

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 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” lecture at KIT in the summer semester of 2023 with Dr. Florian Kurpicz.

2 Algorithms and Data Structures

Predecessor data structures

Predecessor data structures are unique data structures 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, meaning a query only has to look at one “group” of numbers. However, this may result in all numbers being 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 current Elias-Fano data structure 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 straightforwardly. Please see Appendix A on the “minor” Elias-Fano Problem, which I mentioned in person. Otherwise, my implementation of PREDECESSOR is based on the



© 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

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, scanning this range to find the predecessor is straightforward.

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 precisely as in the lecture; namely, I have not implemented the *Cartesian Trees*. I have split the query into three parts for queries that run between block boundaries. 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 ThinkCentre running Ubuntu 22.10 with an Intel(R) Core(TM) i5-8500T CPU @ 2.10GHz processor.

4.2 Elias-Fano

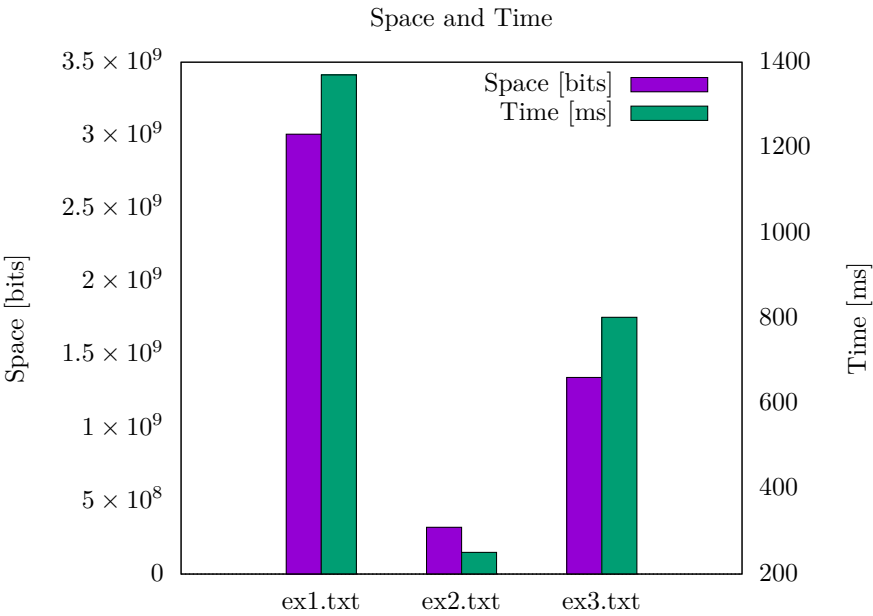
The performance (based on the provided benchmarks) of the Elias-Fano coding can be seen in Figure 1. It can be noticed that both the memory consumption and the query depend on the distribution of the numbers. In the first file *ex1.txt*, 1 000 000 numbers are present in the *total* range of $[0, 2^{64})$, i.e., scattered throughout the universe. On the other hand, in the file *ex2.txt*, only “small” numbers are present, i.e., very one-sidedly distributed in the universe. Many most significant bits are 0 for all numbers, i.e., storing this “unnecessary” information is unnecessary.

4.3 Range-Minimum Queries

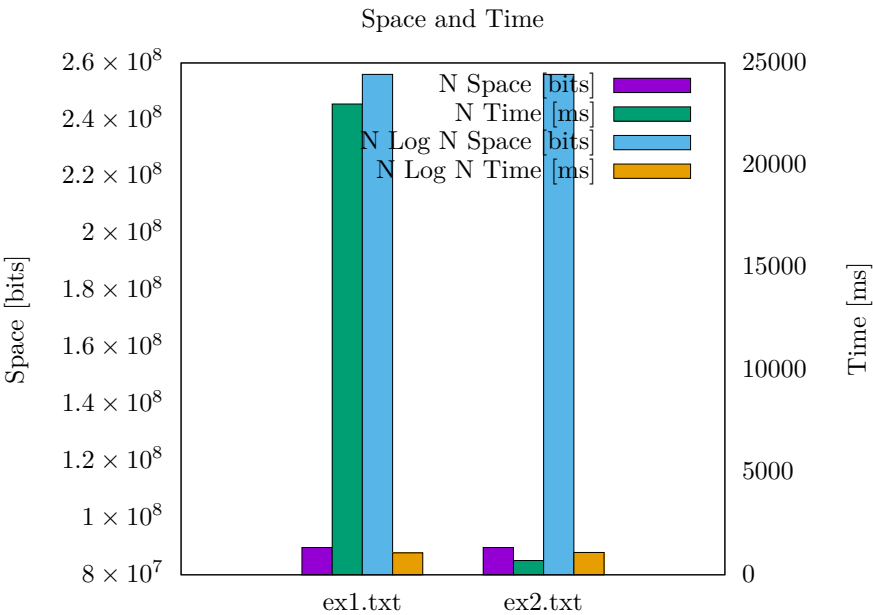
Figure 2 visualises the performance of the RMQ data structures (DS), namely the two with $\mathcal{O}(n \log n)$ and $\mathcal{O}(n)$ additional memory consumption. The data structure with $\mathcal{O}(n^2)$ memory consumption could not be executed due to too high memory consumption. It can be seen that the memory consumption between the two DSs is very different, namely, as expected, with significantly less memory for the DS with linear memory overhead. Interestingly, the query performance varies significantly between the DS with linear memory overhead because the ranges in the second file *ex2.txt* are very close to each other. In my implementation, this only calls a local sequential search, which is the most cache-efficient and, for small ranges, the fastest way. Otherwise, the query performance of the $\mathcal{O}(n \log n)$ DS does not depend on the range since the query performs lookups “independent” of the given range.

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.



■ **Figure 1** Space consumption and query time (for 1 000 000 predecessor queries) for my implementation of Elias-Fano. The three files are the three provided benchmark files.



■ **Figure 2** Shown is the memory consumption and performance of the range-minimum query data structures $N \log N$ and N on the benchmark files. Note: the $\mathcal{O}(n^2)$ data structure could not be executed on the benchmarks due to high memory consumption.

82 2 Dan E. Willard. Log-logarithmic worst-case range queries are possible in space $\theta(n)$. *Informa-*
 83 *tion Processing Letters*, 17(2):81–84, 1983. URL: [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/0020019083900753)
 84 [article/pii/0020019083900753](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).

85 **A A Note on Elias-Fano**

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

88 0, 1, 10, 13, 32, 64, 1926720561250547757, 4354482840745948438, 18446013172173323708

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

$$90 \quad (0000)_2, (0001)_2, (0011)_2, (1111)_2 \quad (1)$$

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