

AARHUS UNIVERSITY

BIOINFORMATICS RESEARCH CENTER

---

**Range Minimum Queries and Interval-Trees for Genome  
Scale String Processing**

---

Pediotidis-Maniatis Dimitrios  
202002459

*Supervisor:*  
Thomas Mailund



Master's Thesis in Bioinformatics

June, 2022

# Acknowledgements

I would like to thank my supervisor Thomas Mailund for his guidance, patience and support throughout this thesis.

Next, I would like to thank my colleagues for the amazing environment they provided. Sharing this experience with them everyday is something I will treasure.

And of course, I would like to thank my friends and family back home. Their unwavering support and encouragement under all circumstances are what eventually made this work a reality.

# Abstract

Suffix Trees have been considered one of the most important data structure in string analysis, since their conception in 1973. Nevertheless, their large memory requirements have pushed for the conception of more succinct data structures. This thesis analyzes the concepts of LCP-Interval trees and Range Minimum Queries as an alternative to suffix trees. More specifically we are concerned with the simulation of a top-down traversal of the suffix tree to perform exact pattern matching. In our approach, we first create the necessary framework. We analyze the construction of the suffix and lcp arrays using the SA-IS and Kasai algorithms respectively. Then, we proceed to explore the LCP-interval trees and different implementations to the Range Minimum Query problem. We end the analysis of Range Minimum Queries with the Fischer-Heun structure, which achieves constant time queries with linear preprocessing of the lcp array. In our experiments we aimed to verify the relationship between the Fischer-Heun structure for Range Minimum Queries and two other implementations. Namely, the sparse table and a hybrid solution which allows queries in logarithmic time. Lastly, we compare two exact pattern matching implementations. The first makes use of the Range Minimum Queries and the LCP-interval tree. The second makes use of the suffix tree. These were implemented in Python programming language.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Suffix Trees and the Suffix Array . . . . .	1
1.2	The LCP array and the LCP-interval tree . . . . .	2
1.3	The Range Minimum Query Problem . . . . .	2
1.4	Contents of this Project . . . . .	3
<b>2</b>	<b>Preliminaries</b>	<b>4</b>
<b>3</b>	<b>SA-IS and Kasai Algorithms</b>	<b>6</b>
3.1	SA-IS Algorithm . . . . .	6
3.2	Kasai Algorithm . . . . .	14
<b>4</b>	<b>LCP-interval trees</b>	<b>17</b>
4.1	Suffix and LCP arrays in tree traversal . . . . .	17
4.2	The LCP-Interval Tree . . . . .	19
4.3	Range Minimum Queries in Tree Traversal . . . . .	22
<b>5</b>	<b>Range Minimum Query</b>	<b>25</b>
5.1	Table Preprocessing RMQ . . . . .	26
5.2	Sparse table Preprocessing RMQ . . . . .	27
5.3	A Hybrid Solution to RMQ . . . . .	30
5.4	The Fischer Heun Structure . . . . .	34
<b>6</b>	<b>Experiments</b>	<b>39</b>
6.1	Running Times . . . . .	39
6.2	Pattern Search . . . . .	40
<b>7</b>	<b>Discussion</b>	<b>42</b>
7.1	Range Minimum Queries Implementations . . . . .	42
7.2	Pattern Search Implementations . . . . .	43

7.3	Remarks . . . . .	43
<b>8</b>	<b>Conclusion</b>	<b>45</b>
8.1	Summary . . . . .	45
8.2	Evaluation . . . . .	45
8.3	Future Work . . . . .	46
	<b>Bibliography</b>	<b>47</b>