

University of Warsaw
Faculty of Mathematics, Informatics and Mechanics

Wiktor Zuba

Efficient enumerations in words
PhD dissertation

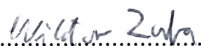
Supervisor:
prof. dr hab. Wojciech Rytter
Institute of Informatics
University of Warsaw

June 2021

Author's declaration:

I hereby declare that this dissertation is my own work.


June 8, 2021


.....
Wiktor Zuba

Supervisor's declaration:

The dissertation is ready to be reviewed

June 8, 2021


.....
Wojciech Rytter

Abstract

The dissertation presented here consists of a set of results in text processing algorithms from four different scientific publications. The papers focus on algorithms and data structures that allow efficient enumeration and counting of certain classes of factors. Most of the results are based on different variations of the same key property - regularity of the set of occurrences of factors of specific types.

In "Efficient Representation and Counting of Antipower Factors in Words" we focused on a class of k -antipowers (words consisting of k pairwise different segments of the same length). The results of the paper consist of effective algorithms for enumerating all occurrences, counting all occurrences and counting different k -antipower factors in a string, working in $O(nk \log k + \text{output})$, $O(nk \log k)$ and $O(nk^4 \log k \log n)$ time, respectively.

The paper "Counting Distinct Patterns in Internal Dictionary Matching" deals with a problem in which we are given an input word T and its internal dictionary D . We want to answer questions of the form: "How many distinct words from D are the factors of $T[l..j]$?" We have shown few algorithms and data structures, that allows us to do that effectively. The algorithms differ in exploited properties and complexities (each one can be more effective from the other depending on the ratio between T and D).

In "The Number of Repetitions in 2D-Strings" we focused on the extensions of the notions of squares and runs to the case of two-dimensional strings. We gave new bounds on the maximal number of their occurrences in strings and effective algorithms for enumerating them.

In the fourth paper titled "Efficient Enumeration of Distinct Factors Using Package Representations", we proposed a new compact representation of string factors. It allows us to obtain new effective algorithms for enumeration and counting of distinct factors it contains. In particular it allowed us to obtain a new simple algorithm for finding distinct squares and a faster algorithm for finding all distinct antipowers in a string (improving one of the results from the first of enclosed papers).

Keywords: algorithm, data structure, enumeration, counting, factor, subword, antipower, runs, two-dimensional string, internal dictionary

Efektywne zliczanie w słowach.

Streszczenie

Przedstawiona rozprawa stanowi zbiór wyników prac na temat algorytmów przetwarzających dane tekstowe. Prace skupiają się na zaprezentowaniu algorytmów i struktur danych pozwalających na efektywne znajdowanie i zliczanie podśłów należących do konkretnych klas podśłów. Innym elementem wspólnym poniższych prac jest to, że wykorzystują one regularności w słowach związane z wystąpieniami podśłów szczególnego typu.

W pracy "Efficient Representation and Counting of Antipower Factors in Words" zajęliśmy się klasą k -antypotę (słów składających się z k parami różnych członów o tej samej długości). Rezultatem pracy są efektywne algorytmy wypisywania wszystkich wystąpień, zliczania wszystkich wystąpień oraz zliczania różnych k -antypotę w słowie, działające w czasach odpowiednio $O(nk \log k + \text{wynik})$, $O(nk \log k)$ i $O(nk^4 \log k \log n)$.

Praca "Counting Distinct Patterns in Internal Dictionary Matching" przedstawia problem w którym mając dane słowo wejściowe T oraz słownik jego podśłów D chcielibyśmy móc efektywnie odpowiadać na pytania "Ile różnych słów ze słownika D jest podśłowami słowa $T[i..j]$?". Przedstawiliśmy kilka algorytmów/struktur danych, które pozwalają na rozwiązanie tego problemu. Algorytmy różnią się wykorzystywanymi założeniami oraz złożonościami (każdy może być lepszy od pozostałych w zależności od zastosowania, oraz stosunku wielkości T i D).

W pracy "The Number of Repetitions in 2D-Strings" zajęliśmy się uogólnieniami pojęć kwadratów i maksymalnych powtórzeń na teksty dwuwymiarowe. Podaliśmy nowe ograniczenia na maksymalną ich ilość w tekstach, oraz efektywne algorytmy ich wyznaczania.

Na końcu przedstawiam pracę "Efficient Enumeration of Distinct Factors Using Package Representations", w której wprowadziliśmy nowy, skompresowany sposób reprezentacji wystąpień wielu podśłów. Pozwala ona na otrzymanie nowych, efektywnych algorytmów wyznaczania i zliczania różnych słów przez nią reprezentowanych. W szczególności pozwoliło nam to na otrzymanie nowego prostego algorytmu wyznaczania kwadratów i szybszego algorytmu wyznaczania antypotę (poprawiając jeden z wyników pierwszej z przedstawionych prac).

Słowa kluczowe: algorytm, struktura danych, wyliczanie, zliczanie, podśło, antypotęga, maksymalne powtórzenie, słowo dwuwymiarowe, słownik wewnętrzny

Contents

1	Extended abstract	1
1.1	Introduction	1
1.1.1	Motivation	1
1.1.2	Contents	1
1.2	Preliminaries	2
1.3	Presentation of the main results	4
1.3.1	Efficient Representation and Counting of Antipower Factors in Words	4
1.3.2	Counting Distinct Patterns in Internal Dictionary Matching	6
1.3.3	The Number of Repetitions in 2D-Strings	8
1.3.4	Efficient Enumeration of Distinct Factors Using Package Representations	11
1.4	Other publications	12
2	Autoreferat	14
2.1	Wstęp	14
2.1.1	Motywacja	14
2.1.2	Skład pracy	14
2.2	Preliminaria	15
2.3	Przegląd wyników	17
2.3.1	Efficient Representation and Counting of Antipower Factors in Words	17
2.3.2	Counting Distinct Patterns in Internal Dictionary Matching	19
2.3.3	The Number of Repetitions in 2D-Strings	21
2.3.4	Efficient Enumeration of Distinct Factors Using Package Representations	25
2.4	Pozostałe publikacje	25
3	Efficient Representation and Counting of Antipower Factors in Words	27
4	Counting Distinct Patterns in Internal Dictionary Matching	58
5	The Number of Repetitions in 2D-Strings	74

6 Efficient Enumeration of Distinct Factors Using Package Representations	93
---	----

Chapter 1

Extended abstract

1.1 Introduction

1.1.1 Motivation

The study of regular words and especially of algorithms finding and counting such words in texts is one of the most popular areas of text algorithms. There are many algorithms, which work differently on periodic and aperiodic parts of strings, as such parts may show distinct properties. For example, high periodicity of a word results in a very good compression rate - we can describe it using only its period and length. On the other hand it may make searching (for it or in it) much more difficult, as an aperiodic word of length m may occur in a word of length n at most $2\frac{n}{m}$ times. That is not true if the word is periodic - for instance a word $aaaaaa$ appears in $aaaaaaaaaaaaaaaaaaaa$ $n - m + 1 = 11$ times. The number of such (distinct) structures in a word can by itself serve as a good measure of the words regularity or complexity.

The presented dissertation is composed of publications in which (together with my coauthors) we have concerned a few classes of such regular (or irregular in a specific sense) subwords. In each of the papers we have shown new, effective algorithms for enumerating and counting (all occurrences or all distinct) factors from a specific class.

1.1.2 Contents

The dissertation is composed of four articles written together with my coauthors during the course of my doctoral studies at the Faculty of Mathematics, Informatics and Mechanics of the University of Warsaw:

1. Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszyński, Tomasz Waleń, Wiktor Zuba: *Efficient Representation and Counting of Antipower Factors in Words*. LATA 2019: 421-433

The paper is included in the version accepted for publication in the Special Issue of Information & Computation derived from LATA 2019 conference.

2. Panagiotis Charalampopoulos, Tomasz Kociumaka, Manal Mohamed, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Waleń, Wiktor Zuba: *Counting Distinct Patterns in Internal Dictionary Matching*. CPM 2020: 8:1-8:15
3. Panagiotis Charalampopoulos, Jakub Radoszewski, Wojciech Rytter, Tomasz Waleń, Wiktor Zuba: *The Number of Repetitions in 2D-Strings*. ESA 2020: 32:1-32:18
4. Panagiotis Charalampopoulos, Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Tomasz Waleń, Wiktor Zuba: *Efficient Enumeration of Distinct Factors Using Package Representations*. SPIRE 2020: 247-261

All those papers contain original results in text algorithms. They were written by the authors representing the University of Warsaw together with scientists from King's College London (Panagiotis Charalampopoulos had a long term stay in Warsaw) and Tomasz Kociumaka, who changed his affiliation from the University of Warsaw to Bar-Ilan University in the meantime.

In section 1.2 of this abstract I introduce the basic notions - regular classes of words and factors, which are used throughout the presented papers.

In section 1.3 I show a detailed overview of our results together with sketches of methods used to obtain them.

In section 1.4 I list the publications I coauthored during the course of my doctoral studies, which however are not included as a part of this dissertation.

Chapter 2 contains a polish version of this extended abstract, while chapters 3, 4, 5 and 6 contain the papers included in the dissertation. The bibliography at the end contains references to papers cited in this abstract.

1.2 Preliminaries

To present the results we need to define the string regularities they use or search for first. Many of them are used by more than just one of the papers included in this dissertation.

Definition 1.1. Let us assume that $x = y_0 y_1 \cdots y_{k-1}$ where $k \geq 2$ and y_i are words of length d . We say that:

- x is a **k -power** if all y_i 's are the same (x is a **square** if in addition $k = 2$);
- x is a **k -antipower** if all y_i 's are pairwise distinct;
- x is a **weak k -power** if it is not a k -antipower, that is, if $y_i = y_j$ for some $i \neq j$;
- x is a **gapped square** if $y_0 = y_{k-1}$.

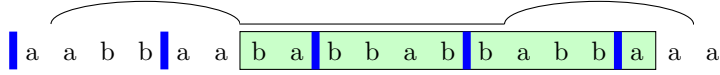


Figure 1.1: Regularities in standard strings. A run (with the period equal to 3) is depicted by a green background. Blue lines mark an occurrence of a 4-antipower. Notice that the words of length 16 starting one or two positions later are not 4-antipowers as they are gapped squares generated by a maximal 3-gapped repeat.

Definition 1.2. We say that:

- p is a **period** of a word w if $w[i] = w[i + p]$ for all $i \in [0 .. |w| - p)$ (by the period of a word we usually mean the shortest one);
- a fragment $w[i..j]$ is a **run** (a maximal repetition) if its length is equal to at least a double of its shortest period and it cannot be extended by a one position to the left or to the right without breaking the period;
- a fragment $w[i..j]$ is a **maximal ρ -gapped repeat** (for $\rho \geq 1$) if it is of a form uvu for $p = |uv| + \rho|u|$ and it cannot be extended by a one position to the left or to the right without breaking the period.

Definition 1.3. For two-dimensional texts we define repetitions analogously.

- A **tandem** is a two-dimensional text of even width, in which its left and right halves are equal as texts (extension of square to the case where instead of letters we have columns of letters).
- A **quartic** is a two-dimensional text Q of even height and width, such that Q and Q are tandems (it is composed of four identical texts in a 2×2 grid).
- A **2D-run** is a $i \times j$ fragment of a text such that its vertical period is smaller or equal to $\frac{i}{2}$ and its horizontal period is smaller or equal to $\frac{j}{2}$. Moreover it cannot be extended by a row or a column from any side (if such exist in the text) without breaking any of those periods.

By a square generated by a run I mean a square fully contained in it with length equal to exactly double of its period. All such squares can be easily obtained from a run, and every square is generated by a run. The generation of gapped-squares by runs and maximal ρ -gapped repeats and the generation of quartics by 2D-runs work analogously.

A square is primitively rooted if its root (half) is not a k -power for any $k > 1$. Analogously a quartic is primitively rooted if its root (quarter) is not a k -power for any $k > 1$ neither horizontally nor diagonally.

Definition 1.4. By an **internal dictionary** we mean a set of factors of a given word represented by the starting and ending positions of a single occurrence of each factor.

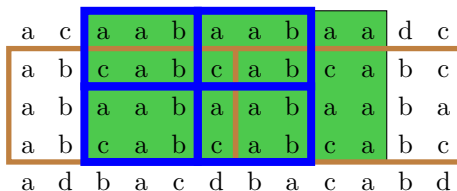


Figure 1.2: Different types of repetitions in 2D-strings. A tandem is marked by brown rims, a quartic by the blue ones, while the green background depicts a 2D-run.

1.3 Presentation of the main results

1.3.1 Efficient Representation and Counting of Antipower Factors in Words

The paper presented here was presented on LATA 2019 conference and published in conference proceedings. Extended version of that paper was later accepted for publication (it is not published yet) in the Special Issue of Information & Computation journal derived from that conference. It is included in this dissertation in the mentioned extended version (in chapter 3).

Introduction

Badkobeh et al. in [5] presented a lower and upper bound $\Theta(n^2/k)$ for the maximal number of occurrences of k -antipowers in a string. Furthermore they presented an $O(n^2/k)$ time algorithm that finds all those occurrences. Due to the lower bound on the maximal size of the output the algorithm is optimal in the pessimistic case. It does not mean however, that nothing more in this field can be done.

In our publication we have focused on the problems that are not directly influenced by those bounds - counting of antipowers and enumerating them with complexity parameterized by the size of the output (more effective if a word contains a small number of antipowers).

More formally - we presented algorithms, which for a word of length n count all the occurrences of k -antipower factors in $O(nk \log k)$ time, find all those occurrences in $O(nk \log k + \text{output})$ time, and count all distinct k -antipowers in a word in $O(nk^4 \log k \log n)$ time. We also presented a data structure that allows us to check efficiently if a given subword is a k -antipower.

In (the most common) case when k is small in comparison to the length of the word (smaller than \bar{n}) our algorithm has a better asymptotic performance time than the one previously known.

Weak powers and a compact representation

Since antipowers itself are very irregular it is not that easy to effectively find them. Due to that, we focused on computing a representation of all the other

words, that is of all weak k -powers. Every weak k -power on the other hand is a consequence of the existence of a gapped square (possibly for a smaller k , but certainly for the same d).

Observation 1.1.

- (a) Each (gapped) square is generated by a run with (not necessarily shortest) period $(q + 1) \cdot d$ or a maximal $(q + 1)$ -gapped repeat with a period equal to $(q + 1) \cdot d$ for some $q \leq k - 2$.
- (b) Each run and each maximal q -gapped repeat generate a single interval of positions in which (gapped) squares of a given length start, moreover this interval can be computed in constant time.
- (c) In a word of length n there are at most $O(n)$ general runs and their representation can be computed in $O(n)$ time ([6]).
- (d) There are only $O(n)$ maximal q -gapped repeats in a word and their representation can be computed in $O(n)$ time ([11, 12]).

With the use of those properties we can compute the representation of weak k -powers for all d 's as $O(nk)$ interval chains (set of intervals of the same lengths and beginnings forming an arithmetic progression) in total. Then, using geometric methods we can determine their set-theoretic sum in $O(nk \log k)$ total time. This immediately gives us the number of all weak k -powers, hence also of all k -antipower factors. Instead of counting the size of the complement of this set we can enumerate through all of the positions it does not contain and thus obtain an algorithm that finds all k -antipowers in $O(nk \log k + \text{output})$ time.

Antipower queries

Paper [5] included descriptions of two data structures, that after being constructed answer queries "Is a subword $w[i..j]$ a k -antipower?". The first structure requires $O(n)$ space and answers those questions in $O(k)$ time, while the other answers them in a constant time, but requires storing information of $O(n^2)$ size.

In our paper we have shown a new structure parameterized by a value $r \in [1, n]$. After being built in $O(n^2/r)$ time our structure stores information of size $O(n^2/r)$ and answers such questions in $O(r)$ time.

If $r < k$ then we know that the base of the power d equals at most $n/r > n/k$. To answer the questions we store for each $d \leq n/r$ separately a data structure based on the range minimum queries structure, that is of $O(n)$ size and answers questions in constant time. Otherwise ($r \geq k$) to obtain the right complexities it is enough to just use the mentioned data structure from [5].

Counting distinct k -antipowers

In the journal version of the paper we have added a new result - an algorithm for counting distinct k -antipowers in a word.

Due to a potentially large number of all occurrences of antipowers, the methods enumerating through those cannot obtain a satisfactory computation time.

To obtain such a result we made use of the connection between antipowers and weak powers once again.

The number of distinct k -antipower factors can be obtained as a subtraction of two other values. The first one is equal to the number of (all) distinct factors of the word of length divisible by k . We can compute it in $O(n)$ time with the use of a suffix tree.

The second value needed is the number of distinct weak k -powers. To obtain that we have strengthened the definition of weak powers, so that each such subword can be generated only in one way. Then with the use of a reduction to a graph problem we have acquired an algorithm computing that number in $O(nk^4 \log k \log n)$ total time.

1.3.2 Counting Distinct Patterns in Internal Dictionary Matching

In [8] my coauthors introduced a problem of pattern matching with internal dictionary - for a given word T and its internal dictionary D they presented a data structure that allows effective answering of questions for subwords $T[i..j]$:

- does $T[i..j]$ contain a word from the dictionary?
- return all subwords of $T[i..j]$ that belong to the dictionary
- return all distinct factors from D contained in $T[i..j]$
- count all subwords of $T[i..j]$ that belong to the dictionary
- count the number of distinct factors from D contained in $T[i..j]$

The structure requires $\tilde{O}(n + d)$ space and construction time (where $n = |T|$ denotes the length of the word, $d = |D|$ the number of factors contained in the dictionary, and \tilde{O} is the asymptotic notation ignoring polylogarithmic factors). It answers the first four questions in $\tilde{O}(1 + \text{output})$ time.

Unfortunately for the last question only a data structure answering those question $O(\log n)$ approximately was shown.

Space	Preprocessing time	Query time	Variant
$\tilde{O}(n + d)$	$\tilde{O}(n + d)$	$\tilde{O}(1)$	2-approximation
$\tilde{O}(n^2/m^2 + d)$	$\tilde{O}(n^2/m + d)$	$\tilde{O}(m)$	exact
$\tilde{O}(nd/m + d)$	$\tilde{O}(nd/m + d)$	$\tilde{O}(m)$	exact
$O(n \log^2 n)$	$O(n \log^2 n)$	$O(\log n)$	$D = \text{squares}$, exact

Table 1.1: Our results for Count Distinct queries ($m \in [1, n]$ is an arbitrarily chosen parameter).

2-approximation

To improve the result from the previous paper we have designed a data structure which stores all the exact results for the base words - that is for all subwords of length $(1 + \frac{1}{q})^p$ for all natural p 's and a fixed value $q = \frac{1}{9}$. To obtain

the result for any given subword $T[i..j]$ we make use of two maximal length base words $T[i..i]$ and $T[j..j]$ contained in $T[i..j]$ to obtain its division into three parts $F_1F_2F_3$ ($F_1F_2 = T[i..i]$, $F_2F_3 = T[j..j]$).

The result is obtained by counting all the factors from the dictionary, which have an occurrence starting in F_1 and ending in F_3 , which at the same time occurs in neither of the two base words. We can do that effectively thanks to a large ratio of the length of F_2 to the lengths of F_1 and F_3 . To the obtained value we add the results stored for the two base words, which however results in only a 2-approximate result.

There are $O(n \log n)$ base words in a word of length n . The exact values for them can be counted in $O(n \log^{1+\epsilon} n + d)$ total time (for any constant $\epsilon > 0$) with the use of the property that we can efficiently update the result after a single position shift. The structure used to count the words described earlier costs us additional $O(n + d \log n)$ space and $O(n \log n / \log \log n + d \log^{3/2} n)$ construction time. It answers the questions in $O(\log^2 n / \log \log n)$ time.

Exact counting

Instead of the base words we can focus on storing the results for all the words of the form $T[c_1m+1..c_2m]$ for a chosen m and all the values $c_1 < c_2$ to obtain a data structure of $O(n^2/m^2 + n + d)$ size constructable in $O((n^2 \log n)/m + n \sqrt{\log n} + d)$ time. We can update the result obtained after extending a subword by a one position in $O(\log n)$ time, and to obtain the exact result for any chosen factor we only require up to $O(m)$ such operations (we extend one of the words for which the result is stored).

We can also approach the problem differently, namely, we can inspect the structure of the dictionary.

If the dictionary does not contain a set of k distinct prefixes of the same word, then at any given position only up to $k - 1$ dictionary words can start. In this case we can store for each of the dictionary words the set of positions in T , where such a word occurs. We can construct such a data structure of $O(nk \log n)$ size in $O(nk \log n)$ time. In return, it allows us to respond to our queries in $O(\log n)$ time.

To obtain a dictionary of such a form we can simply remove from it all such large groups of prefixes (building new dictionaries). There are at most d/k such groups, and for each of those we can answer the queries in $O(\log n)$ time (for any $\epsilon > 0$) using a bounded LCP structure ([14]) of $O(n \sqrt{\log n})$ size.

For $k = \frac{d}{m}$ this gives us a query time equal to $O(m \log n + \log n)$.

Counting distinct squares factors

For the popular case, where we want to count the number of distinct square factors in a given subword we have described a yet different algorithm.

In a word of length n there can be up to $O(n^2)$ squares (only $O(n)$ distinct ones), however thanks to the runs (only $O(n)$ such can appear in the word, and all such can be computed in $O(n)$ time) we can distinguish $m = O(n \log n)$ occurrences which really affect our results. With the use of a geometrical data

structure of $O(m \log m)$ size and construction time from [13] we can answer the queries for the number of distinct squares in a subword in $O(\log m) = O(\log n)$ time.

1.3.3 The Number of Repetitions in 2D-Strings

In the paper from chapter 5 we have focused on the repetitions in two-dimensional texts. The table 1.2 contains a summary of the current knowledge about the maximal possible number of 2D-runs, tandems and quartics (distinct or all occurrences of primitively rooted ones) in a two-dimensional string. It also shows the time complexity of the fastest known algorithms which find those structures in a text. Notice, that even though the maximal number of 2D-runs in a string was previously studied the gap between the known lower and upper bounds was linear. In our paper we have reduced this gap to a polylogarithmic one.

	Bounds on the number in an $n \times n$ string	Computation time
2D-runs	$(n^2), O(n^3)$ [2] $O(n^2 \log^2 n)$	$O(n^2 \log n + \text{output})$ [2] $O(n^2 \log^2 n)$
Occurrences of primitively rooted quartics	$(n^2 \log^2 n)$ [3]	$O(n^2 \log^2 n)$
Distinct quartics	$O(n^2 \log^2 n)$	$O(n^2 \log^2 n)$
Occurrences of primitively rooted tandems	$(n^3 \log n)$ [3]	$O(n^3 \log n)$ [4]
Distinct tandems	$(n^2); O(n^4)$ (trivial) (n^3)	$O(n^3)$

Table 1.2: Overview of the previous knowledge and of our results (written in bold).

2D-runs

To obtain our upper bound on the number of 2D-runs we made use of a lemma from [1], to assign to every such 2D-run a maximal horizontal repetition (no condition on vertical periodicity or maximality) of height equal to a power of two (largest possible not exceeding the height of the 2D-run). Such repetition meets three important conditions as well - its upper left or lower left corner equals the respective corner of the 2D-run, its horizontal period is equal to the horizontal period of the 2D-run, and its width equals at least the width of the 2D-run (see figure 1.3)

With the use of those properties, of the periodicity lemma, and of the three squares lemma we proved that the horizontal repetitions occupying the rows from i to $i + 2^k - 1$ can be assigned to at most $O(n \log n)$ distinct 2D-runs in total. By multiplying this value by the number of feasible choices of i and k we



Figure 1.3: Assignment of a maximal horizontal repetition of height 2^k to a 2D-run.

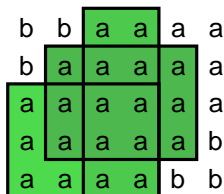


Figure 1.4: Many overlapping 2D-runs

obtain an upper bound on the number of 2D-runs in an $n \times n$ string equal to $O(n^2 \log^2 n)$.

Our bound immediately shows that the algorithm finding all the 2D-runs in a 2D-string from [2] works in $O(n^2 \log^2 n)$ time (their algorithm has time complexity parameterized by the size of the output).

Primitively rooted quartics

In standard strings we can easily extract all the occurrences of primitively rooted squares from the runs - it is enough to take all the subwords of the runs of length equal to the double of their period.

In 2D situation it is somewhat similar - each primitively rooted quartic is a rectangle of width equal to the double of the horizontal period and height equal to the double of the vertical period of a 2D-run, fully contained in it. However, while in one dimension two runs with the same period cannot have a big enough overlap to generate the same square, in 2D many such overlaps can easily appear (see figure 1.4), hence simple adaptation of such a one-dimensional algorithm would result in reporting some of the quartics multiple times.

This problem can be solved with the use of a sweeping line technique ([7]), which allows us to count set-theoretic sums of many families of rectangles on an $n \times n$ grid in $O(n + r + \text{output})$ total time (r denotes the number of all the rectangles).

By dividing all the 2D-runs according to their periods, for each such group of runs we count the set-theoretic sum of the rectangles containing the starting positions of quartics generated by each 2D-run. This way we obtain the desired algorithm.

Due to our bound on r from the previous paragraph (the number of 2D-runs), and the $O(n^2 \log^2 n)$ bound on the number of the occurrences of primitively rooted quartics from [3] our algorithm runs in $O(n^2 \log^2 n)$ total time (optimal

in the pessimistic case by the mentioned bound).

Distinct quartics

In an $n \times n$ unary string there are (n^4) occurrences of general quartics. Due to that in the pessimistic case a trivial algorithm finding all such occurrences in an optimal time. A more interesting problems are the one of bounding the maximal number of distinct quartics (differing as texts and not only by their positioning) and the one of enumerating through all of them.

[3] already gave us a bound on the number of distinct primitively rooted quartics. All the other ones are composed of many adjacent occurrences of the primitively rooted ones forming a larger rectangle. In other words, while a primitively rooted quartic is constructed from occurrences of a primitive text W arranged in a 2×2 grid, in the case of the other ones it is a $k \times l$ grid ($k > 1$ or $l > 1$). We divide such quartics into the thin ones ($k = 1$ or $l = 1$), and the thick ones (both $k, l > 1$).

Then again we perform a yet another division - we group quartics for which their primitive string W has size in $[2^a; 2^{a+1}) \cup [2^b; 2^{b+1})$ for the same values $a, b \rightarrow \log n$.

Due to the periodicity lemma and the three squares lemma any position can be the top left corner of the rightmost occurrence (same quartic does not occur anywhere further to the right) of at most four distinct thin quartics from a single group (it works analogously to a proof of a bound of the number of distinct squares in a standard string [10]).

In the case of the thick quartics the multitude of the grid points in which an occurrence of the W string begins in a large power of W allows us to uniquely assign such points to every smaller power of W . Moreover, due to the periodicity and the three squares lemmas such points cannot belong to a grid for a power of a different string W^0 which belongs to the same group.

Due to this assignments of the points in $[1; n] \times [1; n]$ to the quartics and due to the fact, that there are only $\log^2 n$ choices of a and b we obtain an upper bound on the number of distinct quartics of size $O(n^2 \log^2 n)$.

To compute all the distinct quartics in a text we use our previous algorithm to find all the occurrences of primitively rooted quartics, and then we group them by their roots.

Since other quartics are constructed from the occurrences of the primitive ones (with the same root) we can distinguish the top left corners of all such occurrences, and then find the largest possible grids they form.

This problem can be solved with the sweep line approach to obtain an algorithm working in $O(n^2 \log^2 n)$ time.

Distinct tandems

In the same paper we also considered the related problem of finding all the distinct tandems in a string (as with the quartics, there can be (n^4) occurrences of all the tandems in a 2D-string).

In this problem we can easily identify the tandems occupying the rows from i to j with squares in a standard string by assigning single letters to all the columns of height $j - i + 1$. In a standard word of length n there can be up to $\binom{n}{j-i+1}$ distinct squares, hence the same bound applies to the number of distinct tandems occupying the distinguished rows. Multiplying this value by a number of possible choices of i and j we obtain $O(n^3)$ as the upper bound on the number of distinct tandems.

$\binom{n}{j-i+1}$ as a lower bound can be shown just as easy - it is obtained for a text in which every row is a different unary word.

The assignment of standard letters to the columns of height equal to $j - i + 1$ can be performed constructively. Thanks to that, we can obtain an $O(n^3)$ time algorithm for finding all distinct tandems in a 2D-string, with the use of an algorithm finding distinct squares in a standard string. Due to our lower bound on the maximal size of the output our algorithm is optimal in the pessimistic case.

1.3.4 Efficient Enumeration of Distinct Factors Using Package Representations

In the last of the described papers we introduced a new representation of sets of subwords, which in many cases gives a compact representation, and at the same time allows effective operations on the set.

By a package we mean a set of subwords of the same length, whose starting positions form an interval. Our representation is composed of many such packages represented as triples (starting position, length of words, length of interval). Our greatest interest in this representation lies in the property, that the set of the distinct factors it represents can be found efficiently.

More formally, we are interested in the set

$$\text{Factors}(F) = \{T[j:j+l] : j \in [i:i+k] \text{ and } (i;l;k) \in F\}$$

There are many families of subwords, which can be effectively represented in this way - for example k -powers (when each package corresponds to a run) and k -antipowers (when we can use the representation from section 1.3.1). The obtained representations have size and computation time equal to $O(n)$ and $O(nk^2)$ respectively (each interval chain can be represented by at most k ordinary intervals).

To enumerate through $\text{Factors}(F)$, and compute $|\text{Factors}(F)|$ most effectively we consider two cases. The simpler one, called special, is the one in which if a factor belongs to $\text{Factors}(F)$, then each of its occurrences must be contained in a package from F (it cannot be that some occurrences are represented, and some are not).

In this case for a word of length n and F composed of m packages we make use of a data structure for the longest previous factors ([9]) to obtain efficient enumeration and counting of the result in $O(n + m + \text{output})$ and $O(n + m)$ time respectively. In the case of aforementioned distinct k -powers we obtain a

new, simple algorithm computing them in $O(n)$ time, and in the case of distinct k -antipowers we obtain a new algorithm working in $O(nk^2 + \text{output})$ time. This means, that the newly obtained algorithm has a better performance time than the algorithm described in section 1.3.1, which counts those antipower factors in $O(nk^4 \log k \log n)$ time, and is much more complicated.

The second case, called general does not require this extra assumption, which however makes solving it much more difficult. To do that we adapt a method we used in the paper about k -antipowers (counting of distinct such factors) to obtain algorithms finding and counting $\text{Factor}(F)$ in $O(n \log^2 n + m \log n + \text{output})$ and $O(n \log^2 n + m \log n)$ time respectively.

1.4 Other publications

During the course of my doctoral studies several other papers coauthored by me have been published. I decided not to include them in this dissertation, since they are not directly connected to its topic.

1. Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Faster Recovery of Approximate Periods over Edit Distance. SPIRE 2018: 233-240
2. Wojciech Rytter, Wiktor Zuba: Syntactic View of Sigma-Tau Generation of Permutations. LATA 2019: 447-459
Extended version of the paper was accepted for publication in Theoretical Computer Science journal.
3. Mai Alzamel, Maxime Crochemore, Costas S. Iliopoulos, Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Quasi-Linear-Time Algorithm for Longest Common Circular Factor. CPM 2019: 25:1-25:14
4. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Circular Pattern Matching with k Mismatches. FCT 2019: 213-228
5. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Weighted Shortest Common Supersequence Problem Revisited. SPIRE 2019: 221-238
6. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Shortest Covers of All Cyclic Shifts of a String. WALCOM 2020: 69-80
7. Panagiotis Charalampopoulos, Solon P. Pissis, Jakub Radoszewski, Tomasz Walen, Wiktor Zuba: Unary Words Have the Smallest Levenshtein k -Neighbourhoods. CPM 2020: 10:1-10:12

8. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Internal Quasiperiod Queries. SPIRE 2020: 60-75
9. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Shortest covers of all cyclic shifts of a string. Theor. Comput. Sci. 866: 70-81 (2021)
10. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Circular pattern matching with k mismatches. J. Comput. Syst. Sci. 115: 73-85 (2021)

Chapter 2

Autoreferat

2.1 Wstip

2.1.1 Motywacja

Badanie regularnych fragmentów sów, w szczególności ich znajdowanie i zliczanie jest jedną z najbardziej popularnych części algorytmiki tekstowej. Wiele algorytmów inaczej działa na fragmentach sów które są okresowe niż na takich w których ta własność nie występuje, jako że fragmenty te cechują czysto inne własności. Przykładowo duża okresowość słowa pozwala na łatwiej jego kompresję - opisanie go przy pomocy samego okresu i długości, z drugiej strony może to utrudniać wyszukiwanie takiego słowa - jeśli słowo dane jest nieokresowe, to w tekście dane może się pojawić najwyżej $\frac{2}{m}$ razy, nie jest to jednak prawdą w przypadku sów okresowych (na przykład słowo $aaaaaa$ pojawia się w słowie $aaaaaaaaaaaaaaaaaaaa$ $n = m + 1 = 11$ razy). Sama ilość takich (różnych) struktur w słowie może służyć za miarę jego regularności lub skomplikowania.

W pracach wchodzących w skład doktoratu wraz z współautorami zajęliśmy się kilkoma klasami podów regularnych lub takich, które można przedstawić w zwarty sposób. W każdej z prac przedstawiliśmy nowe, efektywne algorytmy wypisywania i zliczania wszystkich wystąpień, lub wszystkich różnorodnych podów pewnej klasy.

2.1.2 Skład pracy

Rozprawa składa się z czterech artykułów, napisanych podczas moich studiów doktorskich na Wydziale Matematyki, Informatyki i Mechaniki Uniwersytetu Warszawskiego we współpracy z innymi autorami:

1. Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszyński, Tomasz Walczak, Wiktor Zuba: Efficient Representation and Counting of High Power Factors in Words. LATA 2019: 421-433

W wersji przyjętej do druku w Special Issue konferencji LATA 2019, które ma być wydane w czasopiśmie Information & Computation.

2. Panagiotis Charalampopoulos, Tomasz Kociumaka, Manal Mohamed, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszyński, Tomasz Walczak, Wiktor Zuba: Counting Distinct Patterns in Internal Dictionary Matching. CPM 2020: 8:1-8:15
3. Panagiotis Charalampopoulos, Jakub Radoszewski, Wojciech Rytter, Tomasz Walczak, Wiktor Zuba: The Number of Repetitions in 2D-Strings. ESA 2020: 32:1-32:18
4. Panagiotis Charalampopoulos, Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Tomasz Walczak, Wiktor Zuba: Efficient Enumeration of Distinct Factors Using Package Representations. SPIRE 2020: 247-261

Wszystkie prace zawierają oryginalne wyniki z algorytmiki na słowach. Prace zostały napisane przez autorów reprezentujących Uniwersytet Warszawski we współpracy z naukowcami z King's College London (Panagiotis Charalampopoulos przebywał przez dłuższy czas w Warszawie) oraz Tomaszem Kociumaką, który w międzyczasie zmienił afiliację z Uniwersytetu Warszawskiego na Uniwersytet Bar-Ilan.

W sekcji 2.2 niniejszego autoreferatu przedstawiam podstawowe pojęcia - definicje regularnych klas słów i podsłów, które są wykorzystywane w przedstawianych pracach.

W sekcji 2.3 pokazuję dokładny przegląd wyników zakończonych prac oraz skrócone metody użyte do ich uzyskania.

W sekcji 2.4 wymieniam publikacje których jestem współautorem i które zostały opublikowane podczas moich studiów doktoranckich, które jednak nie zostały włączone w skład tej rozprawy.

W rozdziale 1 zamieszczone jest to samo streszczenie w języku angielskim, zaś w rozdziałach 3, 4, 5 oraz 6 kolejne prace wchodzące w skład rozprawy. Bibliografia na końcu zawiera referencje użyte w niniejszym streszczeniu.

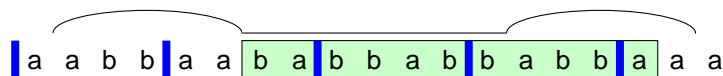
2.2 Preliminaria

Aby przedstawić wyniki zakończonych prac należy zacząć od przedstawienia typów regularności w słowach, przez nie wykorzystywanych. Wiele z tych pojęć pojawia się w więcej niż jednej z tych prac.

Definicja 2.1. Niech $x = y_0y_1 \dots y_{k-1}$ dla $k \geq 2$ i słów y_i o tej samej długości d . Wtedy mówimy, że:

• x jest k -potęgą, jeśli wszystkie y_i są takie same (x jest kwadratem, jeśli dodatkowo $k = 2$);

• x jest k -antypotęgą, jeśli wszystkie y_i są parami różne;



Rysunek 2.1: Regularności w s^aowach jednowymiarowych. Zielonym t^aem zaznaczono maksymalne powtórzenie (o okresie 3). Niebieskie linie pokazują wystąpienie 4-antypotygi. Zauważmy, że jedną i dwie pozycje dalej nie występują 4-antypotygi o tej samej podstawie ponieważ te wyrazy długości 16 są przerywanymi kwadratami generowanymi przez zaznaczone nad s^aowem maksymalne 3-przerywane wystąpienie.

^ x jest s^aab_i k-potygi jeśli nie jest k-antypotygi, to jest jeśli $y_i = y_j$ dla pewnych $i \in j$;

^ x jest przerywanym kwadratem jeśli $y_0 = y_{k-1}$.

Definicja 2.2. Mówimy, że:

^ p jest okresem s^aowaw jeśli $w[i] = w[i + p]$ dla wszystkich $i \in [0 : j - p]$ (przez okres s^aowa najczęściej rozumiemy jego najkrótszy okres);

^ fragment s^aowa $w[i : j]$ jest maksymalnym powtórzeniem jeśli jego długość jest równa co najmniej dwóm jego najkrótszym okresom i nie można go rozszerzyć o jedną pozycję w lewo ani w prawo bez zmiany tego okresu;

^ fragment s^aowaw $[i : j]$ jest maksymalnym -przerywanym powtórzeniem (dla $k > 1$) jeśli jest postaci uv^k dla $p = j - i + 1$ i nie można go rozszerzyć o jedną pozycję w lewo ani w prawo bez zmiany okresu.

Definicja 2.3. Dla tekstów dwuwymiarowych definiujemy analogiczne powtórzenia.

^ Tandemem nazywamy dwuwymiarowy tekst o parzystej długości, w którym lewa i prawa połowa tekstu są takie same (ogólnienie kwadratu na przypadek gdzie pojedynczy liter jest cała kolumna).

^ Kwartyk nazywamy dwuwymiarowy tekst K o parzystej długości i szerokości, taki że zarówno K jak i K^T są tandemami (cztery identyczne teksty ułożone w kratki 2×2).

^ maksymalnym 2D-powtórzeniem nazywamy fragment tekstu dwuwymiarowego o wymiarach $i \times j$, taki że jego okres poziomy jest wielkością co najwyżej $\frac{i}{2}$ a pionowy co najwyżej $\frac{j}{2}$. Co więcej powiększenie tego fragmentu o dodatkową kolumnę lub wiersz z dowolnej strony (o ile takie istnieją w całym tekście) skutkowałyby zmianą tych okresów.

Przez kwadrat generowany przez maksymalne powtórzenie rozumiem kwadrat zawarty w tym powtórzeniu o długości równej dokładnie dwóm okresom powtórzenia.

a	c	a	a	b	a	a	b	a	a	d	c
a	b	c	a	b	c	a	b	c	a	b	c
a	b	a	a	b	a	a	b	a	a	b	a
a	b	c	a	b	c	a	b	c	a	b	c
a	d	b	a	c	d	b	a	c	a	b	d

Rysunek 2.2: Różne rodzaje powtórzeń w tekstach dwuwymiarowych. Kolorem brązowym zaznaczono tandem, niebieskim kwartyk, zaś zielonym tłem maksymalne 2D-powtórzenie.

Wszystkie takie kwadraty można łatwo wyznaczyć z powtórzenia i każdy kwadrat jest generowany przez jakieś maksymalne powtórzenie. Analogicznie wygląda sytuacja generowania przerywanych kwadratów przez powtórzenia i -przerywane powtórzenia oraz generowania kwartyk przez 2D-powtórzenia.

Kwadrat ma pierwiastek pierwotny gdy jego pierwiastek (potłg) nie jest k-potłg dla żadnego $k > 1$. Analogicznie działa to w przypadku kwartyk, dla których jego pierwiastek (kwartyk) nie jest k-potłg dla $k > 1$ ani poziomo ani pionowo.

Definicja 2.4. Słownikiem wewnętrznym nazywamy zbiór pozycji początkowych i końcowych wybranych podśłów zadanego słowa.

2.3 Przegląd wyników

2.3.1 Efficient Representation and Counting of Antipower Factors in Words

Zaprezentowana tutaj praca została zaprezentowana na konferencji LATA 2019 i opublikowana w sprawozdaniu z konferencji. Rozszerzona wersja tej pracy została później przyjęta do druku w Special Issue tej konferencji, wydawanym w czasopiśmie Information & Computation. Praca w tej formie (przyjętej do druku, lecz jeszcze nie opublikowanej) znajduje się do niniejszej rozprawy (w rozdziale 3).

Wstęp

W pracy [5] przedstawione zostało górne i dolne ograniczenie $\pi(n^2=k)$ na maksymalną liczbę wystąpień k-antypotłg w słowie. Praca ta podaje również algorytm wyznaczający te wystąpienia w czasie $O(n^2=k)$. Ze względu na ograniczenia na wielkość wyjścia algorytm ten jest optymalny w przypadku pesymistycznym. Nie oznacza to jednak, że nic w tej dziedzinie nie można poprawić.

W naszej publikacji skupiliśmy się na problemach które nie podlegają temu ograniczeniu - zliczaniu antypotłg oraz ich wyznaczaniu ze złożonością zależną od rozmiaru wyjścia (efektywniejszy gdy rozważane słowo zawiera mało antypotłg).

Piszyc dokladniej - zaprezentowalimy algorytmy, które dla sowa dugości zliczaj wszystkie podsowa bidge k-antypotłgami w czasie $O(nk \log k)$, wyznaczaj wszystkie te podsowa w czasie $O(nk \log k + \text{rozmiar wyniku})$, oraz wyznaczaj liczb! rónych k-antypotłg w sowie w czasie $O(nk^4 \log k \log n)$. Zaprezentowalimy równie» nowi struktur! danych pozwalajici na szybkie sprawdzenie, czy dane podsowo jesk-antypotłgi.

Jak widaç w (najcz!ciej wykorzystywanym) przypadku gdy rozwa»anç jest ma»e w stosunku do dugości sowa (mniejsze ni»n) nasz algorytm uzyskuje lepszy asymptotyczny czas dziaania od poprzednio znanych.

Sabe potłgi i zwarta reprezentacja

Ze wzgl!du na nieregularni natur! antypotłg efektywne ich znajdowanie nie jest proste. Dlatego te», skupilimy si! na wyznaczeniu reprezentacji tych podsów, które nimi nie s!, czyli sabychk-potłg. Jednocze»nie ka»da sabk-potłga jest konsekwencji wyst!pienia przerywanego kwadratu (byç mo»e dla mniejszegok, lecz tego samego).

Obserwacja 2.1.

- Ka»dy (przerywany) kwadrat jest generowany przez maksymalne powtórzenie o (nie koniecznie najkrótszym) okresie $q + 1$ d lub maksymalne $(q + 1)$ -przerywane powtórzenie o okresie $q + 1$ d dla $q = k - 2$.
- Ka»de maksymalne powtórzenie i ka»de przerywane powtórzenie generuje pojedynczy przedzia» pozycji w którym pojawiaj! si! (przerywane) kwadraty okre»lonej dugości, co wi»cej mo»emy ten przedzia» wyznaczyç w czasie sta»ym.
- Wszystkich ogólnych maksymalnych powtórze» w sowie o dugości jest $O(n)$ i ich reprezentacji mo»na wyliczyç w czasie $O(n)$ ([6]).
- Wszystkich maksymalnych -przerywanych powtórze» w sowie jes $O(n)$ i ich reprezentacji mo»na wyliczyç w czasie $O(n)$ ([11, 12]).

Korzystajic z tych w»asno»ci potra my wyliczyç reprezentacji sabychk-potłg dla wszystkich d w postaci »ycznie $O(nk)$ »a»uchów przedziaów (zbiór przedziaów o tej samej dugości i pocz!tkach tworzicych postłp arytmetyczny). Nast!pnie korzystajic z metod geometrycznych potra my wyznaczyç sum! teoriiomnogo»ciow! tych »a»uchów w »icznym czasie $O(nk \log k)$. Daje nam to bezpo»rednio ilo»ç sabychk-potłg dla ka»dego d, a wi»c i ilo»ç k-antypotłg. Zamiast liczyç wielko»ç zbioru pozycji pokrytych przez reprezentacji wyst!pie» sabychk-potłg mo»emy wyznaczyç wszystkie pozycje niepokryte i otrzymaç algorytm wyznaczajicy k-antypotłgi w czasie $O(nk \log k + \text{wynik})$.

Pytania o podsowa

W pracy [5] zaprezentowano równie» dwie struktury danych które po zbudowaniu pozwalaj! odpowiadaç na pytania "Czy dane podsowow[i::j] jest k-antypotłgi?". Pierwsza struktura o rozmiarze $O(n)$ odpowiada na pytania

w czasie $O(k)$, zaś druga robi to w czasie stałym, lecz wymaga zapamiętania informacji rozmiaru $O(n^2)$.

W naszej pracy przedstawiliśmy nową strukturę, parametryzowaną przez wartość $r \geq 2$ [1; n]. Po zbudowaniu w czasie $O(n^2/r)$ struktura ma rozmiar $O(n^2/r)$ i odpowiada na pytania w czasie $O(r)$.

Jeżeli $r < k$, to wiemy, że podstawa potęg d wynosi co najwyżej $n/r > n/k$. Aby odpowiadać na pytania dla $k \geq d \rightarrow n/r$ osobno budujemy strukturę danych opartą na range minimum queries, o rozmiarze $O(n)$ i czasie zapytania $O(1)$. W komplementarnym przypadku aby uzyskać pożądaną rezultat wystarczy skorzystać ze wspomnianej struktury z pracy [5].

Zliczanie różnych k-antypotęg

W wersji do czasopisma dodaliśmy kolejny wynik - zliczanie różnych k-antypotęg w słowie.

Ze względu na potencjalną dużą ilość wszystkich występień antypotęg metody wykorzystujące ich wyznaczenie nie pozwalają na osiągnięcie satysfakcjonującej szybkości. Aby uzyskać taki wynik ponownie skorzystaliśmy z powiązania antypotęg z bardziej regularnymi słabymi potęgami.

Liczba różnych k-antypotęg możemy uzyskać odejmując od siebie dwie inne wartości. Pierwszą wartością jest ilość (wszystkich) różnych podciągów o długości podzielnej przez k . Możemy ją łatwo uzyskać w czasie $O(n)$ przy użyciu drzewa suksowego.

Drugą potrzebną wartością jest liczba różnych słabych k-potęg. Aby ją uzyskać najpierw wzmocniliśmy definicję słabych potęg, tak aby każda była generowana na dokładnie jeden sposób. Następnie przez redukcję do problemu grafowego udało nam się wykonać liczenie różnych takich konstrukcji w łącznym czasie $O(nk^4 \log k \log n)$.

2.3.2 Counting Distinct Patterns in Internal Dictionary Matching

W pracy [8] moi współautorzy wprowadzili problem wyszukiwania przy użyciu słownika wewnętrznego - dla zadanego słowa a i jego słownika wewnętrznego D przedstawili konstrukcję struktury danych pozwalającą efektywnie odpowiadać na pytania dla danego podciągu $T[i:j]$:

- ^ czy $T[i:j]$ zawiera jakieś słowo ze słownika?
- ^ zwróć wszystkie podciągi słowa a $T[i:j]$ które należą do słownika
- ^ zwróć wszystkie słowa ze słownika D które są zawarte w $T[i:j]$
- ^ zwróć liczbę podciągów $T[i:j]$ które należą do słownika
- ^ zwróć liczbę różnych słów ze słownika D które są zawarte w $T[i:j]$

Przedstawiona struktura danych ma rozmiar i czas konstrukcji $O(n+d)$ (gdzie $n = \sum |T_j|$ oznacza długość słowa $a = \sum |D_j|$ liczbę słów w słowniku, zaś O notacji asymptotycznej pomijając czynniki polilogarytmiczne) i odpowiadając na pierwsze cztery pytania w czasie $O(1 + \text{wynik})$.

W przypadku ostatniego pytania tylko $O(\log n)$ aproksymacja wyniku została zaprezentowana.

Rozmiar	Czas konstrukcji	Czas zapytania	Wariant
$O(n + d)$	$O(n + d)$	$O(1)$	2-aproksymacja
$O(n^2 = m^2 + d)$	$O(n^2 = m + d)$	$O(m)$	dokładny
$O(nd = m + d)$	$O(nd = m + d)$	$O(m)$	dokładny
$O(n \log^2 n)$	$O(n \log^2 n)$	$O(\log n)$	D = kwadraty, dokładny

Tabela 2.1: Nasze rezultaty dla zapytania «policz różnicę 2 [1; n] jest dowolnie wybranym parametrem).

2-aproksymacja

Aby poprawić wynik z poprzedniej pracy stworzyliśmy strukturę danych, która pamięta wyniki dla podciągów bazowych, czyli wszystkich podciągów długości $b(1 + \epsilon)^p$ dla wszystkich liczb naturalnych p i ustalonego $\epsilon = \frac{1}{9}$. Aby otrzymać wynik dla dowolnego podciagu $T[i:j]$ korzystamy z maksymalnych sów bazowych postaci $T[i:i^0]$ oraz $T[j^0:j]$ zawartych w $T[i:j]$ i otrzymujemy jego podział na trzy części $F_1 F_2 F_3$ ($F_1 F_2 = T[i:i^0]$, $F_2 F_3 = T[j^0:j]$).

Wynik uzyskujemy poprzez zliczenie sów ze słownika, które mają wystąpienie zaczynające się w F_1 i kończące w F_3 , jednak nie występujące w żadnym z wymienionych sów bazowych, co udaje nam się zrobić efektywnie dzięki dużemu stosunkowi długości F_2 do F_1 i F_3 . Do tego dodajemy ilość sów występujących w obu słowach bazowych, co jednak skutkuje uzyskaniem jedynie 2-aproksymacji ze względu na możliwe powtórzenia.

Wszystkich sów bazowych jest $O(n \log n)$, zaś wartości dla nich możemy policzyć w łącznym czasie $O(n \log^{1+\epsilon} n + d)$ (dla dowolnego stałego $\epsilon > 0$) wykorzystując to, że wynik po przesunięciu o jedną pozycję można zaktualizować niewielkim kosztem. Struktura danych używana do zliczania sów opisanych w poprzednim akapicie kosztuje nas dodatkowo $O(n \log n = \log \log n + d \log^{3-\epsilon} n)$ czasu preprocessingu oraz zajmuje $O(n + d \log n)$ pamięci. Pozwala ona odpowiadać na pytania w czasie $O(\log^2 n = \log \log n)$.

Obliczanie dokładne

Jeśli zamiast dla sów bazowych zapamiętamy wyniki dla wszystkich podciągów postaci $T[c_1 m + 1 : c_2 m]$ dla wybranego m oraz wszystkich $c_1 < c_2$ otrzymamy strukturę danych o rozmiarze $O(n^2 = m^2 + n + d)$ obliczaną w czasie $O((n^2 \log n) = m + n \log n + d)$. Rozszerzając takie słowo po jednej pozycji możemy uzyskać wynik dla dowolnego słowa przy użyciu $O(m)$ operacji o koszcie $O(\log n)$.

Innym podejściem do problemu jest zagłębienie się w strukturę słownika.

Jeśli słownik nie zawiera zbiorów różnicowych preksów tego samego słowa, to na dowolnej pozycji możemy się zacząć co najwyżej od 1 różnicowego słowa z tego

słownika. W takim przypadku możemy dla każdego słowa ze słownika zapamiętać wszystkie pozycje na których ono występuje. Daje nam to strukturę danych wielkości $O(nk \log n)$, obliczalną w czasie $O(nk \log n)$ i pozwalającą odpowiadać na pytania w czasie $O(\log n)$.

Aby uzyskać słownik w takiej postaci wydzielamy z niego (jako nowe słowniki) słowa niespełniające tej własności, to jest duże grupy słów bliźnich prefiksami tego samego słowa. Takich grup (o wielkości co najmniej k) jest co najwyżej d/k , zaś dla każdej takiej grupy potrafiemy efektywnie odpowiadać na pytania przy pomocy struktury do pytań o ograniczone LCP ([14]) o rozmiarze $O(n \log n)$ w czasie $O(\log n)$ dla dowolnego $d > 0$.

Dla $k = d/m$ daje nam to w sumie czas zapytania wielkości $O(m \log n + \log n)$.

Słownik kwadratów

Dla popularnego przypadku obliczania liczby kwadratów w podśowie opisaliśmy odrębny algorytm. Dzięki szczególnym własnościom takich podśów możemy otrzymać lepsze wyniki.

Wszystkich występowanie kwadratów w śowie może być $O(n^2)$ (różnych kwadratów tylko $O(n)$), jednak dzięki skorzystaniu z maksymalnych powtórzeń, których w całym śowie jest $O(n)$ (i które można obliczyć w czasie $O(n)$) potrafiemy wyznaczyć $m = O(n \log n)$ istotnych występowanie kwadratów, które są wystarczająco reprezentatywne dla naszych potrzeb. Dzięki skorzystaniu z geometrycznej struktury danych o rozmiarze i czasie konstrukcji $O(m \log m)$ z pracy [13] potrafiemy odpowiedzieć na pytanie o liczbę różnych takich śów w podśowie w czasie $O(\log m) = O(\log n)$.

2.3.3 The Number of Repetitions in 2D-Strings

W pracy zamieszczonej w rozdziale 5 zajęliśmy się problemem powtórzeń w tekstach dwuwymiarowych. W tabeli 2.2 przedstawiam podsumowanie naszych wyników oraz dotychczasowej wiedzy na temat możliwości maksymalnych 2D-powtórzeń jak i maksymalnej liczby różnych tandemów i kwartyków oraz wszystkich występowanie tych o pierwotnym pierwiastku. Przedstawiam również czasy algorytmów pozwalających wyznaczyć te struktury w tekście. Warto zauważyć, że ilość 2D-powtórzeń była już wcześniej badana, jednak granica pomiędzy dolnym i górnym ograniczeniem pozostawała liniowa. W przedstawianej pracy udało nam się zredukować tę wielkość do polilogarytmicznej.

maksymalne 2D-powtórzenia

Aby uzyskać ograniczenie na liczbę maksymalnych 2D-powtórzeń skorzystaliśmy z lematu z pracy [1] aby przypisać każde takie powtórzenie do maksymalnego powtórzenia poziomego (brak wymagania pionowej okresowości i maksymalności) o wysokości bliźniej potęg dwójki (największy możliwy, nie większy niż wysokość 2D-powtórzenia). Takie powtórzenie poziome spełnia również trzy istotne warunki - jego lewy górny lub lewy dolny narożnik jest równy odpowiadającemu narożnikowi 2D-powtórzenia, jego okres jest równy okresowi poziomemu

	Ograniczenia na ilość w tekście n	Czas wyznaczania
maksymalne 2D-powtórzenia	$(n^2), O(n^3)$ [2] $O(n^2 \log^2 n)$	$O(n^2 \log n + \text{wynik})$ [2] $O(n^2 \log^2 n)$
Wystąpienia kwartyk o pierwotnych pierwiastkach	$(n^2 \log^2 n)$ [3]	$O(n^2 \log^2 n)$
Różne kwartyki	$O(n^2 \log^2 n)$	$O(n^2 \log^2 n)$
Wystąpienia tandemów o pierwotnych pierwiastkach	$(n^3 \log n)$ [3]	$O(n^3 \log n)$ [4]
Różne tandemy	$(n^2); O(n^4)$ (trywialne) (n^3)	$O(n^3)$

Tabela 2.2: Przegląd wcześniejszych wyników i naszych rezultatów (pogrubione).



Rysunek 2.3: Przypisanie maksymalnemu 2D-powtórzeniu maksymalnego powtórzenia poziomego o wysokości k^2

powiązanego 2D-powtórzenia oraz jego długość jest co najmniej taka jak długość 2D-powtórzenia (rysunek 2.3).

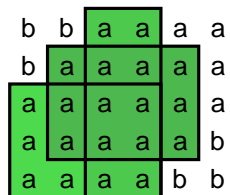
Korzystając z tych własności, lematu o okresowości oraz lematu o trzech kwadratach dowiedzemy, że wszystkim poziomym powtórzeniom okupującym wiersze od $i + 2^k - 1$ może zostać przypisane $\Theta(n \log n)$ różnych maksymalnych 2D-powtórzeń. Mnożąc tę wartość, przez ilość możliwych wyborów i oraz k otrzymujemy ograniczenie $O(n^2 \log^2 n)$ na liczbę maksymalnych 2D-powtórzeń w tekście rozmiaru n .

Nasze ograniczenie na ilość maksymalnych 2D-powtórzeń automatycznie daje nowe ograniczenie na czas działania algorytmu z pracy [2] (algorytm ten działa w czasie zależnym od rozmiaru wyniku).

Kwartyki o pierwotnych pierwiastkach

W przypadku jednowymiarowym wszystkie kwadraty o pierwotnych pierwiastkach można bardzo łatwo wyznaczyć znając maksymalne powtórzenia w słowie - wystarczy wziąć wszystkie słowa o długości dwóch okresów maksymalnego powtórzenia, które są w nim zawarte.

W przypadku dwuwymiarowym jest podobnie - każda kwartyka o pierwotnym pierwiastku jest prostokątem o długości dwóch okresów poziomych i wysokości dwóch okresów pionowych pewnego maksymalnego 2D-powtórzenia, całkowicie w nim zawartym. Pojawia się jednak pewien problem - o ile w przypadku jednowymiarowym dwa maksymalne powtórzenia o tym samym okresie nie mogą



Rysunek 2.4: Wiele przecinających sił maksymalnych 2D-powtórze« o tych samych okresach

mieć du»ego przecięcia i w szczególności nie mogą generować tego samego kwadratu, o tyle w 2D może być bardzo łatwo zachodzić (rysunek 2.4), przez to zwykłe u»ycie takiego algorytmu mogłoby skutkować wielokrotnym zwracaniem tych samych kwartyk.

Rozwiązaniem tego problemu jest metoda zamiatania - algorytm z pracy [7], który pozwala na liczenie sum teoriiomnogościowych kilku różnych rodzin prostokątów na kracie $n \times n$ w czasie $O(n + r + \text{wynik})$, gdzie r to liczba wszystkich prostokątów.

Dzieląc wszystkie maksymalne 2D-powtórzenia na grupy o takich samych okresach i dla każdej takiej biorąc prostokąty wyznaczające lewe górne rogi generowanych kwartyk otrzymujemy po»jdany algorytm.

Dzięki naszemu ograniczeniu na liczbę r z poprzedniej podsekcji oraz ograniczeniu $O(n^2 \log^2 n)$ na liczbę kwartyk o pierwotnym pierwiastku z pracy [3] algorytm działa w czasie $O(n^2 \log^2 n)$ (optymalnym ze względu na dolne ograniczenie z tej samej pracy).

Różne kwartyki

W tekście rozmiar n z»ony z samych litera zawarte jest (n^4) występie« dowolnych kwartyk, przez co trywialny algorytm znajduje wszystkie takie występienia w optymalnym pesymistycznym czasie. Dlatego te» ciekawszym problemem jest wyznaczenie różnych kwartyk (fragmenty różniące sił nie tylko po»o»eniem, ale te» jako teksty).

Wynik z pracy [3] daje ograniczenie na liczbę kwartyk o pierwiastku pierwotnym. Inne kwartyki z»one są z wielu przylegających do siebie występie« takich samych pierwotnych kwartyk tworzących większy prostokąt. Innymi słowy o ile w przypadku tych poprzednich kwartyk występienia pierwotnego tekstu W tworzą k krat 2×2 , to w przypadku tych pozostałych tworzą $2k \times 2l$, gdzie $k > 1$ lub $l > 1$. Kwartyki takie dzielimy na wąskie, czyli takie dla których $k = 1$ lub $l = 1$ i szerokie dla których $k, l > 1$.

Następnie dokonujemy kolejnego podziału - kwartyki dla których słowo pierwotne W ma rozmiary w przedziałach $[2^a; 2^{a+1})$ $[2^b; 2^{b+1})$ dla pewnych a, b - logn rozważamy razem.

Ze względu na lematy o okresowości oraz o trzech kwadratach każdy punkt może być lewym górnym rogiem ostatniego występienia (w tym samym wierszu, bardziej na prawo w tekście nie występuje taka sama kwartyka) co najwy»ej

czterech różnych kwartyk wjskich z jednej grupy (analogicznie dziaa dowód na ograniczenie liczby różnych kwadratów w standardowym tekście [10]).

W przypadku kwartyk szerokich dzięki mnogości punktów kratowych w których zaczyna się wystąpienie sowa W każdej takiej kwartyce bldicej potłgi W możemy przypisać unikalny punkt na tej kracie. Dodatkowo dzięki lematowi o okresowości taki punkt nie może nalee do kraty dla kwartyki bldicej potłgi innego W^0 o rozmiarze naleicym do $[2; 2^{a+1})$ $[2^b; 2^{b+1})$.

Dzięki przypisywaniu punktów z $[1; n]$ $[1; n]$ do kwartyk oraz $\log^2 n$ wyborom $a; b$ otrzymujemy ograniczenia $O(n^2 \log^2 n)$ na ilość różnych kwartyk kaędego z typów.

W celu wyznaczenia wszystkich unikalnych kwartyk w tekście korzystamy z naszego wcześniejszego algorytmu by wyznaczyć wszystkie wystąpienia kwartyk o pierwiastkach pierwotnych i grupujemy je według pierwiastków.

Kwartyki bldice potłgami sowa pierwotnego W powstają z połączenia wielu kwartyk pierwotnych o tym samym pierwiastku. Możemy więc wyznaczyć lewe górne rogi wszystkich wystpie takich kwartyk i wyszukać maksymalne kraty utworzone przez takie punkty.

Problem ten rozwijujemy uęywając ponownie metody zmiatania by uzyskać algorytm działający w czasie $O(n^2 \log^2 n)$.

Róne tandemy

Jako powizywanym problemem zajliemy się z wyznaczeniem różnych tandemów (podobnie jak w przypadku kwartyk wszystkich wystpie tandemów może być (n^4)).

W wypadku tego problemu można aatwo utworzyć tandemy zaczynające się w wierszu i oraz kończące w wierszu j z kwadratami w standardowym sowie (jednowymiarowym), poprzez przypisanie kolumnom wysokość $j - i + 1$ standardowych liter. W standardowym sowie długości n może być (n) różnych kwadratów, stąd te same ograniczenie obowiązuje tandemy okupujące wyznaczone wiersze. Mnożąc tę wartość przez liczbę wyborów wartości i oraz j otrzymujemy ograniczenie $O(n^3)$ na liczbę różnych tandemów.

Dolne ograniczenie (n^3) na maksymalną liczbę takich tandemów otrzymujemy równie prosto - otrzymujemy je na przykład dla tekstu w którym każdy wiersz jest sowem nad alfabetem jednoliterowym, dla każdego wiersza innym.

Przypisanie kolumnom wysokość $j - i + 1$ standardowych liter można zrobić konstrukcyjnie. Dzięki temu i dzięki skorzystaniu z algorytmów wyszukujących różne kwadraty w standardowym sowie otrzymujemy algorytm wyznaczający różne tandemy w sowie dwuwymiarowym w czasie $O(n^3)$. Ze względu na dolne ograniczenie maksymalnego rozmiaru wyniku ten bardziej skomplikowany algorytm nie pozwala na uzyskanie lepszego rezultatu w pesymistycznym przypadku.

2.3.4 Efficient Enumeration of Distinct Factors Using Package Representations

W ostatniej z opisywanych prac zajmiemy się nową reprezentacją podzbiorów podzbiórów która w wielu przypadkach zapewnia zwięzłą reprezentację, która jednocześnie pozwala na efektywne operowanie opisywanym podzbiorem.

Przez pakiet rozumiemy zbiór podzbiórów tej samej długości występujących na kolejnych pozycjach w s-owie. Nasza reprezentacja składa się ze zbioru takich pakietów reprezentowanych jako trójki (pierwsza pozycja początkowa, długość s-ów, długość przedziału). Rzecz, która interesuje nas szczególnie jest zbiór wszystkich podzbiórów które należą do przynajmniej jednego z tych pakietów.

Bardziej formalnie

$$\text{Factor}(F) = \{T[j:j+k] : j \geq 1, i \leq j+k \text{ and } (i,l;k) \in F\}$$

Przykładami rodzin podzbiórów, dla których łatwo można uzyskać taką reprezentację są potęgi (kiedy to kiedyś pakiet reprezentuje bezpośrednio jedno maksymalne powtórzenie) oraz k-antypotęgi (gdzie możemy wykorzystać bezpośrednio ich reprezentację z pracy opisaną w sekcji 2.3.1). Reprezentacje te mają rozmiar (i czas obliczania) równy odpowiednio $O(n)$ i $O(nk^2)$ (kiedyś pakiet przedziałów reprezentuje co najwyżej zwykłych przedziałów).

W celu znalezienia zbioru $\text{Factor}(F)$, czyli różnych podzbiórów należących do reprezentacji oraz $|\text{Factor}(F)|$, czyli ilości takich s-ów rozpatrujemy dwa przypadki. Pierwszy prostszy przypadek, nazywany w pracy specjalnym, to taki w którym jeżeli podzbiórowo należy do $\text{Factor}(F)$, to kiedyś jego wystąpienie w rozwinięciu s-owie należy do któregoś z pakietów należących do

W tym przypadku dla s-owa długości n i F składającej się z m pakietów korzystając ze struktury danych najdłuższych poprzednich podzbiórów (longest previous factors [9]) otrzymujemy wyznaczenie i zliczenie wyniku w czasie odpowiednio $O(n + m + \text{wynik})$ i $O(n + m)$. W przypadku przedstawionych przykładów daje nam to nowy, prosty algorytm wyznaczania różnych potęg w s-owie, działający w czasie $O(n)$, oraz nowy algorytm wyznaczania różnych k-antypotęg działający w czasie $O(nk^2 + \text{wynik})$, a więc lepszy od naszego algorytmu z pracy opisaną w sekcji 2.3.1, który zwraca ilość różnych k-antypotęg, działając w czasie $O(nk^4 \log k \log n)$ i jest dużo bardziej skomplikowany.

Drugi rozwinięty przypadek, nazywany w pracy ogólnym nie wymaga tego dodatkowego założenia, co jednak czyni go istotnie trudniejszym do rozwiązania. Aby go rozwiązać adaptujemy metodę użytą poprzednio w pracy o k-antypotęgach (zliczanie różnych k-antypotęg) i uzyskujemy algorytm wyznaczania i zliczania wyniku $\text{Factor}(F)$, działający w czasie odpowiednio $O(n \log^2 n + m \log n + \text{wynik})$ i $O(n \log^2 n + m \log n)$.

2.4 Pozostałe publikacje

Podczas moich studiów doktorskich zostały opublikowane również inne prace których jestem współautorem. Zdecydowałem się nie dołączać ich do niniejszej

rozprawy ze względu na niedostateczne dopasowanie do jej tematu.

1. Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Faster Recovery of Approximate Periods over Edit Distance. SPIRE 2018: 233-240
2. Wojciech Rytter, Wiktor Zuba: Syntactic View of Sigma-Tau Generation of Permutations. LATA 2019: 447-459
Rozszerzona wersja pracy została przyjęta do druku w czasopiśmie Theoretical Computer Science.
3. Mai Alzamel, Maxime Crochemore, Costas S. Iliopoulos, Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Quasi-Linear-Time Algorithm for Longest Common Circular Factor. CPM 2019: 25:1-25:14
4. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Circular Pattern Matching with k Mismatches. FCT 2019: 213-228
5. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Weighted Shortest Common Supersequence Problem Revisited. SPIRE 2019: 221-238
6. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Shortest Covers of All Cyclic Shifts of a String. WALCOM 2020: 69-80
7. Panagiotis Charalampopoulos, Solon P. Pissis, Jakub Radoszewski, Tomasz Walen, Wiktor Zuba: Unary Words Have the Smallest Levenshtein k -Neighbourhoods. CPM 2020: 10:1-10:12
8. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Internal Quasiperiod Queries. SPIRE 2020: 60-75
9. Maxime Crochemore, Costas S. Iliopoulos, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Shortest covers of all cyclic shifts of a string. Theor. Comput. Sci. 866: 70-81 (2021)
10. Panagiotis Charalampopoulos, Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, Juliusz Straszynski, Tomasz Walen, Wiktor Zuba: Circular pattern matching with k mismatches. J. Comput. Syst. Sci. 115: 73-85 (2021)

Chapter 3

Efficient Representation and Counting of Antipower Factors in Words

Chapter 4

Counting Distinct Patterns in Internal Dictionary Matching

Chapter 5

The Number of Repetitions in 2D-Strings

Chapter 6

Efficient Enumeration of Distinct Factors Using Package Representations

Bibliography

- [1] Amihood Amir, Gad M. Landau, Shoshana Marcus, and Dina Sokol. Two-dimensional maximal repetitions. In 26th Annual European Symposium on Algorithms, ESA 2018, volume 112 of LIPIcs, pages 2:1{2:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018. doi:10.4230/LIPIcs.ESA.2018.2.
- [2] Amihood Amir, Gad M. Landau, Shoshana Marcus, and Dina Sokol. Two-dimensional maximal repetitions. *Theoretical Computer Science* 812:49{61, 2020. doi:10.1016/j.tcs.2019.07.006 .
- [3] Alberto Apostolico and Valentin E. Brimkov. Fibonacci arrays and their two-dimensional repetitions. *Theoretical Computer Science* 237(1-2):263{273, 2000. doi:10.1016/S0304-3975(98)00182-0 .
- [4] Alberto Apostolico and Valentin E. Brimkov. Optimal discovery of repetitions in 2D. *Discrete Applied Mathematics* 151(1-3):5{20, 2005. doi:10.1016/j.dam.2005.02.019 .
- [5] Golnaz Badkobeh, Gabriele Fici, and Simon J. Puglisi. Algorithms for anti-powers in strings. *Information Processing Letters*, 137:57{60, 2018. doi:10.1016/j.ipl.2018.05.003 .
- [6] Hideo Bannai, Tomohiro I, Shunsuke Inenaga, Yuto Nakashima, Masayuki Takeda, and Kazuya Tsuruta. The "runs" theorem. *SIAM Journal on Computing*, 46(5):1501{1514, 2017. doi:10.1137/15M1011032 .
- [7] Jon Louis Bentley. Algorithms for Klee's rectangle problems. Unpublished notes, Computer Science Department, Carnegie Mellon University, 1977.
- [8] Panagiotis Charalampopoulos, Tomasz Kociumaka, Manal Mohamed, Jakub Radoszewski, Wojciech Rytter, and Tomasz Walenczyk. Internal dictionary matching. In 30th International Symposium on Algorithms and Computation, ISAAC 2019, volume 149 of LIPIcs, pages 22:1{22:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019. arXiv:1909.11577 , doi:10.4230/LIPIcs.ISAAC.2019.22 .

- [9] Maxime Crochemore and Lucian Ilie. Computing longest previous factor in linear time and applications. *Inf. Process. Lett.*, 106(2):75–80, 2008. doi : 10.1016/j.ipl.2007.10.006.
- [10] Aviezri S. Fraenkel and Jamie Simpson. How many squares can a string contain? *J. Comb. Theory, Ser. A*, 82(1):112–120, 1998. doi : 10.1006/jcta.1997.2843.
- [11] Paweł Gawrychowski, Tomohiro I, Shunsuke Inenaga, Dominik Köppl, and Florin Manea. Tighter bounds and optimal algorithms for all maximal k -gapped repeats and palindromes - finding all maximal k -gapped repeats and palindromes in optimal worst case time on integer alphabets. *Theory of Computing Systems*, 62(1):162–191, 2018. doi : 10.1007/s00224-017-9794-5.
- [12] Tomohiro I and Dominik Köppl. Improved upper bounds on all maximal k -gapped repeats and palindromes. *Theoretical Computer Science*, 753:1–15, 2019. doi : 10.1016/j.tcs.2018.06.033.
- [13] Haim Kaplan, Natan Rubin, Micha Sharir, and Elad Verbin. Efficient colored orthogonal range counting. *SIAM Journal on Computing*, 38(3):982–1011, 2008. doi : 10.1137/070684483.
- [14] Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. Internal pattern matching queries in a text and applications. In *26th Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015*, pages 532–551. SIAM, 2015. doi : 10.1137/1.9781611973730.36.