

# Advanced Data Structures: Final Report

Steil, Patrick ✉

2513240

## 1 Introduction

In this work, different data structures are implemented that can efficiently answer two types of queries: predecessor data structures and range-minimum queries. For this purpose, the *Elias-Fano coding* and several “Block-Lookup” data structures are implemented and tested. The project is part of the lecture “Advanced Data Structures” at KIT in the summer semester 2023 with Dr. Florian Kurpicz.

## 2 Algorithms and Data Structures

### Predecessor data structures

Predecessor data structures are special data structures that were developed to efficiently solve the problem of predecessor access. The predecessor access problem is to find the predecessor (the element with the highest value that is less than a given element) in a set of elements. Without additional memory consumption, elements can be stored sorted or in a binary search tree, whereby a predecessor query can subsequently be answered with a runtime of  $\mathcal{O}(\log(n))$ . However, this is too slow in practice, which is why techniques such as *x/y-fast-tries* [2] or *Elias-Fano encoding* [1] have been developed. The idea of the last technique is to divide the representation of numbers into two parts, on the one hand the *most significant* bits with a length of  $\lceil \log(n) \rceil$  and on the other hand the remaining bits. This way, all numbers with the same *most significant* bits are “grouped” during a query, which means that a query only has to look at one “group” of numbers. However, this may result in all numbers having to be scanned if the numbers are not favourable.

### Range minimum queries

Range minimum queries (RMQ) refer to a special type of search operation in data structures where the minimum value within a given interval (range) of elements in a set is to be found. The goal is to efficiently find the smallest value among a subset of elements. Analogous to the predecessor data structures, there are trade-offs between additional memory and runtime (the different pre-computation times can be neglected for static data). To be able to answer RMQ quickly, one can use  $\Theta(n^2)$  additional memory to calculate all possible queries in advance and thus obtain a constant query time (a lookup). To reduce the impractical quadratic memory, one can now only precompute queries of certain ranges. If the size of this range is cleverly chosen, one obtains constant query time and down to only linear additional memory.

## 3 Implementation

### 3.1 Elias-Fano

For an efficient Elias-Fano implementation, one needs a bit-vector data structure that can efficiently answer SELECT query. The actual Elias-Fano data structure then consists of two bit vectors: *upperBits* and *lowerBits*. The former encodes in  $2n + 1$  bits the *most significant bits* (MSB) of the numbers, and in *lowerBits* the remaining bits (LSB) of the numbers are stored straightforward. Please see the Appendix A on the “minor” Elias-Fano Problem which I mentioned to you in person.



© P. Steil;  
licensed under Creative Commons License CC-BY 4.0



Leibniz International Proceedings in Informatics  
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Otherwise, my implementation of `PREDECESSOR` is based on the description of <https://arxiv.org/pdf/2003.11835.pdf>. To identify the predecessor of the number  $x$ , use `SELECT` queries to find the “leftmost” and “rightmost” index of numbers for which the MSB is equal to the MSB of  $x$ . Subsequently, it is just a matter of scanning this range to find the predecessor.

## 3.2 Range-Mimimum Queries

For the range-minimum queries I implemented all 3 data structures, i.e. the one with  $\mathcal{O}(n^2)$ ,  $\mathcal{O}(n \log(n))$  and the one with  $\mathcal{O}(n)$  additional memory. I have not implemented the latter data structure exactly as in the lecture, namely I have not implemented the *Cartesian Trees*. I have, for queries that run between block boundaries, split the query into three parts. The initial and final part (from the left index boundary to the first “completely covered” block and analogously from the last covered block to the end) and the covered blocks in between. The initial and final check is usually done with Cartesian trees, I have implemented a linear search over the part of the array.

## 4 Experiments

### 4.1 Hardware & Code

The programming language used was C++, and the compiler was g++ with optimisations enabled. All experiments were conducted on a Macbook Pro with a 2.3 GHz dual-core Intel Core i5 processor.

### 4.2 Elias-Fano

#### References

- 1 Peter Elias. Efficient storage and retrieval by content and address of static files. *J. ACM*, 21(2):246–260, apr 1974. doi:10.1145/321812.321820.
- 2 Dan E. Willard. Log-logarithmic worst-case range queries are possible in space  $\theta(n)$ . *Information Processing Letters*, 17(2):81–84, 1983. URL: <https://www.sciencedirect.com/science/article/pii/0020019083900753>, doi:[https://doi.org/10.1016/0020-0190\(83\)90075-3](https://doi.org/10.1016/0020-0190(83)90075-3).

## A A Note on Elias-Fano

In the script, the first  $\lceil \log(n) \rceil$  bits were defined as MSB, but this definition caused problems with e.g. the following numbers:

0, 1, 10, 13, 32, 64, 1926720561250547757, 4354482840745948438, 18446013172173323708

Now, looking at the various MSBs of length  $\lceil \log(n) \rceil = 4$ , we see:

$$(0000)_2, (0001)_2, (0011)_2, (1111)_2 \tag{1}$$

If *upperBits* has a length of  $2 \cdot 9 + 1 = 19$ , then to decode the last number, you would need to access the index  $(1111)_2 + 8 = 39$ , which is out of bounds for *upperBits*. Therefore, I set the length of the *upperBits* to  $3 \cdot n$ , as discussed in the lecture.