



**FACULTY
OF INFORMATION
TECHNOLOGY
CTU IN PRAGUE**

Assignment of bachelor's thesis

Title: Algorithms related to generalized palindromes in SageMath
Student: Ivan Romanenko
Supervisor: doc. Ing. Štěpán Starosta, Ph.D.
Study program: Informatics
Branch / specialization: Web and Software Engineering, specialization Software Engineering
Department: Department of Software Engineering
Validity: until the end of summer semester 2023/2024





Instructions

A palindrome is a word which reads the same from the left and from the right, in other words, it is fixed by the reversal mapping, which is an antimorphism fixing each word of length 1. Considering also other antimorphisms and their fixed points, we obtain the so-called generalized palindromes (see [1]). The open-source computer algebra system SageMath [2] contains a developed library supporting various algorithms on words and some of those relate to classical palindromes. The goal of the thesis is to extend some of these algorithms to generalized palindromes and, if appropriate, add other related methods.

- 1) Perform a survey on the computer algebra system and its development process, focus on the state of the word combinatorics library.
- 2) Perform a survey on generalized palindromes with a focus on G-palindromic defect [1].
- 3) Design and implement selected suitable improvements of the words combinatorics library, the least required improvement is a method for calculation of G-palindromic defect; implementation in Cython is preferred if it would result in a performance increase.
- 4) Evaluate the complexity of the designed algorithms.
- 5) Start the process of integration of these improvements into SageMath.

[1] E. Pelantová and Š. Starosta, Languages invariant under more symmetries: overlapping factors versus palindromic richness, *Discrete Math.* 313 (2013), 2432-2445, DOI: 10.1016/j.disc.2013.07.002

[2] <https://www.sagemath.org/>

Bachelor's thesis

**ALGORITHMS RELATED
TO GENERALIZED
PALINDROMES IN
SAGEMATH**

Ivan Romanenko

Faculty of Information Technology
Department of Computer Science
Supervisor: doc. Ing. Štěpán Starosta, Ph.D.
January 11, 2024

Czech Technical University in Prague
Faculty of Information Technology

© 2024 Ivan Romanenko. All rights reserved.

This thesis is school work as defined by Copyright Act of the Czech Republic. It has been submitted at Czech Technical University in Prague, Faculty of Information Technology. The thesis is protected by the Copyright Act and its usage without author's permission is prohibited (with exceptions defined by the Copyright Act).

Citation of this thesis: Romanenko Ivan. *Algorithms related to generalized palindromes in SageMath*. Bachelor's thesis. Czech Technical University in Prague, Faculty of Information Technology, 2024.

Contents

Acknowledgments	vi
Declaration	vii
Abstract	viii
Summary	ix
Notations	x
Introduction	1
Aims	3
I Theoretical Background	5
1 SageMath	7
1.1 Brief introduction to SageMath	7
1.2 Conda	7
1.3 Conda-forge	8
1.4 Installing SageMath on Linux using conda	8
1.5 SageMath development workflow	8
1.6 Creating an issue in SageMath	8
1.7 Testing in SageMath	8
2 Mathematical & Algorithmic Background	11
2.1 Palindromic defect	11
2.2 Manacher's algorithm and its generalization	13
2.3 Dynamic level ancestor problem	17
II Practical Implementation	19
3 Analysis of Requirements	21
3.1 Functional Requirements	22
3.2 Non-Functional Requirements	22
4 Design of Algorithms	23
4.1 Computing generalized palindromic defects	24
4.1.1 Θ -defect	24
4.1.2 G -defect	31
4.2 Computing classical palindromic defect	37

5	Implementation in SageMath	39
5.1	SageMath development	39
5.2	My contribution	39
5.3	Overview of changed and new methods	40
5.4	Results of implementation phase	42
6	Testing & Deployment	43
6.1	Hardware and Operating System Specification	43
6.1.1	Hardware info	43
6.1.2	Operating system info	43
6.2	Functional Testing	43
6.3	Non-Functional Testing	44
6.3.1	Performance testing for existing methods	44
6.3.2	Performance testing for G -defect computing method	47
6.4	Deployment	47
	Conclusion	49

List of Tables

6.1	Performance of <code>lengths_unioccurent_lps</code> method	45
6.2	Performance of defect method	45
6.3	Performance of <code>palindromic_complexity</code> method	46
6.4	Performance of lacunas method	46
6.5	Performance of G-defect method based on G	47
6.6	Performance of G-defect method based on the word length	47

List of code listings

2.1	Slow algorithm for finding the longest palindromic substring	14
2.2	$\mathcal{O}(n^2)$ algorithm for finding the longest palindromic substring	14
2.3	Modified $\mathcal{O}(n^2)$ algorithm for finding maximal palindromic substring radiuses	15
2.4	Manacher's algorithm for finding maximal palindromic substring radiuses	15
2.5	Manacher's algorithm for finding maximal Θ -palindromic substring radiuses	16

Thanks to my friends in Prague for support. Thanks to all math professors I studied from in Russia. Special thanks to my first math professor Vladimir Vladimirovich Nyu, Cand.Sc.(Russia) who has always been a huge inspiration for me. Thanks to my parents for financial support during my studies.

Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No. 121/2000 Coll., the Copyright Act, as amended, in particular the fact that the Czech Technical University in Prague has the right to conclude a licence agreement on the utilization of this thesis as a school work pursuant of Section 60 (1) of the Act.

In Prague on January 11, 2024

.....

Abstract

A palindrome is a word which reads the same from the left and from the right. Palindromic defect of word w is the difference between $|w|+1$ and the amount of pairwise distinct palindromic substrings of w . Concepts of palindrome and palindromic defect can be generalized to generalized palindrome, Θ -defect and G -defect, where Θ is an antimorphism, G is a finite group consisting of morphisms and antimorphisms (see [1]). The free open-source mathematics software system SageMath [2] contains a developed library containing numerous algorithms dealing with words. The first goal of the thesis is to present and prove several newly discovered algorithms for computing palindromic defect and its generalizations. As a special case of one of these algorithms, linear time algorithm for computing classical palindromic defect will be shown. The second goal of the thesis is to start adding some of these algorithms into SageMath.

Keywords SageMath, Words, Word algorithms, Palindromic defect, Generalized palindrome, Generalized palindromic defect

Abstrakt

Palindrom je slovo, které se čte stejně zleva doprava i zprava doleva. Palindromický defekt slova w je rozdíl mezi $|w|+1$ a počtem dvojic odlišných palindromických podřetězců slova w . Koncepty palindromu a palindromického defektu lze zobecnit na generalizovaný palindrom, Θ -defekt a G -defekt, kde Θ je antimorfismus, G je konečná grupa skládající se z morfismů a antimorfismů (viz [1]). Open-source matematický softwarový systém SageMath [2] obsahuje vyvinutou knihovnu podporující různé algoritmy na slova. Prvním cílem této práce je prezentovat a dokázat několik nově objevených algoritmů pro výpočet palindromického defektu a jeho zobecnění. Jako speciální případ jednoho z těchto algoritmů bude ukázán lineární časový algoritmus pro výpočet klasického palindromického defektu. Druhým cílem této práce je začít přidávat některé z těchto algoritmů do SageMath.

Klíčová slova SageMath, Řetězce, Algoritmy nad řetězci, Palindromický defekt, Generalizované palindromy, Generalizovaný palindromický defekt

Summary

This thesis aims to enhance the SageMath word library by introducing new functionalities and significantly improving the performance of existing methods.

Additionally, this thesis includes the design of several new algorithms to compute generalized palindromic defects.

Among these algorithms, one computes the Θ -defect in $\mathcal{O}(|w|)$ time. Another algorithm calculates the G -defect in $\mathcal{O}(|G| \cdot |w|)$ time, which currently stands as the fastest known algorithm as of writing this thesis. Besides these, the thesis presents designs for several other algorithms, although they are not included in the subsequent implementation in SageMath.

In this context, w represents a word: $w \in A^*$, Θ signifies an antimorphism, which is a letter

permutation, and G denotes a finite group of morphisms and antimorphisms.

The method for computing Θ -defect already exists in the SageMath word library. Therefore, the algorithm computing Θ -defect in $\mathcal{O}(|w|)$ time is implemented with same signature within SageMath but with rewritten contents. On the other hand, the method for computing G -defect did not exist in SageMath previously. The introduction of this new time-efficient method significantly enhances the SageMath word library, offering researchers new tools to explore palindromic defects behavior.

The SageMath issue on GitHub [2], to which the mentioned changes will be deployed, is numbered 35495. The SageMath fork with these changes is accessible at [3].

Notations

The following notations are assumed, if not stated otherwise:

\mathcal{A}	Alphabet
\mathcal{A}^*	Set of all finite words over \mathcal{A}
$\mathcal{AM}(\mathcal{A}^*)$	Set of all morphisms and antimorphisms on \mathcal{A}^*
$\mathcal{M}(\mathcal{A}^*)$	Set of all morphisms on \mathcal{A}^*
w	Word, $w \in \mathcal{A}^*$
f	Morphism and a letter permutation, $f \in \mathcal{M}(\mathcal{A}^*)$, $\forall a \in \mathcal{A} : f(a) \in \mathcal{A}$
Θ	Antimorphism and a letter permutation, $\Theta \in \mathcal{AM}(\mathcal{A}^*)$, $\forall a \in \mathcal{A} : \Theta(a) \in \mathcal{A}$
G	Finite group of morphisms and antimorphisms

Introduction

SageMath ([4], [5], [6], [7]) is an open-source mathematical software system.

The motivation behind this thesis is based on found inefficient implementations within the existing SageMath codebase related to palindromes and palindromic defects. Additionally, an important palindromic defect computation method was absent, despite its theoretical foundation being already developed. Addressing these issues became the primary objective, aiming to significantly benefit SageMath researchers in the related fields.

This thesis revolves around analyzing and enhancing the SageMath codebase associated with palindromes and palindromic defects. It also introduces novel algorithm designs to compute different palindromic defects. Background information to achieve these objectives includes a concise description of SageMath in chapter 1, followed by algorithms and mathematical background in chapter 2, which were used for design of the new algorithms.

The practical part II adopts a common software development life cycle process. It starts with a summary of performed requirement analysis in chapter 3, followed by algorithm design demonstrations in chapter 4. Subsequently, chapter 5 clarifies details of implementation phase, which included algorithms implementation and additional code refactoring. The thesis culminates with decisive results of functional and non-functional testing in chapter 6, which also includes the current deployment status.

This thesis is strongly related to palindromic defects, so it is only logical to add a small introduction of associated terms here.

It was shown in [8] by Droubay that amount of pairwise distinct palindrome substrings of word w does not exceed $|w| + 1$. This limit is reached, for example, for words formed by pairwise distinct letters (empty word ε is included). Palindromic defect of word w is then defined as difference between the maximal possible amount of pairwise distinct palindrome substrings for words of length $|w|$ (which is equal to $|w| + 1$ for the classical palindromes) and actual amount of pairwise distinct palindrome substrings of w .

After some time, generalization of palindrome to Θ -palindrome was introduced in [9]. Θ -defect is then a generalization of classical palindromic defect, defined the same way as classical palindromic defect, but the $|w| + 1$ term is slightly modified.

Then in [10], definition of G -defect was given in similar manner to Θ -defect.

In this work several algorithms to compute Θ -defect will be presented in section 4.1, which are theoretically backed by chapter 2. Recapitulation of some definitions, lemmas and theorems from [1] is done in section 2.1. Main purpose of section 2.2 is to generalize Manacher's algorithm ([11]) to work with Θ -palindromes instead of classical palindromes. Several already existing solutions of dynamic level ancestor problem when leaves can be added to tree dynamically ([12], [13]) will be mentioned in section 2.3.

Three algorithms for computing Θ -defect will be presented: with $\mathcal{O}(|w| \log |w|)$, $\mathcal{O}(|w|)$ and $\mathcal{O}(|w|)$ time complexities. First two of these algorithms use generalized Manacher's algorithm

from section 2.2 and data structures from section 2.3, while the last algorithm uses only generalized Manacher's algorithm. Readers interested only in this algorithm can safely skip section 2.3.

As for algorithms to compute G -defect, single-threaded algorithm with $\mathcal{O}(|G| \cdot |w|)$ time complexity will be shown. Also multi-threaded algorithm to compute G -defect in $\mathcal{O}(|w|)$ time using $\mathcal{O}(|G|)$ threads will be presented. It is conceptually same as the previous algorithm. It simply integrates classical multi-thread techniques and uses concurrent hash map data structure with minor implementation details. Concurrent hash map can be implemented in different ways and our algorithm has some extra requirements on concurrent hash map implementation to achieve promised time complexity.

Aims

This thesis focuses on two primary objectives:

- Enhancing the performance of the existing method for computing Θ -defect in SageMath.
- Implementing a new, efficient method for computing G -defect in SageMath.

To accomplish these objectives, a significant part of the work involves designing efficient algorithms for computing various palindromic defects.

The theoretical segment (Part I) of this thesis aims to familiarize readers with both SageMath and the existing theory behind palindromic defects.

The steps undertaken to achieve the main goals of this thesis in the practical part (Part II) are as follows:

1. Understanding requirements through discussions with the supervisor and researching the current SageMath codebase.
2. Designing efficient algorithms for computing different palindromic defects.
3. Implementing these algorithms in SageMath while following common development guidelines in SageMath.
4. Conducting comprehensive tests on both functionality and performance of the implemented changes.
5. Proceeding with the deployment process of these changes.

The main outcome of this thesis will be a significant contribution to the development of the SageMath word library. This expansion will be based on recent advancements in palindromic defects theory and will also significantly boost the performance of the existing codebase. These improvements collectively aim to provide researchers in palindromic-related fields with essential tools to conduct studies on a broader range of data and expedite their research processes.

Part I

Theoretical Background



Chapter 1

SageMath

1.1 Brief introduction to SageMath

SageMath ([4], [5], [6], [7]) is an open-source mathematical software system designed to support research and education in mathematics. It provides a unified interface for various mathematical tasks, including algebra, calculus, number theory, discrete mathematics, and more. SageMath combines various open-source mathematical software packages into a single cohesive environment, enabling users to perform computations, visualize data, create interactive notebooks, and collaborate on mathematical projects.

From an end-user perspective, SageMath also offers a user-friendly interface through a web-based notebook system called SageMathCell, where users can write and execute code, visualize results, and document their work using text, images, and mathematical expressions. It allows for integration of different programming languages like Python, R, Cython, and others.

SageMath is commonly used by mathematicians, researchers, educators, and students for:

- **Mathematical Computation:** Performing complex calculations, symbolic manipulation, solving equations, and conducting numerical experiments.
- **Data Visualization:** Creating plots, graphs, and visual representations of mathematical concepts and data.
- **Teaching and Learning:** Supporting educational activities by providing an interactive platform for teaching and learning mathematics.
- **Research Collaboration:** Facilitating collaboration among mathematicians and researchers by sharing notebooks and code.

SageMath can be installed locally for non-development purposes or for development purposes. There are many different ways of how SageMath can be installed. One of the easiest and most practical approaches to installing SageMath locally on Linux is by using conda together with conda-forge. First, both conda and conda-forge need to be installed.

1.2 Conda

Conda [14] is an open-source package management system and environment management system primarily used for Python. It allows users to easily install, manage, and update various software packages and their dependencies within isolated environments. Conda is known for its simplicity in creating and managing environments across different platforms (Windows, macOS, Linux) and for supporting multiple programming languages beyond just Python.

1.3 Conda-forge

Conda-forge [15] is a community-driven collection of conda packages. It is a repository that provides a vast collection of pre-built conda packages contributed and maintained by the community. Conda-forge aims to offer high-quality packages that can be easily installed using the conda package manager.

1.4 Installing SageMath on Linux using conda

After installing conda and conda-forge on our Linux machine, we can now install the non-development version of SageMath by following the steps outlined in the official installation guide in [16] and [17].

To install the development version of SageMath, we should follow the steps provided in sections [16] and [18] of the same guide.

1.5 SageMath development workflow

SageMath development workflow was significantly improved at the start of 2023 by complete transfer from its own ticket tracking system to GitHub [19]. Presently, SageMath maintains its public repository on GitHub, aligning its development workflow closely with the standard process for public repositories on the platform. Here are steps of SageMath development workflow:

1. Create or choose open issue(s) on SageMath repository, which will be fixed by your changes
2. Fork SageMath repository
3. Make changes on your fork: clone, make changes locally, commit and push. This step includes testing your changes
4. Open pull request
5. Go through code review process and finish it successfully
6. Merge your changes into SageMath repository

1.6 Creating an issue in SageMath

SageMath issue system currently uses GitHub Issues for repository [2]. New issues can be created using one of three available templates:

- Bug report
- Template for reporting a build failure
- Feature request

1.7 Testing in SageMath

SageMath uses Python's Doctest [20] approach as base for testing.

Python's Doctest [20] is a testing module included in the Python Standard Library. It allows developers to write tests in the form of examples within docstrings or docfiles (text files containing documentation). These examples resemble interactive Python sessions, including code snippets

and expected outputs. When the tests are run, `doctest` automatically executes the examples and compares the actual output to the expected output specified in the documentation. If there are discrepancies, it raises an error, indicating a test failure. Doctest helps ensure that code examples within documentation remain accurate and functional while serving as executable tests for Python code.

SageMath allows users to configure the specifics of test execution, including the ability to determine which tests are run and how they are executed. Users can specify the number of threads for test execution. Additionally, tests can be marked as 'long,' ensuring they are executed only when a specific parameter for test running is provided. Full documentation for testing in SageMath is located at [21].

Mathematical & Algorithmic Background

2.1 Palindromic defect

An *alphabet* \mathcal{A} is a finite set. Elements of \mathcal{A} are usually called *letters*. A *finite word* w over \mathcal{A} is a finite string $w = w_1w_2 \cdots w_n$ of letters $w_i \in \mathcal{A}$. Its length, denoted by $|w|$, is n . The set of all finite words over \mathcal{A} equipped with the operation of concatenation is the free monoid \mathcal{A}^* . Its neutral element is the empty word ε . A word $v \in \mathcal{A}^*$ is a *factor* of a word $w \in \mathcal{A}^*$ if there exist words $s, t \in \mathcal{A}^*$ such that $w = svt$. If $s = \varepsilon$, then v is a *prefix* of w , if $t = \varepsilon$, then v is a *suffix* of w . Alternatively, a factor of w is called a *substring* of w .

A mapping φ on \mathcal{A}^* is called

- a *morphism* if $\varphi(vw) = \varphi(v)\varphi(w)$ for any $v, w \in \mathcal{A}^*$
- an *antimorphism* if $\varphi(vw) = \varphi(w)\varphi(v)$ for any $v, w \in \mathcal{A}^*$.

We denote the set of all morphisms and antimorphisms on \mathcal{A}^* by $AM(\mathcal{A}^*)$. Together with composition, it forms a monoid with the identity mapping Id as the unit element. The set of all morphisms, denoted by $M(\mathcal{A}^*)$, is a submonoid of $AM(\mathcal{A}^*)$. The reversal mapping R is defined by

$$R(w_1w_2 \cdots w_n) = w_nw_{n-1} \cdots w_2w_1 \quad \text{for all } w = w_1 \cdots w_n \in \mathcal{A}^*$$

It is obvious that any antimorphism is a composition of R and a morphism. Thus

$$AM(\mathcal{A}^*) = M(\mathcal{A}^*) \cup R(M(\mathcal{A}^*))$$

A fixed point of a given antimorphism Θ is called Θ -palindrome, i.e., a word w is a Θ -palindrome if $w = \Theta(w)$. If Θ is the reversal mapping R , we say palindrome or classical palindrome instead of R -palindrome.

In the rest of the thesis, it will be assumed for all antimorphisms that they are letter permutations as well, because this limitation simplifies all results of this thesis and similar results for antimorphisms which are not letter permutations were not researched in this thesis.

► **Definition 2.1** (Θ -defect). *Let Θ be an antimorphism and letter permutation, $w \in A^*$. Then Θ -defect of w written as $D_\Theta(w)$ is defined as:*

$$D_\Theta(w) = |w| + 1 - \#\text{Pal}_\Theta(w) - \gamma_\Theta(w)$$

where $\text{Pal}_\Theta(w)$ is the set of Θ -palindromic factors occurring in w , and

$$\gamma_\Theta(w) := \#\{\{a, \Theta(a)\} \mid a \in \mathcal{A}, a \text{ occurs in } w \text{ and } a \neq \Theta(a)\}$$

► **Theorem 2.2.**

$$D_\Theta(w) \geq 0$$

Proof. See [1]. ◀

In the rest of the thesis, the symbol G stands exclusively for a subset of $AM(\mathcal{A}^*)$ satisfying the two following requirements:

1. G is a finite group
2. G contains at least one antimorphism

The first requirement on G implies the following for an element ν of G . The element ν is non-erasing, i.e., $\nu(a) \neq \varepsilon$ for all $a \in \mathcal{A}$ (otherwise ν has no inverse in G). Moreover, $\nu(a)$ is a letter for all $a \in \mathcal{A}$ (otherwise $\nu^n \neq \text{Id}$ for all $n \geq 1$). We can conclude that ν restricted to \mathcal{A} is a permutation of letters.

The second requirement on G comes from the fact that results of this thesis are based on generalized palindromes and one gets only trivial or no results when dealing with groups consisting of morphisms only.

► **Lemma 2.3.** *The following set of properties apply to G :*

1. every element of G is either a morphism or an antimorphism determined by a permutation of letters of \mathcal{A}
2. G may contain elements of order greater than 2
3. G need not be abelian;
4. the set of antimorphisms of G generates the group G
5. the number of morphism in G equals the number of antimorphisms in G

Proof. See [1]. ◀

We say that finite words $w, v \in \mathcal{A}^*$ are G -equivalent if there exists $\mu \in G$ such that $w = \mu(v)$. The class of equivalence containing a word w is denoted

$$[w] := \{\mu(w) \mid \mu \in G\}$$

As already mentioned, since the group G is finite, any $\mu \in G$ preserves length of words and thus equivalent words have the same length.

A word $w \in \mathcal{A}^*$ is said to be G -palindrome if there exists an antimorphism $\Theta \in G$ such that $w = \Theta(w)$.

► **Definition 2.4** (*G*-defect). Let $w \in A^*$. Then *G*-defect of w written as $D_G(w)$ is defined as:

$$D_G(w) = |w| + 1 - \#\text{Pal}_G(w) - \gamma_G(w)$$

where

$$\text{Pal}_G(w) := \{[v] \mid v \text{ is a factor of } w \text{ and a } G\text{-palindrome}\}$$

and

$$\gamma_G(w) := \#\{[a] \mid a \in \mathcal{A}, a \text{ occurs in } w, \text{ and } a \neq \Theta(a) \text{ for every antimorphism } \Theta \in G\}$$

► **Theorem 2.5.**

$$D_G(w) \geq 0$$

Proof. See [1]. ◀

On top of above theory we will need two additional lemmas in this thesis.

► **Lemma 2.6.** Let $u, v \in [w]$ and u is a *G*-palindrome. Then v is a *G*-palindrome as well.

Proof. By definition of $[w]$, $\exists \varphi \in G : \varphi(u) = v$. By definition of *G*-palindrome, $\exists \psi \in G : \psi(u) = u$ and ψ - antimorphism. Let's notice that in both cases when φ is morphism and antimorphism, $\varphi^{-1} \circ \psi \circ \varphi$ is an antimorphism and $\varphi^{-1} \circ \psi \circ \varphi(v) = v$. This means that v is *G*-palindrome. ◀

It follows from above lemma that either all elements of $[w]$ are *G*-palindromes or none of $[w]$ elements are *G*-palindromes.

► **Lemma 2.7.** Let $u, v \in [w]$ and w is a *G*-palindrome. Then $\exists f \in G : f(u) = v$, f - morphism.

Proof. By definition of $[w]$, $\exists \varphi \in G : \varphi(u) = v$. If φ is morphism, then this lemma is proven. Otherwise, φ is an antimorphism. In such case, let's notice that from the previous lemma follows u is *G*-palindrome so $\exists \psi \in G : \psi(u) = u$ and ψ is an antimorphism. Then $\psi \circ \varphi$ is a morphism and $\psi \circ \varphi(u) = v$. This is exactly what we wanted. ◀

2.2 Manacher's algorithm and its generalization

Manacher's algorithm was first described by Manacher in [11]. This algorithm finds longest palindromic substring of given word in linear time. In reality it does even more. For word w of length n , this is a $\mathcal{O}(n)$ time algorithm which finds maximum radiuses for all positions of w such that for position p and its radius r_p , substring of w starting at $p - r_p$ and ending at $p + r_p$ is a palindrome. Position p can either be position of one of w letters or it could be position in-between two consecutive letters of w (in such case we say that empty string has radius 0, two letter palindrome has radius 1 and so on).

For purposes of this thesis it is not enough to just mention this algorithm. As we are going to build different data structures during run of Manacher's algorithm, it is important to understand how Manacher's algorithm actually functions internally.

Let's start with more simple algorithms. All algorithms will be presented in Python. Naive $\mathcal{O}(n^3)$ time algorithm to find longest palindromic substring works by iterating over all $\mathcal{O}(n^2)$ substrings and checking for each substring if it is palindrome in $\mathcal{O}(n)$ time. Here is an example of algorithm:

■ **Code listing 2.1** Slow algorithm for finding the longest palindromic substring

```

1 def longest_palindromic_substring(s):
2     n = len(s)
3     max_len = 0
4     start = 0
5     for i in range(n):
6         for j in range(i+1, n+1):
7             sub = s[i:j]
8             if sub == sub[::-1] and len(sub) > max_len:
9                 max_len = len(sub)
10                start = i
11    return s[start:start+max_len]

```

It is easy to see that we could actually find all maximum radiuses instead using same idea.

One can notice that we can reuse knowledge about substring starting on a position and ending on b position being a palindrome to find out if substring starting at $a - 1$ and ending at $b + 1$ is a palindrome. With this observation we get $\mathcal{O}(n^2)$ time algorithm to find longest palindromic substring. Here is an example algorithm:

■ **Code listing 2.2** $\mathcal{O}(n^2)$ algorithm for finding the longest palindromic substring

```

1 def longest_palindromic_substring(s):
2     n = len(s)
3     max_len = 1
4     start = 0
5     for i in range(n):
6         # odd-length substrings
7         l, r = i, i
8         while l >= 0 and r < n and s[l] == s[r]:
9             if r - l + 1 > max_len:
10                max_len = r - l + 1
11                start = l
12                l -= 1
13                r += 1
14         # even-length substrings
15         l, r = i, i+1
16         while l >= 0 and r < n and s[l] == s[r]:
17             if r - l + 1 > max_len:
18                max_len = r - l + 1
19                start = l
20                l -= 1
21                r += 1
22    return s[start:start+max_len]

```

For sake of simplification, let's notice that we can add a special character s into our alphabet A and insert it between each pair of consecutive letters of w . Word w did not contain s character before we added it into alphabet, which means that all palindromic substrings are now of odd length (except empty string ε). This idea is important and by default we assume that it is used in all further algorithms.

Here is updated version of above algorithm computing radiuses instead:

■ **Code listing 2.3** Modified $\mathcal{O}(n^2)$ algorithm for finding maximal palindromic substring radiuses

```

1 def maximal_palindromic_substring_radiuses(s):
2     s = '#'.join(s) # insert special character '#'
3     p = [0] * len(s) # stores radiuses of palindromic substrings
4     for i in range(len(s)):
5         l, r = i, i
6         while l >= 0 and r < len(s) and s[l] == s[r]:
7             p[i] += 1
8             l -= 1
9             r += 1
10    odd_pos_results = [(x + 1) // 2 for x in p[::2]]
11    even_pos_results = [x // 2 for x in p[1::2]]
12    p[::2] = odd_pos_results
13    p[1::2] = even_pos_results
14    return p

```

Manacher's algorithm goes even further to achieve $\mathcal{O}(n)$ time complexity. It is based on fact that we can reuse knowledge about found palindromes centered on different positions. More precisely, let's say we have two maximal palindromic substrings w and v . w has center position $center$, start position $left$ and end position $right$. v has center position $oldCenter$ such that $left \leq oldCenter < center$, start position $oldLeft$ and radius $oldRadius$. Then for maximal radius $newRadius$ in center position $newCenter = center + (center - oldCenter)$ the following applies:

- if $oldLeft > left$, then $newRadius = oldRadius$. (otherwise, v is not maximal)
- if $oldLeft < left$, then $newRadius = right - newCenter + 1, newRadius < oldRadius$. (otherwise, w is not maximal)
- if $oldLeft = left$, then $newRadius \geq oldRadius$.

Manacher's algorithm iterates string from left to right and uses above observations. Here is an example code:

■ **Code listing 2.4** Manacher's algorithm for finding maximal palindromic substring radiuses

```

1 def maximal_palindromic_substring_radiuses(s):
2     if len(s) == 0:
3         return []
4     s = '#'.join(s)
5     p = [0] * len(s)
6     center, right = 0, 0
7     for i in range(len(s)):
8         if i < right:
9             p[i] = min(right - i, p[2 * center - i]) # use symmetry
10        l, r = i - p[i], i + p[i]
11        while l >= 0 and r < len(s) and s[l] == s[r]:
12            p[i] += 1
13            l -= 1
14            r += 1
15        if i + p[i] > right:
16            center = i
17            right = i + p[i]
18    odd_pos_results = [(x + 1) // 2 for x in p[::2]]
19    even_pos_results = [x // 2 for x in p[1::2]]
20    p[::2] = odd_pos_results
21    p[1::2] = even_pos_results
22    return p

```

Let's notice using "two pointers technique" that this algorithm indeed has $\mathcal{O}(n)$ time complexity. We can choose i and $right$ as two pointers. They both start with 0 value. On each iteration of outside loop i is increased by 1, while on each iteration of inside loop $right$ is going to be increased by 1 in the following if-statement. Both i and $right$ can not exceed length of modified string. This proves $\mathcal{O}(n)$ time complexity.

For purposes of this thesis we are going to need Manacher's algorithm which works for Θ -palindromes instead of standard palindromes. Actually, there is not much work to do for generalizing Manacher's algorithm on Θ -palindromes.

Firstly, let's say that the special character s which we are adding in our alphabet in Manacher's algorithm works with Θ as follows: $\Theta(s) = s$.

Secondly, let's notice that all observations of Manacher's algorithm on standard palindromes apply on Θ -palindromes as well. Namely, it is not hard to check that previously mentioned list of statements:

- if $oldLeft > left$, then $newRadius = oldRadius$. (otherwise, v is not maximal)
- if $oldLeft < left$, then $newRadius = right - newCenter + 1, newRadius < oldRadius$. (otherwise, w is not maximal)
- if $oldLeft = left$, then $newRadius \geq oldRadius$.

works for Θ -palindromes.

It is important to notice, though, a small technical change for Θ -palindromes compared to standard ones - maximal radius values for modified string are not necessary odd values now, they can be zeros as well.

Here is an example code of Manacher's algorithm which works with Θ -palindromes:

■ **Code listing 2.5** Manacher's algorithm for finding maximal Θ -palindromic substring radiuses

```

1 def maximal_palindromic_substring_radiuses(s, f):
2     if len(s) == 0:
3         return []
4     s = '#' .join(s)
5     p = [0] * len(s)
6     center, right = 0, 0
7     for i in range(len(s)):
8         if i < right:
9             p[i] = min(right-i, p[2*center-i])
10            l, r = i - p[i], i + p[i]
11            while l >= 0 and r < len(s) and f(s[l]) == s[r] and s[l] == f(s[r]):
12                p[i] += 1
13                l -= 1
14                r += 1
15            if i + p[i] > right:
16                center = i
17                right = i + p[i]
18            odd_pos_results = [(x + 1) // 2 for x in p[::2]]
19            even_pos_results = [x // 2 for x in p[1::2]]
20            p[::2] = odd_pos_results
21            p[1::2] = even_pos_results
22            return p

```

2.3 Dynamic level ancestor problem

There are multiple versions of dynamic level ancestor problem and known solutions for them. Let's start from recapitulation of standard level ancestor problem.

Standard level ancestor problem is a problem in which we are given a static tree graph T with n nodes and chosen root node r . The goal is to spend some time preprocessing and then answer queries of type "for node u from T return node which is k nodes higher in T (higher means going into direction of r)".

You can find known solutions for standard level ancestor problem in [12], including the best known $\mathcal{O}(n)$ preprocessing time solution which uses $\mathcal{O}(1)$ time answering queries. This solution is not suitable for our purposes, though. Instead we will use solution (again from [12]) which uses jump-pointers and achieve $\mathcal{O}(\log n)$ query time by $\mathcal{O}(n \log n)$ preprocessing time.

Dynamic version of level ancestor problem which has interest for purposes of this thesis is modification of standard level ancestor problem for non-static trees. Given a tree graph T with n initial nodes and chosen root node r , spend some time for preprocessing and then answer two types of queries: "for node u from T return node which is k nodes higher in T " and "add new leaf node into T ".

It is quite trivial that jump-pointer solution for standard level ancestor problem can be used for the above dynamic version of level ancestor problem. This algorithm is using $\mathcal{O}(n \log n)$ preprocessing time and answers both types of queries in $\mathcal{O}(\log n)$ time (n is node count of T before query).

The best known solution for such dynamic level ancestor problem which we found is solution by Alstrup; Holm described in [13]. They claim their algorithm uses $\mathcal{O}(n)$ preprocessing time and answers both types of our queries in constant time.

Part II

Practical Implementation

Analysis of Requirements

Analysis of Requirements ([22], [23]) refers to the specific needs and constraints identified during the analysis phase of a project or system development. They define the functionalities, behaviors, and constraints necessary for the system to perform as intended.

These requirements are typically split by being:

Functional Requirements: Specifications detailing what the system should do. They describe specific functionalities, features, and operations that the system must perform to meet user needs.

Functional Requirements:

- Specify system tasks.
- Are specific and tangible.
- Directly related to system functionalities.
- Validated by testing individual features and operations.
- Examples: user authentication, data validation, report generation, etc.

Non-functional Requirements: Define system attributes such as performance, security, reliability, usability, scalability, and other quality attributes crucial for the system's success but not directly related to specific functionalities.

Non-functional Requirements:

- Define system behavior or performance.
- More abstract, concerning system qualities.
- Not directly tied to specific functionalities.
- Often measured through qualities like performance, security, usability, etc.
- Examples: system response time, security protocols, user interface aesthetics, etc.

In case of this bachelor's thesis, analysis of requirements was accomplished through my verbal discussions with doc. Ing. Štěpán Starosta, Ph.D., who is supervising this thesis, and by conducting researches on the current SageMath codebase.

3.1 Functional Requirements

Based on initial discussions with supervisor of this thesis, this thesis had only one functional requirement:

- SageMath should contain a method for computing G -defect.

This is due to the fact that SageMath already contained code for computing Θ -defect.

3.2 Non-Functional Requirements

As for non-functional requirements, this thesis had multiple of them:

- Existing method for computing Θ -defect in SageMath should be enhanced by improving its asymptotic time complexity.
- New method for computing G -defect in SageMath should be added with good asymptotic time complexity.
- Adhere to the SageMath coding standards.
- Base implementation of this thesis on existing codebase of word library inside combinatorics library of SageMath.

The first requirement was established only after research of current SageMath code was finished.

We can also see that the second requirement is dependent on the sole functional requirement of this thesis.

Design of Algorithms

Driven by the non-functional requirements, we needed to describe algorithms capable of computing different defects efficiently in terms of time, while maintaining a reasonable level of simplicity for implementation. Moreover, acknowledging the scientific significance of this work, we will also present certain algorithms that, although not utilized for subsequent implementation, still hold considerable scientific value.

In this chapter the following algorithms will be presented in section 4.1:

1. Single-threaded algorithm computing Θ -defect in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm and jump-pointers
2. Single-threaded algorithm computing Θ -defect in $\mathcal{O}(n)$ time by combining Manacher's algorithm with advanced solution of dynamic level ancestor problem
3. Single-threaded algorithm computing Θ -defect in $\mathcal{O}(n)$ time by using Manacher's algorithm and tree traversal
4. Single-threaded algorithm computing G -defect in $\mathcal{O}(n \cdot |G|)$ time by using previous algorithm and tree traversal
5. Multi-threaded algorithm computing G -defect in $\mathcal{O}(n)$ time using $\mathcal{O}(|G|)$ threads

After that in section 4.2, the following algorithms will be presented:

1. Single-threaded algorithm computing classical palindromic defect in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm and jump-pointers
2. Single-threaded algorithm computing classical palindromic defect in $\mathcal{O}(n)$ time by combining Manacher's algorithm with advanced solution of dynamic level ancestor problem
3. Single-threaded algorithm computing classical palindromic defect in $\mathcal{O}(n)$ time by using Manacher's algorithm and tree traversal

These latest algorithms are special cases of algorithms computing Θ -defect from section 4.1 in case when Θ is the reversal mapping R . R is antimorphism and $\forall a \in \mathcal{A} : R(a) = a$.

Important notice is that all mentioned algorithms have promised time complexity only in cases when alphabet size is small compared to n or we are using hashing algorithm on the alphabet which performs queries in $\mathcal{O}(1)$ time. Alternatively, these time complexities can be achieved in case when all letters are enumerated by increasing space complexity of the algorithms.

In this chapter, we will use special notation w' to refer to a word built from word w by adding special character s between each pair of consecutive letters of w as it is done in Manacher's algorithm for Θ -palindromes.

4.1 Computing generalized palindromic defects

4.1.1 Θ -defect

Θ -defect is defined in definition 2.1. Just from observing this definition, we can conclude the following lemma.

► **Lemma 4.1.** *Let S be a single-threaded algorithm which finds amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(f(n))$ time, then there exists a single-threaded algorithm P which computes Θ -defect in $\mathcal{O}(f(n) + n)$ time.*

Proof. Following definition,

$$D_{\Theta}(w) = |w| + 1 - \#\text{Pal}_{\Theta}(w) - \gamma_{\Theta}(w),$$

the $|w|$ and 1 terms are trivial. If we show that a single-thread algorithm K exists which computes the $\gamma_{\Theta}(w)$ term in $\mathcal{O}(n)$ time, then this lemma will be proven, because we can just take P as a consecutive run of S and K , followed by one addition and two subtractions.

By observing definition of the term,

$$\gamma_{\Theta}(w) := \#\{\{a, \Theta(a)\} \mid a \in \mathcal{A}, a \text{ occurs in } w \text{ and } a \neq \Theta(a)\},$$

it is not hard to see that we can simply iterate over all letters of w and for each $a \in w$ check if $a \neq \Theta(a)$. Then we would need to deal with duplicate letters. This can practically be done in several ways.

For example, hash table would give us $\mathcal{O}(1)$ time to add every of $\mathcal{O}(n)$ letters considering alphabet size is small compared to n or chosen hashing algorithm behaves well for our letters.

Another alternative, which is possible if all letters are enumerated, is to just use static array of counters of size $|\mathcal{A}|$. Initially we fill the array with zeros. When iterating over letter a , if $\Theta^2(a) \neq a$, then we increase counter for a . Otherwise, we increase counter for letter with minimal index: a or $\Theta(a)$. In the end, we count the amount of non-zero counters and get $\gamma_{\Theta}(w)$.

So we have built an algorithm K computing the $\gamma_{\Theta}(w)$ term in $\mathcal{O}(n)$ time, which ends proof of this lemma. ◀

If we assume that word w is given to us as sequence of letters, then $f(n)$ in above lemma is asymptotically at least linear.

Also it follows from the above lemma, that we can forget about algorithms for computing Θ -defect and talk about algorithms for counting amount of pairwise distinct Θ -palindromic substrings of a word instead. This is exactly what we are going to do.

We will be presenting algorithms counting pairwise distinct Θ -palindromic substrings of word w . All these algorithms will have following steps in common:

1. Insert a special character s between each pair of consecutive letters of w to get modified word w' . Reasoning here is the same as in section 2.2 when special character is added before main part of Manacher's algorithm for Θ -palindromes. $\Theta(s) = s$, w does not contain s .
2. While applying Manacher's algorithm to w' , build some data structures.
3. By using data structures from previous step, build a tree graph T_w such that set of all nodes of T_w is in bijection with set of all pairwise distinct Θ -palindromic substrings of w' . The amount of nodes of T_w is same as the amount of pairwise distinct Θ -palindromic substrings of w' .

4. By making some simple observations on T_w structure, find amount of pairwise distinct Θ -palindromic substrings of w from it in linear time. This can also be done dynamically while building T_w .

Let's describe structure of T_w . As already mentioned, each node of T_w will represent a Θ -palindromic substring of w' which is different from other such substrings. Here are rules describing T_w structure:

- Nodes of T_w do not contain any data inside them.
- Edges of T_w contain one letter each.
- If there is an edge e with letter a from non-root node h to node l and h is higher in the tree than l and h represents word v , then node l represents word $av\Theta(a)$.
- If there is an edge e with letter a from root node to node l , then l represents word a .
- Root of T_w represents empty string ε .
- Every pairwise distinct Θ -palindromic substring of w' is represented in T_w by exactly one node.

For every possible w there always exists exactly one tree graph to which all these rules apply. This can be seen from the fact that w' contains only odd-length Θ -palindromic substrings (except empty string ε), and the fact that if w' contains Θ -palindromic substring $u = avb$, where $a, b \in \mathcal{A}$, then v is a Θ -palindromic substring of w' as well.

The next lemma helps with finding amount of pairwise distinct Θ -palindromic substrings of w from T_w :

► **Lemma 4.2.** *Amount of pairwise distinct Θ -palindromic substrings of w is equal to amount of nodes of T_w which do not represent word starting with special character s . In other words, it is equal to amount of edges of T_w which do not contain s character plus one.*

Proof. Let's show that there exists a bijection between all nodes of T_w which do not represent word starting with special character s and all pairwise distinct Θ -palindromic substrings of w .

From the way we defined w' by adding special character s between each pair of consecutive letters of w , it follows that a word v is a Θ -palindromic substring of w iff word v' , which does not start with s , is a Θ -palindromic substring of w' , where v' is v with s character added between each pair of consecutive letters. This gives us bijection between all Θ -palindromic substrings of w' not starting with s and all Θ -palindromic substrings of w .

There is an obvious bijection between all nodes of T_w which do not represent word starting with special character s and all Θ -palindromic substrings of w' not starting with s .

Combining these two bijections give us exactly what we wanted. ◀

Now considering above lemma, we can state another important lemma:

► **Lemma 4.3.** *Let S be a single-threaded algorithm which builds T_w in $\mathcal{O}(f(n))$ time, $f(n)$ is asymptotically at least linear. Then there exists a single-threaded algorithm P which finds amount of pairwise distinct Θ -palindromic substrings of w in $\mathcal{O}(f(n))$ time.*

Proof. Let's run S algorithm to build T_w . Then, using lemma 4.2, traverse T_w and find amount of pairwise distinct Θ -palindromic substrings of w in $\mathcal{O}(n)$ time. ◀

So being able to compute T_w fast would provide us with a way to compute amount of pairwise distinct Θ -palindromic substrings of a word fast. The following lemma will help us achieve it by establishing a link between structure of T_w and Manacher's algorithm for Θ -palindromes:

► **Lemma 4.4.** T_w can be built in the following way. At the start T_w contains only root node, which represents empty string. Then Manacher's algorithm for Θ -palindromes is run on w' (special character is not added, w' already has it) and we try to add a new node to T_w each time maximal radius of current position is increased in Manacher's algorithm. The node we are adding to T_w is the one, which represent Θ -palindromic substring in the current position with the currently known maximal radius in this position of running Manacher's algorithm. In case there is already a node in T_w , which represents such Θ -palindrome, we do not add a new node.

Proof. To prove this lemma we need to prove two things.

Firstly, we need to prove that adding a node in the described way is always possible. This means we need to prove that parent node already exists in T_w when we are adding a new node.

Secondly, we need to prove that in the end T_w will contain exactly one node for every pairwise distinct Θ -palindromic substring of w' . We can not have two nodes representing same Θ -palindrome simply by the way we are adding nodes to T_w . So we only need to show that we will try to add to T_w every pairwise distinct Θ -palindromic substring of w' at least once.

The following observation will help us to prove these statements.

By behaviour of Manacher's algorithm (see section 2.2), it can be seen that Manacher's algorithm iterates over w' letter positions from left to right. For each position it first defines initial radius as either zero or by coping it by symmetry from some previous position and possibly reducing copied radius. Then it iterates from the initial radius to actual maximal radius in the position (possibly this iteration has zero steps).

Now we can prove both statements by mathematical induction over the following statements:

- $A(k)$ is "When current position of Manacher's algorithm reaches $k + 1$, then T_w contains nodes representing all Θ -palindromes centered at any of the first k letters of w' "
- $B(k)$ is "When adding all nodes for position $k + 1$ into T_w , parent of each new node will already exist in T_w "

We start numerating letter positions with 1.

To prove that $B(k)$ is true, it is enough to show that parent node for the first node added in position $k + 1$ exists in T_w before we try to add this node. For the rest of the added nodes, parent node already exists because we added it on the previous step of Manacher's algorithm.

Let's prove the statements $A(k)$ and $B(k)$ together by mathematical induction for k from 0 to $|w'|$. $B(|w'|)$ is considered to be true from the start. For $k = 0$, $A(0)$ is trivial, $B(0)$ is true, because T_w contains empty string from the start and initial radius in position 1 is zero.

For step $k \rightarrow k + 1$, let's notice that $A(k + 1)$ follows from $A(k)$ and $B(k)$. Indeed, based on $A(k)$, we only need to prove that all Θ -palindromes centered at $k + 1$ are contained in T_w after Manacher's algorithm finishes iterating from initial radius to actual maximal radius at $k + 1$. From $B(k)$ it follows that all Θ -palindromes centered at $k + 1$ with radius less than initial radius are already represented in T_w before the iteration starts. After iteration is done, the rest of Θ -palindromes centered at $k + 1$ are added to T_w . This proves $A(k + 1)$.

Now we will prove $B(k + 1)$ when we already know all $A(l), l \leq k + 1$. If $k + 1 = |w'|$, then there is nothing to prove. Let's see what happens when current position of Manacher's algorithm is $k + 2 \leq |w'|$. If initial radius at $k + 2$ is zero, then $B(k + 1)$ is proven, because T_w contains node for empty string from the start. Otherwise, initial radius at $k + 2$ was copied from actual maximal radius of some previous position $m \leq k + 1$ and possibly reduced after copy. In both cases, $B(k + 1)$ follows from $A(k + 1)$ and $m \leq k + 1$. This ends induction.

It is not hard to see that statements which we initially needed to prove are followed from statements we just proved by mathematical induction. ◀

Now we will present an algorithm which finds amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(n^2)$ time.

► **Lemma 4.5.** *There exists an algorithm finding amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(n^2)$ time using Manacher's algorithm for Θ -palindromes.*

Proof. It follows from lemma 4.3 that it is enough for us to present an algorithm which builds T_w in $\mathcal{O}(n^2)$ time. We can build T_w using lemma 4.4, but we still need to answer the question of how we will find parent node in T_w of a node we are currently trying to add.

We can just take Θ -palindrome represented by the parent node and find this Θ -palindrome in T_w in $\mathcal{O}(n)$ time by going down from the root of T_w . Doing it every time we try to add a node in T_w would give us an algorithm to build T_w in $\mathcal{O}(n^2)$ time. ◀

It is hard to improve time complexity of the above algorithm, so we need to introduce an additional data structure, which we will be building while running Manacher's algorithm for Θ -palindromes.

For every letter position p in w' we are going to remember two things:

1. Node v_p in T_w
2. Non-negative integer t_p

The limitations on choosing v_p and t_p pair are that v_p currently exists in T_w , depth of v_p in T_w is at least t_p and t_p th ancestor of v_p in T_w is a node representing Θ -palindrome centered at p with maximal possible radius.

► **Lemma 4.6.** *Pairs (v_p, t_p) for all positions p can be computed in $\mathcal{O}(n)$ time while running Manacher's algorithm for Θ -palindromes.*

Proof. We will set (v_p, t_p) pair right after Manacher's algorithm finishes iterating from initial radius to maximal radius in position p .

If initial radius was less than maximal radius in p , then at the last step of iteration we added node into T_w which represents Θ -palindrome centered at p with maximal possible radius, so we can set t_p as 0 and v_p as this node right after we added it.

If initial radius was equal to maximal radius in p , then there are two possibilities. First possibility is that both initial radius and maximal radiuses are 0. In this case we set v_p as root node of T_w and t_p as 0. Second possibility is that initial radius in p was copied from maximal radius in some previous position l and reduced by non-negative integer value d . In this case we can set $v_p = v_l$ and $t_p = t_l + d$. Limitations on pair (v_p, t_p) are not broken here. Limitations on t_p are not broken because maximal radius in position l was at least d simply by definition of d and correctness of Manacher's algorithm. Limitations on v_p are not broken because of behaviour of Manacher's algorithm for Θ -palindromes.

We can notice that in all cases we spend constant time to compute every next pair (v_p, t_p) . And plain Manacher's algorithm for Θ -palindromes runs in $\mathcal{O}(n)$ time. So we find pairs (v_p, t_p) for all positions p in $\mathcal{O}(n)$ time. ◀

We are now ready to present three algorithms which find amount of pairwise distinct Θ -palindromic substrings of a word and have time complexities $\mathcal{O}(n \log n)$, $\mathcal{O}(n)$ and $\mathcal{O}(n)$ respectively.

► **Theorem 4.7.** *There exists an algorithm finding amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm for Θ -palindromes and jump-pointers.*

Proof. Our algorithm will follow the same basic structure as the algorithm from lemma 4.5. On top of that we will be building T_w with jump-pointers (see section 2.3), so adding new nodes into T_w will cost us $\mathcal{O}(\log n)$ time now. And we will be additionally computing (v_p, t_p) pairs while running Manacher's algorithm for Θ -palindromes (see lemma 4.6). Another possibility is to compute all (v_p, t_p) pairs beforehand in $\mathcal{O}(n)$ time.

Now let's show that finding parent node in T_w of a node we are currently trying to add can be done in $\mathcal{O}(\log n)$ time. Then following same logic as in algorithm from lemma 4.5, we will get the $\mathcal{O}(n \log n)$ time algorithm we promised. As can be seen from lemma 4.4, we are trying to add new nodes into T_w only when our Manacher's algorithm is currently in position p , which has initial radius less than maximal radius, and we increase current radius for position p .

Let's show how we can find in $\mathcal{O}(\log n)$ time node representing Θ -palindrome centered in p with radius equal to initial radius in p . If initial radius in p is 0, then we take the root node of T_w in constant time. If initial radius in p is greater than 0, then it was copied from some previous position l and reduced by non-negative integer d . Pair (v_l, t_l) was already computed, so we can take $(t_l + d)$ th ancestor of v_l in $\mathcal{O}(\log n)$ time using jump-pointers. This is exactly the node we are looking for.

So we can find parent nodes for nodes, which are added when current radius in some position p is increased from initial radius, in $\mathcal{O}(\log n)$ time. For the nodes, which are added when current radius in p is increased not from initial radius, their parent node was found or added on the previous step of Manacher's algorithm for Θ -palindromes. So finding the parent node can be done in constant time in such case.

In worst case, we find parent node in T_w of a node we are currently trying to add in $\mathcal{O}(\log n)$ time. As already mentioned above, the promised algorithm with $\mathcal{O}(n \log n)$ time complexity follows from this. ◀

► **Theorem 4.8.** *There exists an algorithm finding amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(n)$ time by combining Manacher's algorithm for Θ -palindromes and advanced solution for dynamic level ancestor problem.*

Proof. Let's notice that algorithm described in theorem 4.7 does not use any properties of jump-pointers except that we use jump-pointers as a solution for dynamic level ancestor problem (see section 2.3).

So if we had a solution for dynamic level ancestor problem, such that we can add new leaf nodes in constant time and find ancestors in constant time, then we would build the promised algorithm based on algorithm from theorem 4.7.

Such solution is mentioned in section 2.3. ◀

► **Theorem 4.9.** *There exists an algorithm finding amount of pairwise distinct Θ -palindromic substrings of a word in $\mathcal{O}(n)$ time by using Manacher's algorithm for Θ -palindromes and tree traversal.*

Proof. From lemma 4.3 it follows that it is enough for us to present an algorithm which builds T_w in $\mathcal{O}(n)$ time.

And from lemma 4.4 follows that T_w can be built in the following way. Initially T_w contains only root node. Then we run Manacher's algorithm for Θ -palindromes on w' , while collecting all pairs (p, r) of positions with radiuses such that Manacher's algorithm increased radius in position p from $r - 1$ to r . And then iterate over all these pairs (p, r) in some order and for each pair add node into T_w which represents palindrome centered in p with radius r . The only limitation on

the order of (p, r) pairs is that when we are adding node into T_w , the parent node of this node must already exist in T_w .

To use the above observation we will first need to build an additional data structure.

While running Manacher's algorithm for Θ -palindromes, we will build an oriented forest graph H which will contain one node h_p for each letter position p of w' .

Every node h_p will contain two non-negative integers $incr_p$ and $decr_p$. $incr_p$ is equal to maximal radius in position p minus initial radius in p as in Manacher's algorithm. If Manacher's algorithm copied initial radius in position p from some previous position and reduced it by non-negative integer d , then $decr_p = d$. Otherwise, initial radius in position p is 0 and we define $decr_p = 0$.

Additionally every node h_p will contain index of position p in w' , so for every node h_p we can go to position p in w' in constant time.

Edges of H will be added based on how Manacher's algorithm copies radiuses of previous positions. An edge e from node h_l to node h_p exists in H iff Manacher's algorithm copied initial radius in position p from maximal radius in previous position l .

It is not hard to see that H can be built in $\mathcal{O}(n)$ time by running Manacher's algorithm for Θ -palindromes, because every $incr_p$ and $decr_p$ can be computed in constant time in all cases and every edge can be added in constant time. Details of implementation are left to readers as an exercise.

The fact, that H is indeed an oriented forest, follows from how edges of H were defined.

Actually, we do not need all the data contained in H , so we will build another oriented forest H' by removing some nodes from H .

H' starts off as a copy of H and then we consecutively remove all nodes h_p from H' such that $incr_p = 0$, while adding $decr_p$ values lost from removal of the nodes to descendant nodes. In other words, when node h_p is removed, then for each direct descendant h_l of h_p we modify $decr_l = decr_l + decr_p$. If h_p is a root node for some tree in H' , then we simply remove node h_p and all of its edges, meaning amount of trees in H' might change after this operation. If h_p is not a root node and has a parent node h_t , then for each direct descendant h_l of h_p we add an edge from h_t to h_p , and finally, we remove node h_p and all of its edges.

The order in which nodes are removed from H' is not arbitrary. At the start we choose an arbitrary order in which trees of H' will be traversed. Then we traverse each tree from the root node using depth-first traversal, while removing the nodes as mentioned above. Here it is important that we traverse the trees in such way that after node h_p was visited, all ascendants of h_p were already visited before and none of them will be removed later.

Let's notice that in the described way of building H' we spend $\mathcal{O}(\deg_{in}(h_p) + \deg_{out}(h_p))$ time to remove node h_p and the value $\deg_{in}(h_l) + \deg_{out}(h_l)$ is not increased for any descendant nodes h_l of h_p . From this observation follows, that the described way of building H' from H has time complexity $\mathcal{O}(\# \text{ nodes}(H) + \# \text{ edges}(H) + \sum_{h_p \in \text{nodes}(H)} \deg(h_p)) = \mathcal{O}(n + n + \# \text{ edges}(H))$ which in fact is $\mathcal{O}(n)$.

In reality we can build H' instead of H from the start during run of Manacher's algorithm for Θ -palindromes. Details of implementation and proof are left to readers as an exercise.

As can be seen from definition of H' , H' contains one node h'_p for each position p of w' such that maximal radius in p is bigger than initial radius in p in Manacher's algorithm. All $incr_p$ values in H' are positive.

H has one property which is important for us. For each node h_p in H , initial radius in position p is equal to sum of all $incr_a - decr_a$ values for all ancestors h_a of h_p minus $decr_p$. H follows this property because of how it is defined based on Manacher's algorithm. Moreover, H' also follows this property. It can be seen from the fact that when we are building H' , it starts off as H , which follows the property, and the property is not broken after each node removal in H' (notice how we modify $decr_p$ values for this).

From now on we will not need H anymore, so any mentions of $incr_p$ and $decr_p$ values will be referring to the values in nodes of H' .

Now let's build T_w using H' . Initially, T_w contains only root node, which represents an empty string ε .

We will present an algorithm S for building part of T_w while traversing one tree F of H' . This algorithm is then applied to every tree of H' consecutively. The order, in which trees are traversed, is arbitrary.

While traversing F , we will be remembering some path $curPath$ of T_w , which start at the root node of T_w . We will also have a dynamic array $nodesLeft$, which has the same size as amount of nodes in $curPath$. $nodesLeft$ will contain as its elements dynamic arrays of references to some nodes of F .

Let's describe how $curPath$ and $nodesLeft$ will be changing during traversal of F and how they will affect traversal of F .

When algorithm S comes to a node $h'_p \in F$ during traversal and it wants to mark h'_p as visited, then it follows the next procedure. By the way S will behave, at this moment $curPath$ will end with a node from T_w which represents Θ -palindrome centered in p with radius equal to initial radius in p . Before marking h'_p as visited, S iterates from initial radius in p plus one to maximal radius in p and tries to add new node into T_w representing Θ -palindrome centered in p with current radius. S also pushes references to these new nodes or already existing nodes into $curPath$, and it fills $nodesLeft$ with empty dynamic arrays so it has the same size as node count in $curPath$. After this part is done, S iterates over all direct descendants h'_l of h'_p in F and it adds a reference to h'_l in $nodesLeft$ array into dynamic array on index $(\text{len}(\text{nodesLeft}) - \text{decr}_l)$. And finally, S marks h'_p as visited and continues. Time complexity of this operation is $\mathcal{O}(\text{incr}_p + \text{deg}_{\text{out}}(h'_p))$, because we get p from F in constant time, we get initial radius from node count of $curPath$, then we do incr_p constant time operations of trying to add new node into T_w , and finally we add reference into $nodesLeft$ for each of $\text{deg}_{\text{out}}(h'_p)$ direct descendants of h'_p , adding each reference costs constant time.

When algorithm S needs to choose node h'_p to apply the above procedure of marking h'_p as visited, it follows the next procedure. If the last dynamic array in $nodesLeft$ is not empty, then S pops the last node from it and takes this node as h'_p . If the last dynamic array in $nodesLeft$ is empty, then S decreases size of $nodesLeft$ by 1 and removes the last node from $curPath$, and after that S repeats this procedure of choosing h'_p . If $nodesLeft$ has size 0, then S finishes.

Initially S starts with procedure of marking root node of F as visited.

As can be seen from definition of S , all nodes of F will indeed be marked as visited at the end of S . It follows from the fact that after we mark some node h'_p as visited, then we add references to all direct descendants of h'_p somewhere into $nodesLeft$ and so we will mark all of them as visited some time later.

From definition of S and properties of H' mentioned above, it follows that for node h'_p from F , before S marked h'_p as visited, S tried to add into T_w nodes representing all Θ -palindromes centered in p with radius bigger than initial radius in p and not bigger than maximal radius in p .

So if we apply S for all trees of H' in some arbitrary order, then we will try to add to T_w nodes representing all Θ -palindromes centered in position p with radius r , which is bigger than initial radius in p and not bigger than maximal radius in p , for all pairs (p, r) .

Then it follows from the observation on how T_w can be built, which was made at the start of this proof, that the above algorithm actually builds whole T_w .

Now it is only left to show that the above algorithm of building T_w from H' has $\mathcal{O}(n)$ time complexity. This will finish the proof, because we can build H in $\mathcal{O}(n)$ time, then build H' from H in $\mathcal{O}(n)$ time, and finally build T_w from H' in $\mathcal{O}(n)$ time. Combining all three algorithms, we get an algorithm building T_w in $\mathcal{O}(n)$ time, which is exactly what we wanted, as mentioned at the very start of this proof.

Let's notice that all calls of procedure of S to choose node h'_p do the same amount of constant-time steps as amount of times $curPath$ is increased by 1 or a some reference is added into $nodesLeft$ somewhere during run of S . Indeed, at the start of S and at the end of S , both

$curPath$ and $nodesLeft$ are empty. Nodes are removed from $curPath$ and references are removed from $nodesLeft$ only during procedure of choosing node h'_p .

It follows that time spent on choosing nodes during run of S is not asymptotically bigger than time spent on marking nodes as visited. So time complexity of S applied to F is equal to $\mathcal{O}\left(\sum_{h'_p \in F} incr_p + \sum_{h'_p \in F} \deg_{out}(h'_p)\right)$. Then time complexity of algorithm for building T_w from H' by using algorithm S is $\mathcal{O}\left(\sum_{h'_p \in H'} incr_p + \sum_{h'_p \in H'} \deg_{out}(h'_p)\right) = \mathcal{O}\left(n + \sum_{h'_p \in H'} incr_p\right)$. Sum $\sum_{h'_p \in H'} incr_p$ is equal to amount of time radius in current position is increased during Manacher's algorithm, so $\mathcal{O}\left(\sum_{h'_p \in H'} incr_p\right) = \mathcal{O}(n)$. Then it follows that $\mathcal{O}\left(n + \sum_{h'_p \in H'} incr_p\right)$ is same as $\mathcal{O}(n)$. This proves that our algorithm to build T_w from H' has $\mathcal{O}(n)$ time complexity, which is exactly what we wanted. ◀

Going back to computing Θ -defect, let's present the initially promised algorithms.

► **Lemma 4.10.** *There exists an algorithm computing Θ -defect in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm for Θ -palindromes and jump-pointers.*

Proof. Combine theorem 4.7 and lemma 4.1. ◀

► **Lemma 4.11.** *There exists an algorithm computing Θ -defect in $\mathcal{O}(n)$ time by combining Manacher's algorithm for Θ -palindromes and advance solution of dynamic level ancestor problem.*

Proof. Combine theorem 4.8 and lemma 4.1. ◀

► **Lemma 4.12.** *There exists an algorithm computing Θ -defect in $\mathcal{O}(n)$ time by using Manacher's algorithm for Θ -palindromes and tree traversal.*

Proof. Combine theorem 4.9 and lemma 4.1. ◀

4.1.2 G -defect

For purposes of multi-threaded algorithms presented in this subsection, we will require an implementation of concurrent hash table/map which has additional properties that any amount of parallel queries to remove the same element from it has constant time complexity and that any amount of parallel queries to add the same element to it has constant time complexity. Implementation details of such version of concurrent hash table/map are left to readers as an exercise.

In this subsection, we will be using symbol r specifically to represent amount of antimorphisms contained in G . It follows from lemma 2.3 that $r = \mathcal{O}(|G|)$.

Any mention of concurrent hash table/map data structure in this subsection will be referring to implementation of concurrent hash table/map mentioned above.

We will also use $T_{w,\Theta}$ notation to refer to tree T_w defined in subsection 4.1.1 built for word w and antimorphism Θ .

G -defect is defined in definition 2.4. Just from observing this definition, we can conclude the following lemma.

► **Lemma 4.13.** *Let S be a single-threaded algorithm which finds $\#\text{Pal}_G(w)$ in $\mathcal{O}(f(n, r))$ time, then there exists a single-threaded algorithm P which computes G -defect in $\mathcal{O}(f(n, r) + r \cdot n)$ time.*

Proof. Following definition,

$$D_G(w) = |w| + 1 - \#\text{Pal}_G(w) - \gamma_G(w)$$

the $|w|$ and 1 terms are trivial. If we show that a single-thread algorithm K exists which computes the $\gamma_G(w)$ term in $\mathcal{O}(r \cdot n)$ time, then this lemma will be proven, because we can just take P as a consecutive run of S and K , followed by one addition and two subtractions.

By observing definition of the term,

$$\gamma_G(w) := \#\{[a] \mid a \in \mathcal{A}, a \text{ occurs in } w, \text{ and } a \neq \Theta(a) \text{ for every antimorphism } \Theta \in G\},$$

it is not hard to see that we can simply iterate over all letters of w and for each $a \in w$ check if $a \neq \Theta(a)$ for every antimorphism $\Theta \in G$ and mark all such letters a . This requires $\mathcal{O}(r \cdot n)$ time. Then we would need to deal with duplicate letters. This can practically be done in several ways.

We will show how to do it using hash table. Adding one letter to hash table would require $\mathcal{O}(1)$ time considering alphabet size is small compared to n or chosen hashing algorithm behaves well for our letters. So we try to add all letters, which we marked above, into hash table. It takes $\mathcal{O}(n)$ and hash table does not contain duplicate letters in the end.

Now when we have the hash table, we need to find how many pairwise distinct classes of equivalence $[a]$ there are among letters a of our hash table. We can do it by taking an arbitrary element a of our hash table and trying to remove elements $\varphi(a)$ from hash table for all $\varphi \in G$. Then repeat this for another arbitrary element of hash table until hash table is empty. Overall time complexity of this is $\mathcal{O}(r \cdot n)$. Amount of iterations of this algorithm is equal to amount of pairwise distinct classes of equivalence $[a]$ in our hash table. This is exactly what we want to compute.

Summing up, we have built an algorithm K computing the $\gamma_G(w)$ term in $\mathcal{O}(r \cdot n)$ time, which ends proof of this lemma. ◀

We can also prove similar lemma for multi-threaded algorithms.

► **Lemma 4.14.** *Let S be a multi-threaded algorithm which uses $k \leq r$ threads and finds $\#\text{Pal}_G(w)$ in $\mathcal{O}(f(n, r, k))$ time, then there exists a multi-threaded algorithm P which computes G -defect in $\mathcal{O}(f(n, r, k) + \frac{r \cdot n}{k})$ time.*

Proof. Following definition,

$$D_G(w) = |w| + 1 - \#\text{Pal}_G(w) - \gamma_G(w)$$

the $|w|$ and 1 terms are trivial. If we show that a multi-thread algorithm K using k threads exists which computes the $\gamma_G(w)$ term in $\mathcal{O}(\frac{r \cdot n}{k})$ time, then this lemma will be proven, because we can just take P as a consecutive run of S and K , followed by one addition and two subtractions.

By observing definition of the term,

$$\gamma_G(w) := \#\{[a] \mid a \in \mathcal{A}, a \text{ occurs in } w, \text{ and } a \neq \Theta(a) \text{ for every antimorphism } \Theta \in G\},$$

it is not hard to see that we can iterate over all letters of w and for each $a \in w$ check if $a \neq \Theta(a)$ for every antimorphism $\Theta \in G$ and mark all such letters a . This can be done in $\mathcal{O}(\frac{r \cdot n}{k})$ time using k threads. Indeed, let's first setup an atomic boolean for each letter of w as true. This setup can be done in $\mathcal{O}(n)$ time using one thread. Then we numerate all pairs (a, Θ) and make thread with index i iterate all pairs with indices j such that $j \equiv i \pmod{k}$. If some thread finds that $a = \Theta(a)$ for some Θ , then it sets the atomic boolean of letter a as false. In the end, we are interested only in letters with atomic boolean with true value.

Now we need to deal with duplicate letters. This can practically be done in several ways.

We will show how to do it using concurrent hash table. Adding one letter to concurrent hash table would require $\mathcal{O}(1)$ time considering alphabet size is small compared to n or chosen hashing algorithm behaves well for our letters. So we try to add all letters, which have atomic boolean with true value, into concurrent hash table. It takes $\mathcal{O}(n)$ using one thread and concurrent hash table does not contain duplicate letters in the end.

Now when we have the concurrent hash table, we need to find how many pairwise distinct classes of equivalence $[a]$ there are among letters a of our concurrent hash table. We can do it by taking an arbitrary element a of our concurrent hash table and trying to remove elements $\varphi(a)$ from concurrent hash table for all $\varphi \in G$. We do it using k threads by numerating all $\varphi_j \in G$ and making thread with index i iterate all φ_j with indices j such that $j \equiv i \pmod{k}$. Then we repeat this for another arbitrary element of concurrent hash table until concurrent hash table is empty. Overall time complexity of this is $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$. Amount of iterations of this algorithm is equal to amount of pairwise distinct classes of equivalence $[a]$ in our concurrent hash table. This is exactly what we want to compute.

Summing up, we have built an algorithm K computing the $\gamma_G(w)$ term in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time, which ends proof of this lemma. ◀

It follows from the above lemmas, that we can forget about algorithms for computing G -defect and talk about algorithms for computing $\#\text{Pal}_G(w)$ instead. This is exactly what we are going to do.

To compute $\#\text{Pal}_G(w)$, our algorithms will be building a special tree graph to which we will refer as $T_{w,G}$. It will be built using trees $T_{w,\Theta}$ built for all antimorphisms $\Theta \in G$.

Before describing structure of $T_{w,G}$, we need to pick an arbitrary order of all morphisms $f \in G$. We will refer to morphism with index i as f_i .

It is also important to say that special character s , which is used to build w' from w , is chosen the same for all antimorphisms $\Theta \in G$, and $\forall \Theta \in G : \Theta(s) = s$. It follows from definition that all G -palindromic substrings of w' have odd length (except empty string ε).

Now let's describe structure of $T_{w,G}$. Every node h of $T_{w,G}$ will represent a G -palindromic class of equivalency $[v']$, v' is a substring of w' , which is different from G -palindromic classes of equivalency represented by other nodes of $T_{w,G}$. We will refer to node representing a class of equivalency $[v']$ as $h_{v'}$. Here are rules describing $T_{w,G}$ structure:

- Edges of $T_{w,G}$ contain two letters each. For edge e , we will refer to the first letter as e_l and to the second letter as e_r . $e_l = e_r$ for all edges from root node of $T_{w,G}$.
- If there is an edge e with $e_l = a$ and $e_r = b$ from non-root node h to node l and h is higher in the tree than l and h represents $[v]$, then node l represents $[avb]$.
- If there is an edge e with $e_l = e_r = a$ from root node to node l , then l represents $[a]$.
- Root of $T_{w,G}$ represents $[\varepsilon]$.
- Every pairwise distinct G -palindromic class of equivalency $[v']$, such that v' is substring of w' , is represented in $T_{w,G}$ by exactly one node.
- Node $h_{v'}$ of $T_{w,G}$ will contain a dynamic array of size r of references to hash maps $\text{mapRef}_{i,v'}$. Two references $\text{mapRef}_{i,v'}$ and $\text{mapRef}_{j,v'}$ might reference the same hash map. We will refer to hash maps by $\text{map}_{i,v'}$. Keys of the hash maps are pairs of letters, and values are references to other such hash maps. All these hash maps are maintained in one dynamic array outside of $T_{w,G}$. We will be saying that hash map $\text{map}_{i,v'}$ represents string $f_i(v')$ and $\text{map}_{i,\varepsilon}$ represents ε .

- For node $h_{v'}$ with parent edge e with $e_l = a$ and $e_r = b$, $mapRef_{i,v'} = mapRef_{j,v'} \iff (f_i(a), f_i(b)) = (f_j(a), f_j(b))$. For root node h_ε of $T_{w,G}$, $\forall i, j : mapRef_{i,\varepsilon} = mapRef_{j,\varepsilon}$.
- There is an edge e with $e_l = a$ and $e_r = b$ from node $h_{v'}$ to node $h_{av'b}$ iff hash map $map_{i,v'}$ in node $h_{v'}$ contains an element with the key $(f_i(a), f_i(b))$ and the value $mapRef_{i,av'b}$.

For every possible w there always exists exactly one tree graph to which all these rules apply. Uniqueness of data inside nodes follows from how data inside nodes is defined. The rest can be seen from the fact that w' contains only odd-length G -palindromic substrings (except empty string ε), and the fact that if w' contains G -palindromic substring $u = avb$, where $a, b \in \mathcal{A}$, then v is a G -palindromic substring of w' as well.

Every hash map map of $T_{w,G}$ is always referenced only in one node h of $T_{w,G}$. We will refer to h as node containing map and put a reference to h alongside map so we can access it in constant time.

The next lemma helps with finding $\# \text{Pal}_G(w)$ from $T_{w,G}$:

► **Lemma 4.15.** $\# \text{Pal}_G(w)$ is equal to amount of nodes of $T_{w,G}$ which do not represent $[v']$ such that v' starts with special character s . In other words, it is equal to amount of edges e of $T_{w,G}$ such that $e_l \neq s$ plus one.

Proof. Let's show that there exists a bijection between all nodes $h_{v'}$ of $T_{w,G}$ such that v' does not start with s and all pairwise distinct G -palindromic classes of equivalency $[v]$, v is a substring of w .

From the way we defined w' by adding special character s between each pair of consecutive letters of w , it follows that a word v is a G -palindromic substring of w iff word v' , which does not start with s , is a G -palindromic substring of w' , where v' is v with s character added between each pair of consecutive letters. This gives us bijection between all G -palindromic substrings of w' not starting with s and all G -palindromic substrings of w .

Also from definition of s it follows that for any v, u two G -palindromic substrings of $w : v \in [u] \iff v' \in [u']$. This gives us bijection between all pairwise distinct G -palindromic classes of equivalency $[v']$, where v' is a substring of w' and v' does not start with s , and all pairwise distinct G -palindromic classes of equivalency $[v]$, v is a substring of w .

By definition of $T_{w,G}$, there is an obvious bijection between all nodes $h_{v'}$ of $T_{w,G}$ such that v' does not start with s and all pairwise distinct G -palindromic classes of equivalency $[v']$, where v' is a substring of w' and v' does not start with s .

Combining the last two bijections give us exactly what we wanted. ◀

Now considering above lemma, we can state several other important lemmas.

► **Lemma 4.16.** Let S be a single-threaded algorithm which builds $T_{w,G}$ in $\mathcal{O}(f(n, r))$ time. Then there exists a single-threaded algorithm P which computes $\# \text{Pal}_G(w)$ in $\mathcal{O}(f(n, r) + n)$ time.

Proof. Let's run S algorithm to build $T_{w,G}$. Then, using lemma 4.15, traverse $T_{w,G}$ and compute $\# \text{Pal}_G(w)$ in $\mathcal{O}(n)$ time. ◀

► **Lemma 4.17.** Let S be a multi-threaded algorithm which uses k threads and builds $T_{w,G}$ in $\mathcal{O}(f(n, r, k))$ time. Then there exists a multi-threaded algorithm P which uses k threads and computes $\# \text{Pal}_G(w)$ in $\mathcal{O}(f(n, r, k) + n)$ time.

Proof. Let's run S algorithm to build $T_{w,G}$. Then, using lemma 4.15, traverse $T_{w,G}$ with one thread and compute $\# \text{Pal}_G(w)$ in $\mathcal{O}(n)$ time. ◀

Now we are ready to present single-threaded algorithms for computing $\#\text{Pal}_G(w)$ and G -defect.

► **Theorem 4.18.** *Let S be a single-threaded algorithm which builds $T_{w,\Theta}$ for arbitrary antimorphism $\Theta \in G$ in $\mathcal{O}(f(n))$ time, then there exists a single-threaded algorithm P which computes $\#\text{Pal}_G(w)$ in $\mathcal{O}(r \cdot f(n) + r \cdot n)$ time.*

Proof. We start by running S for all antimorphisms $\Theta \in G$ and by doing so we build $T_{w,\Theta}$ for all antimorphisms $\Theta \in G$ in $\mathcal{O}(r \cdot f(n))$ time.

We also compute all f_i^{-1} in $\mathcal{O}(r \cdot n)$ time.

If we now present a single-threaded algorithm which builds $T_{w,G}$ in $\mathcal{O}(r \cdot n)$ time using $T_{w,\Theta}$ trees, then by summing up we get a single-threaded algorithm which builds $T_{w,G}$ in $\mathcal{O}(r \cdot f(n) + r \cdot n)$ time. By combining this algorithm with lemma 4.16, we get a single-threaded algorithm which computes $\#\text{Pal}_G(w)$ in $\mathcal{O}(r \cdot f(n) + r \cdot n)$ time, which is exactly what we want.

Now we will present a single-threaded algorithm K which builds $T_{w,G}$ in $\mathcal{O}(r \cdot n)$ time using $T_{w,\Theta}$ trees, which will finish this theorem.

K will consecutively traverse all trees $T_{w,\Theta}$ in arbitrary order by depth-first traversal while trying to add new nodes into $T_{w,G}$.

When K is visiting node in $T_{w,\Theta}$ which represents Θ -palindrome $u' = u_1u_2 \dots u_{2n-1}$, then it will remember hash maps in $T_{w,G}$ which represents $\varepsilon, u_n, \dots, u_2u_3 \dots u_{2n-2}$ and u' .

If K goes up one node in $T_{w,\Theta}$, then it forgets the last of hash maps. This operation is done in constant time.

If K goes down one node in $T_{w,\Theta}$ from node h to node l using edge e with letter a , then there are two cases.

If the last of hash maps contains key $(a, \Theta(a))$ with value being a reference to map , then K adds map as the last hash map it currently remembers. This is done in constant time.

If the last of hash maps does not contain key $(a, \Theta(a))$, then K adds a new node into $T_{w,G}$. $b = \Theta(a)$. This new node $h_{a'v'b'}$ is a direct descendant of node $h_{v'}$ of the last hash map. Edge e from $h_{v'}$ to $h_{a'v'b'}$ has $a' = e_l = f_i^{-1}(a)$ and $b' = e_r = f_i^{-1}(b)$. To choose i in constant time, K needs to additionally remember for every hash map $map_{i,v'}$ one of indices i . After $h_{a'v'b'}$ is added, K needs to fill data for it and create new hash maps. To do this, K builds a hash map with keys being tuples $(mapRef_{i,v'}, f_i(a'), f_i(b'))$ for all i and values being a new hash map $newMap_i$. Then, for all i , K sets $mapRef_{i,a'v'b'} = ref(newMap_i)$ and adds to $map_{i,v'}$ element with key $(f_i(a'), f_i(b'))$ and value $mapRef_{i,a'v'b'}$. Overall time complexity of this operation of K is $\mathcal{O}(r)$.

Let's notice that in total K adds $\mathcal{O}(n)$ nodes into $T_{w,G}$ and operation of K which adds a node into $T_{w,G}$ costs $\mathcal{O}(r)$ time, so total time spent on such operations is $\mathcal{O}(r \cdot n)$. At the same time, all other operations of K have constant time complexity and their amount is $\mathcal{O}(r \cdot n)$. It follows that time complexity of K is $\mathcal{O}(r \cdot n)$.

As can be seen from definition of $K, T_{w,\Theta}$ and $T_{w,G}$, after K finishes traversal of $T_{w,\Theta}$ for some antimorphism Θ , all Θ -palindromic substrings of w' are represented by some hash map in $T_{w,G}$ and so all G -palindromic classes of equivalency $[v']$, where v' is a substring of w' and v' is a Θ -palindrome, are represented by some node in $T_{w,G}$.

Summing up, after K finishes traversal of all $T_{w,\Theta}$, $T_{w,G}$ contains nodes representing all G -palindromic classes of equivalency $[v']$, where v' is a substring of w' . ◀

► **Lemma 4.19.** *There exists an algorithm computing G -defect in $\mathcal{O}(r \cdot n)$ time.*

Proof. By combining algorithm described inside proof of theorem 4.9 and theorem 4.18, we get an algorithm which computes $\#\text{Pal}_G(w)$ in $\mathcal{O}(r \cdot n)$ time. By combining this algorithm with lemma 4.13, we get an algorithm which we want. ◀

Now we will present multi-threaded algorithms for computing $\#\text{Pal}_G(w)$ and G -defect.

► **Theorem 4.20.** *Let S be a single-threaded algorithm which builds $T_{w,\Theta}$ for arbitrary antimorphism $\Theta \in G$ in $\mathcal{O}(f(n))$ time, then there exists a multi-threaded algorithm P using $k \leq r$ threads which computes $\#\text{Pal}_G(w)$ in $\mathcal{O}\left(\frac{r \cdot f(n) + r \cdot n}{k}\right)$ time.*

Proof. We start by running S for all antimorphisms $\Theta \in G$ using k threads, up to k instances of S are running at the same. By doing so we build $T_{w,\Theta}$ for all antimorphisms $\Theta \in G$ in $\mathcal{O}\left(\frac{r \cdot f(n)}{k}\right)$ time.

We also compute all f_i^{-1} in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time.

If we now present a multi-threaded algorithm which uses $k \leq r$ threads and builds $T_{w,G}$ in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time using $T_{w,\Theta}$ trees, then by summing up we get a multi-threaded algorithm which uses $k \leq r$ threads and builds $T_{w,G}$ in $\mathcal{O}\left(\frac{r \cdot f(n) + r \cdot n}{k}\right)$ time. By combining this algorithm with lemma 4.17, we get a multi-threaded algorithm which uses $k \leq r$ threads and computes $\#\text{Pal}_G(w)$ in $\mathcal{O}\left(\frac{r \cdot f(n) + r \cdot n}{k} + n\right) = \mathcal{O}\left(\frac{r \cdot f(n) + r \cdot n}{k}\right)$ time, which is exactly what we want.

Now we will present a multi-threaded algorithm M which uses $k \leq r$ threads and builds $T_{w,G}$ in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time using $T_{w,\Theta}$ trees, which will finish this theorem.

We will be using an algorithm K which builds $T_{w,G}$ in $\mathcal{O}(r \cdot n)$ time described in theorem 4.18. K consecutively traverses all trees $T_{w,\Theta}$ by algorithm L .

M also traverses all trees $T_{w,\Theta}$ by algorithm L , but it runs up to k instances of L .

M uses concurrent hash maps in $T_{w,G}$ and for the additional hash map created when a new node is being added to $T_{w,G}$.

It is important to understand how M synchronizes between threads to achieve $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time complexity.

When a thread detects that it needs to add a new node, then it does not do it by itself, but synchronizes with other threads and reports that a new node needs to be added. This operation of reporting is done in constant time.

After each instance of L in M makes one step of traversal of its $T_{w,\Theta}$, all threads synchronize and check if there is at least one new node to add, reported by one of the threads. If there is a new node to add, then all threads participate in adding this new node, and then all threads rollback their state to state before the last step of traversal. Steps of traversal of $T_{w,\Theta}$ include ascending by one node, descending by one node without new node in $T_{w,G}$, descending by one node with new node in $T_{w,G}$.

Algorithm for adding a new node into $T_{w,G}$ is slightly modified in M compared to how it works in K . Instead of initializing as many as needed new hash maps, M initialized exactly r new concurrent hash maps. Some of these new concurrent hash maps will not be used at all after operation of adding new node to $T_{w,G}$ is finished, it is decided by which thread is faster to add a key into the concurrent hash map with keys $(\text{mapRef}_{i,v'}, f_i(a'), f_i(b'))$.

Synchronized operations in M all have constant time complexity. Initially, amount of synchronized operations is $\mathcal{O}(r \cdot n)$. Considering $\mathcal{O}(n)$ rollbacks of k operations, time spent on synchronized operations using k threads is $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$.

Operation in M to add new node into $T_{w,G}$ has $\mathcal{O}\left(\frac{r}{k}\right)$ time complexity.

It follows that time complexity of M is $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$, which is exactly what we wanted. ◀

► **Lemma 4.21.** *There exists a multi-threaded algorithm which uses $k \leq r$ threads and computes G -defect in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time.*

Proof. By combining algorithm described inside proof of theorem 4.9 and theorem 4.20, we get a multi-threaded algorithm which uses k threads and computes $\#\text{Pal}_G(w)$ in $\mathcal{O}\left(\frac{r \cdot n}{k}\right)$ time. By combining this algorithm with lemma 4.14, we get an algorithm which we want. ◀

4.2 Computing classical palindromic defect

Let's notice that the whole subsection 4.1.1 does not make any actual use from working with Θ -palindromes, except at the very start in lemma 4.1 and at the very end in lemma 4.10, lemma 4.11 and lemma 4.12. So if we skip these parts and read the rest while changing all mentions of Θ -palindromes into classical palindromes, then all statements will work for classical palindromes as well.

We are not going to rewrite whole subsection 4.1.1 here. Instead we will simply present results for classical palindromes as special cases of results, which were already achieved in subsection 4.1.1.

► **Note 4.22.** Classical palindromic defect of word w is defined as $|w|+1$ minus amount of pairwise distinct palindromic substrings of w . So to compute classical palindromic defect in $\mathcal{O}(f(n))$ time, where $f(n)$ is at least linear, it is enough to find amount of pairwise distinct palindromic substrings of w in $\mathcal{O}(f(n))$ time and then perform one addition and one subtraction.

► **Lemma 4.23.** *There exists an algorithm finding amount of pairwise distinct palindromic substrings of a word in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm and jump-pointers.*

Proof. Proof follows as special case of theorem 4.7 when Θ is chosen as the reversal mapping R . R is an antimorphism and $\forall a \in \mathcal{A} : R(a) = a$. ◀

► **Lemma 4.24.** *There exists an algorithm finding amount of pairwise distinct palindromic substrings of a word in $\mathcal{O}(n)$ time by combining Manacher's algorithm and advanced solution for dynamic level ancestor problem.*

Proof. Proof follows as special case of theorem 4.8 when Θ is chosen as the reversal mapping R . R is an antimorphism and $\forall a \in \mathcal{A} : R(a) = a$. ◀

► **Lemma 4.25.** *There exists an algorithm finding amount of pairwise distinct palindromic substrings of a word in $\mathcal{O}(n)$ time by using Manacher's algorithm and tree traversal.*

Proof. Proof follows as special case of theorem 4.9 when Θ is chosen as the reversal mapping R . R is an antimorphism and $\forall a \in \mathcal{A} : R(a) = a$. ◀

► **Lemma 4.26.** *There exists an algorithm computing classical palindromic defect in $\mathcal{O}(n \log n)$ time by combining Manacher's algorithm and jump-pointers.*

Proof. Combine lemma 4.23 and note 4.22. ◀

► **Lemma 4.27.** *There exists an algorithm computing classical palindromic defect in $\mathcal{O}(n)$ time by combining Manacher's algorithm and advance solution of dynamic level ancestor problem.*

Proof. Combine lemma 4.24 and note 4.22. ◀

► **Lemma 4.28.** *There exists an algorithm computing classical palindromic defect in $\mathcal{O}(n)$ time by using Manacher's algorithm and tree traversal.*

Proof. Combine lemma 4.25 and note 4.22. ◀

Implementation in SageMath

5.1 SageMath development

First of all, I started the development with steps described in chapter 1:

- Installed non-development version of SageMath for future testing purposes, and installed development version of SageMath.
- Created an issue #35495 on SageMath, which I intended to fix. During this step, I discovered the existence of another separate issue, #16366, which was also going to be fixed as part of my issue.
- Forked SageMath repository [3] and configured the local environment for development.

To implement a change in SageMath, it is important to first understand SageMath development methodologies and practices.

The software development approach commonly seen in large open-source projects like SageMath doesn't neatly fit into a single category like Agile [24] or Waterfall [25]. Instead, it typically combines elements from both methodologies and often creates a customized approach that suits the specific needs of open-source collaboration.

Rather than strictly following the sequential phases of Waterfall or the specific practices of Agile methodologies like Scrum or Kanban, such projects tend to embrace iterative development, collaboration, and responsiveness to change similar to Agile. Simultaneously, they might include aspects of long-term planning and structured organization similar to Waterfall methodologies, although in a more flexible and adaptable manner.

In essence, big open-source projects usually employ a hybrid approach, drawing from various methodologies while emphasizing adaptability, continuous improvement, and community collaboration as core principles in their software development process.

5.2 My contribution

As a one-time contributor to SageMath, there was no need for me to get much into Waterfall practices used in SageMath development, because the usage of these practices is mostly related to core developers and frequent contributors of SageMath. This left me only with responsibility to apply common Agile coding practices when developing my changes for SageMath.

One of such practices is known as **Collective Code Ownership** (see [26]). It is an important practice for the development of big open-source projects, which involve the work of a significant number of developers. Here is a brief description of this practice taken from [26]:

We are all responsible for all our code.

Collective code ownership means the team shares responsibility for its code. Rather than assigning modules, classes, or stories to specific individuals, the team owns it all. It's the right and responsibility to make improvements to any aspect of your team's code at any time.

In fact, improved code quality is one of the hidden benefits of collective code ownership. Collective ownership allows—no, expects—everyone to fix the problems they find. If you encounter duplication, unclear names, poor automation, or even poorly designed code, it doesn't matter who wrote it. It's your code. Fix it!

Keeping that in mind, my initial steps in implementation involved conducting research to locate the SageMath code which is related to the modifications I intended to make. It turned out that this code was mainly located inside *FiniteWord* class in word library, which is a part combinatorics library in SageMath. Once this part was completed, I proceeded not only with the implementation of my additional code but also with improvements of many lines within the existing codebase in the same scope as my changes. These changes simplified existing code, improved its readability, and enhanced the performance of several already existing methods, which will be mentioned in section 6.3.

5.3 Overview of changed and new methods

There are two key algorithms which were implemented in the following methods of *FiniteWord* class:

- `_get_palindromic_factors_data(self, f=None)`

Description: Private method which returns data about palindromic factors or f-palindromic factors of self.

Input: f - letter permutation on the alphabet of self.

Notes: Implements algorithm described in lemma 4.9 to build T_w tree graph, which is described in subsection 4.1.1. This method is used in many public methods.

- `g_defect(self, morphisms=[], antimorphisms=[])`

Description: Return the G -defect of self.

Input: morphisms - an iterable of letter permutations (default: []) on the alphabet of self. Letter permutations must be callable on letters as well as words (e.g. *WordMorphism*). If the identity morphism is not in morphisms, then it is added automatically.

Input: antimorphisms - an iterable of letter permutations (default: []) on the alphabet of self. If antimorphisms is empty, then antimorphism which acts as identity on letters is added to it. Letter permutations must be callable on letters as well as words (e.g. *WordMorphism*).

Notes: Implements algorithm described in lemma 4.19.

And here is a summary of existing public methods, which were changed during implementation:

- `lps(self, f=None)`

Description: Return the longest palindromic (or f-palindromic) suffix of self.

Input: f - letter permutation on the alphabet of self.

Notes: Refactored implementation to use `lps_lengths` method. Improved readability, decreased code complexity.

- `palindromic_lacunae_study(self, f=None)`

Description: Return interesting statistics about longest (f-)palindromic suffixes and lacunas of self.

Input: f - letter permutation on the alphabet of self.

Notes: Refactored to extract all necessary data from return value of `__get_palindromic_factors_data` method. Improved readability, decreased code complexity.

- `lacunas(self, f=None)`

Description: Return the list of all the lacunas of self.

Input: f - letter permutation on the alphabet of self.

Notes: Previously, it used `palindromic_lacunae_study` method, which was time inefficient. Refactored to use `__get_palindromic_factors_data` method instead. Improved time efficiency

- `length_maximal_palindrome(self, pos, f=None)`

Description: Return the length of the longest palindrome centered at position pos.

Input: f - letter permutation on the alphabet of self.

Input: pos - integer, position of the symmetry axis of the palindrome.

Notes: Refactored this method to PEP 8 and renamed inner variables, which were hard to read, because a lot of common variable names were used. For example, i, j, jj, m. Improved code readability.

- `lengths_maximal_palindromes(self, f=None)`

Description: Return the length of maximal palindromes centered at each position.

Input: f - letter permutation on the alphabet of self.

Notes: Refactored to use Manacher's algorithm instead of a trivial one.

- `lps_lengths(self, f=None)`

Description: Return the length of the longest palindromic suffix of each prefix.

Input: f - letter permutation on the alphabet of self.

Notes: Refactored to extract needed data from results of `lengths_maximal_palindromes` method. Improved code readability.

- `palindromic_complexity(self, n, f=None)`

Description: Return the number of distinct palindromic factors of length n of self.

Input: f - letter permutation on the alphabet of self.

Notes: Previously, it used `palindromes` method, which was time inefficient. Refactored to extract result data from output of `__get_palindromic_factors_data` method.

- `defect(self, f=None)`

Description: Return the defect of self.

Input: f - letter permutation on the alphabet of self.

Notes: Refactored to use algorithm described in lemma 4.12.

5.4 Results of implementation phase

In the end of implementation phase, all functional and non-functional requirements from chapter 3 were completed and ready for testing during the testing phase. At this point of time, expected asymptotic time complexity for method computing Θ -defect was $\mathcal{O}(n)$, because the algorithm was implemented as it is described in chapter 4, lemma 4.12. And expected asymptotic time complexity for method computing G -defect was $\mathcal{O}(n \cdot |G|)$, the implemented algorithm was as it is described in chapter 4, lemma 4.19.

Testing & Deployment

6.1 Hardware and Operating System Specification

Part of testing phase for this thesis included performance testing. Therefore, it is essential to provide a description of the hardware and operating system setup on which the testing was conducted.

6.1.1 Hardware info

- CPU: 6-core AMD Ryzen 5 5625U with Radeon Graphics
- RAM: 15.0 GiB
- GPU: AMD Barcelo
- Storage: 512 GB SSD x1

6.1.2 Operating system info

- Operating system: Linux Mint 21.2 Cinnamon
- Cinnamon Version: 5.8.4
- Linux Kernel: 5.15.0-88-generic

6.2 Functional Testing

Functional testing is a type of software testing that focuses on verifying that each function of the software application operates in accordance with the functional specifications. It validates whether the software performs the specific tasks and functionalities as intended, ensuring that the system meets the functional requirements outlined in the project or system documentation, which are also mentioned in section 3.1 in case of this thesis. Functional testing involves testing individual features by providing input and verifying the output against expected results. This type of testing aims to ensure that the software works as expected and delivers the intended functionalities.

More information about different types of tests and testing techniques can be found in book [27].

Functional testing for this thesis was done by adding automated tests using SageMath approach for implementing tests, which was described in section 1.7. Since existing methods were already covered by functional tests, the primary focus was on integrating new testing for computing the G -defect method. To achieve this, I adapted tests from already existing method for computing Θ -defect as tests for method computing G -defect, when G consists of identity morphism and a single antimorphic involution. On top of that, I added several non-trivial tests when G consists of several morphisms and antimorphisms.

After all that, I run all the affected existing and new tests. All tests passed, which concluded functional testing of this thesis.

6.3 Non-Functional Testing

Non-functional testing evaluates aspects of a software system that are not related to specific functionalities but are essential for the overall quality of the system. It focuses on attributes such as performance, reliability, usability, security, scalability, and other quality factors. Non-functional testing aims to assess how well the system behaves under various conditions, measures its response times, evaluates its security protocols, and checks its usability, among other aspects. Unlike functional testing, which verifies what the system does, non-functional testing evaluates how well the system performs.

In case of this thesis, crucial part of non-functional testing can be completed by verifying that the actual execution time of algorithms mentioned in non-functional requirements section 3.2 aligns with the final expected time complexities of these algorithms, as described in section 5.4.

For collecting execution times both on non-development version of SageMath and development version of SageMath I used Python *timeit* library inside interactive session of SageMath. More info about this library can be found, for example, in book [28].

6.3.1 Performance testing for existing methods

For existing methods with affected performance I chose to gather not only the actual execution times of the methods on my development branch but also to collect the same execution times on the current master branch of SageMath. In this approach, these times can be compared with each other to see achieved performance improvements. These methods include not only method for computing Θ -defect but also several others, as highlighted in chapter 5. The list of these affected methods includes:

- *lengths_unioccurent_lps*
- *defect*
- *palindromic_complexity*
- *lacunas*

For further insight into the functionalities of these methods, readers can refer to the SageMath documentation or examine these methods within the existing SageMath codebase.

I set required parameter n in method *palindromic_complexity* as 15 for all performance test runs.

The next problem was to select test data for taking performance measurement. Given that all affected methods are in some way related to the value of defect, I decided to use the following datasets:

- The first n characters of Thue-Morse word, where n equals 10^4 , 10^5 , 10^6 and 10^7 . For Thue-Morse word prefixes, defect equals approximately $\frac{1}{5}$ of the length of the prefix.
- The first n characters of Fibonacci word, where n equals 10^4 , 10^5 , 10^6 and 10^7 . Defect equals 0 for Fibonacci word prefixes.
- Random word from 3 different characters with lengths n , where n equals 10^4 , 10^5 , 10^6 and 10^7 . For such random words, the defect almost matches the length of the word, exceeding $\frac{19}{20}$ of the word's length across all test data.

Here are tables of execution times, collected as described above:

lengths_unioccurent_lps	SageMath master	My branch
Thue-Morse 10^4	0.12 s	0.08 s
Thue-Morse 10^5	2.94 s	0.75 s
Thue-Morse 10^6	528 s	7.72 s
Thue-Morse 10^7	>15 minutes	77.4 s
Fibonacci 10^4	0.14 s	0.13 s
Fibonacci 10^5	1.26 s	0.91 s
Fibonacci 10^6	57 s	7.81 s
Fibonacci 10^7	>15 minutes	80.4 s
Random 10^4	0.33 s	0.16 s
Random 10^5	2.96 s	1.56 s
Random 10^6	41.2 s	15.7 s
Random 10^7	>15 minutes	151 s

■ **Table 6.1** Performance of lengths_unioccurent_lps method

defect	SageMath master	My branch
Thue-Morse 10^4	0.05 s	0.11 s
Thue-Morse 10^5	1.11 s	0.71 s
Thue-Morse 10^6	308 s	6.8 s
Thue-Morse 10^7	>15 minutes	67.4 s
Fibonacci 10^4	0.08 s	0.11 s
Fibonacci 10^5	0.49 s	0.79 s
Fibonacci 10^6	34.6 s	6.83 s
Fibonacci 10^7	>15 minutes	68.4 s
Random 10^4	0.11 s	0.15 s
Random 10^5	0.71 s	1.48 s
Random 10^6	6.9 s	14.6 s
Random 10^7	70.3 s	146 s

■ **Table 6.2** Performance of defect method

palindromic_complexity	SageMath master	My branch
Thue-Morse 10^4	0.08 s	0.11 s
Thue-Morse 10^5	1.19 s	0.69 s
Thue-Morse 10^6	313 s	6.97 s
Thue-Morse 10^7	>15 minutes	68.2 s
Fibonacci 10^4	0.04 s	0.11 s
Fibonacci 10^5	0.65 s	0.75 s
Fibonacci 10^6	36 s	7.39 s
Fibonacci 10^7	>15 minutes	75.1 s
Random 10^4	0.12 s	0.15 s
Random 10^5	0.72 s	1.52 s
Random 10^6	6.8 s	15 s
Random 10^7	68.4 s	150 s

■ **Table 6.3** Performance of palindromic_complexity method

lacunas	SageMath master	My branch
Thue-Morse 10^4	0.12 s	0.08 s
Thue-Morse 10^5	2.49 s	0.69 s
Thue-Morse 10^6	522 s	6.69 s
Thue-Morse 10^7	>15 minutes	67.2 s
Fibonacci 10^4	0.11 s	0.08 s
Fibonacci 10^5	0.95 s	0.74 s
Fibonacci 10^6	53.9 s	7.46 s
Fibonacci 10^7	>15 minutes	71.8 s
Random 10^4	0.27 s	0.19 s
Random 10^5	2.4 s	1.52 s
Random 10^6	35.6 s	14.9 s
Random 10^7	>15 minutes	151 s

■ **Table 6.4** Performance of lacunas method

From above tables we can make three following observations.

Firstly, the methods *defect* and *palindromic_complexity* showed slightly better performance than the new algorithms on certain test data, although this advantage was not consistent across all possible inputs. On some test data, old algorithms declined in speed much more rapidly compared to the new algorithms. By analyzing the old algorithms codebase and observing the performance data, we can state that the old algorithms performed well on words with defect value near the length of the word.

Secondly, the results of performance testing confirm that expected asymptotic time complexity for method *defect* is $\mathcal{O}(n)$, as it was assumed in the end of implementation phase in section 5.4.

Lastly, performance measurements make it clear that the asymptotic time complexities for all four methods were improved. Now these asymptotic time complexities are either linear or close to being linear.

6.3.2 Performance testing for G -defect computing method

Another part of performance testing was to clarify that the actual time execution for G -defect computing method aligns with the expected asymptotic time complexity of this method, as outlined in the conclusion of the implementation phase in section 5.4. This expected asymptotic time complexity equals $\mathcal{O}(n \cdot |G|)$, where n is length of the word.

To validate that the actual asymptotic time complexity corresponds to a formula involving two parameters, we need to select two different datasets. One dataset maintains the first parameter fixed while varying the second parameter, and vice versa. For all datasets, I decided to generate random words using 4 different letters. The first dataset has $n = 10^6$ and $|G| = 2, 4, 8, 48$. The second dataset has $|G| = 48$ and $n = 10^3, 10^4, 10^5, 10^6$.

Here are tables of execution times, collected as described above:

G-defect	$n = 10^6$
$ G = 2$	12.9 s
$ G = 4$	27.6 s
$ G = 8$	52.3 s
$ G = 48$	274 s

■ **Table 6.5** Performance of G-defect method based on G

G-defect	$ G = 48$
$n = 10^3$	0.37 s
$n = 10^4$	2.84 s
$n = 10^5$	28.5 s
$n = 10^6$	274 s

■ **Table 6.6** Performance of G-defect method based on the word length

These performance testing results align with the expected time complexity of G -defect computing method, as they show a linear growth pattern based on either n or $|G|$ when the other parameter is constant.

6.4 Deployment

Upon the successful completion of all testing phase steps with all tests passing, I initiated the deployment process for my SageMath fork. The following steps outline the procedure I followed:

- Merged all updates from the SageMath master into my branch. I encountered and resolved some issues during this step due to additional syntax checks introduced by SageMath while I have been working on other phases of this thesis.
- Run all automated functional tests and checked that they all passed.
- Pushed my local branch into my fork on GitHub.
- Created a pull request for my fork into SageMath.

Currently, I am awaiting the review of my pull request, anticipating the final deployment of my changes into the SageMath repository.

Conclusion

This thesis significantly boosted palindrome library inside SageMath by speeding up existing methods and adding a highly efficient new algorithm.

Through Agile methods, robust testing, and optimization, part of SageMath related to palindromes underwent significant improvements. These advances mean faster execution times and expanded functionalities for researchers.

All functional and non-functional requirements of this thesis from chapter 3 were fulfilled.

Other than that, this thesis has contributed to the growing body of knowledge on the topic of algorithms dealing with generalized palindromes and generalized palindromic defects, and SageMath implementation of these time efficient algorithms offers valuable insights for future researchers in this field.

Although, there are still some open theoretical questions worth mentioning left at the time of writing this thesis. Here is a list of some of these questions:

1. Let's take the algorithm described in theorem 4.7, but instead of using jump-pointers just go m nodes up the T_w in $\mathcal{O}(m)$ time when getting m th ancestor of a node in T_w . This algorithm finds amount of pairwise distinct Θ -palindromic substrings of w . For sure it has time complexity $\mathcal{O}(n^2)$. Does this algorithm actually have $\Theta(n^2)$ time complexity or is its actual time complexity asymptotically faster than $\Theta(n^2)$? If it turned out that actual time complexity of this algorithm is $\mathcal{O}(n)$, then there would be a much more simple version of $\mathcal{O}(n)$ time algorithm computing Θ -defect compared to algorithms presented in this thesis.
2. What is the best achievable time complexity of a single-threaded algorithm computing G -defect? The fastest algorithm presented in this thesis has $\mathcal{O}(|w| \cdot |G|)$ time complexity.
3. What time complexities can be achieved for multi-threaded algorithms computing Θ -defect? Based on these algorithms, what promising multi-threaded algorithms computing G -defect can be presented and what time complexities will they have?
4. Can algorithms for computing Θ -defect be generalized for antimorphisms which are not letter permutations? Should definition of Θ -defect be modified for such antimorphisms compared to definition 2.1?
5. During the writing of the thesis we discovered an article by Rubinchik and Shur [29] which describes a rather sophisticated algorithm counting palindromes in a word in $\mathcal{O}(n \log n)$ time. As this article is very technical, we have not included any details from it. It is a natural question whether this approach can be generalized to count G -palindromes, or if it can be used to improve our approach.

Bibliography

1. PELANTOVA, Edita; STAROSTA, Štěpán. Palindromic richness for languages invariant under more symmetries. *Theoretical computer science*. 2014, vol. 518, pp. 42–63. ISBN 0304-3975.
2. THE SAGEMATH DEVELOPERS. *SageMath*. 2022. Version 9.5. Available from DOI: 10.5281/zenodo.593563.
3. ROMANENKO, I. *SageMath improvements based on this thesis*. 2023. Available also from: <https://github.com/giovanni-romanenko/sage-35495>.
4. SAGEMATH DEVELOPMENT TEAM. *SageMath: Open Source Mathematics Software*. The Mathematical Association of America, 2005. Available also from: <https://www.maa.org/press/maa-reviews/sage-open-source-mathematics-software>.
5. SAGEMATH DEVELOPMENT TEAM. *SageMath*. 2023. Available also from: <https://www.sagemath.org/>.
6. SAGEMATH DEVELOPMENT TEAM. *SageMath Documentation*. 2023. Available also from: <https://doc.sagemath.org/>.
7. ZIMMERMANN, Paul; CASAMAYOU, Alexandre; COHEN, Nathann; CONNAN, Guillaume; DUMONT, Thierry; FOUSSE, Laurent; MALTEY, François; MEULIEN, Matthias; MEZZAROBBA, Marc; PERNET, Clément; THIÉRY, Nicolas M.; BRAY, Erik; CREMONA, John; FORETS, Marcelo; GHITZA, Alexandru; THOMAS, Hugh. *Computational Mathematics with SageMath*. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2018. Available from DOI: 10.1137/1.9781611975468.
8. DROUBAY, Xavier; JUSTIN, Jacques; PIRILLO, Giuseppe. Episturmian words and some constructions of de Luca and Rauzy. *Theoretical computer science*. 2001, vol. 255, no. 1, pp. 539–553. ISBN 0304-3975.
9. GARZON, Max H.; YAN, Hao. Watson-Crick Conjugate and Commutative Words. In: Germany: Springer Berlin / Heidelberg, 2008, vol. 4848. DNA Computing. ISBN 9783540779612; 3540779612;
10. PELANTOVA, Edita; STAROSTA, Štěpán. Languages invariant under more symmetries: Overlapping factors versus palindromic richness. *Discrete mathematics*. 2013, vol. 313, no. 21, pp. 2432–2445. ISBN 0012-365X.
11. MANACHER, Glenn. A New Linear-Time “On-Line” Algorithm for Finding the Smallest Initial Palindrome of a String. *Journal of the ACM*. 1975, vol. 22, no. 3, pp. 346–351. ISBN 0004-5411.
12. BENDER, Michael A.; FARACH-COLTON, Martín. The Level Ancestor Problem simplified. *Theoretical computer science*. 2004, vol. 321, no. 1, pp. 5–12. ISBN 0304-3975.

13. ALSTRUP, Stephen; HOLM, Jacob. Improved Algorithms for Finding Level Ancestors in Dynamic Trees. In: ed. by MONTANARI, U.; ROLIM, JDP; WELZL, E. Berlin, Heidelberg: Springer Berlin Heidelberg, 2000, vol. 1853, pp. 73–84. AUTOMATA LANGUAGES AND PROGRAMMING. ISBN 0302-9743.
14. ANACONDA, INC. *Conda*. 2023. Available also from: <https://conda.io>.
15. CONDA-FORGE CONTRIBUTORS. *conda-forge*. 2023. Available also from: <https://conda-forge.org>.
16. SAGEMATH DEVELOPMENT TEAM. *SageMath Installation Documentation: Conda-forge*. 2023. Available also from: <https://doc.sagemath.org/html/en/installation/conda.html#sec-installation-conda>.
17. SAGEMATH DEVELOPMENT TEAM. *SageMath Installation Documentation: Non-development version*. 2023. Available also from: <https://doc.sagemath.org/html/en/installation/conda.html#installing-all-of-sagemath-from-conda-not-for-development>.
18. SAGEMATH DEVELOPMENT TEAM. *SageMath Installation Documentation: Development version*. 2023. Available also from: <https://doc.sagemath.org/html/en/installation/conda.html#sec-installation-conda-develop>.
19. GITHUB, INC. *GitHub*. 2023. Available also from: <https://github.com>.
20. WIKIPEDIA CONTRIBUTORS. *Doctest — Wikipedia, The Free Encyclopedia*. 2023. Available also from: <https://en.wikipedia.org/wiki/Doctest>.
21. SAGEMATH DEVELOPMENT TEAM. *SageMath Documentation: Running Sage’s Doctests*. 2023. Available also from: <https://doc.sagemath.org/html/en/developer/doctesting.html>.
22. SOMMERVILLE, I. *Software Engineering*. Pearson Education, 2015. ISBN 9780133943276.
23. PRESSMAN, R.S.; BRUCE R. MAXIM, D. *Software Engineering: A Practitioner’s Approach*. McGraw-Hill Education, 2014. ISBN 9780078022128.
24. RIGBY, D.; ELK, S.; BEREZ, S. *Doing Agile Right: Transformation Without Chaos*. Harvard Business Review Press, 2020. ISBN 9781633698710.
25. BLOKDYK, G. *Waterfall Model: Master the Art of Design Patterns*. CreateSpace Independent Publishing Platform, 2017. ISBN 9781979444217.
26. SHORE, J.; WARDEN, S.; SAFARI, an O’Reilly Media Company. *The Art of Agile Development, 2nd Edition*. O’Reilly Media, Incorporated, 2021. ISBN 9781492080695.
27. JORGENSEN, P.C. *Software Testing: A Craftsman’s Approach, Second Edition*. Taylor & Francis, 2002. Software engineering series. ISBN 9780849308093.
28. BEAZLEY, David; JONES, Brian K. *Python Cookbook*. O’Reilly Media, Inc., 2013. ISBN 1449340377.
29. RUBINCHIK, Mikhail; SHUR, Arseny M. Counting Palindromes in Substrings. In: FICI, Gabriele; SCIORTINO, Marinella; VENTURINI, Rossano (eds.). *String Processing and Information Retrieval*. Cham: Springer International Publishing, 2017, pp. 290–303. ISBN 978-3-319-67428-5.