

Master's Thesis

Fast and Space Efficient Wavelet Tree Construction

Overview

The *FM-index* combines two data structures: the *Burrows-Wheeler transform* and *wavelet trees*. It is a very prominent full-text index and used in most DNA read aligners [4] and in Bioinformatics in general. In this Master's thesis, we focus on the efficient construction of the second data structure—wavelet trees.

The wavelet tree is a binary tree data structure that can be used to answer *rank* and *select* queries on texts of size n over an alphabet of size σ in $O(\lg \sigma)$ time. Here, $rank_{\alpha}(i)$ queries ask for the number of occurrences of the symbol α before the position i and $select_{\alpha}(i)$ queries return the text position of the i -th occurrence of the symbol α .

Let T be a text of length n over an alphabet of size σ . The corresponding wavelet tree consists of $\lceil \lg \sigma \rceil$ bit vectors of size n , see Fig. 1. Even though all $n \lceil \lg \sigma \rceil$ entries in the bit vectors have to be looked at during construction, the wavelet tree can be computed in $O(n \lg \sigma / \sqrt{\lg n})$ time using broadword programming [1, 5]. There exists an implementation of such an algorithm by Kaneta [3], which heavily relies on specialized CPU instructions like *parallel bit extract* and *packed shuffle bytes*. The reported construction times are faster than the previously fastest sequential WT construction algorithm [2]. However, the algorithm has one significant disadvantage—it requires twice as much memory.

Objective

The main objective of this Master's thesis is to develop a fast *and* space efficient wavelet tree construction algorithm that computes the wavelet tree in $O(n \lg \sigma / \sqrt{\ln n})$ time using specialized CPU instructions.

The parallelization of this algorithm is another (minor) goal of this Master's thesis. There exists a very fast meta-algorithm that can be used to parallelize all wavelet tree construction algorithms, e.g., [2] and it should be experimentally evaluated whether a parallelization of the algorithm described above is faster than using the meta-algorithm for the parallelization.

Requirements

- Excellent C++ programming skills
- Interest in string algorithms and compact data structures

Contact

Florian Kurpicz (kurpicz@kit.edu)

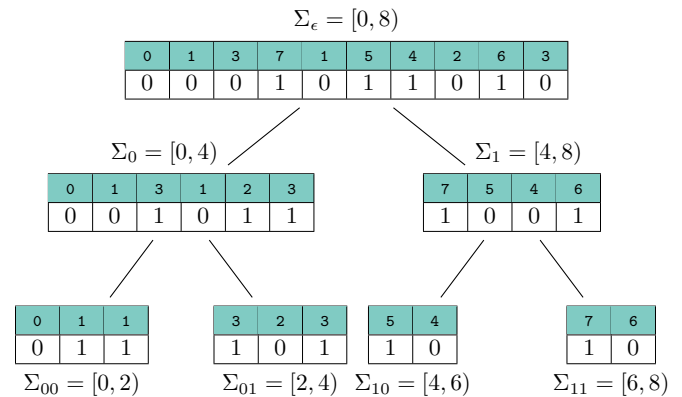


Figure 1: The wavelet tree of $T = [0, 1, 3, 7, 1, 5, 4, 2, 6, 3]$. The light teal (●) arrays contain the characters represented at the corresponding position in the bit vector and are not a part of the wavelet tree. Note that all bit vectors on the same depth can be concatenated to a single bit vector, while retaining the same functionality. Σ_{α} denotes the characters that are represented by the bit vector for $\alpha \in \{\epsilon, 0, 1, 00, 01, 10, 11\}$. All this auxiliary information is not stored explicitly.

References

- [1] Maxim A. Babenko, Munro, Pawel Gawrychowski, Tomasz Kociumaka, and Tatiana Starikovskaya. Wavelet trees meet suffix trees. In *SODA*, pages 572–591. SIAM, 2015.
- [2] Johannes Fischer, Florian Kurpicz, and Marvin Löbel. Simple, fast and lightweight parallel wavelet tree construction. In *ALENEX*, pages 9–20. SIAM, 2018.
- [3] Yusaku Kaneta. Fast wavelet tree construction in practice. In *SPIRE*, volume 11147 of *Lecture Notes in Computer Science*, pages 218–232. Springer, 2018.
- [4] Ben Langmead and Steven L. Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*, 9(4):357, 2012.
- [5] J. Ian Munro, Yakov Nekrich, and Jeffrey Scott Vitter. Fast construction of wavelet trees. *Theor. Comput. Sci.*, 638:91–97, 2016.