\mathcal{O}(\sqrt[4]{\Delta})$ to the total time, and there are $\mathcal{O}(m/\sqrt{\Delta})$ of them, which adds up to $\mathcal{O}(m/\sqrt[4]{\Delta})$. In every phase, this is $\mathcal{O}(n/\sqrt[4]{\Delta})$, so $\mathcal{O}(n)$ overall by Corollary 3.

5 Computing Runs

Now we adapt the algorithm such that it computes all runs. We start with the algorithm from Sections 3.2 and 4 without the final improvement from Section 4.2. First, note that the key properties of the Δ -approximate LZ factorisation, in particular Lemmas 6 and 7, also hold for the computation of runs. This is expressed by the lemmas below.

Lemma 16. Let $b_1b_2...b_z$ be a Δ -approximate LZ factorisation of a string T. For every run $\langle s, e, p \rangle$ of length $e - s + 1 \ge 8\Delta$, there is at least one phrase b_i with $|\mathsf{tail}(b_i)| \ge \frac{e-s+1}{8} \ge \Delta$ such that $\mathsf{tail}(b_i)$ and the right-hand side $T[s + \lceil \frac{e-s+1}{2} \rceil ...e]$ of the run intersect.

Proof. Let $\ell = \frac{e-s+1}{2}$ and note that $\frac{\ell}{4} \ge \Delta$ and $e = s + 2\ell - 1$. Assume that all tails that intersect $T[s + \lceil \ell \rceil ..e]$ are of length less than $\frac{\ell}{4}$, then the respective phrases of these tails are of length at most $\frac{\ell}{4} + \Delta - 1 \le \frac{\ell}{2} - 1$ (because each head is of length less than Δ). This means that $T[s + \lceil \ell \rceil ..e]$ (of length $\lfloor \ell \rfloor$) intersects at least $\lceil \lfloor \ell \rfloor / (\frac{\ell}{2} - 1) \rceil \ge 3$ phrases (the inequality holds for $\ell \ge 4$, which is implied by $\Delta \ge 1$). Thus there is some phrase $b_i = T[x..y]$ properly contained in $T[s + \lceil \ell \rceil ..e]$, formally $s + \lceil \ell \rceil < x \le y < e$. However, this contradicts the definition of the Δ -approximate LZ factorisation because T[x..e + 1] is the prefix of a standard LZ phrase (due to T[x..e] = T[x - p..e - p]). The contradiction implies that $T[s + \lceil \ell \rceil ..e]$ intersects a tail of length at least $\frac{\ell}{4}$.

Before we show how to algorithmically apply Lemma 16, we need to explain how Lemma 1 extends to computing runs, and then how this implies that the approach of Main and Lorentz [61] easily extends to computing all runs. We do not claim this to be a new result, but the original paper only talks about finding a representation of all squares, and we need to find runs, and hence include a description for completeness.

Lemma 17. Given two strings x and y over a general alphabet, we can compute all runs in xy that include either the last character of x or the first character of y using O(|x| + |y|) time and comparisons.

Proof. Consider a run $\langle s, e, p \rangle$ in t = xy that includes either the last character of x or the first character of y, meaning that $s \leq |x| + 1$ and $e \geq |x|$. Let $\ell = \lfloor \frac{e-s+1}{2} \rfloor \geq p$. We separately compute all runs with $s + \ell \leq |x| + 1$ and $s + \ell > |x| + 1$. Below we describe the former, and the latter is symmetric.

Due to $s + \ell \leq |x| + 1$, the length-*p* substring x[|x| - p + 1.. |x|] is fully within the run. This suggests the following strategy to generate all runs with $s + \ell \leq |x| + 1$. We iterate over the possible values of p = 1, 2, ..., |x|. For a given *p*, we calculate the length of the longest common prefix of x[|x| - p + 1.. |x|]y and *y*, denoted pref, and the length of the longest common suffix of x[1.. |x| - p] and *x*, denoted suf. It is easy to see that $t[|x| - p + 1 - \mathfrak{suf}. |x| + \mathfrak{pref}]$ is a lengthwise maximal *p*-periodic substring, and its length is $\ell' = p + \mathfrak{suf} + \mathfrak{pref}$. If $\mathfrak{pref} + \mathfrak{suf} \geq p$ and $s + \lfloor \ell'/2 \rfloor \leq |x| + 1$, then we report the substring as a run. (The latter condition ensures that each run gets reported by exactly one of the two symmetric cases.)

We use a prefix table to compute the longest common prefixes. For a given string, this table contains at position *i* the length of the longest substring starting at position *i* that is also a prefix of the string. For computing the values **pref**, we use the prefix table of y\$xy (where \$ is a new character that does not match any character in *x* nor *y*). Similarly, for computing the values **suf**, we use the prefix table of the reversal of a new string x\$x. The tables can be computed in $\mathcal{O}(|x| + |y|)$ time and comparisons (see, e.g., computation of table *lppattern* in [61]). Then, each value of *p* can be checked in constant time.

Lemma 18. Computing all runs in a length-n string over a general unordered alphabet can be implemented in $\mathcal{O}(n \log n)$ time and comparisons.

Proof. Let the input string be T[1..n]. We apply divide-and-conquer. Let $x = T[1..\lfloor n/2 \rfloor]$ and $y = T[\lfloor n/2 \rfloor + 1..n]$. First, we recursively compute all runs in x and y. Of the reported runs, we filter out all the ones that contain either the last character of x or the first character of y, which takes $\mathcal{O}(|x| + |y|)$ time. In this way, if some reported run is a run with respect to x (or y), but not with respect to xy, then it will be filtered out. We have generated all runs except for the ones that contain the last character of x or the first character of y (or both). Thus we simply invoke

Lemma 17 on xy, which will output exactly the missing runs in $\mathcal{O}(|x|+|y|)$ time and comparisons. There are $\mathcal{O}(\log n)$ levels of recursion, and each level takes $\mathcal{O}(n)$ time and comparisons in total. \Box

Lemma 19. Let $T = b_1 b_2 \dots b_z$ be a Δ -approximate LZ factorisation, and $\chi = \sum_{|\mathsf{tail}(b_i)| \geq \Delta} |\mathsf{tail}(b_i)|$. We can compute in $\mathcal{O}(\chi + z)$ time and $\mathcal{O}(\chi)$ comparisons a multiset R of size $\mathcal{O}(\chi)$ of runs with the property that a run T[s..e] is possibly not in R only if $e - s + 1 < 8\Delta$ or there is some tail $\mathsf{tail}(b_i) = T[a_2..a_3]$ with $a_2 < s$ and $e < a_3$.

Proof. Let n = |T|. We consider each phrase $b_i = T[a_1..a_3]$ with $\mathsf{head}(b_i) = T[a_1..a_2 - 1]$ and $\mathsf{tail}(b_i) = T[a_2..a_3]$ separately. Let $k = |\mathsf{tail}(b_i)|$. If $k \ge \Delta$, we apply Lemma 17 to $x_1 = T[a_2 - 8k..a_2 - 1]$ and $y_1 = T[a_2..a_3 + 4k]$, as well as $x_2 = T[a_2 - 8k..a_3 - 1]$ and $y_2 = T[a_3..a_3 + 4k]$ trimmed to T[1..n]. This takes $\mathcal{O}(|\mathsf{tail}(b_i)|)$ time and comparisons and reports $\mathcal{O}(|\mathsf{tail}(b_i)|)$ runs with respect to $x_1y_1 = x_2y_2 = T[a_2 - 8k..a_3 + 4k]$ (trimmed to T[1..n]). Of these runs, we filter out the ones that contain any of the positions $a_2 - 8k$ (only if $a_2 - 8k > 1$) and $a_3 + 4k$ (only if $a_3 + 4k < n$), which takes $\mathcal{O}(|\mathsf{tail}(b_i)|)$ time. This way, each reported run is not only a run with respect to x_1y_1 , but also a run with respect to T. In total, we report $\mathcal{O}(\chi)$ runs (including possible duplicates) and spend $\mathcal{O}(\chi)$ time and comparisons when applying Lemma 17. Additional $\mathcal{O}(z)$ time is needed to check if $|\mathsf{tail}(b_i)| \ge \Delta$ for each phrase.

Now we show that the described strategy computes all runs of length at least 8Δ , except for the ones that are properly contained in a tail. Let $\langle s, e, p \rangle$ be a run of length 2ℓ , where $\ell \ge 4\Delta$ is a multiple of $\frac{1}{2}$. Due to Lemma 16, the right-hand side $T[s + \lceil \ell \rceil ...e]$ of this run intersects some tail $tail(b_i) = T[a_2..a_3]$ of length $k = |tail(b_i)| \ge \frac{\ell}{4} \ge \Delta$. Due to the intersection, we have $a_2 \le e$ and $a_3 \ge s + \lceil \ell \rceil$. Thus, when processing b_i and applying Lemma 17, the starting position of x_1 and x_2 satisfies $a_2 - 8k \le e - 8\frac{\ell}{4} < s$, while the end position of y_1 and y_2 satisfies $a_3 + 4k \ge s + \lceil \ell \rceil + 4\frac{\ell}{4} > e$. Therefore, the run is contained in the fragment $T[a_2 - 8k..a_3 + 4k]$ (trimmed to T[1..n]) corresponding to x_1y_1 and x_2y_2 , and the run does not contain positions $a_2 - 8k$ and $a_3 + 4k$. If $s \le a_2 \le e$, we find the run when applying Lemma 17 to x_1 and y_1 . If $s \le a_3 \le e$, we find the run when applying Lemma 17 to x_2 and y_2 . Otherwise, T[s..e] is entirely contained in $T[a_2 + 1..a_3 - 1]$ and we do not have to report the run.

Now we describe how to compute all runs using $\mathcal{O}(n \log \sigma)$ comparisons and $\mathcal{O}(n \log \sigma + n \log^* n)$ time. We again use the sequence $\sigma_t = 2^{2^{\lceil \log \log n \rceil - t}}$, for $t = 0, 1, \dots, \lceil \log \log n \rceil$. We observe that $\sigma_{t-1} = (\sigma_t)^2$, and proceed in phases corresponding to the values of t. In the tth phase we aim to compute runs of length at least σ_t and less than $(\sigma_t)^2$. We stress that this condition depends on the length of the run and not on its period. We partition the whole T[1..n] into blocks of length $(\sigma_t)^2$, and denote the kth block by B_k . A run of length less than $(\sigma_t)^2$ is fully contained within some two consecutive blocks $B_i B_{i+1}$, and there is always a pair of consecutive blocks such that the run contains neither the first nor the last position of the pair (unless the first position is T[1] or the last position is T[n] respectively). Hence we consider each pair B_1B_2 , B_2B_3 , and so on. We first apply Lemma 13 with $\Delta = \sigma_t/8$ and $\tilde{\sigma} = (\sigma_t)^{1/4}/\log(\sigma_t)$ to find an $(\sigma_t/8)$ -approximate LZ factorisation of the corresponding fragment of T[1..n], and then use Lemma 19 to compute all runs of length at least σ_t , apart from possibly the ones that are properly contained in a tail. Of the computed runs, we discard the ones that contain the first or last position of the block pair (unless the first position is T[1] or the last position is T[n] respectively). This way, each reported run is a run not only with respect to the block pair, but with respect to the entire T[1..n]. If we do not report some run of length at least σ_t and less than $(\sigma_t)^2$ in this way, then it is properly contained in one of the tails.

We cannot always afford to apply Lemmas 13 and 19 to all block pairs. Thus, we have to deactivate some of the blocks. During the current phase t, for each tail T[s..e] of length at least Δ ,

we deactivate all block pairs in phase t + 3 that are contained in T[s + 1..e - 1]. By similar logic as in Section 3.2, if a tail contributes e - s + 1 comparisons and time to the application of Lemma 19, then it permanently deactivates $\Omega(e - s + 1)$ positions of the string, and thus the total time and comparisons needed for all invocations of Lemmas 13 and 19 are bounded by $\mathcal{O}(n)$ (apart from the additional $\mathcal{O}(n \log^* n)$ total time for Lemma 13). Whenever we apply Lemma 13, we add all the tails of length at least Δ to a list \mathcal{L} , where each tail is annotated with the position of its previous occurrence. After the algorithm terminates, \mathcal{L} contains all sufficiently long tails from all phases. We have already shown that the total time needed for Lemma 19 is bounded by $\mathcal{O}(n)$, and thus the total length of the tails in \mathcal{L} is at most $\mathcal{O}(n)$.

If any of the calls to Lemma 13 in the current phase detects that $\sigma > \tilde{\sigma}$, or if $\tilde{\sigma} < 256$, we immediately switch to applying Lemma 18 on every pair of blocks $B_i B_{i+1}$ of the current phase, which takes $\mathcal{O}(n \log \sigma)$ time (because the length of a block pair is polynomial in $\tilde{\sigma}$). Again, after applying Lemma 13 to $B_i B_{i+1}$, we discard all runs that contain the first or last position of $B_i B_{i+1}$ (unless the first position is T[1] or the last position is T[n], respectively). After this procedure terminates, we have computed all runs, except for possibly some of the runs that were properly contained in a tail in list \mathcal{L} . We may have reported some duplicate runs, which we filter out as follows. The number of runs reported so far is $r = \mathcal{O}(n \log \sigma)^4$. We sort them in additional $\mathcal{O}(n+r) = \mathcal{O}(n \log \sigma)$ time, e.g., by using radix sort, and remove duplicates. The running time so far is $\mathcal{O}(n \log \sigma)$.

5.1 Copying Runs From Previous Occurrences

Lastly, we have to compute the runs that were properly contained in a tail in \mathcal{L} . Consider such a run $\langle r_s, r_e, p \rangle$, and let T[s..e] be a tail in \mathcal{L} with $s < r_s$ and $r_e < e$. If multiple tails match this criterion, let T[s..e] be the one that maximizes e. In \mathcal{L} , we annotated T[s..e] with its previous occurrence T[s - d..e - d]. Note that $\langle r_s - d, r_e - d, p \rangle$ is also a run. Thus, if we compute the runs in an appropriate order, we can simply copy the missing runs from their respective previous occurrences. For this sake, we annotate each position $i \in [1, n]$ with:

- a list of all the runs $\langle i, e, p \rangle$ that we already computed, arranged in increasing order of end position e. We already sorted the runs for duplicate elimination, and can annotate all position in $\mathcal{O}(n)$ time.
- a pair (e^*, d^*) , where $e^* = d^* = 0$ if there is no tail T[s..e] such that s < i < e. Otherwise, among all tails T[s..e] with s < i < e, we choose the one that maximizes e. Let T[s-d..e-d] be its previous occurrence, then we use $e^* = e$ and $d^* = d$. As explained earlier, the total length of all tails in \mathcal{L} is $\mathcal{O}(n)$, and thus we can simply scan each tail and update the annotation pair of each contained position whenever necessary.

Observe that, if a position is annotated with (0,0), then none of the runs starting at position i is fully contained in a tail, and thus we have already annotated position i with the complete list of the runs starting at i. Now we process the positions $i \in [1, n]$ one at a time and in increasing order. We inductively assume that, at the time at which we process i, we have already annotated each j < i with the complete list of runs starting at j. Hence our goal is to complete the list of i such that it contains all runs starting at i. If i is annotated with (0,0), then the list is already complete. Otherwise, i is annotated with (e,d), every missing run $\langle i, e_r, p \rangle$ satisfies $e_r < e$, and the annotation list of i-d already contains the run $\langle i-d, e_r - d, p \rangle$ (due to $T[i-1..e_r+1] = T[i-d-1..e_r-d+1]$

⁴a more careful analysis would reveal that it is $\mathcal{O}(n)$, but this is not necessary for the proof

and the inductive assumption). For each run $\langle i - d, r_e - d, p \rangle$ in the annotation list of position i - d, we insert the run $\langle i, e_r, p \rangle$ into the annotation list of i. We perform this step in a merging fashion, starting with the shortest runs of both lists and zipping them together. As soon as we are about to insert a run $\langle i, e_r, p \rangle$ with $e_r \ge e$, we do not insert it and abort. Thus, the time needed for processing i is linear in the number of runs starting at position i. By the runs theorem [8], the total number of runs is less than n, making the total time for this step $\mathcal{O}(n)$.

Apart from the new steps in Section 5.1, the complexity analysis works exactly like in Section 3.2. Hence we have shown:

Theorem 5.1. Computing all runs in a length-*n* string that contains σ distinct symbols from a general unordered alphabet can be implemented in $\mathcal{O}(n \log \sigma)$ comparisons and $\mathcal{O}(n \log \sigma + n \log^* n)$ time.

5.2 Final Improvement for Computing Runs

The goal is now to adapt the final algorithm to detect all runs. We can no longer stop as soon as we detect a square, and we cannot simply deactivate pairs of blocks that occur earlier. However, Theorem 2.2 is actually capable of reporting all runs in $T[i \cdot \tau ..i \cdot \tau + \ell_i)$ and $T[(i+1) \cdot \tau - r_i ..(i+1) \cdot \tau)$ in $\mathcal{O}(\ell_i + r_i)$ time, and we do not need to terminate the algorithm if these fragments are not squarefree. Thus, we can indeed deactivate any block pair fully contained in $T[i \cdot \tau ..i \cdot \tau + \ell_i)$ and $T[(i+1) \cdot \tau - r_i ..(i+1) \cdot \tau)$. Next, we also deactivate block pairs fully contained in the longest prefix of $T[i_j + \tau ..i_{j+1} + \tau)$ equal to $T[i_{j'} + \tau ..i_{j'+1} + \tau)$, for some j' < j. Denoting the length of this prefix by ℓ , we treat $T[i_j + \tau ..i_j + \ell]$ as a tail and add it to the list \mathcal{L} (annotated with $i_{j'}$). The total length of all fragments added to \mathcal{L} is still $\mathcal{O}(n)$.

Theorem 5.2. Computing all runs in a length-n string that contains σ distinct symbols from a general unordered alphabet can be implemented in $\mathcal{O}(n \log \sigma)$ comparisons and $\mathcal{O}(n \log \sigma)$ time.

References

- J. Allouche and J. O. Shallit. The ubiquitous Prouhet-Thue-Morse Sequence. In C. Ding, T. Helleseth, and H. Niederreiter, editors, Sequences and their Applications - Proceedings of SETA 1998, Singapore, December 14-17, 1998, Discrete Mathematics and Theoretical Computer Science, pages 1–16. Springer, 1998.
- [2] A. Amir, E. Eisenberg, and A. Levy. Approximate periodicity. In O. Cheong, K. Chwa, and K. Park, editors, Algorithms and Computation - 21st International Symposium, ISAAC 2010, Jeju Island, Korea, December 15-17, 2010, Proceedings, Part I, volume 6506 of Lecture Notes in Computer Science, pages 25–36. Springer, 2010.
- [3] A. Amir, E. Eisenberg, and A. Levy. Approximate periodicity. Inf. Comput., 241:215–226, 2015.
- [4] A. Amir and A. Levy. Approximate period detection and correction. In L. Calderón-Benavides, C. N. González-Caro, E. Chávez, and N. Ziviani, editors, String Processing and Information Retrieval - 19th International Symposium, SPIRE 2012, Cartagena de Indias, Colombia, October 21-25, 2012. Proceedings, volume 7608 of Lecture Notes in Computer Science, pages 1–15. Springer, 2012.

- [5] A. Apostolico and D. Breslauer. An optimal $\mathcal{O}(\log \log n)$ -time parallel algorithm for detecting all squares in a string. SIAM J. Comput., 25(6):1318–1331, 1996.
- [6] A. Apostolico and R. Giancarlo. Periodicity and repetitions in parameterized strings. *Discret.* Appl. Math., 156(9):1389–1398, 2008.
- [7] A. Apostolico and F. P. Preparata. Optimal off-line detection of repetitions in a string. *Theor. Comput. Sci.*, 22:297–315, 1983.
- [8] H. Bannai, T. I, S. Inenaga, Y. Nakashima, M. Takeda, and K. Tsuruta. The "runs" theorem. SIAM J. Comput., 46(5):1501–1514, 2017.
- D. R. Bean, A. Ehrenfeucht, and G. F. McNulty. Avoidable patterns in strings of symbols. *Pacific Journal of Mathematics*, 85(2):261 – 294, 1979.
- [10] M. A. Bender and M. Farach-Colton. The LCA problem revisited. In G. H. Gonnet, D. Panario, and A. Viola, editors, LATIN 2000: Theoretical Informatics, 4th Latin American Symposium, Punta del Este, Uruguay, April 10-14, 2000, Proceedings, volume 1776 of Lecture Notes in Computer Science, pages 88–94. Springer, 2000.
- [11] J. Berstel and L. Boasson. Partial words and a theorem of Fine and Wilf. Theor. Comput. Sci., 218(1):135–141, 1999.
- [12] F. Blanchet-Sadri, D. Bal, and G. Sisodia. Graph connectivity, partial words, and a theorem of Fine and Wilf. Inf. Comput., 206(5):676–693, 2008.
- [13] F. Blanchet-Sadri and R. A. Hegstrom. Partial words and a theorem of Fine and Wilf revisited. *Theor. Comput. Sci.*, 270(1-2):401–419, 2002.
- [14] F. Blanchet-Sadri, S. Simmons, A. Tebbe, and A. Veprauskas. Abelian periods, partial words, and an extension of a theorem of Fine and Wilf. *RAIRO Theor. Informatics Appl.*, 47(3):215– 234, 2013.
- [15] D. Breslauer. Efficient String Algorithmics. PhD thesis, Columbia University, 1992.
- [16] S. Brlek and S. Li. On the number of squares in a finite word. CoRR, abs/2204.10204, 2022.
- [17] M. G. Castelli, F. Mignosi, and A. Restivo. Fine and Wilf's theorem for three periods and a generalization of Sturmian words. *Theor. Comput. Sci.*, 218(1):83–94, 1999.
- [18] R. Cole and R. Hariharan. Tighter upper bounds on the exact complexity of string matching. SIAM J. Comput., 26(3):803–856, 1997.
- [19] R. Cole, R. Hariharan, M. Paterson, and U. Zwick. Tighter lower bounds on the exact complexity of string matching. SIAM J. Comput., 24(1):30–45, 1995.
- [20] S. Constantinescu and L. Ilie. Fine and Wilf's theorem for abelian periods. Bulletin of the EATCS, 89:167–170, 01 2006.
- [21] T. M. Cover and J. A. Thomas. *Elements of information theory*. Wiley, 2nd edition, 2006.
- [22] M. Crochemore. An optimal algorithm for computing the repetitions in a word. Inf. Process. Lett., 12(5):244–250, 1981.

- [23] M. Crochemore. Transducers and repetitions. Theor. Comput. Sci., 45(1):63–86, 1986.
- [24] M. Crochemore and L. Ilie. Maximal repetitions in strings. J. Comput. Syst. Sci., 74(5):796– 807, 2008.
- [25] M. Crochemore, L. Ilie, and L. Tinta. The "runs" conjecture. Theor. Comput. Sci., 412(27):2931–2941, 2011.
- [26] M. Crochemore, C. S. Iliopoulos, T. Kociumaka, R. Kundu, S. P. Pissis, J. Radoszewski, W. Rytter, and T. Walen. Near-optimal computation of runs over general alphabet via noncrossing LCE queries. In S. Inenaga, K. Sadakane, and T. Sakai, editors, *String Processing and Information Retrieval - 23rd International Symposium, SPIRE 2016, Beppu, Japan, October* 18-20, 2016, Proceedings, volume 9954 of Lecture Notes in Computer Science, pages 22–34, 2016.
- [27] J. D. Currie. Pattern avoidance: themes and variations. Theor. Comput. Sci., 339(1):7–18, 2005.
- [28] A. Deza, F. Franek, and A. Thierry. How many double squares can a string contain? Discret. Appl. Math., 180:52–69, 2015.
- [29] J. Duval, T. Lecroq, and A. Lefebvre. Linear computation of unbordered conjugate on unordered alphabet. *Theor. Comput. Sci.*, 522:77–84, 2014.
- [30] J. Ellert and J. Fischer. Linear Time Runs Over General Ordered Alphabets. In N. Bansal, E. Merelli, and J. Worrell, editors, 48th International Colloquium on Automata, Languages, and Programming (ICALP 2021), volume 198 of Leibniz International Proceedings in Informatics (LIPIcs), pages 63:1–63:16, Dagstuhl, Germany, 2021. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- [31] M. Farach. Optimal suffix tree construction with large alphabets. In 38th Annual Symposium on Foundations of Computer Science, FOCS '97, Miami Beach, Florida, USA, October 19-22, 1997, pages 137–143. IEEE Computer Society, 1997.
- [32] N. J. Fine and H. S. Wilf. Uniqueness theorems for periodic functions. In Proceedings of the American Mathematical Society, volume 16, pages 109–114, 1965.
- [33] J. Fischer, S. Holub, T. I, and M. Lewenstein. Beyond the runs theorem. In C. S. Iliopoulos, S. J. Puglisi, and E. Yilmaz, editors, *String Processing and Information Retrieval - 22nd International Symposium, SPIRE 2015, London, UK, September 1-4, 2015, Proceedings*, volume 9309 of *Lecture Notes in Computer Science*, pages 277–286. Springer, 2015.
- [34] J. Fischer, T. I, D. Köppl, and K. Sadakane. Lempel-Ziv factorization powered by space efficient suffix trees. *Algorithmica*, 80(7):2048–2081, 2018.
- [35] A. S. Fraenkel and J. Simpson. How many squares can a string contain? J. Comb. Theory, Ser. A, 82(1):112–120, 1998.
- [36] F. Franek and Q. Yang. An asymptotic lower bound for the maximal number of runs in a string. Int. J. Found. Comput. Sci., 19(1):195–203, 2008.

- [37] H. N. Gabow. Data structures for weighted matching and nearest common ancestors with linking. In D. S. Johnson, editor, *Proceedings of the First Annual ACM-SIAM Symposium* on Discrete Algorithms, 22-24 January 1990, San Francisco, California, USA, pages 434–443. SIAM, 1990.
- [38] Z. Galil and J. I. Seiferas. Time-space-optimal string matching. J. Comput. Syst. Sci., 26(3):280–294, 1983.
- [39] P. Gawrychowski, T. Kociumaka, W. Rytter, and T. Walen. Faster longest common extension queries in strings over general alphabets. In R. Grossi and M. Lewenstein, editors, 27th Annual Symposium on Combinatorial Pattern Matching, CPM 2016, June 27-29, 2016, Tel Aviv, Israel, volume 54 of LIPIcs, pages 5:1–5:13. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2016.
- [40] M. Giraud. Not so many runs in strings. In C. Martín-Vide, F. Otto, and H. Fernau, editors, Language and Automata Theory and Applications, Second International Conference, LATA 2008, Tarragona, Spain, March 13-19, 2008. Revised Papers, volume 5196 of Lecture Notes in Computer Science, pages 232–239. Springer, 2008.
- [41] M. Giraud. Asymptotic behavior of the numbers of runs and microruns. *Inf. Comput.*, 207(11):1221–1228, 2009.
- [42] G. Gourdel, T. Kociumaka, J. Radoszewski, W. Rytter, A. M. Shur, and T. Walen. String periods in the order-preserving model. *Inf. Comput.*, 270, 2020.
- [43] D. Gusfield. Algorithms on Strings, Trees, and Sequences Computer Science and Computational Biology. Cambridge University Press, 1997.
- [44] Y. Han and M. Thorup. Integer sorting in $\mathcal{O}(n\sqrt{\log \log n})$ expected time and linear space. In *FOCS*, pages 135–144. IEEE Computer Society, 2002.
- [45] S. Holub. Prefix frequency of lost positions. Theor. Comput. Sci., 684:43–52, 2017.
- [46] J. Hong and G. Chen. Efficient on-line repetition detection. Theor. Comput. Sci., 407(1-3):554–563, 2008.
- [47] L. A. Idiatulina and A. M. Shur. Periodic partial words and random bipartite graphs. Fundam. Informaticae, 132(1):15–31, 2014.
- [48] L. Ilie. A note on the number of squares in a word. Theor. Comput. Sci., 380(3):373–376, 2007.
- [49] J. Justin. On a paper by Castelli, Mignosi, Restivo. RAIRO Theor. Informatics Appl., 34(5):373–377, 2000.
- [50] T. Kociumaka, J. Radoszewski, W. Rytter, and T. Walen. A periodicity lemma for partial words. Inf. Comput., 283:104677, 2022.
- [51] R. M. Kolpakov, G. Bana, and G. Kucherov. mreps: efficient and flexible detection of tandem repeats in DNA. *Nucleic Acids Res.*, 31(13):3672–3678, 2003.
- [52] R. M. Kolpakov and G. Kucherov. Finding maximal repetitions in a word in linear time. In 40th Annual Symposium on Foundations of Computer Science, FOCS '99, 17-18 October, 1999, New York, NY, USA, pages 596–604. IEEE Computer Society, 1999.

- [53] S. R. Kosaraju. Computation of squares in a string (preliminary version). In M. Crochemore and D. Gusfield, editors, *Combinatorial Pattern Matching, 5th Annual Symposium, CPM 94, Asilomar, California, USA, June 5-8, 1994, Proceedings, volume 807 of Lecture Notes in Computer Science*, pages 146–150. Springer, 1994.
- [54] D. Kosolobov. Online square detection. CoRR, abs/1411.2022, 2014.
- [55] D. Kosolobov. Lempel-Ziv Factorization May Be Harder Than Computing All Runs. In E. W. Mayr and N. Ollinger, editors, 32nd International Symposium on Theoretical Aspects of Computer Science (STACS 2015), volume 30 of Leibniz International Proceedings in Informatics (LIPIcs), pages 582–593, Dagstuhl, Germany, 2015. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [56] D. Kosolobov. Online detection of repetitions with backtracking. In F. Cicalese, E. Porat, and U. Vaccaro, editors, *Combinatorial Pattern Matching - 26th Annual Symposium, CPM* 2015, Ischia Island, Italy, June 29 - July 1, 2015, Proceedings, volume 9133 of Lecture Notes in Computer Science, pages 295–306. Springer, 2015.
- [57] D. Kosolobov. Computing runs on a general alphabet. Inf. Process. Lett., 116(3):241–244, 2016.
- [58] D. Kosolobov. Finding the leftmost critical factorization on unordered alphabet. Theor. Comput. Sci., 636:56–65, 2016.
- [59] S. Li, J. Pachocki, and J. Radoszewski. A note on the maximum number of k-powers in a finite word. CoRR, abs/2205.10156, 2022.
- [60] M. Maekawa. A \sqrt{N} algorithm for mutual exclusion in decentralized systems. ACM Trans. Comput. Syst., 3(2):145–159, 1985.
- [61] M. G. Main and R. J. Lorentz. An $\mathcal{O}(n \log n)$ algorithm for finding all repetitions in a string. J. Algorithms, 5(3):422–432, 1984.
- [62] W. Matsubara, K. Kusano, H. Bannai, and A. Shinohara. A series of run-rich strings. In A. Dediu, A. Ionescu, and C. Martín-Vide, editors, *Language and Automata Theory and Applications, Third International Conference, LATA 2009, Tarragona, Spain, April 2-8, 2009. Proceedings*, volume 5457 of *Lecture Notes in Computer Science*, pages 578–587. Springer, 2009.
- [63] W. Matsubara, K. Kusano, A. Ishino, H. Bannai, and A. Shinohara. New lower bounds for the maximum number of runs in a string. In J. Holub and J. Zdárek, editors, *Proceedings of* the Prague Stringology Conference 2008, Prague, Czech Republic, September 1-3, 2008, pages 140–145. Prague Stringology Club, Department of Computer Science and Engineering, Faculty of Electrical Engineering, Czech Technical University in Prague, 2008.
- [64] Y. Matsuoka, T. Aoki, S. Inenaga, H. Bannai, and M. Takeda. Generalized pattern matching and periodicity under substring consistent equivalence relations. *Theor. Comput. Sci.*, 656:225– 233, 2016.
- [65] O. Merkurev and A. M. Shur. Searching runs in streams. In N. R. Brisaboa and S. J. Puglisi, editors, String Processing and Information Retrieval 26th International Symposium, SPIRE 2019, Segovia, Spain, October 7-9, 2019, Proceedings, volume 11811 of Lecture Notes in Computer Science, pages 203–220. Springer, 2019.

- [66] O. Merkurev and A. M. Shur. Computing the maximum exponent in a stream. Algorithmica, 84(3):742–756, 2022.
- [67] S. J. Puglisi, J. Simpson, and W. F. Smyth. How many runs can a string contain? Theor. Comput. Sci., 401(1-3):165–171, 2008.
- [68] W. Rytter. The number of runs in a string: Improved analysis of the linear upper bound. In B. Durand and W. Thomas, editors, STACS 2006, 23rd Annual Symposium on Theoretical Aspects of Computer Science, Marseille, France, February 23-25, 2006, Proceedings, volume 3884 of Lecture Notes in Computer Science, pages 184–195. Springer, 2006.
- [69] A. M. Shur and Y. V. Gamzova. Partial words and the interaction property of periods. *Izvestiya: Mathematics*, 68(2):405–428, apr 2004.
- [70] A. M. Shur and Y. V. Konovalova. On the periods of partial words. In J. Sgall, A. Pultr, and P. Kolman, editors, *Mathematical Foundations of Computer Science 2001, 26th International Symposium*, MFCS 2001 Marianske Lazne, Czech Republic, August 27-31, 2001, Proceedings, volume 2136 of Lecture Notes in Computer Science, pages 657–665. Springer, 2001.
- [71] J. Simpson. Modified Padovan words and the maximum number of runs in a word. Australas. J Comb., 46:129–146, 2010.
- [72] A. Thierry. A proof that a word of length n has less than 1.5n distinct squares. CoRR, abs/2001.02996, 2020.
- [73] A. Thue. Über unendliche Zeichenreihen. Norske Vid. Selsk. Skr., I Mat.-Nat. Kl., Christiania, 7:1–22, 1906.
- [74] R. Tijdeman and L. Zamboni. Fine and Wilf words for any periods. Indagationes Mathematicae, 14(1):135–147, 2003.
- [75] J. Vuillemin. A unifying look at data structures. Commun. ACM, 23(4):229–239, 1980.