



Università
Ca' Foscari
Venezia

Artificial Intelligence and Data Engineering

Master Degree Thesis

A New Compression Technique for Repetitive Tries

Supervisor

Nicola Prezza

Co-supervisor

Alessio Campanelli

Author

Davide Tonetto

Student ID 884585

Academic year

2024/2025

Abstract

Contents

Abstract	ii
1 Introduction	1
1.1 Challenges and Contributions	1
1.2 Structure of The Thesis	3
2 Preliminary Concepts	5
2.1 Notation	5
2.2 Labeled Trees and Tries	7
2.2.1 Introduction and Motivation	7
2.2.2 Applications	7
2.2.3 Indexing	8
2.2.4 Information-Theoretic Lower Bound	9
2.2.5 State of The Art	10
2.2.6 Conclusion	12
2.3 The Extended Burrows-Wheeler Transform	12
2.3.1 Introduction and Motivation	12
2.3.2 Definition	13
2.3.3 Properties	14
2.3.4 Construction	16
2.3.5 Inversion	19
2.3.6 Compressing Labeled Trees	21
2.3.7 Indexing a Compressed Labeled Tree	22
2.3.8 Implementation	26
2.3.9 Experiments	28
2.4 Finite State Automata	30
2.4.1 Introduction and Motivation	30
2.4.2 DFA Minimization	31
2.4.3 Revuz' Minimization Algorithm	32
2.5 Wheeler and p -sortable Graphs	36
2.5.1 Introduction and Motivation	36
2.5.2 Orders	36
2.5.3 Wheeler Graphs	37
2.5.4 The Co-lex Width of an Automaton	38
2.5.5 Indexing Finite State Automata	41
2.6 Min-Weight Perfect Bipartite Matching	41
2.6.1 Introduction and Motivation	41
2.6.2 Bipartite Graphs	42
2.6.3 Problem Definition	42
2.6.4 Hall's Marriage Theorem	44
2.6.5 State of the Art	45
2.6.6 Implementation	45

3	Tree Compression Scheme	47
3.1	Compression Scheme Pipeline	47
3.2	Reducing the String Partitioning Problem to the MWPBM Problem .	48
3.2.1	Bipartite Graph Construction	49
3.2.2	Proof	52
3.3	Collapsing Nodes in Chains	58
3.3.1	How to Collapse Nodes	58
3.3.2	Language Equivalence	60
4	Experiments	64
	Bibliography	66
	Web bibliography	69

Chapter 1

Introduction

The problem of compressing large sets of strings, or finite languages, is a fundamental challenge in computer science with applications in areas like bioinformatics, natural language processing, and data indexing. A finite language can be naturally represented by an acyclic deterministic finite automaton (ADFA), commonly known as a trie. In this representation, each string in the language corresponds to a unique path from the root to a final state. Compressing the language is therefore equivalent to compressing its corresponding trie structure.

Traditional compression algorithms often fail to exploit the inherent structural properties of tries. To address this, specialized techniques have been developed. Among the most prominent is the *Extended Burrows-Wheeler Transform (XBWT)* [12], which extends the classical Burrows-Wheeler Transform [4] to labeled trees and can be applied to tries to achieve significant compression by capturing their structural regularities.

However, existing techniques may not be optimal when dealing with tries that exhibit a high degree of repetitiveness. Such is the case for languages containing many strings with shared substrings, leading to tries with large, identical subtrees. Many real-world datasets, such as genomic databases or dictionaries of related terms, generate such highly repetitive structures. This thesis introduces and analyzes a novel compression technique specifically designed to exploit these repetitions. The core idea is to identify and merge identical subtrees by reducing the trie to its minimal deterministic finite automaton (DFA) representation. We implement this method and evaluate its performance against state-of-the-art approaches like XBWT, assessing its effectiveness on various datasets.

1.1 Challenges and Contributions

The primary goal of this thesis is to develop a data structure that both compresses a given finite language and efficiently supports indexing queries, such as navigational and subpath queries (see Definition 8). This challenge involves navigating a fundamental trade-off between compression and indexability, which we explore in detail in Section 2.5. Two straightforward approaches highlight the extremes of this spectrum:

- **Full Compression, Difficult Indexing:** One could minimize the input trie into the smallest possible equivalent DFA using an algorithm like Revuz’s [33]. While this yields optimal compression through DAG compression (see Section 2.1), indexing the resulting general DFA is a notoriously difficult problem. As shown by Equi et al. [11], pattern matching on general DFAs requires super-polynomial time unless the Strong Exponential Time Hypothesis (SETH) fails,

making this approach unsuitable for most indexing purposes.

- **Full Indexability, No Compression:** At the other extreme, the input trie itself can be used as an index. Tries are Wheeler graphs [14], specifically 1-sortable automata (Definition 18), a property that makes them highly amenable to efficient indexing [7]. While this provides excellent query performance through the co-lexicographic ordering of states (see Definition 20), it offers no compression, as even highly repetitive subtrees are stored explicitly.

This thesis proposes a novel algorithm that finds a sweet spot in this trade-off, which we develop throughout Section 2.3. The central idea is to partially minimize the input trie while ensuring the resulting automaton remains efficiently indexable. We achieve this by leveraging the theory of p -sortable graphs (see Section 2.5), developing a method that strategically increases the sortability parameter p just enough to enable significant compression. The motivation for this approach is rooted in the observation that a small increase in p can lead to substantial compression. As noted by Policriti et al. [28], there are cases where increasing p from 1 to 2 allows for an exponential reduction in the automaton's size, a phenomenon we explore in detail in Section 2.5.

Our compression scheme, presented in Section 2.3, works by partitioning the trie's nodes into a predefined number p of chains and then optimizing this partition to merge the maximum number of equivalent states while ensuring p -sortability. To make this optimization problem more concrete, we can frame it as a string partitioning problem. Consider the sequence of nodes in the trie, when read in co-lexicographic order (see Definition 20), as a single long string. The "character" corresponding to each node is its Myhill-Nerode equivalence class (see Theorem 3), which determines if it can be merged with other nodes. The task is to partition this string of nodes into p subsequences such that the number of runs (Definition 1) is minimized. Minimizing the number of runs directly corresponds to maximizing the number of merged states, yielding a compact p -sortable automaton. We call this problem the STRING PARTITIONING PROBLEM (Definition 2).

Definition 1 (Run). *Let $S = s_1s_2 \dots s_n$ be a string over an alphabet Σ . A substring $S[i..j] = s_is_{i+1} \dots s_j$ (where $1 \leq i \leq j \leq n$) is a run if it satisfies the following conditions:*

1. **Homogeneity:** *All characters in the substring are identical, i.e., $s_k = s_i$ for all k such that $i \leq k \leq j$.*
2. **Maximality:** *The substring cannot be extended to the left or right without violating homogeneity. That is:*
 - *If $i > 1$, then $s_{i-1} \neq s_i$.*
 - *If $j < n$, then $s_{j+1} \neq s_j$.*

Example 1.1: Run

For example, in the string $S = AAABBC$, the runs are AAA , BB , and C .

Definition 2 (String Partitioning Problem). *Let S be a string of length n over an alphabet Σ , and let $p \in [2, n]$ be a positive integer. A p -partition of S is a set of*

p subsequences S_1, S_2, \dots, S_p such that every character of S belongs to exactly one subsequence. The total number of runs in a p -partition is the sum of the number of runs in each subsequence S_i . The goal of the **STRING PARTITIONING PROBLEM** is to find a p -partition of S that minimizes the total number of runs.

This thesis shows that the **STRING PARTITIONING PROBLEM** can be reduced to the problem of finding a minimum weight perfect matching on a bipartite graph (MWPBM), allowing us to use efficient, well-studied algorithms to find the optimal solution. The result is a compressed automaton that is p -sortable by construction and thus supports efficient queries using the data structure developed by Cotumacchio et al. [7]. As our experimental results in ?? will show, this method is particularly effective for the highly repetitive datasets common in real-world applications, achieving a balance of compression and indexability that prior methods could not attain.

1.2 Structure of The Thesis

Chapter 2

Preliminary Concepts

2.1 Notation

The following definitions and notation will be used throughout the entire thesis

We first recall the formal definition of a labeled tree.

Definition 3. A **labeled tree** is a rooted structure in which every edge is assigned a label from a finite alphabet Σ . Formally, a labeled tree is a quadruple $T = (V, E, \lambda, r)$, where:

- V is the set of vertices,
- $E \subseteq V \times V$ is the set of directed edges such that (V, E) forms a rooted tree,
- $\lambda : E \rightarrow \Sigma$ is an edge-labeling function,
- $r \in V$ is the root vertex.

In the case of ordered labeled trees, the children of each node are totally ordered; the degree and shape of the tree, as well as the size of the alphabet Σ , are otherwise unconstrained.

We now give the formal definition of a trie.

Definition 4 (Trie). Let Σ be a finite alphabet and $L \subseteq \Sigma^*$ be a finite set of strings. A **trie** for L is a 5-tuple $T = (V, E, \lambda, r, F)$ where:

- V is a finite set of vertices,
- $E \subseteq V \times V$ is a set of edges such that (V, E) forms a rooted tree,
- $\lambda : E \rightarrow \Sigma$ is an edge-labeling function,
- $r \in V$ is the root vertex,
- $F \subseteq V$ is the set of terminal (final) vertices.

The following properties hold:

1. **Determinism:** For every vertex $v \in V$ and every symbol $a \in \Sigma$, there is at most one edge $(v, u) \in E$ such that $\lambda((v, u)) = a$.
2. **String Representation:** For any vertex $v \in V$, let $\text{str}(v)$ be the string obtained by concatenating the labels on the unique path from the root r to v . We set $\text{str}(r) = \varepsilon$ (the empty string).
3. **Language Correspondence:** The language represented by the trie is exactly L , i.e., $L = \{ \text{str}(v) \mid v \in F \}$.

4. **Prefix Property:** The set of all prefixes of words in L coincides with $\{ \text{str}(v) \mid v \in V \}$.

The determinism property ensures that each string prefix identifies a unique path from the root, which makes tries well suited to our compression techniques. A trie can be viewed as a DFA (see Definition 5), allowing the application of automata minimization to compress the tree structure (DAG compression). **DAG compression** is the process of transforming a tree into a directed acyclic graph (DAG) by identifying and merging all isomorphic subtrees. In the context of tries, this is equivalent to finding and merging states with the same right language, which is the standard procedure in DFA minimization. The result is the smallest possible automaton that recognizes the same language.

We now recall the definition of a deterministic finite automaton (DFA), a fundamental computational model.

Definition 5 (Deterministic Finite Automaton). *A deterministic finite automaton (DFA) is a 5-tuple $M = (Q, \Sigma, \delta, q_0, F)$ where:*

- Q is a finite set of states,
- Σ is a finite input alphabet,
- $\delta : Q \times \Sigma \rightarrow Q$ is the transition function,
- $q_0 \in Q$ is the initial state,
- $F \subseteq Q$ is the set of final (accepting) states.

A DFA processes an input string $s \in \Sigma^*$ symbol by symbol, starting from q_0 and following transitions prescribed by δ . The string s is accepted if the automaton halts in a state of F after consuming all symbols; otherwise, it is rejected. The language recognized by a DFA is the set of all strings it accepts.

Also, we define the notation I_q to denote the set of strings reaching state q from the initial state: $I_q = \{ \alpha \in \Sigma^* \mid q = \delta(s, \alpha) \}$

We next define a non-deterministic finite automaton (NFA).

Definition 6 (Non-deterministic Finite Automaton). *A non-deterministic finite automaton is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$ where:*

- Q is a finite set of states,
- Σ is a finite alphabet,
- $\delta : Q \times (\Sigma \cup \{ \varepsilon \}) \rightarrow \mathcal{P}(Q)$ is the transition function (with $\mathcal{P}(Q)$ the power set of Q),
- $q_0 \in Q$ is the initial state,
- $F \subseteq Q$ is the set of accept states.

The principal difference between DFAs and NFAs is that NFAs allow (i) multiple outgoing transitions from a state on the same input symbol and (ii) ε -transitions, whereas DFAs have exactly one outgoing transition per state and input symbol and admit no ε -transitions.

For an automaton $\mathcal{N} = (Q, \Sigma, \delta, q_0, F)$, we extend the transition function to strings as follows, for all $u \in Q$, $a \in \Sigma$, and $\alpha \in \Sigma^*$:

$$\delta(u, \varepsilon) = \{u\}, \quad \delta(u, \alpha a) = \bigcup_{v \in \delta(u, \alpha)} \delta(v, a).$$

Thus, $\delta(q_0, \alpha)$ denotes the set of states reachable from the start state q_0 by reading α . A string α is *accepted* if $\delta(q_0, \alpha) \cap F \neq \emptyset$. The language recognized by \mathcal{N} is denoted by $\mathcal{L}(\mathcal{N})$ and is defined as

$$\mathcal{L}(\mathcal{N}) = \{\alpha \in \Sigma^* \mid \delta(q_0, \alpha) \cap F \neq \emptyset\}.$$

We introduce now the concept of regular language.

Definition 7 (Regular Language). *A language $L \subseteq \Sigma^*$ is called a **regular language** if it is the language recognized by some finite automaton.*

2.2 Labeled Trees and Tries

In this chapter, we will discuss labeled trees, a fundamental data structure used to represent hierarchical data. We will define the concept of a labeled tree, explore its common applications, and introduce the core concepts behind its compression and indexing.

2.2.1 Introduction and Motivation

Before delving into specific compression techniques, it is essential to establish a solid theoretical foundation regarding labeled trees and tries. These structures are fundamental for representing hierarchical data across diverse fields, from bioinformatics to document processing. This chapter provides the necessary background, defining labeled trees, exploring their common applications, and introducing the core concepts behind their compression and indexing.

While labeled trees encompass a broad class of hierarchical structures, this thesis focuses specifically on tries. This focus is motivated by our compression strategy, which relies on identifying and merging identical subtrees. To formally identify these subtrees, we compute the Myhill-Nerode equivalence classes of the trie's states (see Theorem 3), a task that can be accomplished using algorithms for DFA minimization. A fundamental prerequisite for these algorithms is determinism, as they require that from any given state, each symbol corresponds to at most one transition. Tries inherently provide this guarantee, since from any node, there is at most one outgoing edge for each symbol in the alphabet. This deterministic nature is therefore what allows us to apply powerful automata minimization techniques to compress the tree structure.

2.2.2 Applications

Labeled trees are widely used in computer science and data representation due to their hierarchical structure and flexibility in modeling relationships. Prominent applications include:

1. **XML Data Representation:** XML documents are often modeled as labeled trees, where each element is a node labeled by its tag, and hierarchical nesting represents parent-child relationships.
2. **JSON Data Representation:** JSON documents can be viewed as labeled trees, with keys as labels and values as children.
3. **Bioinformatics:** Labeled trees are used to represent phylogenetic trees, genome annotations, and hierarchical clustering.
4. **Compiler Design:** Abstract Syntax Trees (ASTs) for programming languages are labeled trees that capture the structure of code.
5. **File Systems:** The directory structure of file systems can be viewed as a labeled tree.

Efficient representation, navigation, and querying of labeled trees are essential for many applications, motivating the development of specialized data structures and algorithms.

2.2.3 Indexing

The goal of compressing and indexing labeled trees is to design a compressed storage scheme for a labeled tree T with t nodes that allows for efficient navigation operations in T , as well as fast search and retrieval of subtrees or paths within T . To be effective, the compressed representation should minimize the space required to store the tree while supporting a wide range of operations in optimal ($O(1)$) or (near-)optimal time.

We define the following navigation operations on T :

Definition 8 (Tree Operations). *Let T be a labeled tree with node set V , $u \in V$ be a node, and $c \in \Sigma$ be a symbol. We define the following fundamental operations on T :*

- **Navigational queries:** ask for the parent of u , the i -th child of u , or the label of u . The last two operations might be restricted to the children of u with a specific label c .
- **Visualization queries:** retrieve the nodes in the subtree rooted at u (any possible order should be implemented).
- **Subpath queries:** Let moreover $P(\alpha)$ be the set of vertices of N which are reached by a path labeled with a given word $\alpha \in \Sigma^*$. Efficient subpath queries that, given a query word $\alpha \in \Sigma^*$, solve:
 - *Existential queries:* Determine whether $P(\alpha) \neq \emptyset$, i.e., whether α matches a substring of some string in the language of T .
 - *Count queries:* Compute the cardinality of $P(\alpha)$.
 - *Locate queries:* Return a representation for all the vertices in $P(\alpha)$.

A naive solution to index labeled trees is to store the tree as a list of nodes with their labels and parent-child relationships using pointers in $O(t \log t)$. However, this

representation is not space-efficient and does not support fast navigation or query operations.

Many data structures have been proposed to compress and index labeled trees, each with its trade-offs in terms of space usage, query performance, and supported operations. One of the most successful approaches is the Extended Burrows-Wheeler Transform, which extends the classical Burrows-Wheeler Transform (BWT) to handle labeled trees efficiently (Subsection 2.2.5).

2.2.4 Information-Theoretic Lower Bound

Before computing the information-theoretic lower bound for labeled trees, it is essential to define the concept of worst-case entropy, which provides a formal measure of the minimum number of bits required to represent any object from a given set. As detailed by Navarro [30], this is a fundamental concept in data structure design.

Definition 9 (Worst-case entropy). *Let U be a universe of combinatorial objects. The worst-case entropy of U is*

$$H_{wc}(|U|) = \lceil \log_2(|U|) \rceil$$

This definition establishes that the minimum number of bits to uniquely identify any object in a set U is the logarithm of the size of U , rounded up to the nearest integer. We now apply this principle to determine the lower bound for labeled trees.

Lemma 1. *The information-theoretic lower bound for storing a labeled tree T with t nodes over an alphabet Σ is $2t + t \log_2 |\Sigma| - \Theta(\log t)$ bits.*

Proof. The total information required to store a labeled tree can be decomposed into two components: the space needed to encode the tree's structure and the space needed to encode the labels on its nodes.

1. Structural Information (Unlabeled Tree): The number of distinct unlabeled binary trees with t nodes is given by the t -th Catalan number, $C_t = \frac{1}{t+1} \binom{2t}{t}$. Using Stirling's approximation for factorials, the Catalan number can be approximated as:

$$C_t \approx \frac{4^t}{t^{3/2} \sqrt{\pi}}$$

Then, the worst-case entropy (or the information-theoretic minimum number of bits to encode the structure of the tree) is:

$$\log_2 C_t \approx \log_2 \left(\frac{4^t}{t^{3/2} \sqrt{\pi}} \right) = \log_2(4^t) - \log_2(t^{3/2} \sqrt{\pi}) = 2t - \frac{3}{2} \log_2 t - \frac{1}{2} \log_2 \pi$$

The lower-order terms can be expressed using Big Theta notation as $\Theta(\log t)$. Therefore, the space required for the structure is $2t - \Theta(\log t)$ bits.

2. Labeling Information: For a tree with t nodes and an alphabet Σ , each node must be assigned a label. To distinguish between $|\Sigma|$ possible labels, a minimum

of $\log_2 |\Sigma|$ bits is required for each node. Consequently, the total space required to store the labels for all t nodes is:

$$t \log_2 |\Sigma| \text{ bits}$$

Finally, by adding the space required for the structure and the labels, the total information-theoretic lower bound for storing a labeled tree is the sum of the two components:

$$(2t - \Theta(\log t)) + (t \log_2 |\Sigma|) = 2t + t \log_2 |\Sigma| - \Theta(\log t) \text{ bits}$$

This completes the proof. \square

2.2.5 State of The Art

The field of tree indexing and compression has evolved through two main paradigms: succinct data structures that achieve space-optimal representations, and compression techniques that exploit structural repetitions.

In the realm of succinct tree structures, early work by Kosaraju [23] proposed a method to index labeled trees by extending the concept of prefix sorting from strings to labeled trees using trie structures. He introduced the idea of constructing a suffix tree for a reversed trie, enabling subpath queries in $O(|P| \log |\Sigma| + occ)$ time, where occ is the number of occurrences of P in T . However, this approach still required $O(t \log t)$ space (where t is the number of nodes of the tree) and thus was not compressed.

A significant advancement in this direction came with the Extended Burrows-Wheeler Transform (XBWT) [12], a data structure designed for efficient compression and indexing of ordered node-labeled trees. The XBWT works by linearizing a labeled tree into two arrays capturing the structural properties of the tree and its labels. This transformation allows for efficient representation, navigation, and querying of the tree. The key advantage of the XBWT lies in its ability to compress labeled trees while supporting a wide range of operations, such as navigation, visualization and subpath queries (see Definition 8), within (near-)optimal time bounds and entropy-bounded space. The XBWT provides significant improvements in both compression ratio and query performance compared to traditional compression schemes, making it a valuable resource for intensive applications.

Complementing succinct approaches, tree compression has been extensively studied through different paradigms that exploit structural repetitions in distinct ways. One of the classical approaches is *DAG compression*, which represents a tree as a minimal directed acyclic graph (DAG) by identifying and merging identical rooted subtrees. Concretely, whenever two identical subtrees occur, only one copy is stored and all occurrences point to it. The resulting structure can be exponentially smaller than the original tree and can be computed in linear time. DAG compression has been widely used in programming languages, binary decision diagrams, and XML representations [3].

Another line of research extends the well-known LZ77 factorisation from strings to trees. Here, the tree is decomposed into edge-disjoint fragments, each being

either a single node or a copy of a fragment that appeared earlier in a breadth-first traversal. Each fragment is thus defined by pointers to earlier occurrences, much like in the string version of LZ77. This factorisation uniquely determines the tree, and by minimising the number of fragments one obtains a compressed representation. Importantly, such factorizations can be computed in polynomial time (and in linear time for restricted variants), and they yield representations no larger than the smallest tree grammar, thus bridging block compression and grammar-based compression [15].

More recently, top tree compression has been proposed as a method that combines the advantages of subtree sharing and grammar-like approaches. The key idea is to build a hierarchical top tree decomposition, where the input tree is recursively partitioned into clusters that capture connected patterns. These clusters are then merged following a restricted set of operations, producing a binary decomposition tree whose internal repetitions are turned into subtree repeats. Finally, this decomposition is compressed using standard DAG compression, resulting in a so-called top DAG. This approach achieves close-to-optimal worst-case bounds, can be exponentially more succinct than DAG compression, and crucially, supports a wide range of navigational queries (e.g., parent, child, depth, nearest common ancestor) in logarithmic time directly on the compressed representation [3].

Another notable approach is Tree Re-Pair [26], a grammar-based compression technique adapted for tree structures. It extends the principles of the original Re-Pair algorithm [25] to handle the hierarchical nature of trees by identifying and compactly representing frequently occurring patterns. The core idea of the tool is to identify frequently occurring patterns within the tree and represent them more compactly. The process involves the linearization of the tree (e.g., using a specific traversal order) and then the application of the Re-Pair logic. In this way, it finds the most frequent pair of adjacent elements (which could represent nodes, labels, or structural components, depending on the linearization) in the sequence. The pair is then replaced by a new non-terminal symbol, and the corresponding production rule is added to a grammar. All this process is then repeated until no more pairs occur frequently enough or some other stopping criterion is met. The final output is a relatively small grammar (a set of production rules) and a sequence of symbols (including the newly introduced non-terminals) that can be used to reconstruct the original tree. An application of Tree Re-Pair to XML documents can be found in [27].

In summary, DAG compression is efficient but limited to subtree repeats, LZ77 factorisation captures more general repetitions while relating closely to grammar-based methods such as Tree Re-Pair, and top tree compression strikes a balance by exploiting both subtree and pattern repeats while still enabling efficient query support.

This thesis focuses on developing a novel technique specifically tailored for highly repetitive tries. Our approach leverages the structural properties unique to tries and their repetitive patterns. Therefore, we use the XBWT as our primary benchmark for comparison, as it represents a well-established and high-performance baseline specifically designed for trie compression in the field.

2.2.6 Conclusion

This section has established the initial theoretical groundwork for the thesis, beginning with the formal definition of labeled trees and narrowing the focus to tries. We have underscored the deterministic nature of tries as a critical property that makes them suitable for DAG compression techniques. Furthermore, we defined a standard set of navigational and query operations, providing a benchmark for the functionality that an efficient compressed structure must support. By establishing the information-theoretic lower bound for storing labeled trees, we have set a clear goal for compression effectiveness.

The primary objective of this thesis is to develop and analyze a novel approach that leverages DAG compression to reduce the space footprint of tries, particularly those with highly repetitive substructures. Crucially, our goal is not merely to compress but to do so while preserving the trie's indexability. The subsequent chapters will detail the methods for achieving this, aiming to create a compressed representation that approaches the theoretical space limits while supporting essential query operations efficiently.

In the next section, we present a detailed explanation of the Extended Burrows-Wheeler Transform (XBWT), examining its construction, properties, and how it enables efficient compression and indexing of labeled trees.

2.3 The Extended Burrows-Wheeler Transform

The fundamental definitions, properties, and algorithms related to the Extended Burrows-Wheeler Transform presented in this chapter are based on the work introduced by Ferragina et al. 'Compressing and Indexing Labeled Trees, with Applications' [12].

2.3.1 Introduction and Motivation

This chapter explores the Extended Burrows-Wheeler Transform (XBWT), introduced by Ferragina et al. [12], a state of the art technique for labeled tree compression. Understanding the principles and performance of the XBWT is crucial as it will serve as the primary benchmark against which we will evaluate the novel compression scheme proposed in this thesis. By establishing a baseline with a well-regarded method like the XBWT, we can effectively demonstrate the potential advantages and contributions of our new approach, particularly for trees exhibiting high repetitiveness.

In 2005, Ferragina et al. [12] introduced an innovative approach to labeled tree compression by transforming it into a more tractable string compression problem. Their key contribution, the Extended Burrows-Wheeler Transform (XBWT), is a sophisticated data structure that achieves highly efficient compression by combining entropy-compressed edge labels with a succinct representation of the tree topology. This elegant solution not only simplifies the compression process but also maintains the structural relationships essential for tree operations.

The XBWT works by linearizing a labeled tree into two coordinated arrays capturing the structural properties of the tree and storing its labels. This transformation allows for efficient representation, navigation, and querying of the tree. The key advantage of the XBWT lies in its ability to compress labeled trees while supporting a wide range of operations, such as navigation, visualization and subpath queries (see Definition 8), within (near-)optimal time bounds and entropy-bounded space.

One of the primary applications of the XBWT is in compressing and indexing hierarchical data formats, such as XML documents. It provides significant improvements in both compression ratio and query performance compared to traditional tools, making it an invaluable resource for data-intensive applications in fields like bioinformatics, information retrieval, and big data analytics.

This chapter aims to explore the XBWT data structure and its applications in the context of labeled trees. We will start by providing an overview of the theoretical foundations of the XBWT. Finally, we will describe and compare the algorithms for constructing the XBWT and demonstrate its use in compressing and indexing labeled trees.

Key Aspects

The XBWT has several key properties that make it an effective tool for labeled tree compression and indexing:

- **Succinctness:** The XBWT representation of a labeled tree uses space close to the worst-case entropy (Lemma 1), which is $2t - \Theta(\log t) + t \log |\Sigma|$ bits for a tree with t nodes and an alphabet of size $|\Sigma|$.
- **Efficient Querying:** The XBWT supports the navigational queries (Definition 8) in optimal time $O(1)$ if $|\Sigma| = O(\text{polylog}(t))$, otherwise in $O(\log \log^{1+\epsilon} |\Sigma|)$ time. Whereas, given $s \in \Sigma^*$, subpath queries (Definition 8) are supported in $O(|s|)$ if $|\Sigma| = O(\text{polylog}(t))$, otherwise in $O(|s| \log \log^{1+\epsilon} |\Sigma|)$ time.
- **Scalability:** The XBWT is particularly useful for large-scale hierarchical data, such as XML documents or phylogenetic trees, where both compression and fast querying are critical.

2.3.2 Definition

The **Extended Burrows-Wheeler Transform** is a data structure designed to efficiently compress and index *ordered node-labeled trees*. Inspired by the classical Burrows-Wheeler Transform (BWT) [4] for strings, the XBWT extends these principles to hierarchical structures, enabling efficient storage, navigation, and querying of trees. It is particularly effective for trees where each node has a label drawn from an alphabet Σ and the tree structure has an arbitrary shape and degree.

Definition 10. *Let T be a totally ordered node-labeled tree of arbitrary fan-out, depth, and shape, with n internal nodes and l leaves (t nodes in total) and alphabet Σ . Let u be a node in T , we define the following information:*

- *$\text{last}(u)$: a binary value that is 1 if u is the last (rightmost) child of its parent in the total order, and 0 otherwise.*

- $\alpha(u)$: denotes the label of node u concatenated with one bit that is 1 if u is a leaf and 0 otherwise.
- $\pi(u)$: the string obtained by concatenating the labels of the nodes on the UPWARD PATH from u 's parent to the root of T . Formally, we define $\pi(u)$ recursively as:

$$\pi(u) = \begin{cases} \varepsilon & \text{if } u \text{ is the root} \\ \text{label}(\text{parent}(u)) \circ \pi(\text{parent}(u)) & \text{otherwise} \end{cases}$$

where \circ is the concatenation operator and ε is the empty string.

If the leaf-label alphabet Σ_L is disjoint from the internal-node label alphabet Σ_N , the additional indicator bit in $\alpha(u)$ is redundant and may be omitted. In that case, we identify $\alpha(u)$ with the node label, i.e., $\alpha(u) := \text{label}(u)$.

The definition of the XBWT relies on a sequence S , which contains a triplet $(\text{last}(u), \alpha(u), \pi(u))$ for each node u in the tree T .

The construction of S is a two-step process. First, an intermediate vector of triplets is created by traversing the tree T in pre-order and generating a triplet $(\text{last}(u), \alpha(u), \pi(u))$ for each node. Second, it is stably sorted according to the lexicographical order of the ' π ' component to produce the final sequence S (see Table 2.1 for an example).

Theorem 1. *The XBWT of a labeled tree T consists of two arrays, S_{last} and S_{α} . These are constructed from the sequence S of triplets. Specifically, for each i from 1 to t , $S_{\text{last}}[i]$ is the 'last' component of the i -th triplet in S , and $S_{\alpha}[i]$ is the ' α ' component. The total space required is $2t + t \log |\Sigma|$ bits.*

S_{π} (for each i from 1 to t , $S_{\pi}[i]$ is the ' π ' component of the i -th triplet in S), therefore is not needed after the construction of the XBWT. However, in the following discussion, we will still refer to it as it possesses some important properties. See Figure 2.1 for an example of the tree T and its corresponding sequence S .

2.3.3 Properties

The XBWT's effectiveness as an indexing structure stems from a key property of the sequence S . This property, along with its consequences, arises directly from the transform's definition and the sorting process.

Key Property: Grouping by Parent

The fundamental property of the XBWT is that the children of any node u in the tree T form a contiguous block in the sequence S . Let u_1, \dots, u_z be the children of node u in their original order. Their corresponding triplets will appear consecutively in S in that same order.

Example 2.1:

Consider the node u in Figure 2.1. Looking at Table 2.1, we can see that its children form a contiguous block in positions $[5, 6, 7]$ of the sequence S .

This grouping provides several important consequent properties:

Unary Degree Encoding: The subarray S_{last} for the block of children $[u_1, \dots, u_z]$ encodes the degree of u in unary. Specifically, $S_{\text{last}}[u_z] = 1$ and $S_{\text{last}}[u_i] = 0$ for $1 \leq i < z$.

Example 2.2:

Consider the node u in Figure 2.1. Looking at Table 2.1, we can observe that $S_{\text{last}}[5 \dots 7] = \{0, 0, 1\}$, which encodes the degree 3 in unary notation, matching the number of children of node u .

Preservation of Sibling Order: If two nodes u and v have the same label, and the triplet for u precedes the triplet for v in S , then the entire block of children of u will also precede the block of children of v .

Example 2.3:

Consider nodes u and v in Figure 2.1. In Table 2.1, node u appears at index 2 while node v appears at index 4 in the sequence S . Following the preservation of sibling order property, all children of u (occupying positions $[5, 6, 7]$) appear before the child of v (at position 8).

Path-based Indexing: This property extends to entire paths. For any label $c \in \Sigma$, all triplets whose π -components are prefixed by c form a contiguous block in S . If u is the i -th node with label c in S_α , its children's block is located within the larger block of all nodes with paths prefixed by c . This block is delimited by the $(i - 1)$ -th and i -th '1's in the corresponding section of S_{last} .

Example 2.4:

Let's examine nodes u and v in Figure 2.1, both labeled 'B'. In the sequence S shown in Table 2.1, u is the first node with label 'B' (at index 2), and v is the second (at index 4).

The children of all nodes labeled 'B' form a contiguous block in S . In this case, the children of both u and v are located in the range $[5, 8]$. We can distinguish between the children of u and the children of v using the S_{last} array:

- The block of children for u (the first 'B' node) starts at the beginning of the range (index 5) and ends at the position of the first '1' in $S_{\text{last}}[5 \dots 8]$.
- The block of children for v (the second 'B' node) starts after the first '1' and ends at the position of the second '1' in $S_{\text{last}}[5 \dots 8]$.

Other Properties

Additional properties of the XBWT components include:

- The first triplet in S always corresponds to the root of the tree T .
- S_{last} contains n ones (for internal nodes) and l zeros (for leaves).
- S_α is a permutation of the node labels in T .

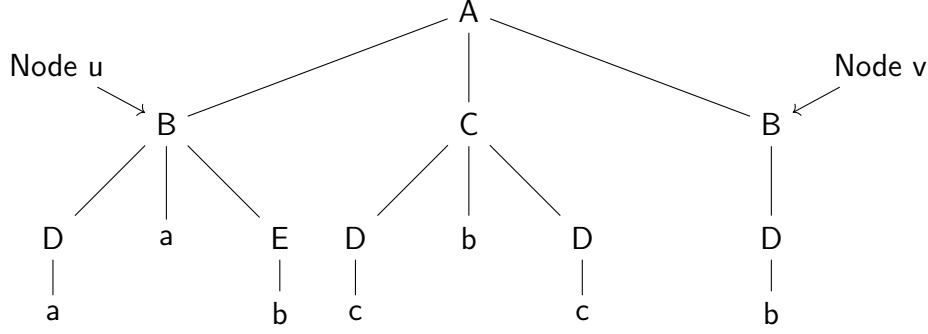


Figure 2.1: A labeled tree T where $\Sigma_N = \{A, B, C, D, E\}$ and $\Sigma_L = \{a, b, c\}$. Notice that $\alpha(u) = \alpha(v) = B$ and $\pi(u) = \pi(v) = A$.

	S_{last}	S_α	S_π
1	0	A	ϵ
2	0	B	A
3	0	C	A
4	1	B	A
5	0	D	BA
6	0	a	BA
7	1	E	BA
8	1	D	BA
9	0	D	CA
10	0	b	CA
11	1	D	CA
12	1	a	DBA
13	1	b	DBA
14	1	c	DCA
15	1	c	DCA
16	1	b	EBA

Table 2.1: The sequence S for the tree shown in Figure 2.1, obtained by stably sorting triplets according to their ' π ' components. In this representation, nodes u and v from the original tree T appear at indices 2 and 4, respectively. The children's block of node u occupies positions 5 through 7, while node v 's single child is located at index 8.

2.3.4 Construction

A naive approach to build the XBWT would be to explicitly construct S through the concretization of π -strings and then sort it using a stable sorting algorithm. However, this approach would require $\Theta(t^2)$ space in the worst case, which is not feasible for deep trees. To overcome this issue, Ferragina et al. [12] proposed a more efficient algorithm that builds S in linear time and $O(t \log t)$ bits of working space.

The linear time algorithm is called **pathSort**, it is based on a generalization of the Skew algorithm for suffix array construction of strings [21]. Let's see briefly how the Skew algorithm works.

Skew Algorithm

The Skew algorithm is an efficient method for constructing the suffix array of a string in linear time. A suffix array is a data structure that lists the starting indices of all the suffixes of a string in lexicographical order, and it is widely used in various string processing algorithms.

Algorithm Overview

Step 0: Construct a Sample. For $k \in \{0, 1, 2\}$, define the index sets

$$B_k = \{i \in [0, n] \mid i \bmod 3 = k\}.$$

Let $C = B_1 \cup B_2$ be the set of sample positions, and let S_C denote the set of sample suffixes.

Step 1: Sort Sample Suffixes. For $k \in \{1, 2\}$, construct the string R_k whose characters are the triples $[t_i t_{i+1} t_{i+2}]$ for $i \in B_k$ in increasing order of i . The last character of R_k is unique because we pad with sentinels so that $t_{n+1} = t_{n+2} = 0$. Let $R = R_1 R_2$ be the concatenation of R_1 and R_2 . Then the nonempty suffixes of R correspond to the sample suffixes in S_C in an order-preserving way: sorting the suffixes of R yields the order of S_C . To sort the suffixes of R , first radix sort the characters of R (the triples) and rename them by their ranks to obtain a string R' . If all characters are distinct, their order directly gives the order of suffixes. Otherwise, sort the suffixes of R' recursively using the same DC3 procedure. Once S_C is sorted, assign ranks to sample suffixes: for $i \in C$, let $\text{rank}(S_i)$ be the rank of S_i within S_C . Additionally, define $\text{rank}(S_{n+1}) = \text{rank}(S_{n+2}) = 0$. For $i \in B_0$, $\text{rank}(S_i)$ is undefined.

Step 2: Sort Nonsample Suffixes. Represent each nonsample suffix S_i with $i \in B_0$ by the pair

$$(t_i, \text{rank}(S_{i+1})).$$

Note that $\text{rank}(S_{i+1})$ is always defined for $i \in B_0$ by the previous step. Then, radix sort these pairs to obtain the order of S_{B_0} .

Step 3: Merge. Merge the two sorted sets S_C and S_{B_0} using a standard comparison-based merging. To compare a suffix $S_i \in S_C$ against a suffix $S_j \in S_{B_0}$, distinguish two cases:

$$\begin{aligned} \text{if } i \in B_1 : S_i \leq S_j &\iff (t_i, \text{rank}(S_{i+1})) \leq_{\text{lex}} (t_j, \text{rank}(S_{j+1})), \\ \text{if } i \in B_2 : S_i \leq S_j &\iff (t_i, t_{i+1}, \text{rank}(S_{i+2})) \leq_{\text{lex}} (t_j, t_{j+1}, \text{rank}(S_{j+2})). \end{aligned}$$

The ranks used above are defined in all cases by Step 1. Each comparison inspects $O(1)$ characters/ranks, and the two-way merge advances each pointer at most once, so the merge runs in linear time overall.

Time Complexity. Excluding the recursive call, all steps are linear-time via radix sorting and a single-pass merge. The recursion operates on a string of length $2n/3$, so the overall time satisfies

$$T(n) = T(2n/3) + O(n) = O(n).$$

PathSort Algorithm

The pseudocode of the pathSort algorithm is shown in Algorithm 1. As we can see, the algorithm is based on the Skew algorithm, but it is adapted to work on labeled trees. Given a value $j \in \{0, 1, 2\}$, the main idea is to recursively sort the upward subpaths of the tree starting at nodes in levels $\not\equiv j \pmod{3}$, then sort the upward subpaths starting at nodes in levels $\equiv j \pmod{3}$ using the result of the previous step, and finally merge the two sets of sorted subpaths by exploiting their lexicographic names. j is chosen in such a way that the number of nodes of the shrunk tree whose level is $\equiv j \pmod{3}$ is at least $t/3$ so that a constant fraction of upward paths are ensured to be dropped at each recursive step. It is important to note that:

1. the height of the new (contracted) tree shrinks by a factor of three, hence the node naming requires the radix sort over triples of names;
2. given the choice of j , the number of nodes of the new (contracted) tree will be at most $2t/3$, thus ensuring that the running time of the algorithm satisfies the recurrence $R(t) = R(2t/3) + \Theta(t) = \Theta(t)$;
3. following an argument similar to [21], the names of the dropped subpaths can be computed in $O(t)$ time from the names of the non-dropped subpaths, by radix sorting.

Algorithm 1 PATHSORT(T)

- 1: Initialize the array of triplets `IntNodes`[1.. t].
 - 2: Visit the internal nodes of T in pre-order. For the i -th visited internal node u , set `IntNodes`[i] $\leftarrow (\alpha(u), \text{level}(u), \text{parent}(u))$.
 - 3: Let $j \in \{0, 1, 2\}$ be such that the number of nodes in `IntNodes` whose level is $\equiv j \pmod{3}$ is at least $t/3$. Sort recursively the upward subpaths starting at nodes in levels $\not\equiv j \pmod{3}$.
 - 4: Sort the upward subpaths starting at nodes in levels $\equiv j \pmod{3}$ using the result of Step 3.
 - 5: Merge the two sets of sorted subpaths by exploiting their lexicographic names.
-

Recursive Step of PathSort

At each recursive step, the algorithm constructs the array `IntNodes`, which stores the triplets $(\alpha(u), \text{level}(u), \text{parent}(u))$ for every internal node u in the given tree T .

Next, the algorithm selects a value j such that the number of nodes in `IntNodes` with depth $\equiv j \pmod{3}$ is at least $t/3$. Based on this choice, two separate arrays are created:

- `IntNodesAtPosJ`, containing nodes at levels $\equiv j \pmod{3}$,
- `IntNodesNotAtPosJ`, containing nodes at levels $\not\equiv j \pmod{3}$

For each node u in `IntNodesNotAtPosJ`, the algorithm extracts the upward path consisting of the first three ancestors of u . These paths are then sorted using radix sort. If the sorted upward paths contain duplicates, the algorithm recursively calls

the PathSort function on a new contracted tree, where nodes are renamed according to their sorted paths. Otherwise, if all upward paths are unique, the nodes in `IntNodesAtPosJ` are sorted and subsequently merged with `IntNodesNotAtPosJ` using lexicographic ordering, following the same merging strategy as in the Skew algorithm.

2.3.5 Inversion

The ability to invert the XBWT is fundamental to its utility as a compression technique. Invertibility guarantees that the original tree can be perfectly reconstructed from its transformed representation (S_{last} and S_α). This ensures that the compression is lossless, meaning no information is lost during the process, which is a critical requirement for most applications.

Property 'Path-based Indexing' (Subsection 2.3.3) ensures that the two arrays S_{last} and S_α of the XBWT can be used to reconstruct the original tree T . The algorithm to invert the XBWT is linear in time and requires $O(t \log t)$ bits of space.

Algorithm 2 operates in three main steps. First, it constructs two auxiliary arrays, F and J , which are crucial for navigating the tree structure within the compressed format.

- **The F array:** This array maps each character $c \in \Sigma$ to the index of the first occurrence in S of a triplet whose π -component is prefixed by c . It essentially marks the starting points of blocks of nodes that share the same initial path label.
- **The J array:** For each entry i in S , $J[i]$ stores the index in S corresponding to the first child of the node represented by $S[i]$. If $S[i]$ represents a leaf, $J[i]$ is set to a sentinel value (e.g., -1).

Example 2.5: F and J arrays

Considering the XBWT in Table 2.1, the F array would map 'A' to index 2 (for node r), 'B' to index 5 (for the children of nodes with label 'B'), and so on. For the J array, let's take the node u at index 2 in S . Its first child is at index 5. Therefore, $J[2]$ would be 5.

Finally, the algorithm employs the array J to simulate a depth-first visit of T , creates its labeled nodes, and properly connects them to their parents.

Algorithm 2 RebuildTree(XBWT[T])

```

1:  $F = \text{BuildF}(\text{XBWT}[T]);$ 
2:  $J = \text{BuildJ}(\text{XBWT}[T], F);$ 
3: Create node  $r$  and set  $Q = \{\langle 1, r \rangle\};$  ▷  $Q$  is a stack
4: while  $Q \neq \emptyset$  do ▷ We still have nodes to create in  $T$ 
5:    $\langle i, u \rangle = \text{pop}(Q);$ 
6:    $j = J[i];$  ▷ Take the block of  $u$ 's children in  $S$ 
7:   if  $j = -1$  then ▷  $u$  is a leaf of  $T$ 
8:     continue;
9:   end if
10:  Find first  $j' \geq j$  such that  $S_{\text{last}}[j'] = 1;$  ▷  $S[j, j']$  are the children of  $u$  in  $T$ 
11:  for  $h = j'$  downto  $j$  do
12:    Create the node  $v$  labeled  $S_\alpha[h];$ 
13:    Attach  $v$  as first child of  $u;$ 
14:     $\text{push}(\langle h, v \rangle, Q);$ 
15:  end for
16: end while
17: return node  $r.$ 

```

Algorithm 3 BuildF(XBWT[T])

```

1: for  $i = 1, \dots, |\Sigma_N|$  do
2:    $C[S_\alpha[i]] \leftarrow C[S_\alpha[i]] + 1;$  ▷ Count the occurrences of node labels
3: end for
4:  $F[1] = 2;$  ▷  $S_\pi[1]$  is the empty string
5: for  $i \in \{1, \dots, |\Sigma_N| - 1\}$  do ▷ Consider just the internal-node labels
6:    $s = 0; j = F[i];$ 
7:   while  $s \neq C[i]$  do ▷ Not all blocks of children have been passed
8:      $j = j + 1;$ 
9:     if  $S_{\text{last}}[j] = 1$  then ▷ One further block of children has passed
10:       $s = s + 1;$ 
11:     end if
12:   end while
13:    $F[i + 1] = j;$ 
14: end for
15: return  $F.$ 

```

Algorithm 4 BuildJ(XBWT[T], F)

```

1: for  $i = 1, \dots, t$  do
2:   if  $S_\alpha[i] \in \Sigma_L$  then
3:      $J[i] = -1$ ;  $\triangleright S_\alpha[i]$  is a leaf label
4:   else
5:      $z = J[S_\alpha[i]]$ ;
6:     while  $S_{\text{last}}[z] \neq 1$  do  $\triangleright$  Reach the last child of  $S_\alpha[i]$ 
7:        $z = z + 1$ ;
8:     end while
9:      $F[S_\alpha[i]] = z + 1$ ;
10:  end if
11: end for
12: return  $J$ .

```

2.3.6 Compressing Labeled Trees

The XBWT[T] exhibits a local homogeneity property on S_α : specifically, the labels $\alpha(u)$ of nodes u whose upward paths $\pi(u)$ share long common prefixes appear in S_α in contiguous (or tightly bounded) clusters. This phenomenon can be formalized via the notion of k -contexts on trees. This property mirrors the strong local homogeneity exhibited by strings under the Burrows-Wheeler Transform [4] when applied to labeled trees.

To illustrate this, let us consider two arbitrary nodes u and v in T , and examine their contexts $\pi(u)$ and $\pi(v)$. Given the sorting of S , the greater the length of the shared prefix between $\pi(u)$ and $\pi(v)$, the closer the corresponding labels $\alpha(u)$ and $\alpha(v)$ will be in the string S_α . These closely spaced labels are expected to be few in number, resulting in S_α exhibiting local homogeneity. As a consequence, we can leverage the advanced algorithmic techniques developed for BWT-based compression methods to achieve efficient compression.

At the end, the XBWT is used for turning the labeled tree compression problem into a string compression problem. To this aim, two string compressors C_α and C_{last} are used to compress the two strings that compose XBWT[T], by exploiting their fine specialties. Of course, many choices are possible for C_α and C_{last} , each having implications on the algorithmic time and compression bounds.

In general, the following theorem holds:

Theorem 2 ([12], Theorem 4). *let C_α be a k -th order string compressor that compresses any string w into $|w|H_k(w) + |w| + o(|w|)$ bits, taking $O(|w|)$ time; and let C_{last} be an algorithm that stores S_{last} without compression. With this simple instantiation, the labeled tree T can be compressed within $tH_k(S_\alpha) + 2t + o(t)$ bits and takes $O(t)$ optimal time.*

Since $H_k(S_\alpha) \leq (\log |\Sigma|) + 1$ (Recall that when a unique alphabet is used to label both internal nodes and leaves, one needs a further bit to distinguish between them. The additional term $+1$ in the bound takes this into account), the above bound is at most $t(\log |\Sigma| + 3) + o(t)$ bits, and can be significantly better than the information-theoretic lower bound and the plain storage of XBWT[T] (both taking $2t + t \log |\Sigma|$

bits), depending on the distribution of the labels among its nodes.

2.3.7 Indexing a Compressed Labeled Tree

In order to implement the efficient operations listed in Subsection 2.2.3 using the compressed arrays S_{last} and S_α of XBWT, we need the chosen compressors C_α and C_{last} to support the following operations:

Given a string $S[1, t]$ over alphabet Σ

- $\text{rank}_c(S, q)$: gives the number of times the symbol $c \in \Sigma$ appears in $S[1, q]$.
- $\text{select}_c(S, i)$: gives the position of the i -th occurrence of the symbol $c \in \Sigma$ in S .

The compressed indexing of $\text{XBWT}[T]$ will be based on three compressed data structures that support rank and select queries over the two strings S_α and S_{last} , and over an auxiliary binary array $A[1, t]$ defined as: $A[1] = 0$, $A[j] = 1$ if and only if the first symbol of $S_\pi[j]$ differs from the first symbol of $S_\pi[j - 1]$. Hence, A contains at most $|\Sigma| + 1$ bits set to 1 out of t positions. It is also easy to see that, through rank and select operations over A , we can succinctly implement the array F employed in Algorithms 2 and 3.

In this section, let $S := (S_{\text{last}}, S_\alpha)$ denote the $\text{XBWT}[T]$ obtained after the construction phase. The following methods are supported by the compressed index:

GetRankedChild(i, k): Returns the position in S of the k -th child of the node at index i . If the child does not exist, it returns -1.

Example 2.6:

In Table 2.2, **GetRankedChild**(2, 2) returns 6.

GetCharRankedChild(i, c, k): Returns the position in S of the k -th child labeled c of the node at index i . If the child does not exist, it returns -1.

Example 2.7:

In Table 2.2, **GetCharRankedChild**(1, B, 2) returns 4.

GetDegree(i): Returns the total number of children of the node at index i in S .

GetCharDegree(i, c): Returns the number of children of the node at index i in S that have the label c .

GetParent(i): Returns the position in S of the parent of the node at index i . If the node is the root (at index 1), it returns -1.

Example 2.8:

In Table 2.2, **GetParent**(8) returns 4.

GetSubtree(i): Retrieves the labels of all nodes in the subtree rooted at the node at index i in S . The labels can be returned in any standard traversal order (e.g., pre-order, in-order, or post-order).

	A	S_{last}	S_{α}	S_{π}
1	0	0	A	ϵ
2	1	0	B	A
3	0	0	C	A
4	0	1	B	A
5	1	0	D	BA
6	0	0	a	BA
7	0	1	E	BA
8	0	1	D	BA
9	1	0	D	CA
10	0	0	b	CA
11	0	1	D	CA
12	1	1	a	DBA
13	0	1	b	DBA
14	0	1	c	DCA
15	0	1	c	DCA
16	1	1	b	EBA

Table 2.2: The sequence S for the tree shown in Figure 2.1, obtained by stably sorting triplets according to their ' π ' components. In this representation, nodes u and v from the original tree T appear at indices 2 and 4, respectively. The children's block of node u occupies positions 5 through 7, while node v 's single child is located at index 8. Also, the auxiliary binary array A is shown.

SubPathSearch(P): For a given labeled path $P = c_1c_2 \cdots c_k$, this function finds the range $S[\text{First} \dots \text{Last}]$ such that all strings in $S_{\pi}[\text{First} \dots \text{Last}]$ are prefixed by the reversed path $P^R = c_k \cdots c_2c_1$.

Example 2.9:

In Table 2.2, **SubPathSearch(BD)** results in the range $[12, 13]$, and **SubPathSearch(AB)** gives the range $[5, 8]$.

It is important to note that their time complexity is dependent on the specific implementation for rank and select over the compressed strings S_{α} and S_{last} .

Let's now see how to implement some of the above methods (from which the others can be derived) using the rank and select operations over the compressed strings S_{α} and S_{last} .

GetChildren(i)

Algorithm 5 exploits directly the properties described before, in particular Property 'Path-based Indexing' (Subsection 2.3.3). The rank operation at line 5 is used to get the number r of nodes labeled c up to position i in S_{α} . Then, the position $F[c]$ is obtained through a select operation on A (line 6). By Property 'Path-based Indexing', the children of $S[i]$ are located at the r -th block of children following position $F[c]$. Lines 8 – 9 identify this block.

Example 2.10:

Let's walk through an example using Table 2.2. Consider the node u at index 2 labeled with B . To find its children:

1. First, we compute $r = 1$ since this is the first occurrence of B in S_α up to position 2.
2. Next, we find $y = F[B] = 5$, which marks the start of the block containing children of all nodes labeled B .
3. Then, we count $z = 1$ ones in S_{last} up to position $y - 1$.
4. Finally, the children block is delimited by the $z + r - 1 = 1\text{st}$ and $z + r = 2\text{nd}$ ones in S_{last} , giving us the range $[5, 7]$.

This range $[5, 7]$ indeed contains the three children of the node at index 2, as we can verify from the tree structure in Figure 2.1.

Algorithm 5 GetChildren(i)

```

1: if  $S_\alpha[i] \in \Sigma_L$  then
2:   return  $-1$                                  $\triangleright S[i]$  is a leaf
3: end if
4:  $c \leftarrow S_\alpha[i]$                              $\triangleright S[i]$  is labeled  $c$ 
5:  $r \leftarrow \text{rank}_c(S_\alpha, i)$ 
6:  $y \leftarrow \text{select}_1(A, c)$                      $\triangleright y = F[c]$ 
7:  $z \leftarrow \text{rank}_1(S_{\text{last}}, y - 1)$ 
8: First  $\leftarrow \text{select}_1(S_{\text{last}}, z + r - 1) + 1$ 
9: Last  $\leftarrow \text{select}_1(S_{\text{last}}, z + r)$ 
10: return (First, Last)
    
```

GetParent(i)

Algorithm 6 is based on Property ‘Path-based Indexing’ (Subsection 2.3.3) and it is the inverse of the GetChildren method. In line 4, the algorithm computes the label c of the parent of $S[i]$ that prefixes the upward path leading to $S[i]$. Then, the parent of $S[i]$ is searched among the nodes labeled c in S_α by exploiting Property ‘Path-based Indexing’ in a reverse manner. Namely, the number k of children-blocks in the range $S[y, i]$ is computed; these are children of nodes labeled c and preceding i in the stable sort of S . Then, the k -th occurrence of c in S_α is selected, which is indeed the parent of $S[i]$.

Example 2.11:

Let’s illustrate how to find a node’s parent using Table 2.2. Consider node v located at index 4 with label B . The process to find its parent involves:

1. Computing $c = \text{rank}_1(A, 4) = 1$, which tells us the parent has label ‘A’ (as A contains exactly one 1 up to position 4).
2. Locating $y = F[A] = 2$, which indicates where the block of children for nodes labeled ‘A’ begins.
3. Calculating $k = \text{rank}_1(S_{\text{last}}, 4 - 1) - \text{rank}_1(S_{\text{last}}, 2 - 1) = 0$, meaning no complete child blocks appear before position 4.
4. Therefore, v ’s parent is the first $((k + 1)\text{-th})$ occurrence of ‘A’ in S_α , corresponding to index 1 (the root of \mathcal{T}).

This example demonstrates how the XBWT structure efficiently encodes parent-child relationships using just the S_{last} and S_α arrays.

Algorithm 6 GetParent(i)

```

1: if  $i = 1$  then
2:   return  $-1$   $\triangleright S[i]$  is the root of  $\mathcal{T}$ 
3: end if
4:  $c \leftarrow \text{rank}_1(A, i)$ 
5:  $y \leftarrow \text{select}_1(A, c)$ 
6:  $k \leftarrow \text{rank}_1(S_{\text{last}}, i - 1) - \text{rank}_1(S_{\text{last}}, y - 1)$ 
7:  $p \leftarrow \text{select}_c(S_\alpha, k + 1)$ 
8: return  $p$ 
    
```

SubPathSearch(P)

We assume that $P = c_1 c_2 \dots c_k$ algorithm SubPathSearch computes the range $[First, Last]$ in $|P| = l$ phases, each one preserving the following invariant:

- Invariant of Phase i . At the end of the phase, $S_\pi[First]$ is the first entry prefixed by $P[1, i]^R$, and $S_\pi[Last]$ is the last entry prefixed by $P[1, i]^R$, where s^R is the reversal of string s .

At the beginning (i.e., $i = 1$), First and Last are easily determined via the entries $F[c_1]$ and $F[c_1 + 1] - 1$, which point to the first and last entry of S_π prefixed by c_1 (by definition of array F). Since we do not have the F array, we implement these operations via rank and select queries over array A . Let us assume that the invariant holds for Phase $i - 1$, and prove that the i -th iteration of the for-loop in algorithm SubPathSearch preserves the invariant. More precisely, let $S_\pi[First, Last]$ be all entries prefixed by $P[1, i - 1]^R$. So $S[First, Last]$ contains all nodes descending from $P[1, i - 1]$. SubPathSearch determines $S[z_1]$ (respectively $S[z_2]$) as the first (respectively last) node in $S[First, Last]$ that descends from $P[1, i - 1]$ and is labeled c_i , if any. Then it jumps to the first child of $S[z_1]$ and the last child of $S[z_2]$. From Property 2 (item 2) and the correctness of algorithms GetChildren and GetDegree, we infer that the positions of these two children are exactly the first (respectively last) entry in S whose π -component is prefixed by $P[1, i]^R$.

The time complexity of the SubPathSearch algorithm is $O(l)$, where l is the length of the input path P .

Example 2.12:

Consider the tree in Figure 2.1, and let $P = BD$. The algorithm SUBPATH-SEARCH(P) returns the range $[12, 13]$ through the following steps:

1. Initially, $First = F[B] = 5$ and $Last = F[C] - 1 = 8$. The range $S[5, 8]$ contains all nodes descending from paths prefixed by B .
2. For $c_2 = D$:
 - Compute $k_1 = 0$ and $k_2 = 2$
 - This yields $z_1 = 5$ and $z_2 = 8$

- The first child of $S[5]$ is at position 12
- The last (and only) child of $S[8]$ is at position 13

3. Therefore, the algorithm returns the range $[12, 13]$

Note that the number of occurrences of subpath P is 2, as evidenced by the two occurrences of 1 in range $S_{\text{last}}[12, 13]$.

Algorithm 7 SubPathSearch(P)

```

1:  $First \leftarrow F(c_1)$ ;  $Last \leftarrow F(c_1 + 1) - 1$ 
2: if  $First > Last$  then
3:   return “ $P$  is not a subpath of  $T$ ”
4: end if
5: for  $i \leftarrow 2, \dots, k$  do
6:    $k_1 \leftarrow \text{rank}_{c_i}(S_\alpha, First - 1)$ ;  $z_1 \leftarrow \text{select}_{c_i}(S_\alpha, k_1 + 1)$   $\triangleright$  first entry in
      $S_\alpha[First, t]$  labeled  $c_i$ 
7:    $k_2 \leftarrow \text{rank}_{c_i}(S_\alpha, Last)$ ;  $z_2 \leftarrow \text{select}_{c_i}(S_\alpha, k_2)$   $\triangleright$  last entry in  $S_\alpha[1, Last]$ 
     labeled  $c_i$ 
8:   if  $z_1 > z_2$  then
9:     return “ $P$  is not a subpath of  $T$ ”
10:  end if
11:   $First \leftarrow \text{GetRankedChild}(z_1, 1)$   $\triangleright$  get the first child of  $S[z_1]$ 
12:   $Last \leftarrow \text{GetRankedChild}(z_2, \text{GetDegree}(z_2))$   $\triangleright$  get the last child of  $S[z_2]$ 
13: end for
14: return  $(First, Last)$ 

```

2.3.8 Implementation

The XBWT data structure has been implemented in C++ using the Succinct Data Structure Library 2.0 (SDSL) for efficient representation and manipulation of compressed data structures. We will develop two algorithms for constructing the XBWT: one efficient linear-time recursive algorithm and one more straightforward iterative algorithm. Also, we will implement the necessary data structures and algorithms for navigating and querying the XBWT, such as parent-child navigation and path-based searches.

The implementation of the XBWT is based on the descriptions provided in the previous sections. Also, it is available on GitHub at the following link: <https://github.com/davide-tonetto-884585/XBWT>.

Implementation Choices

Follows a list of the main choices made during the implementation of the XBWT:

- The implementation is not focused on a specific kind of data, such as XML documents or JSON files, but it is designed to work with any kind of labeled tree.
- The construction method takes as input a labeled tree. It constructs directly a

compressed indexing scheme based on the Extended Burrows-Wheeler Transform of the tree as described in the previous sections.

- The implementation is based on the Succinct Data Structure Library (SDSL) to handle the compressed data structures generated by the XBWT. The SDSL library provides efficient implementations of various compressed data structures and algorithms, which are essential for representing and querying the XBWT efficiently.
- The labels of the alphabet are encoded as integers, starting from 0 to $|\Sigma| - 1$, where $|\Sigma|$ is the cardinality of the alphabet. This encoding respects the order of the labels in the alphabet and allows simplifying and reducing the space needed to store the labels in the compressed data structure. For this reason, the constructor of the XBWT class takes as input a generic labeled tree.

Succinct Data Structures

The implementation of the XBWT relies heavily on succinct data structures to achieve space efficiency while maintaining fast query operations. In particular, we use succinct data structures to compress the two main arrays of the XBWT: S_α and S_{last} . These arrays, which can be quite large for substantial trees, benefit significantly from compression.

The compression is achieved through the Succinct Data Structure Library (SDSL), which provides efficient implementations of various compressed data structures. For S_{last} , which is a binary sequence, we utilize a compressed bit vector that supports fast rank and select operations. For S_α , which contains labels from a potentially large alphabet, we employ a wavelet tree structure that provides both compression and efficient query capabilities.

The SDSL is a C++ library that provides efficient implementations of various compressed data structures and algorithms. It is used in this project to handle the compressed data structures composing the XBWT. The SDSL library provides a wide range of succinct data structures, such as bit vectors, wavelet trees, and compressed suffix arrays, which are essential for representing and querying the XBWT efficiently. The library is available at <https://github.com/simongog/sdsl-lite> [17]. Let's see the implementation details of the SDSL data structures used in the XBWT implementation.

RRR Bit Vector

The RRR bit vector is designed to provide space-efficient representations of bit vectors while supporting efficient rank and select operations. This data structure implements the RRR (Raman, Raman, and Rao) encoding method, which compresses bit vectors by partitioning them into fixed-size blocks and encoding each block based on its population count (the number of 1s) and specific configuration [32].

The space needed for an RRR bit vector of length n with m set bits is $nH_0 + o(n)$ ($\approx \lceil \log \binom{n}{m} \rceil$). The rank support is provided by `sdsl::rank_support_rrr`, adding 80 bits and requiring $O(\log k)$ time for rank queries, where k is the number of set

bits. The select support is provided by `sdsl::select_support_rrr`, adding 64 bits and requiring $O(\log n)$ time for select queries.

This data structure is used to represent S_{last} , a dedicated binary array B_α that stores the additional information associated with each entry of S_α (i.e., $B_\alpha[i] = 1$ if the i -th symbol in S_α corresponds to a leaf and 0 otherwise), and the A array of the XBWT.

Wavelet Tree

The Wavelet tree is designed to efficiently handle sequences over large alphabets, such as integer sequences. It provides a space-efficient representation while supporting fast access, rank, and select operations. The wavelet tree is a balanced binary tree that recursively partitions the alphabet into two equal-sized subsets and encodes the sequence based on the partitioning [18]. The `sdsl::wt_int` uses the RRR bit vectors or other succinct representations for storing the bit vectors in each node of the wavelet tree. This makes the structure space-efficient.

If RRR-compressed bitvectors are used for the internal bitmaps, a wavelet tree over a sequence $S \in \Sigma^n$ (with $|\Sigma| = \sigma$) occupies $nH_0(S) + o(n \log \sigma) + \Theta(\sigma \log n)$ bits of space, where $H_0(S)$ is the zero-order empirical entropy of S , and it supports access, rank, and select queries in $O(\log \sigma)$ time.

This data structure is used to represent the S_α array of the XBWT.

2.3.9 Experiments

Davide T.: Questa sezione andrà riadattata una volta che avremo deciso quali esperimenti fare con l'altro algoritmo

The experiments have been run on a machine with an AMD Ryzen 9 5600Hs CPU with 24 GB of RAM. The results are shown in Table Table 2.3 and Table Table 2.4. The source code for the experiments can be found in the `experiments.cpp` file.

Construction Performance

To evaluate the performance of the implemented algorithms, we conducted a series of experiments on randomly generated trees created using the Python library `networkx`. The trees were generated with sizes ranging from 100 to 900,000 nodes. For each tree, we executed the construction algorithms 10 times, measuring the average execution time for both the linear *PathSort* (P.S.) algorithm and the naive *UpwardStableSort* (N.S.) algorithm used for constructing the XBWT. This approach allowed us to compare their performance across different tree sizes and assess their scalability.

The results are shown in Table Table 2.3. **Alessio:** Spiega subito i risultati.

Alessio: I numeri vanno sempre allineati a destra, così si capisce meglio chi è più grande di chi. Inoltre, anche il numero di cifre dopo la virgola deve essere sempre lo stesso, così la virgola è fissa e si possono leggere meglio i dati. Su questo faccio io, tu fai sulla prossima :)

Nodes	Depth	P.S. Time (s)	N.S. Time (s)
100	22	0.002	0.001
500	45	0.004	0.002
1000	74	0.006	0.003
5000	175	0.028	0.015
10000	288	0.056	0.053
50000	486	0.310	0.350
100000	754	0.690	1.250
500000	2246	4.700	16.460
900000	2658	8.510	34.200

Table 2.3: Performance comparison between PathSort (P.S.) and Naive Sort (N.S.) algorithms.

Space Analysis

To evaluate the space savings achieved through XBWT compression, we conducted experiments on the same set of randomly generated trees used for the construction performance tests. For each tree, we compared the memory usage (in bytes) of three representations: the plain tree, the uncompressed XBWT, and the compressed XBWT.

The plain tree representation consists of the simple balanced parenthesis encoding of the tree structure combined with the edge labels. For example for tree in Figure ??, the plain tree representation would be:

(A(B(D(a))(a)(E(b)))(C(D(c))(b)(D(c)))(B(D(b))))).

By *uncompressed XBWT*, we refer to the XBWT arrays S_{last} and S_{α} (including the additional bit) stored without any compression. Specifically, S_{last} is represented as a plain bitvector (`sds1::bit_vector`), and S_{α} is stored as a wavelet tree (`sds1::wt_int`) with plain bitvectors (`sds1::bit_vector`). In contrast, the *compressed XBWT* representation stores S_{last} and S_A as compressed RRR bitvectors (`sds1::rrr_vector`), and S_{α} as a wavelet tree with RRR bitvectors, as described in the previous chapter.

Table Table 2.4 reports the sizes (in bytes) for each representation of the trees across different sizes. The last column highlights the space savings achieved by the compressed XBWT compared to the plain tree representation, expressed as a percentage. These results illustrate the substantial space reductions achieved through compression, especially as the tree size increases.

Alessio: Oltre ai punti di prima, metti la percentuale anche per UXBWT, magari non come un'altra colonna ma metti tra parentesi. Te lo faccio sulle prime righe per la C.XBWT. Se ti piace, ricorda di spiegare cosa sono i numeri tra parentesi nella descrizione.

Conclusions

From the results shown in Table Table 2.3, we can draw several conclusions about the performance of the PathSort (P.S.) algorithm compared to the Naive Sort (N.S.) algorithm and the space savings achieved by compressing the XBWT.

Firstly, the PathSort algorithm consistently outperforms the Naive Sort algorithm

Nodes	Plain tree (B)	U. XBWT (B)	C. XBWT (B)	Saving (%)
100	390	424	496 (-27.18%)	
500	2390	1112	1136 (52.47%)	
1000	4890	2242	2056	57.96
5000	28890	12911	10400	64
10000	58890	45625	21848	62.90
50000	338890	175146	123216	63.64
100000	688890	349478	259376	62.35
500000	3888890	1850850	1451570	62.67
900000	7088890	3480190	2718570	61.65

Table 2.4: Space analysis of the XBWT. Plain tree is the size in bytes of the tree in the simple balanced parenthesis representation plus the edge labels, U. XBWT is the size in bytes of the tree in the uncompressed XBWT, and C. XBWT is the size in bytes of the tree in the compressed XBWT. The last column shows the space-saving percentage between plain tree and compressed XBWT.

in terms of execution time, especially as the number of nodes increases. For smaller trees, the difference in execution time between the two algorithms is minimal. However, as the number of nodes grows, the PathSort algorithm demonstrates significantly better scalability. For instance, with 900,000 nodes, the PathSort algorithm takes 8.51 seconds, whereas the Naive Sort algorithm takes 34.2 seconds , giving speedup of more than $4\times$.

Secondly, the depth of the tree appears to increase with the number of nodes, which is expected in randomly generated trees. This increase in depth does not seem to adversely affect the performance of the PathSort algorithm as much as it does the Naive Sort algorithm.

For small trees, the compressed XBWT does not always provide immediate savings due to the overhead of succinct data structures. For instance, for 100 nodes, the compressed representation is larger than the plain tree, showing a -27.18% increase in space. However, as the number of nodes increases, the compression becomes more effective, achieving savings of over 60% for large trees.

The space reduction becomes particularly evident for trees with more than 500 nodes. These results confirm that the compressed XBWT provides a scalable and space-efficient alternative for storing and indexing labeled trees. The efficiency gains are particularly beneficial for applications requiring large-scale tree processing, such as bioinformatics and text indexing.

In conclusion, the PathSort algorithm is a more efficient choice for constructing the XBWT, especially for larger trees, and the compression method provides significant space savings, making the overall process more efficient in terms of both time and space.

2.4 Finite State Automata

2.4.1 Introduction and Motivation

Tree compression schemes that effectively exploit repetitive structures require efficient techniques for identifying and representing such repetitions compactly. A powerful approach to this problem is to view a tree as a finite language, where each

path from the root to a leaf represents a word. Such a language can be recognized by a Deterministic Finite Automaton (DFA). More specifically, since trees are inherently acyclic, they can be represented by Acyclic Deterministic Finite Automata (ADFAs).

The problem of finding and compressing identical subtrees is thus equivalent to minimizing the corresponding DFA. DFA minimization ensures that equivalent substructures are merged efficiently, leading to a more compact encoding. The minimized DFA provides a canonical representation of the repetitive structures, which can then be leveraged in our compression pipeline. This theoretical foundation enables us to identify and encode tree patterns systematically, ultimately improving the compression efficiency.

While general-purpose minimization algorithms like Hopcroft's are highly efficient for any DFA, the specific structure of ADFAs allows for even faster, linear-time algorithms. In this context, we focus on Revuz's algorithm, which is specifically designed to minimize ADFAs and is therefore particularly well-suited for compressing tree structures.

This chapter provides the necessary theoretical background on DFA minimization. We will first introduce the concepts of DFAs and their minimization, followed by a detailed look at both Hopcroft's algorithm as a general solution and Revuz's algorithm as a specialized, linear-time solution for acyclic graphs, which is central to our tree compression methodology.

2.4.2 DFA Minimization

The process of automata minimization consists of reducing the number of states in a DFA while preserving its accepted language. The minimization of DFA is crucial for a variety of applications, such as model checking, hardware design, and compilers, as it produces a more effective and compact representation of the automaton, allowing for faster processing and reduced memory usage.

The minimization of DFA is a well-studied problem in automata theory, and there are several algorithms available for this purpose. One of the most popular algorithms for DFA minimization is Hopcroft's algorithm, which was proposed by John Hopcroft in 1971 [20]. Hopcroft's algorithm is an efficient and simple algorithm that can minimize a DFA in $O(n \log n)$ time, where n is the number of states in the DFA.

The algorithm enables computing equivalence classes of nodes, in particular, the Myhill-Nerode equivalence classes [31, 29]. The Myhill-Nerode theorem states that a language is regular if and only if it has a finite number of Myhill-Nerode equivalence classes. This theorem provides a powerful tool for determining the regularity of languages and is a cornerstone of automata theory. Let's formalize the concept of equivalence classes and the Myhill-Nerode theorem.

Definition 11 (Myhill-Nerode Equivalence Relation). *For a language $L \subseteq \Sigma^*$ and any strings $x, y \in \Sigma^*$, we say x is equivalent to y with respect to L (written as $x \approx_L y$) if and only if for all strings $z \in \Sigma^*$:*

$$xz \in L \Leftrightarrow yz \in L$$

That is, strings x and y are equivalent if they have the same behavior with respect to the language L : either they both lead to acceptance or both lead to rejection when any suffix z is appended.

Theorem 3 (Myhill-Nerode theorem [31, 29]). *Let L be a language over an alphabet Σ . Then L is regular if and only if there exists a finite number of Myhill-Nerode equivalence classes for L . Specifically, the number of equivalence classes is equal to the number of states in the minimal DFA recognizing L .*

Throughout this section let $M = (Q, \Sigma, \delta, q_0, F)$ be a DFA. For $q \in Q$ and $a \in \Sigma$, we adopt the shorthand $q.a := \delta(q, a)$. We extend δ to words by the usual recursion:

$$\delta^*(q, \varepsilon) = q, \quad \delta^*(q, wa) = \delta(\delta^*(q, w), a) \quad \text{for } w \in \Sigma^*, a \in \Sigma.$$

For a word $w = w_1w_2 \dots w_n \in \Sigma^*$, we then write $q.w := \delta^*(q, w)$ for the (unique) state reached from q by reading w . A word w is accepted by M iff $\delta^*(q_0, w) \in F$.

Also, Let $M = (Q, \Sigma, \delta, q_0, F)$ be a DFA recognizing L . For states (nodes) $u, v \in Q$, we say that u and v are MN-equivalent iff

$$\forall \alpha \in \Sigma^* : u.\alpha \in F \iff v.\alpha \in F.$$

2.4.3 Revuz' Minimization Algorithm

For our purpose, we will focus on a specific type of finite automaton: an acyclic deterministic finite automaton. An ADFA is a DFA where the transition graph contains no cycles. The acyclic property is key, as it simplifies the minimization process significantly.

In this section, we will discuss an efficient algorithm for minimizing acyclic deterministic finite automata in linear time on the number of edges [33].

Let's begin by providing some definitions needed to understand the algorithm.

Definition 12 (Height function). *For a state s in an automaton, the height $h(s)$ is defined as the length of the longest path starting at s and going to a final state.*

$$h(s) = \max\{|w| : s.w \text{ is final}\}$$

This height function induces a partition Π_i of Q , where Π_i denotes the set of states of height i .

Lets now define the canonical label of each state that will be necessary to identify MN-equivalent states. For $s \in Q$, let $l_1 < \dots < l_k$ be the symbols of the outgoing transitions defined at s (listed in increasing order of Σ). With $b = F$ if $s \in F$ and $b = NF$ otherwise, we set

$$l(s) := (b, l_1, s.l_1, l_2, s.l_2, \dots, l_k, s.l_k).$$

Also, the algorithm uses a function R to map the labels of states to a new signature. This function is defined as follows:

Definition 13 (Signature map R). Let $N[\cdot]$ be the current renaming array that assigns to each state its equivalence class identifier Myhill-Nerode. For a state s labeled

$$l(s) = (b, l_1, s.l_1, l_2, s.l_2, \dots, l_k, s.l_k),$$

where $b \in \{F, NF\}$, $l_i \in \Sigma$ (listed in increasing order), and $nl_i \in Q$, define

$$R(l(s)) = (b, l_1, N[s.l_1], l_2, N[s.l_2], \dots, l_k, N[s.l_k]).$$

It is important to notice that, since the automaton is acyclic, every transition $s \xrightarrow{a} t$ strictly decreases the height: $h(t) < h(s)$. The main loop of Algorithm 8 processes levels in increasing order $i = 0, 1, \dots$, so by the time we handle a state $s \in \Pi_i$, all its targets t lie in $\bigcup_{j < i} \Pi_j$ and have already been assigned a Myhill-Nerode equivalence class.

Algorithm 8 RevuzMinimization(ADFA)

Require: ADFA $M = (Q, \Sigma, \delta, q_0, F)$

Ensure: Minimal DFA $M' = (\{1, \dots, n\}, \Sigma, \delta', N[q_0], F')$ with $F' = \{N[q] \mid q \in F\}$
and $\delta'(N[q], a) = N[\delta(q, a)]$

- 1: Calculate height $h(s)$ for every state s .
 - 2: Create partitions $\Pi_i = \{s \in Q \mid h(s) = i\}$.
 - 3: $N[1, |Q|] = \{1, 2, \dots, |Q|\};$ \triangleright Renaming array
 - 4: $n = 0;$
 - 5: **for** $i := 0$ to $h(q_0)$ **do** $\triangleright q_0$ is the initial state
 - 6: Sort states in Π_i based on $R(l(q)), q \in \Pi_i$.
 - 7: $n = n + 1;$
 - 8: $N[\Pi_i[1]] = n;$
 - 9: **for** $j := 2$ to $|\Pi_i|$ **do**
 - 10: **if** $R(l(\Pi_i[j])) \neq R(l(\Pi_i[j - 1]))$
 - 11: $n = n + 1;$
 - 12: **end if**
 - 13: $N[\Pi_i[j]] = n;$
 - 14: **end for**
 - 15: **end for**
-

The algorithm proceeds level by level, from $i = 0$ up to the maximum height, ensuring that states at each level are correctly partitioned into Myhill-Nerode equivalence classes. For each level i , it groups the states in Π_i based on their signatures computed by the function R (see Definition 13). As explained before, when processing level i , the equivalence classes for all states in lower levels ($j < i$) have already been finalized. The signature $R(l(s))$ for a state s depends on its finality and the equivalence classes of its immediate successors. Therefore, two states $s, t \in \Pi_i$ have the same signature if and only if they are MN-equivalent. The algorithm assigns a unique class identifier to each group of states with the same signature.

The whole algorithm can be implemented to run in time $O(m)$ for an acyclic automaton with m edges. Heights may be computed in linear time by a bottom-up traversal. The lists of states of a given height are collected during this traversal.

The signature of a state is easy to compute provided the edges starting in a state have been sorted (by a bucket sort for instance to remain within the linear time constraint). Sorting states by their signature again is done by a lexicographic sort [2].

Example 2.13:

Now we are going to see an example of reduction for a given ADFA. The ADFA is represented in figure Figure 2.2 and, as we can notice, it is also a valid ordered rooted tree with $n = 11$ nodes, $e = 10$ edges, and the following alphabet: $\Sigma = \{0, 1\}$. The node a is the root of the tree and the initial state of the automaton, while the leaf nodes e, g, h, i, l, m are final states. It is important to note that while the algorithm applies to any ADFA, our focus is on those that are also trees, as this is the specific case relevant to our work.

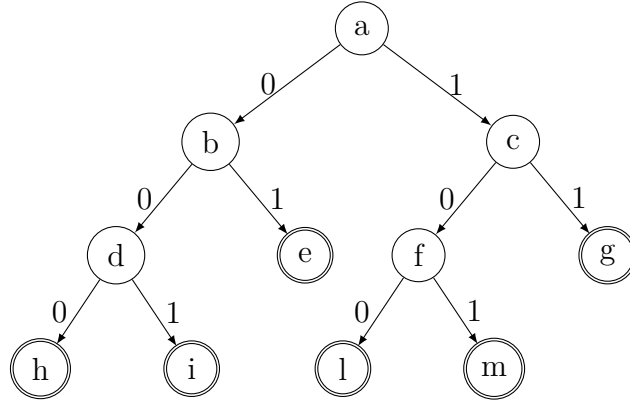


Figure 2.2: Example ADFA to be minimized

Now, let's apply the minimization algorithm step by step:

1. **Height Computation:** First, we compute the height of each state. The height is the length of the longest path to a final state. The final states (e, g, h, i, l, m) have a height of 0. For the other states, the height is calculated as follows:

- $h(d) = 1 + \max(h(h), h(i)) = 1 + 0 = 1$
- $h(f) = 1 + \max(h(l), h(m)) = 1 + 0 = 1$
- $h(b) = 1 + \max(h(d), h(e)) = 1 + \max(1, 0) = 2$
- $h(c) = 1 + \max(h(f), h(g)) = 1 + \max(1, 0) = 2$
- $h(a) = 1 + \max(h(b), h(c)) = 1 + \max(2, 2) = 3$

This gives us the following partitions based on height:

- $\Pi_0 = \{e, g, h, i, l, m\}$
- $\Pi_1 = \{d, f\}$
- $\Pi_2 = \{b, c\}$
- $\Pi_3 = \{a\}$

2. **Processing Π_0 :** All states in Π_0 are final and have no outgoing transitions, so they are all equivalent. We merge them into a single class, let's call it $D = \{e, g, h, i, l, m\}$. After this step, we have a new state D which is final.

3. **Processing Π_1 :** Now we examine the states in Π_1 : d and f . We check their transitions:

- State d : $\delta(d, 0) = h \in D$ and $\delta(d, 1) = i \in D$.
- State f : $\delta(f, 0) = l \in D$ and $\delta(f, 1) = m \in D$.

Since both states transition to the same equivalence class (D) for both symbols 0 and 1, they are equivalent. We merge them into a new class, $C = \{d, f\}$.

4. **Processing Π_2 :** Next, we process the states in Π_2 : b and c .

- State b : $\delta(b, 0) = d \in C$ and $\delta(b, 1) = e \in D$.
- State c : $\delta(c, 0) = f \in C$ and $\delta(c, 1) = g \in D$.

Both states have transitions to class C on symbol 0 and to class D on symbol 1. Therefore, b and c are equivalent. We merge them into a new class, $B = \{b, c\}$.

5. **Processing Π_3 :** Finally, we process Π_3 , which contains only state a . There is nothing to compare it with, so it forms its class, $A = \{a\}$.

After applying the algorithm, we obtain the minimized ADFA represented in figure Figure 2.3. Each node of the original ADFA is represented by a node in the minimized ADFA (equivalence classes). The edges represent transitions between these nodes. The root node A is the initial state of the minimized ADFA, while the node D is the final state.

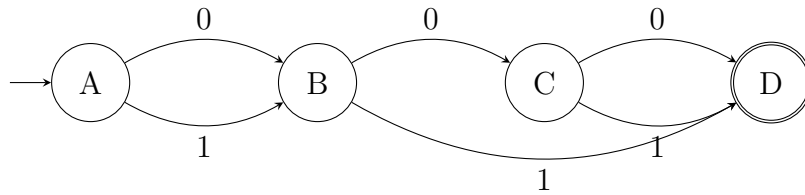


Figure 2.3: Minimized ADFA

The equivalence classes of the nodes are listed in table Table 2.5.

Class	States
A	a
B	b, c
C	d, f
D	e, g, h, i, l, m

Table 2.5: Equivalence classes of the nodes

2.5 Wheeler and p -sortable Graphs

2.5.1 Introduction and Motivation

As established in the introduction, the primary goal of this thesis is to find an effective balance between compressing a finite language and preserving its indexability. The two extremes—full DFA minimization and the raw input trie—are inadequate, as one sacrifices indexing for compression and the other sacrifices compression for indexing. The solution to this problem lies in a specific class of graphs that are structured enough to be indexed efficiently yet flexible enough to allow for significant compression. This chapter introduces the theoretical framework that underpins our approach: Wheeler graphs and their generalization, p -sortable graphs.

A crucial observation is that the input trie representing our language is already a highly structured object. It is a Wheeler graph, a type of graph that admits a special ordering on its nodes and edges, making it exceptionally well-suited for indexing. In formal terms, a trie is a 1-sortable graph. This property explains both its powerful indexing capabilities and its inherent lack of compression.

The concept of p -sortability offers a way to navigate the trade-off. By controllably increasing the sortability parameter p , we can begin to merge MN-equivalent states (see Theorem 3), thereby compressing the graph. The resulting automaton is no longer a simple trie but a more general p -sortable graph that retains strong indexing properties. This chapter will formally define these concepts, which are the foundation of our algorithm for achieving a practical compromise between compression and indexability.

2.5.2 Orders

The core property that makes Wheeler and p -sortable graphs efficiently indexable is the existence of a specific ordering on their states. This ordering provides the necessary structure to navigate the automaton and answer queries quickly, a task that is computationally hard on general graphs. The fundamental ordering used in this context is the *co-lexicographic order* (co-lex), which compares states based on the labels of the paths that reach them. This section formally defines this order and the related concepts that are essential for understanding the structure of indexable automata.

Definition 14 (Co-lexicographic Order on Σ^*). *The co-lex order \preceq is defined as follows. Given two strings $\alpha, \beta \in \Sigma^*$, we say that $\alpha \preceq \beta$ if and only if either:*

- α is a suffix of β , or
- there exist strings $\alpha', \beta', \gamma \in \Sigma^*$ and symbols $a, b \in \Sigma$, such that $\alpha = \alpha'a\gamma$, $\beta = \beta'b\gamma$, and $a \prec b$.

Now, let's define the formal concept of partial order and the width of a partial order.

Definition 15 (Partial Order). *A partial order is a binary relation \leq over a set S that is reflexive, antisymmetric, and transitive. That is, for all $a, b, c \in S$:*

- $a \leq a$ (reflexivity)

- if $a \leq b$ and $b \leq a$, then $a = b$ (antisymmetry)
- if $a \leq b$ and $b \leq c$, then $a \leq c$ (transitivity)

A partial order (S, \leq) can be visualized using a *Hasse diagram*. In a Hasse diagram, each element of S is represented by a node, and given $a, b, c \in S$ there is a line segment or curve going upward from a to b if $a \leq b$ and there is no element c such that $a \leq c \leq b$. The direction of the relation is implicitly understood to be upwards, so arrows are not needed. Also, if two elements of S are incomparable under \leq they are displayed at the same level in the diagram. An example is given in Example 2.14.

Example 2.14:

In the example shown in Figure 2.4, the set is composed of the divisors of 12, and the relation is divisibility. An edge is drawn from a to b if a divides b and there is no other element c in the set such that a/c and c/b . For instance, there is an edge from 2 to 4 because 2 divides 4, and no other element in the set is a multiple of 2 and a divisor of 4. There is no direct edge from 2 to 12 because the relationship is captured transitively through other elements, such as $2/4/12$ or $2/6/12$.

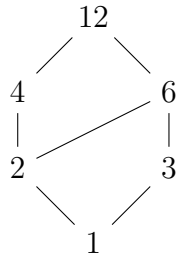


Figure 2.4: Hasse diagram for the set $\{1, 2, 3, 4, 6, 12\}$ with the "divides" relation.

Now we can define the concept of width of a partial order.

Definition 16 (Antichain). *An antichain of a partially ordered set (S, \leq) is a subset of S where any two distinct elements are incomparable. That is, for any two distinct elements a, b in the antichain, neither $a \leq b$ nor $b \leq a$ holds.*

Definition 17 (Width of a partial order, [8]). *The width of a partially ordered set is the size of the largest possible antichain.*

By Dilworth's Theorem [8], the width of a partially ordered set (S, \leq) is equal to the cardinality of its largest antichain; this can be equivalently defined as the minimum number of chains needed to partition S , where each chain is a totally ordered subset of S under the relation \leq .

2.5.3 Wheeler Graphs

With the concept of co-lex order established, we can now define the class of graphs that form the starting point of our work. A Wheeler automaton is an automaton where the states can be arranged in a strict, total order.

Definition 18 (Wheeler automaton, [14]). *A finite state automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ is a Wheeler automaton if there exists a total order \leq on its set of states Q that satisfies the following axioms:*

1. The initial state precede all other states in the order.

For any two transitions $u \in \delta(u', a)$ and $v \in \delta(v', b)$:

2. $a < b \implies u \leq v$,
3. $a = b \wedge u' < v' \implies u \leq v$.

The order \leq is called a *Wheeler order*.

Consequently, we define the concept of Wheeler language.

Definition 19 (Wheeler language). *A Wheeler language L is a language accepted by a deterministic Wheeler automaton.*

The most important example for Wheeler automaton in this thesis is the trie. Any trie representing a finite language is a Wheeler automaton. The co-lexicographic order of the strings spelling the paths from the root to each node provides the required total ordering of the states. This is why tries are inherently indexable. However, this rigid structure also means they are uncompressed, as every unique path must be stored explicitly, even if it corresponds to a substring that appears many times in the language. Our work begins with this observation: we start with a Wheeler automaton (the trie) and seek to compress it while preserving efficient indexability.

2.5.4 The Co-lex Width of an Automaton

Now that we have the concepts of co-lex order and width, we can combine them to formally define the class of indexable automata that are central to this thesis. The width of the co-lex partial order on an automaton's states is the critical measure of its structural complexity from an indexing perspective.

The co-lex order can be extended to the set of states of an automaton. The idea of co-lex order on the states of an automaton was first introduced with the notion of Wheeler graphs by Gagie et al. [14] and was later generalized to arbitrary finite automata by Cotumaccio and Prezra [6], where a partial order replaces the total order. The definition of co-lex order on an automaton is as follows:

Definition 20 ([7]). *Let $N = (Q, \Sigma, \delta, q_0, F)$ be an NFA. A co-lex order on N is a partial order \leq on Q that satisfies the following two axioms:*

1. For every $u, v \in Q$, if $u < v$, then $\max \lambda(u) < \min \lambda(v)$
2. For every $a \in \Sigma$ and $u, v, u', v' \in Q$, if $u \in \delta(u', a)$, $v \in \delta(v', a)$ and $u < v$, then $u' \leq v'$

Where $\lambda(q)$ denotes the set of labels of transitions entering state q , and $\min \lambda(q)$ and $\max \lambda(q)$ represent the minimum and maximum element of the set, respectively.

The two axioms in Definition 20 allow for pair of states of a finite automaton to be compared. When \leq is total, we say that the co-lex order is a Wheeler order (introduced in [14]).

Consequently, we can introduce the concept of co-lex width of an automaton.

Definition 21 ([7]). *The co-lex width of an NFA N is the minimum width of a co-lex order on N .*

$$\text{width}(N) = \min\{|\text{width}(\leq)| \mid \leq \text{ is a co-lex order on } N\}$$

The requirement of a Wheeler order is powerful but restrictive. Many automata, especially those resulting from DAG compression, may not satisfy it (see Example 2.15). This introduces a fundamental trade-off: while DAG compression minimizes an automaton's size, it can destroy the very structure that enables efficient indexing. In fact, it has been shown that indexing general graphs—and thus, highly compressed automata—to support fast string matching is computationally expensive, as showed in [11]. The second axiom of Definition 20 does not always enforce an ordering between any two states, leading to a partial order instead of a total one. This gives rise to the more general notion of a *p-sortable automaton*, where p is the co-lex width of the automaton (see Example 2.16). Under this definitions, a Wheeler automaton is a **1-sortable automaton**, as a total order has a width of 1 (the largest antichain is a single element).

Example 2.15:

State incomparability can arise in several situations. For example, consider two states u and u' .

- As illustrated in Figure 2.5-(a), if we have two same-labeled transitions where the target states u and v are incomparable, we have no information about the relative order of the source states u' and v' . Consequently, they may also be incomparable.
- Conflicting constraints from different labels can force incomparability, as shown in Figure 2.5-(b). An existing order on the sources of a -transitions (e.g., $u' < v'$) may require $u < v$ to satisfy the Wheeler axioms, while an order on the sources of b -transitions (e.g., $v'' < u''$) may require the opposite, $v < u$. Since both cannot be true, the targets u and v must be incomparable.

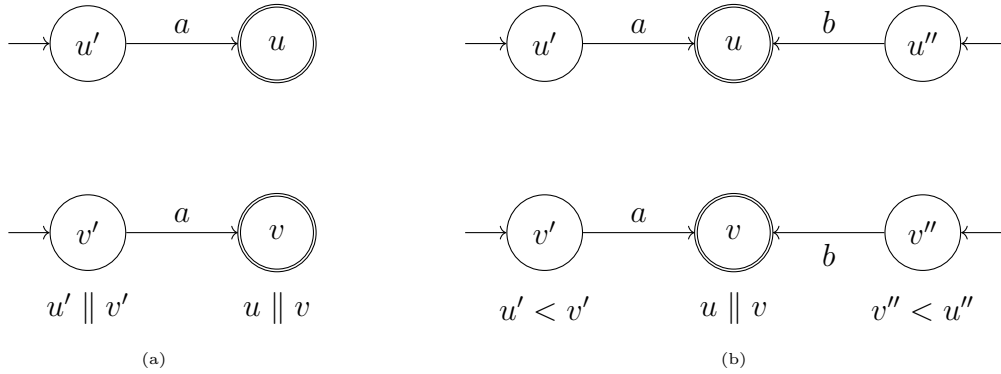


Figure 2.5: Examples of state incomparability in automata.

Example 2.16:

Consider the following DFA D of Figure 2.6-(a) and its partial co-lex order Figure 2.6-(b). The DFA is not Wheeler because states v_3 and v_5 and states v_4 and v_7 are incomparable (as shown in the Hasse diagram). D admits a partition into 2 \leq -chains, for example one possible \leq -chain partition is given by:

- Chain 1: $v_1 \leq v_2 \leq v_3 \leq v_4$
- Chain 2: $v_5 \leq v_6 \leq v_7$

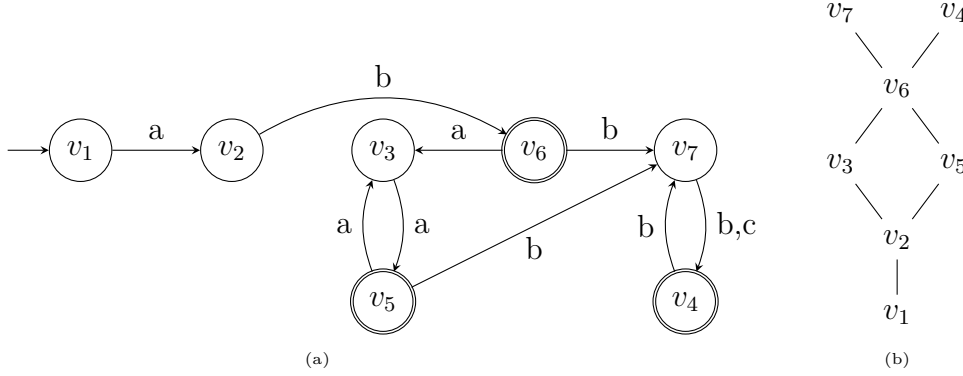


Figure 2.6: An example of a 2-sortable DFA and its corresponding Hasse diagram of the partial order.

We now introduce another important result by [28]. Let $D = (Q, \Sigma, \delta, s, F)$ denote the minimal DFA accepting a Wheeler language L , and let $D^w = (Q^w, \Sigma, \delta^w, s^w, F^w)$ denote the minimal Wheeler DFA (WDFA) accepting L . Since D^w is Wheeler, the rational embedding ensures that for any two distinct states $q, q' \in Q^w$, the associated intervals $I_q, I_{q'}$ are disjoint. This property does not generally hold for D , where states may correspond to overlapping sets of prefixes. As a consequence, when transforming D into D^w , certain states of D may need to be *split* into several states in D^w , potentially leading to an exponential blow-up in the number of states.

Example 2.17: [28]

We now provide a simple example of an automaton D with width n that accepts a Wheeler language, yet its minimum equivalent Wheeler DFA, D^w , is exponentially larger. Let D be the automaton depicted in Figure 2.7.

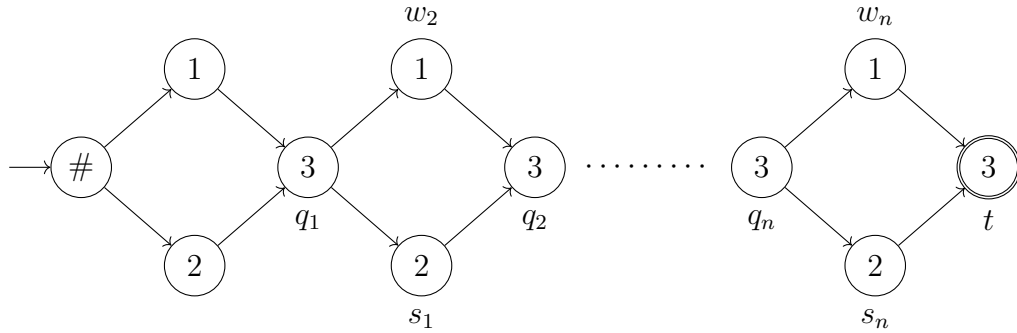


Figure 2.7: A DFA accepting a finite (and thus Wheeler) language, for which the minimal equivalent Wheeler DFA is exponentially larger.

The language $L = \mathcal{L}(D)$, being finite, is a Wheeler language [1]. However, any Wheeler automaton accepting L must have a number of states that is exponential in n . This is because for any pair of distinct strings $\alpha, \gamma \in I_t$, it is possible to find another string β such that the co-lexicographic order of the strings is $\alpha < \beta < \gamma$. Since there are exponentially many such pairwise distinct strings leading to t , a Wheeler automaton must partition the set I_t into an exponential number of sub-intervals. This forces the state t to be "split" into exponentially

many copies, leading to an exponential blow-up in the size of the minimal Wheeler DFA, D_w .

The previous example highlights a crucial trade-off: enforcing the strict ordering of a Wheeler DFA can lead to an exponential increase in the number of states compared to a minimal DFA. Theorem 4 formalizes this observation, showing that the size of a minimal Wheeler DFA cannot be bounded by a polynomial in the size of the minimal DFA.

Theorem 4 ([28], Theorem 29). *Let $L = \mathcal{L}(D) = \mathcal{L}(D^w)$, where L is Wheeler, D is minimal, D^w is minimal Wheeler, and let $f(\cdot, \cdot)$ be such that $|D^w| = \mathcal{O}(f(|D|, \text{width}(D)))$. Then, for any $k, p \in \mathbb{N}$, $f(n, p) \notin \mathcal{O}(n^k + 2^p)$.*

Since this work aims to transform a trie (a 1-sortable automaton) into a more general p -sortable graph (with width $p > 1$) in a way that introduces DAG compression while maintaining efficient indexability, it is motivated by the powerful result of Theorem 4 leading to the fact that even a small increase in sortability (for example from $p = 1$ to $p = 2$) can yield exponential compression. This highlights the potential of exploring the trade-off between sortability and size, which is the central theme of this thesis.

2.5.5 Indexing Finite State Automata

Now we introduce the current state of the art in indexing finite state automata. In 2023, Cotumaccio et al. [6] introduced a compressed data structure for automata whose performance and space complexity are directly tied to the automaton's colex width, p . This structure generalizes the famous Burros-Wheeler transform [4] and supports subpath queries (Definition 8) on a query word α of length m in $\mathcal{O}(mp^2 \log(p|\Sigma|))$ time. The space required is $\log(|\Sigma|) + \log p + 2$ bits per edge for DFAs and $\log(|\Sigma|) + 2 \log p + 2$ bits per edge for NFAs. This highlights a direct trade-off: both query time and space per edge depend on the width parameter p , which governs the automaton's compressibility.

To highlight the importance of this data structure, we recall that the final output of our compression pipeline is a p -sortable DAG compressed automaton with a controlled colex width p . This allows us to leverage these advanced indexing capabilities on the compressed automata produced by our method.

2.6 Min-Weight Perfect Bipartite Matching

2.6.1 Introduction and Motivation

The fundamental goal of our compression scheme is to transform an input trie into a compressed, p -sortable automaton by partitioning its nodes into p chains in an optimal way.

To make this optimization problem more concrete, we can frame it as a string partitioning problem. Imagine the sequence of nodes in the trie, when read in colexicographic order, as a single long string. The "character" corresponding to each node is its Myhill-Nerode equivalence class, which determines if it can be merged

with other nodes. The task is to partition this string of nodes into p subsequences such that the number of runs is minimized (consecutive nodes of the same equivalence class are merged into a single "run"). For instance, a subsequence $AAABBA$ contains three runs. Minimizing the number of runs directly corresponds to maximizing the number of merged states, yielding a compact p -sortable automaton.

In this sections, we will provide the necessary background on the Minimum Weight Perfect Bipartite Matching (MWPBM) problem, a fundamental challenge in combinatorial optimization. We will then demonstrate in Section 2.3 a formal reduction from our partitioning problem to MWPBM. This reduction is the key to our method, as it allows us to model our problem as a bipartite graph and leverage well-known, efficient algorithms to find the optimal solution for our compression task.

2.6.2 Bipartite Graphs

Definition 22. A graph $G = (V, E)$ is called *bipartite* if its vertex set V can be partitioned into two disjoint subsets $V = V_1 \cup V_2$ such that every edge in E has the form (v_1, v_2) where $v_1 \in V_1$ and $v_2 \in V_2$.

In other words, the vertices of the graph can be divided into two separate groups such that all edges connect a vertex from the first group to a vertex from the second group. An example of a bipartite graph is shown in Figure 2.8.

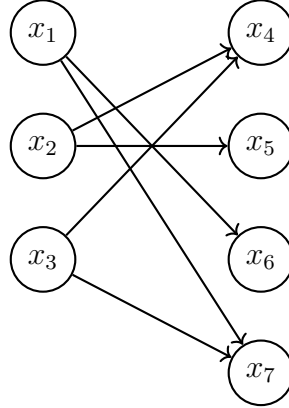


Figure 2.8: Example of a bipartite graph $G = (V, E)$ where $V_1 = \{x_1, x_2, x_3\}$, $V_2 = \{x_4, x_5, x_6, x_7\}$ and $E = \{(x_1, x_4), (x_1, x_6), (x_1, x_7), (x_2, x_4), (x_2, x_5), (x_2, x_7), (x_3, x_4), (x_3, x_6), (x_3, x_7)\}$.

Definition 23 (Weighted Bipartite Graph). A *weighted bipartite graph* is a bipartite graph $G = (V, E, w)$, where w is a weight function $w : E \rightarrow \mathbb{R}$ that assigns a real-valued weight to each edge.

In a weighted bipartite graph, each edge has a numerical value, or "weight", associated with it.

2.6.3 Problem Definition

Given a bipartite graph $G = (V, E)$ (Definition 22), let's define the concept of a matching.

Definition 24 (Matching). Given a bipartite graph $G = (V, E)$, a *matching* $M \subseteq E$ is a collection of edges such that every vertex of V is incident to at most one edge.

of M .

In other words, a matching is a set of edges such that no two edges share a common vertex. If a vertex v has no edge of M incident to it, then v is said to be exposed (or unmatched). A matching is **perfect** if no vertex is exposed; in other words, a matching is perfect if its cardinality is equal to $|V_1| = |V_2|$ [16].

Example 2.18:

In Figure 2.9, we illustrate three distinct scenarios. Subfigure (a) depicts a set of edges that does not constitute a valid matching, as vertex u_1 is incident to more than one edge, namely (u_1, v_1) and (u_1, v_2) , violating the definition of a matching. Subfigure (b) presents a valid, yet non-perfect matching; here, vertices u_3 and v_3 are exposed, meaning they are not incident to any edge in the matching. Finally, subfigure (c) shows a perfect matching, where every vertex in the graph is incident to exactly one edge in the matching, satisfying the condition $|M| = |V_1| = |V_2| = 3$.

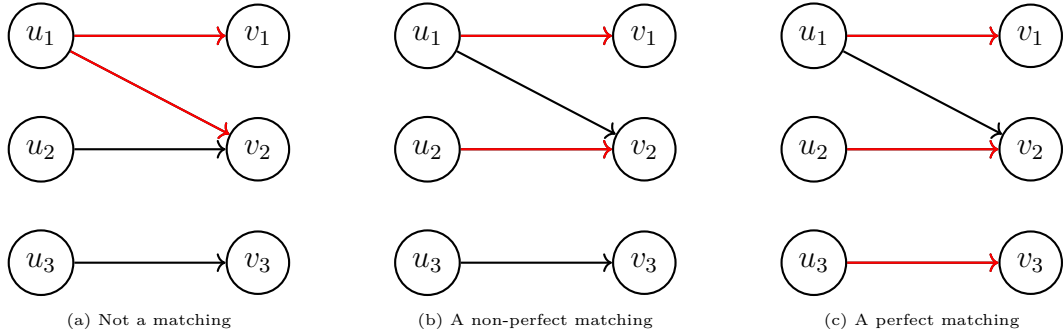


Figure 2.9: Examples of a non-matching (a), a non-perfect matching (b), and a perfect matching (c) in a bipartite graph. The edges in the set M are highlighted in red.

The problem of finding a minimum weight perfect matching in a weighted bipartite graph (Definition 23) is a well-known problem in combinatorial optimization. The problem can be formulated as follows:

Definition 25 (Minimum weight perfect matching in weighted bipartite graphs). *Given a weighted bipartite graph $G = (V, E, w)$, find a perfect matching M such that the sum of the weights of the edges in M is minimized.*

The weight of a matching is the sum of the weights of the edges in the matching. We define the weight of a matching M as follows:

$$w(M) = \sum_{e \in M} w(e) \quad (2.1)$$

Example 2.19:

Consider the weighted bipartite graph in Figure 2.10. The goal is to find a perfect matching with the minimum possible total weight. Both subfigures show a valid perfect matching; however, only one of them has the minimum weight.

- Subfigure (a) shows the perfect matching $M_a = \{(u_1, v_2), (u_2, v_1), (u_3, v_3)\}$. Its total weight is $w(M_a) = w(u_1, v_2) + w(u_2, v_1) + w(u_3, v_3) = 2 + 1 + 1 = 4$. This is a valid perfect matching, but it is not optimal.

- Subfigure (b) shows the perfect matching $M_b = \{(u_1, v_1), (u_2, v_2), (u_3, v_3)\}$. Its total weight is $w(M_b) = w(u_1, v_1) + w(u_2, v_2) + w(u_3, v_3) = 1 + 1 + 1 = 3$.

Since $w(M_b) < w(M_a)$, the matching in (b) is a minimum weight perfect matching for this graph, while the matching in (a) is a non-minimum perfect matching.

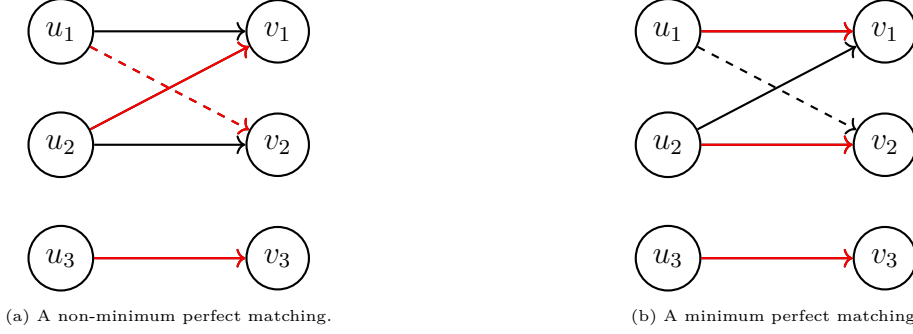


Figure 2.10: Example of a non-minimum perfect matching (a) and minimum perfect matching (b) in a weighted bipartite graph. Dashed edges have weight 2 while solid edges 1. The edges in a matching are highlighted in red.

2.6.4 Hall's Marriage Theorem

Hall's Marriage Theorem [19] provides a necessary and sufficient condition for the existence of a matching that saturates an entire partition of a bipartite graph. Specifically, it allows us to determine if there is a matching that includes every vertex from the smaller of the two partitions. In the special case where the partitions are of equal size ($|V_1| = |V_2|$), the theorem gives a condition for the existence of a perfect matching.

This theorem will be crucial in proving the correctness of our reduction from the tree partitioning problem to the minimum weight perfect bipartite matching problem. Specifically, we will use it to show that our constructed bipartite graph always admits a perfect matching, ensuring that our reduction is valid.

Definition 26 (Neighborhood). *For a subset of vertices $W \subseteq V_1$, the **neighborhood** of W , denoted by $N(W)$, is the subset of all vertices in V_2 that are adjacent to at least one vertex in W .*

$$N(W) = \{v \in V_2 \mid \exists u \in W \wedge \{u, v\} \in E\}$$

Theorem 5 (Hall's Marriage Theorem [19]). *Let $G = (V, E)$ be a bipartite graph with $V = V_1 \cup V_2$. There exists a matching M in G that covers V_1 if and only if for every subset $W \subseteq V_1$, the following condition holds:*

$$|N(W)| \geq |W|$$

This condition is known as **Hall's condition**.

In simpler terms, a matching that covers all vertices in V_1 exists if and only if every group of vertices chosen from V_1 collectively has at least as many neighbors in V_2 as there are vertices in the chosen group.

2.6.5 State of the Art

There are several algorithms to solve the problem of finding a minimum weight perfect matching in a bipartite graph. The first algorithm to solve this problem was proposed by Kuhn in 1955 [24]. The algorithm is based on the Hungarian method, which is a combinatorial optimization algorithm that solves the MWPBM problem in polynomial time. In the original paper the complexity of the algorithm was $O(n^4)$ where n is the number of nodes in the bipartite graph. Later Dinic and Kronrod [9] showed that the algorithm can be implemented in $O(n^3)$ time.

The Hungarian method is a powerful algorithm; however, it is not very intuitive and can be difficult to implement. In recent years, several other algorithms have been proposed to solve this problem. In 1970, Edmonds and Karp [10] proposed an algorithm that solves the problem in $O(nm + n^2 \log n)$ time, where m is the number of edges. In 1989 Gabow and Tarjan [13] proposed an algorithm that solves the problem in $O(\sqrt{nm} \log(nW))$ time, where W denote the highest edge weight in the graph; costs are assumed to be integral. The algorithms work by scaling. Lastly, in 2009, Sankowski and Piotr [34] introduced a randomized algorithm that solves the problem in $O(Wn^\omega)$ time, where ω is the exponent of matrix multiplication, and W is the highest edge weight in the graph.

In 2022, Chen, Li, et al. [5] proposed a nearly linear time algorithm for the Minimum Cost Flow (MCF) problem, running in $O(m'^{1+o(1)})$ on a network with m' edges. This is highly relevant as the MWPBM problem can be reduced to MCF. Specifically, an MWPBM instance on a bipartite graph with n vertices and m edges can be transformed into an MCF problem on a network with $m' = n^2 + m$ edges.

Davide T.: move in experiment chapter (?)

2.6.6 Implementation

To run experiments for the proposed compression scheme, we used a C++ implementation by Vladimir Kolmogorov of the minimum cost perfect matching algorithm described in [22]. The implementation is available at <https://pub.ista.ac.at/~vnk/software.html>.

In brief, paper [22] presents Blossom V, a practical implementation of Edmonds' blossom algorithm for computing minimum-cost perfect matchings in undirected weighted graphs. While the theoretical worst-case bounds for the blossom family have steadily improved since Edmonds' original $O(n^2m)$ algorithm, Blossom V is designed for strong empirical performance rather than new asymptotic guarantees. It combines two ingredients that had proven effective separately in prior work: the variable δ (variable dual updates) strategy popularized by Blossom IV, and systematic use of priority queues to efficiently select minimum-slack edges.

Blossom V targets general (not necessarily bipartite) graphs, and thus directly applies to our bipartite instances as a special case. In our experiments we use the publicly available Blossom V implementation as a black-box solver to compute minimum-cost perfect matchings for the graphs generated by our compression pipeline.

Chapter 3

Tree Compression Scheme

As introduced in the first chapter, the primary goal of this thesis is to develop a novel tree compression scheme that effectively leverages repetitive structures within the input trie. The proposed algorithm is designed to identify and compactly represent these recurring patterns, thereby improving compression performance, particularly for highly repetitive tries. This chapter provides an overview of the proposed compression scheme.

3.1 Compression Scheme Pipeline

The overall pipeline of our proposed method is outlined in Algorithm 9. It takes an input trie T and a width parameter p and produces a compressed, p -sortable automaton.

Algorithm 9 *CompressTrie*(T, p)

Require: Input trie T , width integer parameter p

Ensure: A compressed, p -sortable automaton \mathcal{A}

- 1: $V_{sorted} \leftarrow \text{PathSort}(T)$
 - 2: $N[1, \dots, |Q|] \leftarrow \text{ComputeEquivalenceClasses}(T)$
 - 3: $G_{bipartite} \leftarrow \text{ConstructMWPBMInstance}(V_{sorted}, N, p)$
 - 4: $M \leftarrow \text{SolveMWPBM}(G_{bipartite})$
 - 5: $C[1, \dots, p] \leftarrow \text{ExtractChainsFromMatching}(M, V_{sorted})$
 - 6: $\mathcal{A} \leftarrow \text{CollapseChains}(C, N)$
 - 7: **return** \mathcal{A}
-

The first step of the pipeline (line 2) is to establish a total order on the nodes of the trie. This is achieved by sorting the nodes co-lexicographically using the **Path Sort** algorithm, which we detail in Subsection 2.3.4. This sorting is fundamental, as it arranges the nodes in the order required for a Wheeler automaton.

Next, the algorithm identifies which nodes are candidates for merging (line 3). This is done by partitioning the nodes into equivalence classes based on the structure of the subtrees rooted at each node. Two nodes are in the same class if and only if their subtrees are isomorphic. This is equivalent to computing the Myhill-Nerode equivalence classes for the finite language accepted by the trie, a process we adapt from Revuz’s algorithm for minimizing acyclic DFAs (Subsection 2.4.3).

The core of the algorithm lies in lines 5-7, where the problem of optimally partitioning the sorted nodes into p chains is solved. As detailed in ??, we reduce this problem to finding a Minimum Weight Perfect Bipartite Matching (MWPBM). A

bipartite graph is constructed where the weights on the edges correspond to the "cost" of placing two nodes in the same chain. By finding a perfect matching with minimum weight using standard algorithms (Subsection 2.6.5), we can reconstruct a set of p chains that minimizes the total number of runs (Definition 1), thereby maximizing compression.

Finally, the compressed automaton is constructed (line 8). The algorithm iterates through each of the p chains and "collapses" any consecutive sequence of nodes belonging to the same equivalence class into a single state. This run-length compression, which we describe in Section 3.3, produces the final p -sortable automaton, which can then be indexed for efficient querying (??).

Example 3.1:

In our running example, we begin with the tree shown in Figure 3.1. Each node in this tree is labeled with its corresponding Myhill-Nerode equivalence class, as determined in Example 2.13. By traversing the nodes in co-lexicographic order, we construct a string S where each character represents the Myhill-Nerode equivalence class of a node:

$$S = \text{ABCDDCBDDDD}$$

This string S then becomes the input to the STRING PARTITIONING PROBLEM (see ??), where the objective is to partition it into p subsequences while minimizing the total number of runs.

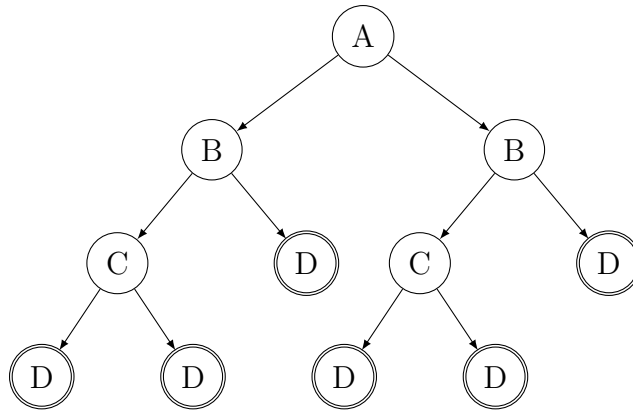


Figure 3.1: Tree ADFA of Figure 2.2. Each node is labeled with its equivalence class.

3.2 Reducing the String Partitioning Problem to the MWFBM Problem

In the previous chapters, we modeled the task of partitioning trie nodes into p chains as the String Partitioning Problem (see Definition 2). We now demonstrate that this problem can be solved in polynomial time by reducing it to the MWFBM problem. This section will first detail the construction of a bipartite graph from an instance of the String Partitioning Problem. Then, we will prove that a minimum-weight perfect matching in this graph directly corresponds to an optimal solution for the partitioning problem. Let's start with an example.

Example 3.2: String Partitioning Problem

Let $S = \text{AABACABB}$ be the input string and let $p = 2$ be the desired number of subsequences. Our goal is to partition the characters of S into two subsequences, S_1 and S_2 , such that the total number of runs (Definition 1) is minimized.

Consider the following partition:

- S_1 is formed by taking the 1st, 2nd, 4th, and 6th characters of S : AAAA.
- S_2 is formed by the remaining characters (3rd, 5th, 7th, and 8th): BCB.

The number of runs for each subsequence is:

- $\text{Runs}(S_1) = 1$ (the run is "AAAA").
- $\text{Runs}(S_2) = 3$ (the runs are "B", "C", "BB").

The total number of runs for this partition is $1 + 3 = 4$. An optimal solution to the String Partitioning Problem would be a partition that achieves the minimum possible total number of runs. In this case, 4 is indeed the optimal value.

3.2.1 Bipartite Graph Construction

Now, we will show how to construct a bipartite graph that allows us to solve the String Partitioning problem.

Definition 27 (Bipartite graph construction). *Let $A = a_1a_2 \dots a_n$ be a string from an alphabet Σ and let p be the number of subsequences we want to partition S in. We can construct a weighted bipartite graph $G = (V, E, w)$ such that vertices are divided in two disjoint sets $V = V_1 \cup V_2$ in the following way:*

- V_1 contains $n + p$ nodes composed by p source nodes s_1, s_2, \dots, s_p (referred to collectively as \mathcal{S}) followed by the n characters a_1, a_2, \dots, a_n (referred to collectively as \mathcal{T}_1) of A . The nodes in V_1 follow a strict ordering $s_1 \prec s_2 \prec \dots \prec s_p \prec a_1 \prec a_2 \prec \dots \prec a_n$.
- V_2 contains $n + p$ nodes composed by the n characters a_1, a_2, \dots, a_n (referred to collectively as \mathcal{T}_2) of A followed by p destination nodes d_1, d_2, \dots, d_p (referred to collectively as \mathcal{D}). The nodes in V_2 follow a strict ordering $a_1 \prec a_2 \prec \dots \prec a_n \prec d_1 \prec d_2 \prec \dots \prec d_p$.

Then the edges of the graph G are constructed in the following way:

- **Source Edges:** For each source node $s \in \mathcal{S}$ and each node $v_j \in \mathcal{T}_2$, an edge (s, v_j) is created with weight $w(s, v_j) = 1$. These edges represent the start of a new subsequence.
- **Internal nodes Edges:** For each pair of indices i, j such that $1 \leq i < j \leq n$, an edge is created between $u_i \in \mathcal{T}_1$ and $v_j \in \mathcal{T}_2$. The weight of this edge, $w(u_i, v_j)$, is 0 if the characters $A[i]$ and $A[j]$ are the same, and 1 otherwise. Formally:

$$w(u_i, v_j) = \begin{cases} 0 & \text{if } A[i] = A[j] \\ 1 & \text{if } A[i] \neq A[j] \end{cases}$$

This weight represents the cost of adding character $A[j]$ after $A[i]$ in a subsequence, where a cost of 1 is incurred if a new run is started.

- **Destination Edges:** For each node $u_i \in \mathcal{T}_1$ and each destination node $d \in \mathcal{D}$, an edge (u_i, d) is created with weight $w(u_i, d) = 1$. These edges represent the end of a subsequence.

Example 3.3: Vertices

Let's apply the reduction to the string $S = \text{ABCDDCBDDDD}$ from our running example (Example 3.1), with a target of $p = 2$ subsequences. Following the construction rules, we build a bipartite graph. The vertices of this graph are structured as shown in Figure 3.2.

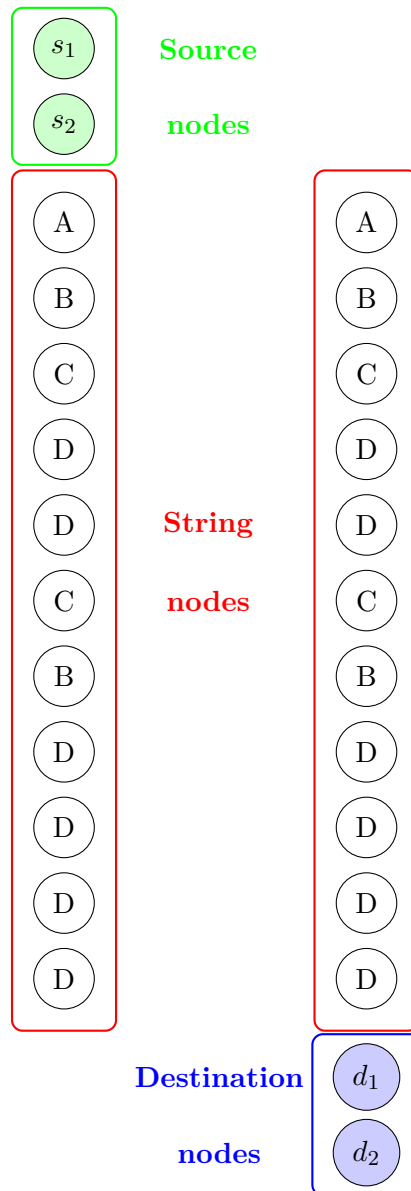


Figure 3.2: Bipartite graph vertices structure for string $S = \text{ABCDDCBDDDD}$ with $p = 2$. The nodes are ordered from top to bottom.

Example 3.4: Edges

Let's see a small example for each case of Definition 27. Consider $p = 2$. In Figure 3.3-(a), there is an example for the sources' edges. As stated before, each source is connected with weight 1 to all nodes in \mathcal{T}_2 .

In Figure 3.3-(b), we illustrate the edges from \mathcal{T}_1 to \mathcal{T}_2 . These edges model the cost of appending a character to a subsequence. An edge from u_i to v_j (for $j > i$) has weight 1 if $S[i] \neq S[j]$ (starting a new run) and weight 0 if $S[i] = S[j]$ (extending an existing run). For instance, the node for the first 'A' connects to the nodes for 'B' and 'C' with weight 1, but connects to the node for 'A' with weight 0.

Lastly, Figure 3.3-(c) shows the destination edges. These edges terminate a subsequence. An edge from any node $u_i \in \mathcal{T}_1$ to any destination node $d_k \in \mathcal{D}$ has a weight of 0, ensuring that ending a chain does not increase the run count. For example, if the node for 'B' is the last element of a subsequence, it is matched with a destination node, and this edge (u_B, d_k) contributes 0 to the total weight.

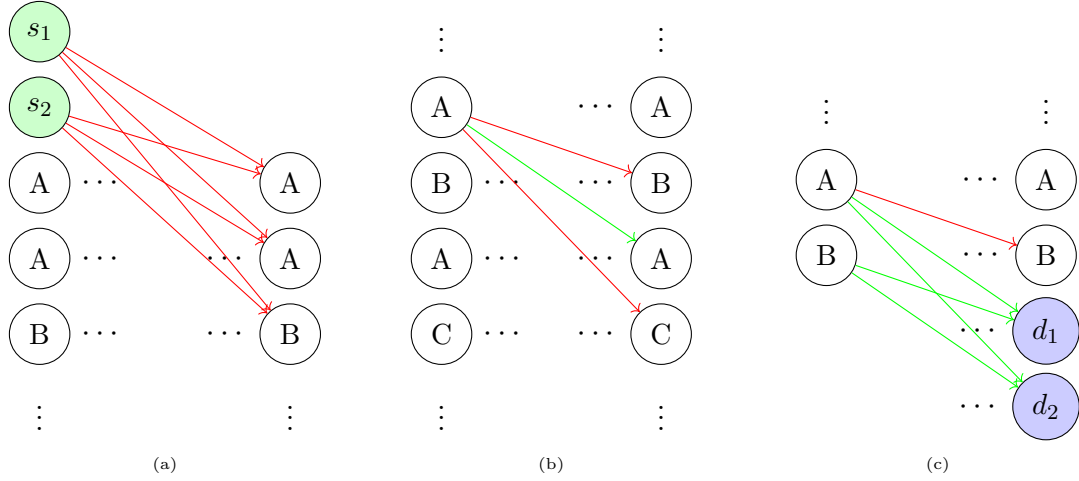


Figure 3.3: Examples of the connection construction in the bipartite graph for $p = 2$, showing the cases for source nodes \mathcal{S} (a), internal nodes \mathcal{T}_1 and \mathcal{T}_2 (b), and destination nodes \mathcal{D} (c). Red arrows indicate edges with weight 1, while green arrows indicate edges with weight 0.

let now state the following theorem regarding the number of edges in the bipartite graph resulting from Definition 27. This theorem is essential for understanding the complexity of the final algorithm employed to solve the MWPBM problem and so, the String Partitioning problem.

Theorem 6 (Bipartite graph properties). *The bipartite graph G constructed as stated in Definition 27 has $2n + 2p$ nodes and $2np + \frac{n(n-1)}{2}$ edges.*

Proof. The total number of edges in the graph G is the sum of the edges from the three categories defined in the construction:

- **Source Edges:** There are p source nodes in \mathcal{S} and n tree nodes in \mathcal{T}_2 . Each source node connects to every node in \mathcal{T}_2 , resulting in $p \times n = np$ edges.
- **Destination Edges:** There are n tree nodes in \mathcal{T}_1 and p destination nodes in \mathcal{D} . Each node in \mathcal{T}_1 connects to every node in \mathcal{D} , resulting in $n \times p = np$ edges.

- **Internal nodes Edges:** For each node $u_i \in \mathcal{T}_1$, edges are created to all nodes $v_j \in \mathcal{T}_2$ where $j > i$. The number of such edges is $\sum_{i=1}^{n-1} (i) = \frac{n(n-1)}{2}$.

Summing these up, the total number of edges is $2np + \frac{n(n-1)}{2}$. \square

3.2.2 Proof

In this section, we will prove that a minimum-weight perfect matching in a bipartite graph constructed as stated in Definition 27 directly corresponds to an optimal solution for the String Partitioning problem.

Let's start by defining the concept of valid partition:

Definition 28 (Valid Partition). *A partition of a string S into p subsequences, denoted $\Pi = \{S_1, S_2, \dots, S_p\}$, is **valid** if it satisfies two conditions:*

1. **Character Conservation:** *The multiset of characters in S is the disjoint union of the multisets of characters in each subsequence S_i . This ensures every character from S is used exactly once.*
2. **Order Preservation:** *For any subsequence $S_k \in \Pi$, if S_k contains characters corresponding to $S[i]$ and $S[j]$ from the original string with $i < j$, then the character from $S[i]$ must appear before the character from $S[j]$ in S_k .*

We also define $\text{runs}(S)$ as the number of runs in the string S .

Lemma 2. *Let $\Sigma_S \subseteq \Sigma$ be the set of unique characters that appear in the string S . The cost of the optimal solution of an instance $\mathcal{I} = (S, p)$ of the String Partitioning Problem is always greater than or equal to $|\Sigma_S|$, for any p .*

Proof. Let $\Pi = \{S_1, S_2, \dots, S_p\}$ be any valid partition of the string S into p subsequences. The cost of this partition is the total number of runs, given by $\sum_{i=1}^p \text{runs}(S_i)$.

For any character $c \in \Sigma_S$, there must be at least one occurrence of c in S . This occurrence must be placed in at least one of the subsequences, say S_j . The first occurrence of character c in any subsequence S_j will start a new run. For each distinct character in Σ_S , its occurrences must form at least one run in the partition. If a character appears in multiple subsequences, it will contribute at least one run for each subsequence it is in. If it appears only in one subsequence, it will contribute at least one run to that subsequence.

Therefore, every unique character $c \in \Sigma_S$ must contribute at least one run to the total cost. Summing over all unique characters, the total number of runs must be at least $|\Sigma_S|$. \square

Claim 1. *For the String Partitioning Problem, the optimal cost for an instance $\mathcal{I} = (S, p)$ where $p > |\Sigma_S|$ is not lower than the optimal cost for the instance $\mathcal{M} = (S, |\Sigma_S|)$, where Σ_S is the set of unique characters in S .*

Proof. The proof builds upon Lemma 2. If we use a number of chains $p > |\Sigma_S|$, we would have at least $p - |\Sigma_S|$ empty chains, since there are only $|\Sigma_S|$ non-empty equivalence classes of nodes to partition. As the minimum cost for any chain is 1, these empty chains contribute to the total cost. An optimal arrangement would

involve $|\Sigma_S|$ chains, each containing nodes from a single equivalence class, costing $|\Sigma_S|$, and $p - |\Sigma_S|$ empty chains, each costing 1. The total cost would be $|\Sigma_S| + (p - |\Sigma_S|) = p$. Since $p > |\Sigma_S|$, this cost is greater than the optimal cost of $|\Sigma_S|$ achievable with $p = |\Sigma_S|$ chains. Therefore, any solution with $p > |\Sigma_S|$ is suboptimal. \square

Therefore, we will only consider instances of the problem where $p < |\Sigma_S|$, as they do not present a trivial solution.

Let $r : \mathcal{I}_{StringPartitioning} \rightarrow \mathcal{M}_{MWPBM}$ be the reduction function that maps an instance $\mathcal{I} = (S, p)$ of the String Partitioning Problem to an instance $\mathcal{M} = (G)$ of the MWPBM Problem, where G is the bipartite graph constructed as stated in Definition 27.

Claim 2. *Let $\mathcal{I} = (S, p)$ be an instance of the String Partitioning Problem. In any perfect matching on $r(\mathcal{I})$, a node $u_i \in \mathcal{T}_1$ can only be matched with a node $v_j \in \mathcal{T}_2$ if $j > i$.*

Proof. This property is a direct consequence of the definition of internal edges in the graph construction. Edges between the sets \mathcal{T}_1 and \mathcal{T}_2 are explicitly defined only for pairs of nodes (u_i, v_j) where the index j is strictly greater than i . No edges are created for any case where $j \leq i$.

Since a matching is a subset of the graph's edges, it is impossible for a node u_i to be matched with a node v_j with an equal or smaller index. This structural constraint is fundamental to ensuring that the subsequences in the partition respect the original order of characters in the string S leading to a valid partition of the string. \square

Lemma 3. *Let $\mathcal{I} = (S, p)$ be an instance of the String Partitioning Problem. There always exists a perfect matching on $r(\mathcal{I})$.*

Proof. Let G be the bipartite graph of the MWPBM instance $r(\mathcal{I})$. The proof comes from the construction of the bipartite graph G and from Theorem 5. We are going to prove that G satisfies Hall's condition (see Theorem 5) and so, since by construction $|V_1| = |V_2|$, a perfect matching for G exists.

To verify Hall's condition, we need to prove that for any subset $W \subseteq V_1$ we have that $|N(W)| \geq |W|$, where $N(W)$ is the neighborhood of W (Definition 26). We have the following cases:

- If $W \subseteq \mathcal{S}$: Let $|W| = k$, where $1 \leq k \leq p$. By construction, every source node in \mathcal{S} is connected to every node in \mathcal{T}_2 . Therefore, the neighborhood of any non-empty subset $W \subseteq \mathcal{S}$ is the entire set \mathcal{T}_2 , so $N(W) = \mathcal{T}_2$. The size of the neighborhood is $|N(W)| = |\mathcal{T}_2| = n$. We need to show that $n \geq k$. From the problem definition, we can assume $p \leq n$, since having more partitions than characters offers no advantage (Claim 1). As W is a subset of \mathcal{S} , we have $k \leq p$. Therefore, $k \leq p \leq n$, which confirms that $|N(W)| \geq |W|$.
- If $W \subseteq \mathcal{T}_1$, Let $|W| = k$. Let the indices of the nodes in W be $\{i_1, i_2, \dots, i_k\}$, sorted such that $i_1 < i_2 < \dots < i_k$.

By the construction of the graph, every node $u_i \in \mathcal{T}_1$ is connected to every destination node in \mathcal{D} . Therefore, the entire set \mathcal{D} is part of the neighborhood of W , so $\mathcal{D} \subseteq N(W)$. This contributes p nodes to the neighborhood.

Additionally, each node $u_{i_r} \in W$ is connected to all nodes $v_j \in \mathcal{T}_2$ where $j > i_r$. The union of these neighbors in \mathcal{T}_2 is $\bigcup_{r=1}^k \{v_j \mid j > i_r\} = \{v_j \mid j > i_1\}$, since for any $r > 1$, the set $\{v_j \mid j > i_r\}$ is a subset of $\{v_j \mid j > i_1\}$. The size of this set of neighbors in \mathcal{T}_2 is $n - i_1$.

Combining these, the total size of the neighborhood is $|N(W)| = |\{v_j \mid j > i_1\} \cup \mathcal{D}| = (n - i_1) + p$.

We need to prove that $|N(W)| \geq |W|$, which means we must show $n - i_1 + p \geq k$.

Since the k indices in W are distinct and sorted, we know that $i_k \geq i_1 + k - 1$. As $i_k \leq n$, it follows that $i_1 + k - 1 \leq n$, which can be rearranged to $k \leq n - i_1 + 1$.

Since $p \geq 2$ (as there is at least two partitions), we have $n - i_1 + p \geq n - i_1 + 2$. Combining these inequalities, we get $|N(W)| = n - i_1 + p \geq n - i_1 + 2 \geq k$.

Thus, $|N(W)| \geq |W|$ holds for this case.

- If $W = W_S \cup W_U$, where $W_S \subseteq \mathcal{S}, W_U \subseteq \mathcal{T}_1$. By construction:
 - Since W_S is a non-empty subset of \mathcal{S} , every node in W_S is connected to every node in \mathcal{T}_2 . Thus, the neighborhood of W_S is the entire set \mathcal{T}_2 , so $N(W_S) = \mathcal{T}_2$.
 - Since W_U is a non-empty subset of \mathcal{T}_1 , every node in W_U is connected to every node in \mathcal{D} . Thus, the neighborhood of W_U includes the entire set \mathcal{D} , so $N(W_U) \supseteq \mathcal{D}$.

Combining these, the neighborhood of the mixed set W contains both \mathcal{T}_2 and \mathcal{D} : $N(W) = N(W_S) \cup N(W_U) = \mathcal{T}_2 \cup \mathcal{D}$.

Since the entire right side of the bipartition is $V_2 = \mathcal{T}_2 \cup \mathcal{D}$, this means the neighborhood of W is the entire set V_2 . Thus, $|N(W)| \geq |W|$ holds for this case.

Since Hall's condition is satisfied for all possible non-empty subsets of V_1 , a perfect matching always exists in the graph G . \square

Before proving the main theorem, we need to define how to retrieve an optimal partitions for an instance of the String Partitioning Problem $\mathcal{I} = (S, p)$ from an optimal perfect matching in $r(\mathcal{I})$.

Definition 29 (Partition decomposition). *Given an instance of the String Partitioning Problem $\mathcal{I} = (S, p)$ and an optimal perfect matching M in the bipartite graph G of $r(\mathcal{I})$, a **partition decomposition** of M is a way to partition the nodes of G into p non-empty subsets, where each subset represents a path from a source to a destination, thus, leading to a partition of S . The partition decomposition proceed as follows:*

Since M is a perfect matching, each source node $s_k \in \mathcal{S}$ begins exactly one path. This path is a sequence of characters from S determined by following the edges of M . A path starting from s_k corresponds to a sequence of indices (j_1, j_2, \dots, j_m) such that:

1. The first edge is $(s_k, v_{j_1}) \in M$.
2. Then, $(u_{j_i}, v_{j_{i+1}}) \in M \forall i = 1, \dots, m-1$.
3. The final edge in the path is $(u_{j_m}, d_q) \in M$ for some destination node d_q .

The subsequence S_k is then constructed from the characters at these indices: $S_k = S[j_1]S[j_2] \dots S[j_m]$. This process uniquely defines a valid partition of S into p subsequences.

Example 3.5:

Consider the example in Figure 3.4, which shows a perfect matching for the instance $r(\mathcal{I})$ with $\mathcal{I} = (S = \text{AAB}, p = 2)$. The solid arrows represent the edges of the matching M . The dashed arrows are a visual aid showing the correspondence between a character's representation in \mathcal{T}_2 (on the right) and its representation in \mathcal{T}_1 (on the left), which is essential for tracing the paths.

Applying the partition decomposition procedure (Definition 29), we extract two subsequences:

- **Subsequence 1 (red):**

1. Start at source s_1 . The matching edge is (s_1, v_1) , where v_1 corresponds to the first character, $S[1] = \text{'A'}$. The subsequence is now "A".
2. Following the conceptual link to u_1 , we find the matching edge (u_1, v_2) , where v_2 corresponds to the second character, $S[2] = \text{'A'}$. The subsequence is now "AA".
3. Following the link to u_2 , we find the matching edge (u_2, d_1) . Since d_1 is a destination node, the path terminates.
4. The final subsequence is $S_1 = \text{"AA"}$.

- **Subsequence 2 (blue):**

1. Start at source s_2 . The matching edge is (s_2, v_3) , where v_3 corresponds to the third character, $S[3] = \text{'B'}$. The subsequence is "B".
2. Following the link to u_3 , we find the matching edge (u_3, d_2) . Since d_2 is a destination, the path terminates.
3. The final subsequence is $S_2 = \text{"B"}$.

This example illustrates how the partition decomposition procedure correctly reconstructs the subsequences from the perfect matching, yielding the partition $\Pi = \{\text{"AA"}, \text{"B"}\}$.

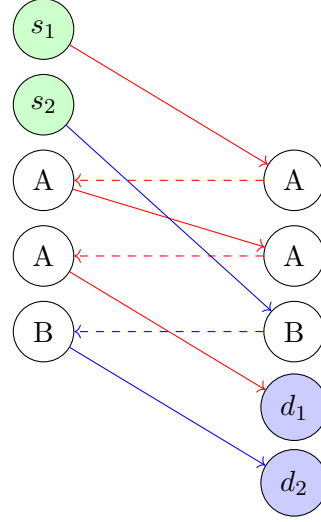


Figure 3.4: An example of a perfect matching (solid lines) in the constructed bipartite graph. The matching defines a partition into two paths (red and blue), which are traced by following the solid and dashed arrows.

Now we introduce the main theorem for the reduction.

Theorem 7. *Let $\mathcal{I} = (S, p)$ be an instance of the String Partition problem with $S = a_1 a_2 \dots a_n$. An optimal solution of $r(\mathcal{I})$ can be used to compute an optimal solution of \mathcal{I} in time $O(2np + \frac{n(n-1)}{2})$.*

Proof. Let $\mathcal{I} = (S, p)$ be an instance of the String Partitioning problem. We assume $p \leq |\Sigma_S|$ where Σ_S is the set of unique characters in S , per Claim 1. Let G be the bipartite graph of instance $r(\mathcal{I})$. We will demonstrate a bijection between the set of valid partitions of S (see Definition 28) and the set of perfect matchings in G , such that the cost of a partition equals the weight of its corresponding perfect matching.

First, we establish the existence of a perfect matching. By construction, the graph G is bipartite with partitions V_1 and V_2 such that $|V_1| = |V_2| = t + p$. Lemma 3 ensures that a perfect matching exists in G .

Let $\Pi = \{S_1, \dots, S_p\}$ be a valid partition of S . We construct a perfect matching M_Π by converting each subsequence $S_k \in \Pi$ into a path in G . Let the ordered indices of S_k be $i_{k,1} < i_{k,2} < \dots < i_{k,l_k}$. The corresponding path in G is formed by the following edges:

- A **start edge** $(s_k, v_{i_{k,1}})$. This edge exists because all source nodes are connected to all nodes in \mathcal{T}_2 .
- A set of **middle edges** $(u_{i_{k,j}}, v_{i_{k,j+1}})$ for $j = 1, \dots, l_k - 1$. These edges exist because the order preservation property of a valid partition ensures $i_{k,j} < i_{k,j+1}$, and our graph construction connects $u_i \in \mathcal{T}_1$ to $v_j \in \mathcal{T}_2$ for all $i < j$.
- An **end edge** $(u_{i_{k,l_k}}, d_k)$. This edge exists because all nodes in \mathcal{T}_1 are connected to all destination nodes.

Since every character of S belongs to exactly one subsequence, this process uses every node in \mathcal{T}_1 and \mathcal{T}_2 exactly once. By assigning each of the p paths to a unique source s_k and a unique destination d_k , we ensure that all source and destination nodes are

also used exactly once. Therefore, the resulting set of edges M_Π constitutes a perfect matching. The weight of a matching is given by:

$$\begin{aligned} W(M) &= \sum_{(u,v) \in M} w(u,v) \\ &= p + |\{(u_i, u_j) \in M \mid u_i \in \mathcal{T}_1, u_j \in \mathcal{T}_2, S[i] \neq S[j]\}| \end{aligned}$$

where p represents the contribution from the source nodes s_i , as each source node must be connected with weight 1 to start a chain. A class change occurs exactly when a path in the matching uses a weight-1 edge between tree nodes. Therefore, $W(M_\Pi)$ is exactly equal to the total number of runs of the partition defined by the matching M_Π .

Conversely, let M be a perfect matching in G . The structure of G ensures that M consists of p disjoint paths starting from source nodes $\{s_1, \dots, s_p\}$ and ending at destination nodes $\{d_1, \dots, d_p\}$. From Definition 29 we know that each path defines a valid subsequence S_k in S . By Claim 2, the node order within these chains is consistent with characters order in S . Thus, M maps to a valid partition of S . The cost of this partition is equal to $W(M)$.

The construction of the graph G takes time proportional to its number of edges. As established in Theorem 6, G has $O(2np + \frac{n(n-1)}{2})$ edges. Therefore, finding a minimum weight perfect matching in G allows us to solve the String Partitioning Problem for \mathcal{I} in time polynomial in the input size. \square

Example 3.6:

Let's apply the reduction to the string $S = \text{ABCDDCBDDDD}$ from our running example (Example 3.1), with a target of $p = 2$ subsequences. In Figure 3.5 we have one of the possible minimum perfect matchings for the instance having weight 5.

By applying the partition decomposition procedure (Definition 29) to this matching, we can trace the two paths from the source nodes to the destination nodes to obtain the following optimal partition:

- **Subsequence 1:** The path starting from s_1 traces the characters corresponding to indices (1, 3, 6, 7), yielding the subsequence $S_1 = \text{"ACCB"}$. The number of runs is $\text{runs}(S_1) = 3$.
- **Subsequence 2:** The path starting from s_2 traces the characters for indices (2, 4, 5, 8, 9, 10, 11), yielding the subsequence $S_2 = \text{"BDDDDDD"}$. The number of runs is $\text{runs}(S_2) = 2$.

The total cost of this partition is the sum of the runs, $3 + 2 = 5$, which matches the weight of the perfect matching. This demonstrates how the reduction finds an optimal solution for the String Partitioning Problem.

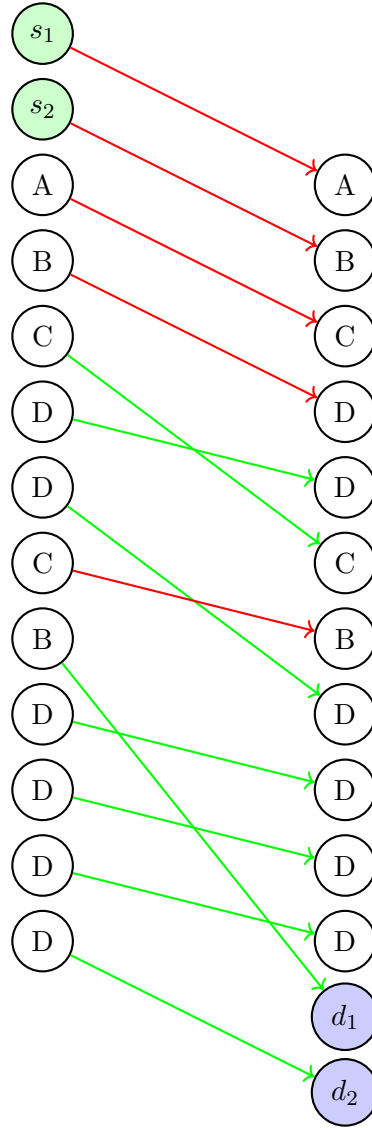


Figure 3.5: Example of an optimal perfect matching for the graph G of the MWPBM instance $r((S = ABCDDCBDDDD, p = 2))$. Green edges weigh 0, while red edges weigh 1.

3.3 Collapsing Nodes in Chains

The next crucial step of our compression scheme is to reduce the space required by each chain by collapsing equivalent nodes. Specifically, any sequence of consecutive nodes within the same chain that belong to the same equivalence class is merged into a single representative node. This new node preserves the connectivity of the original structure by inheriting all distinct outgoing and ingoing edges from the nodes it replaces. In the following subsections we will introduce the concept of non-deterministic finite automaton and define how nodes are collapsed in a chain. Then, we prove that this transformation is language-preserving (see Lemma 6).

3.3.1 How to Collapse Nodes

Now, we define the concept of collapsing consecutive equivalent nodes in a chain.

Definition 30 (Collapsing consecutive equivalent nodes). *Let V be the set of nodes*

of the tree, let Σ be the alphabet, and let $E \subseteq V \times \Sigma \times V$ be the set of labeled edges. Let $\mathcal{P} = \{C_1, C_2, \dots, C_m\}$ be the set of all chains partitioning V .

For each chain $C_i = (u_1^{(i)}, u_2^{(i)}, \dots, u_{n_i}^{(i)}) \in \mathcal{P}$, partition it into maximal consecutive blocks $B_1^{(i)}, \dots, B_{k_i}^{(i)}$, where each block $B_t^{(i)} = (u_j^{(i)}, \dots, u_\ell^{(i)})$ satisfies $\mathcal{C}(u_r^{(i)}) = \mathcal{C}(u_s^{(i)})$ for all $r, s \in \{j, \dots, \ell\}$, and the block is maximal (cannot be extended).

The collapsed chain is $C'_i = (v_1^{(i)}, \dots, v_{k_i}^{(i)})$, where each block $B_t^{(i)}$ is replaced by a single node $v_t^{(i)}$.

Define the global collapse map $\Phi : V \rightarrow V'$ where $V' = \bigcup_{i=1}^m C'_i$, such that $\Phi(u) = v_t^{(i)}$ if $u \in B_t^{(i)}$ for some chain C_i and block $B_t^{(i)}$.

Then:

- Two consecutive nodes $u_j^{(i)}, u_{j+1}^{(i)} \in C_i$ are collapsed into the same node if and only if $\mathcal{C}(u_j^{(i)}) = \mathcal{C}(u_{j+1}^{(i)})$.
- The edge set after collapsing all chains is

$$E' = \{(\Phi(x), a, \Phi(y)) : (x, a, y) \in E\},$$

where parallel duplicates are removed (i.e., E' is treated as a set).

Example 3.7:

Consider Example 3.6 where we obtained the chains $C_1 = \{A, C, C, B\}$ and $C_2 = \{B, D, D, D, D, D, D\}$ for the tree ADFA in Figure 2.2 by applying the reduction from CHAIN-DIVISION to MWPBM. The nodes inside each chain are the following:

- $C_1 = \{a, d, f, c\}$
- $C_2 = \{b, h, l, e, i, m, g\}$

Applying the collapsing operation from Definition 30:

- For $C_1 = (a, d, f, c)$ with classes (A, C, C, B) :
 - Block $B_1 = \{a\}$ (class A) \rightarrow collapsed node v_1 . The node a is the initial state. It has two outgoing edges: $a \xrightarrow{0} b$ and $a \xrightarrow{1} c$. Since b and c collapse to w_1 and v_3 respectively, we obtain $v_1 \xrightarrow{0} w_1$ and $v_1 \xrightarrow{1} v_3$.
 - Block $B_2 = \{d, f\}$ (class C) \rightarrow collapsed node v_2 . The outgoing edges of d and f are:

$$d \xrightarrow{0} h, d \xrightarrow{1} i, \quad f \xrightarrow{0} l, f \xrightarrow{1} m.$$

After collapsing, we obtain:

$$v_2 \xrightarrow{0} w_2, v_2 \xrightarrow{1} w_2, \quad v_2 \xrightarrow{0} w_2, v_2 \xrightarrow{1} w_2.$$

Since we have two identical edges we can keep only one of each.

- Block $B_3 = \{c\}$ (class B) \rightarrow collapsed node v_3 . It has two outgoing edges: $c \xrightarrow{0} f$ and $c \xrightarrow{1} g$. Since f and g collapse to v_2 and w_2 respectively, we obtain $v_3 \xrightarrow{0} v_2$ and $v_3 \xrightarrow{1} w_2$.

Result: $C'_1 = (v_1, v_2, v_3)$ with classes (A, C, B) . Here, v_1 is the initial state.

- For $C_2 = (b, h, l, e, i, m, g)$ with classes (B, D, D, D, D, D, D) :
 - Block $B_1 = \{b\}$ (class B) \rightarrow collapsed node w_1 . It has two outgoing edges: $b \xrightarrow{0} d$ and $b \xrightarrow{1} e$. As d, e collapse to v_2 and w_2 respectively, we obtain $w_1 \xrightarrow{0} v_2$ and $w_1 \xrightarrow{1} w_2$.
 - Block $B_2 = \{h, l, e, i, m, g\}$ (all class D) \rightarrow collapsed node w_2 . The node w_2 collects all incoming edges formerly targeting any of b, h, l, e, i, m, g , and it is accepting.

Result: $C'_2 = (w_1, w_2)$ with classes (B, D) , and w_2 is the unique accepting state for this example.

The collapsed chains preserve all distinct outgoing and incoming edges through the collapse map Φ , significantly reducing the space complexity from 11 nodes to 5 nodes total. The resulting NFA is shown in Figure 3.6.

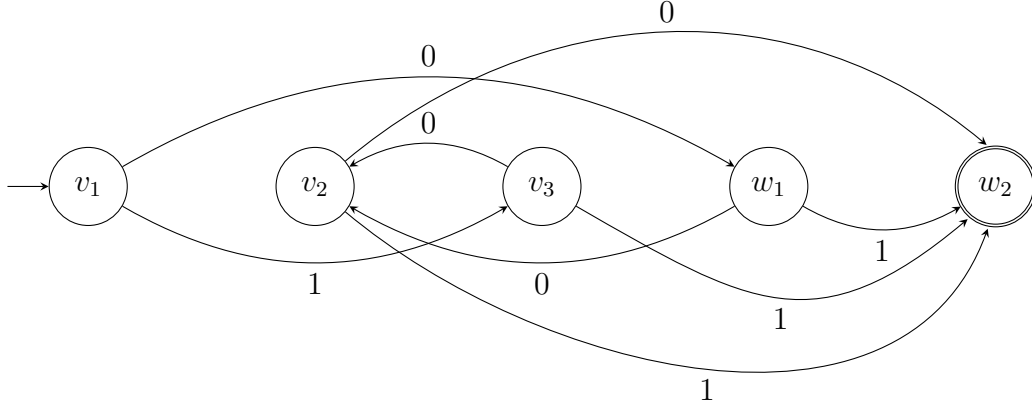


Figure 3.6: NFA obtained after collapsing equivalent nodes in chains C_1 and C_2 .

3.3.2 Language Equivalence

Now, we need to prove that the language recognized by the NFA obtained after collapsing equivalent nodes in chains following Definition 30 is equivalent to the language of the original tree ADFA. We define the language of an NFA or ADFA F as $L(F)$.

Lemma 4 (Root forms a singleton class under Revuz). *In a tree ADFA minimized by Revuz's algorithm (Subsection 2.4.3), the root forms a singleton equivalence class.*

Proof. Let T be a tree ADFA, let r be its root, and let $D = h(r)$ be its height as in Definition 12. For every other node $u \neq r$ at depth $k \geq 1$, we have $h(u) \leq D - k < D$. Hence r is the unique state in the level Π_D . Since Revuz's algorithm partitions by height and refines within each level, r cannot be merged with any other state and thus forms its own equivalence class. \square

Lemma 5 (All leaves fall into the same class under Revuz). *In a tree ADFA minimized by Revuz's algorithm, all leaves belong to the same equivalence class.*

Proof. Every leaf ℓ is final and has no outgoing transitions, so $h(\ell) = 0$. In the base step on level Π_0 , the label of a state depends only on its finality and on transitions to already distinguished classes. All leaves share the same label (final, no outgoing transitions), hence they are merged into a single equivalence class in Π_0 . \square

Lemma 6. *Let T be a tree ADFA with alphabet Σ . Let C_1, C_2, \dots, C_m be the chains partitioning T as defined in ???. Let N be the NFA obtained after collapsing equivalent nodes in chains C_1, C_2, \dots, C_m as defined in Definition 30. Then, $L(N) = L(T)$.*

Proof. We prove the equality $L(N) = L(T)$ by showing both inclusions $L(T) \subseteq L(N)$ and $L(N) \subseteq L(T)$.

$(L(T) \subseteq L(N))$: Let $w \in L(T)$. Then there exists an accepting path in T from the root to some accepting state:

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_n} q_n$$

where q_0 is the root, q_n is accepting, and $w = a_1 a_2 \dots a_n$.

By Definition 30, each node q_i in the original tree is mapped to a collapsed node $\Phi(q_i)$ in N . Since the collapse map preserves all edges (by the definition of E'), there exists a corresponding path in N :

$$\Phi(q_0) \xrightarrow{a_1} \Phi(q_1) \xrightarrow{a_2} \Phi(q_2) \xrightarrow{a_3} \dots \xrightarrow{a_n} \Phi(q_n)$$

Since q_0 is the root of T , by Lemma 4 its image $\Phi(q_0)$ is the unique initial state of N . Since q_n is accepting in T , by Lemma 5 its image $\Phi(q_n)$ lies in the unique leaf class, which is marked accepting in N . Hence, $w \in L(N)$.

$(L(N) \subseteq L(T))$: Let $w \in L(N)$. Then there exists an accepting path in N :

$$v_0 \xrightarrow{a_1} v_1 \xrightarrow{a_2} v_2 \xrightarrow{a_3} \dots \xrightarrow{a_n} v_n$$

where v_0 is the initial state, v_n is accepting, and $w = a_1 a_2 \dots a_n$.

Each collapsed node v_i corresponds to some block $B_t^{(j)}$ in the original tree. By Definition 30, every edge (v_{i-1}, a_i, v_i) in N corresponds to at least one edge (u, a_i, u') in the original tree T , where $\Phi(u) = v_{i-1}$ and $\Phi(u') = v_i$.

Since the chains preserve the ordering from the original tree structure, we can construct a valid path in T by selecting appropriate representatives from each collapsed block. Specifically, we can choose nodes u_0, u_1, \dots, u_n in T such that $\Phi(u_i) = v_i$ and $(u_{i-1}, a_i, u_i) \in E$ for all $i = 1, \dots, n$.

Since v_0 corresponds to the root block and v_n is accepting, for Lemmas 4 and 5 we have u_0 as the root and u_n as an accepting state in T respectively. Therefore, $w \in L(T)$.

Thus, $L(N) = L(T)$. \square

Collapsing the nodes as in Definition 30 preserves the language of the original tree ADFA and the resulting chains inherit a total order. This enables the application of the NFA indexing scheme of Cotumaccio et al. [7], which we present in the next chapter.

Theorem 8 (Myhill–Nerode). *Let $L \subseteq \Sigma^*$ and define the indistinguishability relation*

$$x \sim_L y \iff \forall z \in \Sigma^*, xz \in L \iff yz \in L.$$

Then:

1. \sim_L is an equivalence relation on Σ^* .
2. L is regular iff \sim_L has finitely many equivalence classes.
3. The states of the minimal DFA for L are in one-to-one correspondence with the equivalence classes of \sim_L .
4. L is exactly the union of those equivalence classes that intersect L :

$$L = \bigcup \{[x]_{\sim_L} \mid x \in L\}.$$

Lemma 7. *Let $M = (Q, \Sigma, \delta, q_0, F)$ be a DFA recognizing $L \subseteq \Sigma^*$. If two states $p, q \in Q$ correspond to the same Myhill–Nerode class for L (i.e., for all $w \in \Sigma^*$ we have $\delta(p, w) \in F \iff \delta(q, w) \in F$), then merging p and q into a single state yields an automaton (possibly nondeterministic) that still recognizes exactly L .*

Proof. By the Myhill–Nerode theorem, every state of M corresponds to a unique equivalence class of \sim_L , and L is exactly the union of those classes that intersect L . If p and q belong to the same equivalence class, then for every continuation $z \in \Sigma^*$ we have

$$\delta(p, z) \in F \iff \delta(q, z) \in F.$$

Thus replacing p with q (or vice versa) in any path does not affect whether the run ends in an accepting state. Therefore merging p and q does not alter the set of accepted strings, i.e. the recognized language remains L . \square

Chapter 4

Experiments

Bibliography

- [1] Jarno Alanko et al. “Wheeler languages”. In: *Information and Computation* 281 (2021), p. 104820 (cit. on p. 40).
- [2] Jean Berstel et al. “Minimization of automata”. In: *arXiv preprint arXiv:1010.5318* (2010) (cit. on p. 34).
- [3] Philip Bille et al. “Tree Compression with Top Trees”. In: *Information and Computation* 243 (Aug. 2015), pp. 166–177. ISSN: 08905401. DOI: 10.1016/j.ic.2014.12.012. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0890540114001643> (visited on 08/18/2025) (cit. on pp. 10, 11).
- [4] M. Burrows and D. Wheeler. *A block-sorting lossless data compression algorithm*. Tech. rep. 124. Digital Equipment Corporation, 1994 (cit. on pp. 1, 13, 21, 41).
- [5] Li Chen et al. “Maximum flow and minimum-cost flow in almost-linear time”. In: *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE. 2022, pp. 612–623 (cit. on p. 45).
- [6] Nicola Cotumaccio and Nicola Prezza. “On Indexing and Compressing Finite Automata”. In: *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*. Society for Industrial and Applied Mathematics, Jan. 2021, pp. 2585–2599. ISBN: 9781611976465. DOI: 10.1137/1.9781611976465.153 (cit. on pp. 38, 41).
- [7] Nicola Cotumaccio et al. “Co-lexicographically ordering automata and regular languages-Part I”. In: *Journal of the ACM* 70.4 (2023), pp. 1–73 (cit. on pp. 2, 3, 38, 39, 61).
- [8] Robert P Dilworth. “A decomposition theorem for partially ordered sets”. In: *The Dilworth Theorems: Selected Papers of Robert P. Dilworth*. Springer, 1990, pp. 7–12 (cit. on p. 37).
- [9] EA Dinic and MA Kronrod. “An algorithm for the solution of the assignment problem”. In: *Soviet Math. Dokl.* Vol. 10. 6. 1969, pp. 1324–1326 (cit. on p. 45).
- [10] Jack Edmonds and Richard M Karp. “Theoretical improvements in algorithmic efficiency for network flow problems”. In: *Journal of the ACM (JACM)* 19.2 (1972), pp. 248–264 (cit. on p. 45).
- [11] Massimo Equi, Veli Mäkinen, and Alexandru I. Tomescu. “Graphs Cannot Be Indexed in Polynomial Time for Sub-Quadratic Time String Matching, Unless SETH Fails”. In: *Theoretical Computer Science* 975 (Oct. 2023), p. 114128. ISSN: 03043975. DOI: 10.1016/j.tcs.2023.114128. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0304397523004413> (visited on 08/22/2025) (cit. on pp. 1, 39).
- [12] Paolo Ferragina et al. “Compressing and Indexing Labeled Trees, with Applications”. In: *Journal of the ACM* 57 (2009). DOI: 10.1145/1613676.1613680 (cit. on pp. 1, 10, 12, 16, 21).
- [13] Harold N Gabow and Robert E Tarjan. “Faster scaling algorithms for network problems”. In: *SIAM Journal on Computing* 18.5 (1989), pp. 1013–1036 (cit. on p. 45).

- [14] Travis Gagie, Giovanni Manzini, and Jouni Sirén. “Wheeler graphs: A framework for BWT-based data structures”. In: *Theoretical computer science* 698 (2017), pp. 67–78 (cit. on pp. 2, 37, 38).
- [15] Pawel Gawrychowski and Artur Jez. “LZ77 Factorisation of Trees”. In: *LIPICs, Volume 65, FSTTCS 2016* 65 (2016). Ed. by Akash Lal et al., 35:1–35:15. ISSN: 1868-8969. DOI: 10.4230/LIPICs.FSTTCS.2016.35. URL: <https://drops.dagstuhl.de/entities/document/10.4230/LIPICs.FSTTCS.2016.35> (visited on 08/18/2025) (cit. on p. 11).
- [16] Michel Goemans. *Advanced Algorithms: Matching Notes*. Lecture notes, Massachusetts Institute of Technology. 2009. URL: <https://math.mit.edu/~goemans/18433S09/matching-notes.pdf> (cit. on p. 43).
- [17] Simon Gog et al. “From Theory to Practice: Plug and Play with Succinct Data Structures”. In: *13th International Symposium on Experimental Algorithms, (SEA 2014)*. 2014, pp. 326–337 (cit. on p. 27).
- [18] Roberto Grossi, Ankur Gupta, and Jeffrey Scott Vitter. “High-order entropy-compressed text indexes”. In: *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. USA: Society for Industrial and Applied Mathematics, 2003, pp. 841–850 (cit. on p. 28).
- [19] Philip Hall. “On representatives of subsets”. In: *Journal of the London Mathematical Society* 10.1 (1935), pp. 26–30 (cit. on p. 44).
- [20] John Hopcroft. “AN $n \log n$ ALGORITHM FOR MINIMIZING STATES IN A FINITE AUTOMATON”. In: *Theory of Machines and Computations*. Ed. by Zvi Kohavi and Azaria Paz. Academic Press, 1971, pp. 189–196. ISBN: 978-0-12-417750-5. DOI: <https://doi.org/10.1016/B978-0-12-417750-5.50022-1>. URL: <https://www.sciencedirect.com/science/article/pii/B9780124177505500221> (cit. on p. 31).
- [21] Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. “Linear work suffix array construction”. In: *Journal of the ACM* 53.6 (2006), pp. 918–936. DOI: 10.1145/1217856.1217858 (cit. on pp. 16, 18).
- [22] Vladimir Kolmogorov. “Blossom V: A New Implementation of a Minimum Cost Perfect Matching Algorithm”. In: *Mathematical Programming Computation* 1.1 (July 2009), pp. 43–67. ISSN: 1867-2949, 1867-2957. DOI: 10.1007/s12532-009-0002-8. URL: <http://link.springer.com/10.1007/s12532-009-0002-8> (visited on 08/13/2025) (cit. on p. 45).
- [23] S. R. Kosaraju. “Efficient tree pattern matching”. In: *Proceedings of the 20th IEEE Foundations of Computer Science (FOCS)*. 1989, pp. 178–183 (cit. on p. 10).
- [24] Harold W Kuhn. “The Hungarian method for the assignment problem”. In: *Naval research logistics quarterly* 2.1-2 (1955), pp. 83–97 (cit. on p. 45).
- [25] N. Jesper Larsson and Alistair Moffat. “Off-line dictionary-based compression”. In: *Proceedings DCC 2000. Data Compression Conference*. IEEE. 2000, pp. 296–305 (cit. on p. 11).
- [26] Markus Lohrey, Sebastian Maneth, and Roy Mennicke. “Tree structure compression with repair”. In: *2011 Data Compression Conference*. IEEE. 2011, pp. 353–362 (cit. on p. 11).
- [27] Markus Lohrey, Sebastian Maneth, and Roy Mennicke. “XML tree structure compression using RePair”. In: *Information Systems* 38.8 (2013), pp. 1150–1167 (cit. on p. 11).

- [28] Giovanni Manzini et al. “The Rational Construction of a Wheeler DFA”. In: *LIPICs, Volume 296, CPM 2024* 296 (2024). Ed. by Shunsuke Inenaga and Simon J. Puglisi, 23:1–23:15. ISSN: 1868-8969. DOI: 10.4230/LIPICS.CPM.2024.23. URL: <https://drops.dagstuhl.de/entities/document/10.4230/LIPICS.CPM.2024.23> (visited on 08/23/2025) (cit. on pp. 2, 40, 41).
- [29] John Myhill. *Finite automata and the representation of events*. Tech. rep. WADC Technical Report 57-624. Wright Air Development Center, 1957 (cit. on pp. 31, 32).
- [30] Gonzalo Navarro. *Compact Data Structures: A Practical Approach*. Cambridge University Press, 2016 (cit. on p. 9).
- [31] Anil Nerode. “Linear automaton transformations”. In: *Proceedings of the American Mathematical Society* 9.4 (1958), pp. 541–544 (cit. on pp. 31, 32).
- [32] Rajeev Raman, V. Raman, and S. Srinivasa Rao. “Succinct Indexable Dictionaries with Applications to representations of k-ary trees and multi-sets”. In: *Proceedings of the 13th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 2002 (cit. on p. 27).
- [33] Dominique Revuz. “Minimisation of acyclic deterministic automata in linear time”. In: *Theoretical Computer Science* 92.1 (1992), pp. 181–189 (cit. on pp. 1, 32).
- [34] Piotr Sankowski. “Maximum weight bipartite matching in matrix multiplication time”. In: *Theoretical Computer Science* 410.44 (2009), pp. 4480–4488 (cit. on p. 45).

Web bibliography

- [16] Michel Goemans. *Advanced Algorithms: Matching Notes*. Lecture notes, Massachusetts Institute of Technology. 2009. URL: <https://math.mit.edu/~goemans/18433S09/matching-notes.pdf> (cit. on p. 43).