



**HAL**  
open science

## Optimal Parsing for dictionary text compression

Alessio Langiu

► **To cite this version:**

Alessio Langiu. Optimal Parsing for dictionary text compression. Other [cs.OH]. Université Paris-Est; Università degli studi (Palerme, Italie), 2012. English. NNT : 2012PEST1091 . tel-00804215

**HAL Id: tel-00804215**

**<https://theses.hal.science/tel-00804215v1>**

Submitted on 25 Mar 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

UNIVERSITÀ DEGLI STUDI DI PALERMO

*Dipartimento di Matematica e Applicazioni*

*Dottorato di Ricerca in Matematica e Informatica*

*XXIII<sup>e</sup> Ciclo - S.S.D. Inf/01*

---

UNIVERSITÉ PARIS-EST

*École doctorale MSTIC*

*Thèse de doctorat en Informatique*

***Optimal Parsing for Dictionary  
Text Compression***

*Author: LANGIU ALESSIO*

*Ph.D. Thesis in Computer Science*

*April 3, 2012*

PhD Commission:

---

Thesis Directors:	<i>Prof. CROCHEMORE MAXIME</i>	<i>University of Paris-Est</i>
	<i>Prof. RESTIVO ANTONIO</i>	<i>University of Palermo</i>
Examiners:	<i>Prof. ILIOPOULOS COSTAS</i>	<i>King's College London</i>
	<i>Prof. LECROQ THIERRY</i>	<i>University of Rouen</i>
	<i>Prof. MIGNOSI FILIPPO</i>	<i>University of L'Aquila</i>
Referees:	<i>Prof. GROSSI ROBERTO</i>	<i>University of Pisa</i>
	<i>Prof. ILIOPOULOS COSTAS</i>	<i>King's College London</i>
	<i>Prof. LECROQ THIERRY</i>	<i>University of Rouen</i>



# Summary

Dictionary-based compression algorithms include a parsing strategy to transform the input text into a sequence of dictionary phrases. Given a text, such process usually is not unique and, for compression purposes, it makes sense to find one of the possible parsing that minimize the final compression ratio. This is the parsing problem. An optimal parsing is a parsing strategy or a parsing algorithm that solve the parsing problem taking into account all the constraints of a compression algorithm or of a class of homogeneous compression algorithms. Compression algorithm constraints are, for instance, the dictionary itself, i.e. the dynamic set of available phrases, and how much a phrase weights on the compressed text, i.e. the number of bits of which the codeword representing such phrase is composed, also denoted as the encoding cost of a dictionary pointer.

In more than 30 years of history of dictionary-based text compression, despite plenty of algorithms, variants and extensions have appeared and while dictionary approach to text compression became one of the most appreciated and utilized in almost all the storage and communication processes, only few optimal parsing algorithms were presented. Many compression algorithms still lack optimality of their parsing or, at least, proof of optimality. This happens because there is not a general model of the parsing problem including all the dictionary-based algorithms and because the existing optimal parsing algorithms work under too restrictive hypotheses.

This work focuses on the parsing problem and presents both a general model for dictionary-based text compression called Dictionary-Symbolwise Text Compression and a general parsing algorithm that is proved to be optimal under some realistic hypotheses. This algorithm is called Dictionary-

Symbolwise Flexible Parsing and covers almost all of the known cases of dictionary-based text compression algorithms together with the large class of their variants where the text is decomposed in a sequence of symbols and dictionary phrases.

In this work we further consider the case of a free mixture of a dictionary compressor and a symbolwise compressor. Our Dictionary-Symbolwise Flexible Parsing covers also this case. We have indeed an optimal parsing algorithm in the case of dictionary-symbolwise compression where the dictionary is prefix closed and the cost of encoding dictionary pointer is *variable*. The symbolwise compressor is one the classic variable-length codes that works in linear time. Our algorithm works under the assumption that a special graph that will be described in the following, is well defined. Even if this condition is not satisfied, it is possible to use the same method to obtain almost optimal parses. In detail, when the dictionary is LZ78-like, we show how to implement our algorithm in linear time. When the dictionary is LZ77-like our algorithm can be implemented in time  $O(n \log n)$ . Both have  $O(n)$  space complexity.

Even if the main purpose of this work is of theoretical nature, some experimental results will be introduced to underline some practical effects of the parsing optimality in terms of compression performance and to show how to improve the compression ratio by building extensions Dictionary-Symbolwise of known algorithms. A specific appendix reports some more detailed experiments.

# Résumé

Les algorithmes de compression de données basés sur les dictionnaires incluent une stratégie de parsing pour transformer le texte d'entrée en une séquence de phrases du dictionnaire. Étant donné un texte, un tel processus n'est généralement pas unique et, pour comprimer, il est logique de trouver, parmi les parsing possibles, celui qui minimise le plus le taux de compression finale.

C'est ce qu'on appelle le problème du parsing. Un parsing optimal est une stratégie de parsing ou un algorithme de parsing qui résout ce problème en tenant compte de toutes les contraintes d'un algorithme de compression ou d'une classe d'algorithmes de compression homogène.

Les contraintes de l'algorithme de compression sont, par exemple, le dictionnaire lui-même, c'est-à-dire l'ensemble dynamique de phrases disponibles, et combien une phrase pèse sur le texte comprimé, c'est-à-dire quelle est la longueur du mot de code qui représente la phrase, appelée aussi le coût du codage d'un pointeur de dictionnaire.

En plus de 30 ans d'histoire de la compression de texte par dictionnaire, une grande quantité d'algorithmes, de variantes et d'extensions sont apparus. Cependant, alors qu'une telle approche de la compression du texte est devenue l'une des plus appréciées et utilisées dans presque tous les processus de stockage et de communication, seuls quelques algorithmes de parsing optimaux ont été présentés.

Beaucoup d'algorithmes de compression manquent encore d'optimalité pour leur parsing, ou du moins de la preuve de l'optimalité. Cela se produit parce qu'il n'y a pas un modèle général pour le problème de parsing qui inclut tous les algorithmes par dictionnaire et parce que les parsing optimaux

existants travaillent sous des hypothèses trop restrictives.

Ce travail focalise sur le problème de parsing et présente à la fois un modèle général pour la compression des textes basée sur les dictionnaires appelé la théorie Dictionary-Symbolwise et un algorithme général de parsing qui a été prouvé être optimal sous certaines hypothèses réalistes. Cet algorithme est appelé Dictionary-Symbolwise Flexible Parsing et couvre pratiquement tous les cas des algorithmes de compression de texte basés sur dictionnaire ainsi que la grande classe de leurs variantes où le texte est décomposé en une séquence de symboles et de phrases du dictionnaire.

Dans ce travail, nous avons aussi considéré le cas d'un mélange libre d'un compresseur par dictionnaire et d'un compresseur symbolwise. Notre Dictionary-Symbolwise Flexible Parsing couvre également ce cas-ci. Nous avons bien un algorithme de parsing optimal dans le cas de compression Dictionary-Symbolwise où le dictionnaire est fermé par préfixe et le coût d'encodage des pointeurs du dictionnaire est variable. Le compresseur symbolwise est un compresseur symbolwise classique qui fonctionne en temps linéaire, comme le sont de nombreux codeurs communs à longueur variable.

Notre algorithme fonctionne sous l'hypothèse qu'un graphe spécial, qui sera décrit par la suite, soit bien défini. Même si cette condition n'est pas remplie, il est possible d'utiliser la même méthode pour obtenir des parsing presque optimaux. Dans le détail, lorsque le dictionnaire est comme LZ78, nous montrons comment mettre en œuvre notre algorithme en temps linéaire. Lorsque le dictionnaire est comme LZ77 notre algorithme peut être mis en œuvre en temps  $O(n \log n)$  où  $n$  est la longueur du texte. Dans les deux cas, la complexité en espace est  $O(n)$ . Même si l'objectif principal de ce travail est de nature théorique, des résultats expérimentaux seront présentés pour souligner certains effets pratiques de l'optimalité du parsing sur les performances de compression et quelques résultats expérimentaux plus détaillés sont mis dans une annexe appropriée.

# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Background</b>	<b>7</b>
1.1 Self-Information and Entropy . . . . .	7
1.2 Entropy Encoding . . . . .	9
1.3 Encoding of Numbers and Commas . . . . .	11
1.4 Dictionary Methods . . . . .	11
<b>2 Dictionary-Symbolwise Text Compression</b>	<b>15</b>
2.1 Dictionary Compression . . . . .	15
2.2 Dictionary-Symbolwise Compression . . . . .	18
2.3 The Graph-Based Model . . . . .	23
2.4 On Parsing Optimality . . . . .	26
2.5 Dictionary-Symbolwise Can Be Better . . . . .	29
<b>3 About the Parsing Problem</b>	<b>37</b>
3.1 The Optimal Parsing Problem . . . . .	37
3.2 An Overview about Theory and Practice . . . . .	38
3.3 A Generalization of the Cohn's Theorem . . . . .	40
<b>4 Dictionary-Symbolwise Flexible Parsing</b>	<b>47</b>
4.1 The <i>c-supermaximal</i> Edges . . . . .	48
4.2 The Subgraph $G'_{\mathcal{A},T}$ . . . . .	55
4.3 The Dictionary-Symbolwise Flexible Parsing . . . . .	58
4.4 Time and Space Analyses . . . . .	61



<b>5</b>	<b>The Multilayer Suffix Tree</b>	<b>67</b>
5.1	Preliminary Definitions . . . . .	67
5.2	The Data Structure . . . . .	69
5.3	On the Query 1 . . . . .	73
5.4	On <i>Weighted-Exist</i> and <i>Weighted-Extend</i> . . . . .	75
<b>6</b>	<b>Conclusion</b>	<b>85</b>
<b>A</b>	<b>Experiments</b>	<b>93</b>
<b>B</b>	<b>Time Analysis of the Constant Case</b>	<b>97</b>

# Introduction

Data compression concerns with transformations through a more concise data representation. When such transformation is perfectly invertible we have a lossless data compression, otherwise, a lossy compression. Since data preservation is usually required for textual data, lossless data compression is often called text compression. On the opposite, usually working on visual data, such as the images or video, on sound data and on data from many other domains, a certain degree of approximation is allowed to the compression-decompression process in favour of a stronger compression, i.e. a smaller compression ratio.

Roughly speaking, compression ratios greater than a certain threshold given by the percentage of information contained in the data, are reachable by text compression techniques as they strip just redundancy in the text. Stronger compressions imply data approximation because part of their information is lost along the compression process. The quantity of information in a certain data or, more precisely, the average information inside the data provided by a source, is called entropy. The entropy ratio is then a limit for text compression, i.e. it is a lower bound for the compression ratio.

Entropy, data complexity and data compression are therefore bidden all together. Indeed, fundamental and seminal methods for dictionary-based compression, such as the Lempel' and Ziv's methods, were firstly introduced as text complexity measures.

Lempel' and Ziv's methods are still the basis of almost all the recent dictionary compression algorithms. More in detail, they are the LZ77 and the LZ78 compression methods, i.e. the Lempel and Ziv compression methods presented in 1977 and 1978 years. They are the first relevant dictionary

methods that use dynamic dictionaries. Static dictionary compression was already known as it is a side effect of the code and the transducer theories. Static dictionary compression was the topic of many works around '70, as the text substitution methods in Schuegraf' and Heaps's work (1974) or in the Wagner's work (1973).

Dictionary-based compression include, more or less explicitly, a parsing strategy that transforms the input text into a sequence of dictionary phrases. Since that usually the parsing of a text is not unique, for compression purpose it makes sense to find one of the possible parsing that minimizes the final compression ratio. This is the parsing problem.

In the foundational methods (such as the work of Lempel and Ziv), the parsing problem was not immediately clear as it was confused with the dictionary building strategy. The overall compression algorithms have strictly imposed the parsing of the text. As soon as many variants of such methods appeared along the sequent years, like the Storer' and Szymanski's variant (1982) or the Welch's variant (1984), the maintenance of the dynamic dictionary was clearly divided from the text parsing strategy and, in the meantime, the importance of coupling a kind of compression on the symbols different from the compression for the dictionary phrases taken place. This last feature was initially undervalued in the theoretical model of the compression processes.

One of the first parsing problem models is due to Schuegraf et al. (see [33]). They associated a graph with as many nodes as the characters that form the text and one edge for each dictionary phrase. In this model, the optimal parsing is obtained by using shortest path algorithms on the associated graph. But this approach was not recommended for practical purpose as it was considered too time consuming. Indeed, the graph can have quadratic size with respect to the text length.

A classic formalization of a general dictionary compression algorithm was proposed by Bell et al. in the late 1990, focusing on just three points: the dictionary definition, the dictionary phrases encoding method and the parsing strategy. This model does not acquire all the richness of many advanced dictionary-based compression algorithms as it does not take account of the symbolwise compression.

Recently, in chronological order, [12], [25], [8] and [9] revised both a more general dictionary compression algorithms definition and the graph-based model for the parsing problem and they also introduced a new optimal parsing algorithm. A similar result for the LZ77-like dictionary case, were independently presented in [17], where the symbolwise feature is still not considered.

The study of free mixtures of two compressors is quite involved and it represents a new theoretical challenge. Free mixture has been implicitly or explicitly used for a long time in many fast and effective compressors such as the gzip compression utility (see [30, Sect. 3.23]), the PkZip Archiving Tool (see [30, Sect. 3.23]), the Rolz Compressor<sup>1</sup>, and the MsZip cabinet archiving software (see [30, Sect. 3.7]), also known as CabArc. In order to glance at compression performances see the web page of Mahoney's challenge<sup>2</sup> about large text compression. In detail, there are two famous compression methods that can work together: the *dictionary* encoding and the *statistical* encoding, which are also called *parsing* (or macro) encoding and *symbolwise* encoding, respectively. The fact that these methods can work together is commonly accepted in practice even if the first theory of Dictionary-Symbolwise methods started in [12].

This work focus on the parsing problem and introduce a twofold result; a general model for dictionary-based text compression called Dictionary-Symbolwise theory and a general parsing algorithm that is proved to be optimal under some realistic hypothesis. The Dictionary-Symbolwise model extend both the Bell dictionary compression formalization and the Schuegraf parsing model based on graphs to fit better to the wide class of common compression algorithms.

The parsing algorithm we present is called Dictionary-Symbolwise Flexible Parsing and it covers almost all the cases of the dictionary-based text compression algorithms together with the large class of their variants where

---

<sup>1</sup>For an example see the RZM Order-1 ROLZ Compressor by Christian Martelock (2008) web site: <http://encode.ru/threads/1036>. Last verified on March 2012.

<sup>2</sup>Matt Mahoney's Large Text Compression Benchmark is a competition between lossless data compression programs. See the web page: <http://mattmahoney.net/dc/text.html>. Last verified on March 2012.

the text is parsed as a sequence of symbols and dictionary phrases. It exploits the prefix closed property of common dictionaries, i.e. both the LZ77 and LZ78-like dictionaries. It works for *dynamic* dictionaries and *variable* costs of dictionary phrases and symbols. His main part concerns with the construction of a smallest subgraph that guarantees parsing optimality preservation, and then a shortest path is found by using a classic single source shortest path approach.

The symbolwise encoding can be any classical one that works in linear time, as many common variable-length encoders do. Our algorithm works under the assumption that a special graph that will be described in the following is well defined. Even if this condition is not satisfied it is possible to use the same method to obtain almost optimal parses. In detail, when the dictionary is LZ78-like, we show that our algorithm has  $O(n)$  complexity, where  $n$  is the size of the text. When the dictionary is LZ77-like our algorithm can be implemented in time  $O(n \log n)$ . Both above solutions have  $O(n)$  space complexity.

This thesis is organized as follows. Chapter 1 is devoted to background notions for data compression. Chapter 2 defines and explores the dictionary-symbolwise model for dictionary compression as well as the graph-based model for the parsing problem. Chapter 3 is an overview about historic parsing solutions and contains the generalization to the dynamic case of the classic Cohn's theorem on greedy parsing optimality for suffix-closed dictionaries. Chapter 4 concerns with the new optimal parsing algorithm called dictionary-symbolwise flexible parsing. Chapter 5 presents a new indexing data structure that solves efficiently the rightmost position query of a pattern over a text. This problem is involved in the graph building process for LZ77-like algorithms, where edge labels, i.e. dictionary pointers costs, are not uniquely determined. Chapter 6 contains the conclusions of this thesis and some open problem.

Even if the main aim of this work is theoretical, some experimental results are introduced in the Appendix A to underline some practical effects of the parsing optimality in compression performance. We have experimental evidence that many of the most relevant LZ77-like commercial compressors

use an optimal parsing. Therefore this thesis contains both a good model for many of the commercial dictionary compression algorithms and a general parsing algorithm with proof of optimality. This fills the gap between theory and best practice about text compression.



# Chapter 1

## Background

This chapter concerns with some well known concepts from the field of the Information Theory, that are fundamental to deal with data compression. Information Theory literature is quite large by now. We remand to [30], [31] and [32] books for a comprehensive look on background notions and standard techniques of data compression. We report here just few prerequisites to make readers comfortable with notation and concepts we will use in the rest of this thesis.

### 1.1 Self-Information and Entropy

A foundational concept for Information Theory is the Shannon's self-information definition. It is a quantitative measure of information. Let  $A$  be a probabilistic event, i.e.  $A$  is the set of outcomes of some random experiment. If  $P(A)$  is the probability that the event  $A$  will occur, then the self-information associated with  $A$  is given by:  $i(A) = -\log_2 P(A)$  bits.

If we have a set of independent events  $A_i$ , which are sets of outcomes of some experiment  $\mathcal{S}$ , which sample space is  $S = \cup A_i$ , then the average self-information associated with the random experiment  $\mathcal{S}$  is given by  $H(\mathcal{S}) = \sum P(A_i) i(A_i) = -\sum P(A_i) \log_2 P(A_i)$  bits. This quantity is called the *entropy* associated with the experiment.

Now, if the experiment is a source  $\mathcal{S}$  that emits a string  $S$  of symbols over the alphabet  $\Sigma = \{1, \dots, m\}$ , i.e.  $S = s_1 s_2 s_3 \dots$  with  $s_i \in \Sigma$ , then the



sample space is the set of all the strings the source can produce, i.e. the set of all the possible sequences of alphabet symbols of any length. The *entropy* of the source  $\mathcal{S}$  is given by

$$H(\mathcal{S}) = \lim_{n \rightarrow \infty} \frac{1}{n} G_n$$

with

$$G_n = - \sum_{i_1=1}^m \cdots \sum_{i_n=1}^m P(s_1 = i_1, \dots, s_n = i_n) \log P(s_1 = i_1, \dots, s_n = i_n).$$

If each symbol in the string is independent and identically distributed (*iid*), then we have that

$$G_n = -n \sum_{i=1}^m P(i) \log P(i) \quad \text{and} \quad H(\mathcal{S}) = - \sum_{i=1}^m P(i) \log P(i).$$

When the symbol probabilities are not independent from each other, the distribution follow an intrinsic model of probability of the source. In this case, the above two entropy equations are not equal and we distinguish them calling the latter *first order entropy*.

The probability distribution over the symbols of a source is not usually *a priori* known and the best we can do is to infer the distribution looking inside some sample strings. Obviously, the underlay assumption is that the source is an *ergodic* source, i.e. its output at any time has the same statistical properties.

The Markov process is the common way to model the source distribution when symbols are not independent each other. In this case we have that each new outcome depends on all the previous one. A *discrete time Markov chain* is a special type of Markov model for those experiments where each observation depends on just the  $k$  previous one, i.e.

$$P(s_n | s_{n-1}, s_{n-2}, \dots) = P(s_n | s_{n-1}, s_{n-2}, \dots, s_{n-k})$$

where the set  $\{s_{n-1}, s_{n-2}, \dots, s_{n-k}\}$  is the state of the  $k$ -order Markov process. The entropy of a Markov process is defined as the average value of the entropy at each state, i.e.

$$H(\mathcal{M}_k) = - \sum_{s_{n-k}} P(s_{n-k}) \sum_{s_{n-k+1}} P(s_{n-k+1} | s_{n-k}) \sum_{s_{n-k+2}} P(s_{n-k+2} | s_{n-k+1}, s_{n-k}) \cdots$$

$$\cdots \sum_{s_{n-1}} P(s_{n-1}|s_{n-2}, \dots, s_{n-k}) \sum_{s_n} P(s_n|s_{n-1}, \dots, s_{n-k}) \log P(s_n|s_{n-1}, \dots, s_{n-k})$$

where  $s_i \in \Sigma$ . In the data compression field is common to refer to the state  $\{s_{n-1}, \dots, s_{n-k}\}$  of previous symbols by using the string  $s_{n-k} \dots s_{n-1}$  called the context of length  $k$  of  $s_n$ .

### Empirical Entropy

The  $k$ -order empirical entropy (see [16]) is the measure of information of a text  $T$  based on the number of repetitions in  $T$  of any substring  $w$  of length  $k$ . Let be

$$H_k(T) = -\frac{1}{n} \sum_{w \in \Sigma^k} n_w \left[ \sum_{\sigma \in \Sigma} \frac{n_{w\sigma}}{n_w} \log \left( \frac{n_{w\sigma}}{n_w} \right) \right]$$

where  $n = |T|$ ,  $\Sigma$  is the alphabet,  $w \in \Sigma^k$  is a string over  $\Sigma$  of length  $k$ ,  $w\sigma$  is the string  $w$  followed by the symbol  $\sigma$  and  $n_w$  is the number of occurrences of  $w$  in  $T$ .

This quantity does not refer to a source or to a probabilistic model, but it only depends from the text  $T$ . The empirical entropy is used to measure the performance of compression algorithms as a function of the string structure, without any assumption on the input source.

## 1.2 Entropy Encoding

Entropy encoding, statistical codes or symbolwise codes, as they are also called, are those compression methods that use the expectation value to reduce the symbol representation. There are static model as well as adaptive or dynamic models. They are usually coupled with a probabilistic model that is in charge of providing symbol probability to the encoder. Common models use symbol frequencies or the symbol context to predict the next symbol.

The most simple statistical encoder is the 0-order arithmetic encoding. It considers all the symbols as if they are independent each other. The adaptive version use to estimate symbol probability with the frequency of occurrence of any symbol in the already seen text.

Huffman code keeps count of the symbol frequencies while reads the input text or by preprocessing it, and then assigns shorter codewords of a prefix-free code to the most occurring symbols accordingly with the Huffman tree. Notice that the notation  $k$ -order used in this thesis refers to models where  $k$  is the length of the context, i.e. a 0-order model is a model where the symbol probabilities just depend on the symbol itself.

## Arithmetic coding

The basic idea of arithmetic coding is to represent the entire input with an interval of real numbers between 0 and 1. The initial interval is  $[0, 1)$  and then it is divided in slots accordingly to the symbol probability. Once that a symbol is encoded, the corresponding slot of the interval is divided again accordingly with the adapted symbol distribution. While the active slots becomes finer and finer, its internal points bit representation grows. As soon as the extremal points of the slot have an equal upper part in their bit representation, these bits are outputted and the slot is scaled to be maintained under the finite precision of the representation of real values inside the machine. As any point of a slot represents an infinite set of infinite strings, all having the same prefix, one of them is chosen when the input string terminate to be outputted. The termination ambiguity is usually handled by using a special terminal symbol or by explicitly giving the text length at the beginning.

The output length of arithmetic codes can be accurately estimated by using the Markov process entropy or the empirical entropy. Moreover, it is proved that their compression ratio converges in probability to the entropy of any i.i.d. source. A similar result can be stated in the case of Markov chains. In practice, when the source is unknown, better results are obtained when higher order models are used, because the models get “closer” to the real source. On the other side, higher order models need more time and space to be handled.

### 1.3 Encoding of Numbers and Commas

Encoding is a fundamental stage of many compression algorithms which consists of uniquely representing a sequence of integers as a binary sequence. In the most simple case the encoder makes use of a code, that is a mapping of the positive integers onto binary strings (codewords), in order to replace each value in input with its corresponding codeword. Codewords can be of variable-lengths as long as the resulting code is uniquely decodable, e.g. the prefix-free codes. Prefix-free property requires that no codeword can be equal to a prefix of another codeword. Several codes have been proposed that achieve small average codeword-length whenever the frequencies of the input integers are monotonically distributed, such that smaller values occur more frequently than larger values.

The unary encoding of an integer  $n$  is simply a sequence of  $n$  1s followed by a 0. Unary encoding is rarely used as stand-alone tool and it is often component of more complex codes. It achieves optimality when integer frequencies decrease exponentially as  $p(i + 1) \leq p(i)/2$ .

The Elias codes is a family of codes where codewords have two parts. The first one is devoted to states the codeword length and the second one is the standard binary representation of the integer, without the most significant bit. The first Elias encoder is the well-known  $\gamma$ -code, which stores the prefix-part in unary. Elias  $\delta$ -code differs from  $\gamma$ -code because it encodes also first-part of the codewords with a  $\gamma$ -code, rather than using the unary code. This is an *asymptotically optimal* code as the ratio between the codeword length and the binary representation length asymptotically tends to 1. The Elias  $\omega$ -code can be seen as the iteration of as many  $\delta$ -code nested encodings until a length of two or three bits is reached.

### 1.4 Dictionary Methods

Dictionary compression methods are based on the substitution of phrases in the text with references to dictionary entries. A dictionary is an ordered collection of phrases, and a reference to a dictionary phrase is usually called dictionary pointer. The idea is that if the encoder and the decoder share

the same dictionary and, for most of the dictionary phrases, the size of the representation in output of a dictionary pointer is less than the size of the phrase itself, then a shorter representation of the input text is obtained replacing phrases with pointers. In order to proceed to the phrase substitution, the text has to be divided into a sequence of dictionary phrases. Such decomposition is called *parsing* and is not usually unique. For compression purpose it makes sense to find one of the possible parsing that minimizes the final compression ratio. This is the parsing problem.

The foundational methods in dictionary compression class are Lempel' and Ziv's LZ77 and LZ78 algorithms that will be extensively considered along this thesis. Lempel' and Ziv's methods are the basis of almost all the dictionary compression algorithms. They are the first relevant dictionary methods that use dynamic dictionaries.

The LZ77 method consider the already seen text as the dictionary, i.e. it uses a dynamic dictionary that is the set of all the substrings of the text up to the current position. Dictionary pointers refer to occurrences of the pointed phrase in the text by using the couple  $(length, offset)$ , where the *offset* stand for the backward offset w.r.t. the current position. Since a phrase is usually repeated more than once along the text and since pointers with smaller offset are usually smaller, the occurrence close to the current position is preferred. Notice that this dictionary is both prefix and suffix closed. The parsing strategy use the greedy approach to find the longest phrase in the dictionary equal to a prefix of the rest of the text.

The LZ78 dictionary is a subset of the LZ77 one. It is prefix-closed but it is not suffix-closed. Each dictionary phrases is equal to another dictionary phrase with a symbol appended at the end. Exploiting this property, dictionary is implemented as an ordered collection of couples  $(dictionary\ pointer, symbol)$ , where the dictionary pointer refers to a previous dictionary phrase or to the empty string. As long as the input text is analyzed, the longest match between the dictionary and the text is selected to form a new dictionary phrase. Indeed, a new couple is formed by this selected dictionary phrase and the symbol in the text that follows the occurrence of this phrase. This new dictionary phrase is added to the dynamic dictionary and it is chosen also to be part of the parsing of the text accordingly with the greedy

parsing.

More detail about these method will be reported in next chapters.



## Chapter 2

# Dictionary-Symbolwise Text Compression

Many dictionary-based compression algorithms and their practical variants use to parse the text as a sequence of both dictionary phrases and symbols. Different encoding are used for those two kinds of parse segments. Indeed, many variants of the classic Lempel and Ziv algorithms allow to parse the text as a free mixture of dictionary phrases and symbols. This twofold nature of the parsing segments was not caught in classic formulation of the dictionary-based compression theory. In this chapter we recall the classical dictionary compression algorithm formulation and the classic model of the parsing problem before presenting the more general framework for Dictionary-Symbolwise compression that better fits to almost all the dictionary-based algorithms.

### 2.1 Dictionary Compression

In [4] it is possible to find a survey on Dictionary and Symbolwise methods and a description of the deep relationship among them (see also [3, 11, 30, 31]).

**Definition 2.1.** A dictionary compression algorithm, as noticed in [4], can be fully described by:

1. The dictionary description, i.e. a static collection of phrases or a



complete algorithmic description on how the dynamic dictionary is built and updated.

2. The encoding of dictionary pointers in the compressed data.
3. The parsing method, i.e. the algorithm that splits the uncompressed data in dictionary phrases.

We notice that any of the above three points can depend on each other, i.e. they can be mutually interdependent.

As the reader can notice, above three points are general enough to describe both static and dynamic dictionary and both static and variable costs for the dictionary phrase representation in the output data. We want now to focus on its third point where the parsing is defined as just a sequence of dictionary pointers. The drawback of this constraint is to lead to an overuse of formalism as it is not easy to describe the role played by symbols. Let us show this effect by examples. The characterization of the classic LZ77 and LZ78 algorithms according to the above Definition 2.1 are stated in what follows.

### **LZ77 characterization**

Given an input text  $T \in \Sigma^*$ , it is processed left to right. At time  $i$ ,  $D_i$  is the current state of the dictionary and  $T_i = T[1 : i]$  is the prefix of  $T$  of length  $i$  that has been already parsed.  $T[i - P : i]$  is called the *search buffer* and  $T[i + 1 : i + Q]$  is called the *look-ahead buffer*, where  $P$  is the maximum offset for text factors,  $Q$  is the maximum length for dictionary phrases.

1. Let be  $D_i = \{wa, \text{ such that } w \in \text{Fact}(T[i - P : i]), a \in \Sigma \text{ and } |wa| \leq Q\}$ , where  $\text{Fact}(x)$  is the set of all the factors (or substrings) of the text  $x$ . Let us notice that this dictionary is essentially composed by the factors having length less than or equal to  $Q$  that appear inside a sliding window of size  $P$  ending at position  $i$  over the text.
2. The dictionary phrase  $wa = T[i - p : i - p + q]a$  is represented by the vector  $(p, q, a)$  where  $p$  is the backward offset in the search buffer of an occurrence of  $w$  and  $q = |w|$ . The vector  $(p, q, a)$  is coded by using

three fixed length sequences of bits where  $p$  has length  $\log_2(P)$ ,  $q$  has length  $\log_2(Q)$  and  $a$  is represented with 8 bits by using the ascii code for symbols.

3. An online greedy parsing is used. At time  $i$ ,  $i$  is the position in the text at which the parsing ends up. At this point, the longest match between a dictionary phrase  $wa \in D_i$  and a prefix of the *look-ahead buffer*  $T[i + 1 : i + Q]$  is added to the parsing. The new parsing covers now the text up to the position  $i + |wa|$ . The next parsing phrase will be chosen at position  $i + |wa|$ . For instance if the parsing of the text up to the position  $i$  is the sequence of phrases  $w_1a_1 w_2a_2 \dots w_ja_j$ , then the parsing up to the position  $i + |wa|$  is  $w_1a_1 w_2a_2 \dots w_ja_j w_{j+1}a_{j+1}$ , with  $w_{j+1}a_{j+1} = wa$ .

### LZ78 characterization

Let us suppose that we have a text  $T \in \Sigma^*$  and that we are processing it left to right. We also suppose that at time  $i$  the text up to the  $i$ th character has been encoded. The algorithm maintains a dynamic table  $M_i$  of phrases, initialized with the empty word  $M_0 = [\epsilon]$ .

1. The dictionary  $D_i$  is defined as  $D_i = \{wa \text{ such that } w \in M_i \text{ and } a \in \Sigma\}$ . At time  $i$ , the longest dictionary phrase  $wa \in D_i$  that matches with the text at position  $i$  is chosen to be part of the set  $M_{i+|wa|}$ , while  $M_i = M_j$  with  $i \leq j < i + |wa|$ . Then,  $wa$  is added at the first empty row of  $M_i$  that becomes  $M_{i+|wa|} = M_i \cup wa$ . Consequently,  $D_{i+|wa|} = D_i \cup wa\Sigma$  and  $D_i = D_j$  with  $i \leq j < i + |wa|$ .  $D_i$  is prefix closed at any time by construction. Many practical implementations use a bounded size dictionary by keeping fixed the dictionary once it gets full or by using a prune strategy that preserves the prefix-closed property.
2. The dictionary phrase  $wa \in D_i$  is represented by the couple  $(x, a)$  where  $x$  is the index of  $w$  over  $M_i$ . The couple  $(x, a)$  is encoded by using a fixed length encodings for the integer  $x$  followed by the ascii value of  $a$ .

3. The parsing is the greedy parsing. It is the sequence of the longest matches between the dictionary and a prefix of the uncompressed part of the text. The parsing phrases are equal to the dictionary phrases  $wa$  used as support for dictionary extension. For instance if the parsing of the text up to the position  $i$  is the sequence of dictionary phrases  $w_1a_1 w_2a_2 \dots w_p a_p = T[1 : i]$ , then the parsing up to the position  $i + |wa|$  is  $w_1a_1 w_2a_2 \dots w_p a_p wa$  with  $wa \in D_i$ .

## 2.2 Dictionary-Symbolwise Compression

We propose a new definition for the class of dictionary-based compression algorithms that takes account of the presence of single characters beside to dictionary phrases. For this reason we chose to name them dictionary-symbolwise algorithms. The following definition is an extension of the above Definition 2.1 due to Bell et al. (see [4]) and it refines what was presented in [8, 12, 25].

**Definition 2.2.** A dictionary-symbolwise compression algorithm is specified by:

1. The dictionary description.
2. The encoding of dictionary pointers.
3. The symbolwise encoding method.
4. The encoding of the flag information.
5. The parsing method.

A *dictionary-symbolwise* algorithm is a compression algorithm that uses both dictionary and symbolwise compression methods. Such compressors may parse the text as a free mixture of dictionary phrases and literal characters, which are substituted by the corresponding pointers or literal codes, respectively. Therefore, the description of a dictionary-symbolwise algorithm also includes the so called *flag information*, that is the technique used to distinguish the actual compression method (dictionary or symbolwise) used for each phrase or factor of the parsed text. Often, as in the case of

LZSS (see [36]), an extra bit is added either to each pointer or encoded character to distinguish between them. Encoded information flag can require less space than one bit according to the encoding used.

For instance, a dictionary-symbolwise compression algorithm with a fixed dictionary  $D = \{ab, cbb, ca, bcb, abc\}$  and the static symbolwise codeword assignment  $[a = 1, b = 2, c = 3]$  could compress the text  $abccacbbabbcbcb$  as  $F_d1F_s3F_d3F_d2F_d1F_d4F_d2$ , where  $F_d$  is the flag information for dictionary pointers and  $F_s$  is the flag information for the symbolwise code.

More formally, a parsing of a text  $T$  in a dictionary-symbolwise algorithm is a pair  $(parse, Fl)$  where  $parse$  is a sequence  $(u_1, \dots, u_s)$  of words such that  $T = u_1 \dots u_s$  and where  $Fl$  is a boolean function that, for  $i = 1, \dots, s$  indicates whether the word  $u_i$  has to be encoded as a dictionary pointer or as a symbol. See Table 2.1 for an example of dictionary-symbolwise compression.

### LZ77 characterization

Given a text  $T \in \Sigma^*$  and processing it left to right, at time  $i$  the text up to the  $i$ th character has been read.

1. Let be  $D_i = \{w, \text{ such that } w \in \text{Fact}(T[i - P : i]) \text{ and } |w| < Q\}$ , where  $P$  is the maximum offset for text factors,  $Q$  is the maximum length for dictionary phrases. Let  $T[i - P : i]$  be called the *search buffer* and  $T[i + 1 : i + Q]$  be called the *look-ahead buffer*.
2. The dictionary phrase  $w = T[i - p : i - p + q]$  is represented by the vector  $(p, q)$  where  $p$  is the backward offset in the search buffer at

<b>Input</b>	$ab$	$c$	$ca$	$cbb$	$ab$	$bcb$	$cbb$
<b>Output</b>	$F_d1$	$F_s3$	$F_d3$	$F_d2$	$F_d1$	$F_d4$	$F_d2$

Table 2.1: Example of compression for the text  $abccacbbabbcbcb$  by a simple Dictionary-Symbolwise algorithm that use  $D = \{ab, cbb, ca, bcb, abc\}$  as a static dictionary, the identity as a dictionary encoding and the mapping  $[a = 1, b = 2, c = 3]$  as a symbolwise encoding.

which the phrase  $w$  appears. The vector  $(p, q)$ , also called dictionary pointer, is coded by using two fixed length sequences of bits, where  $p$  has length  $\log_2(P)$  and  $q$  has length  $\log_2(Q)$ .

3. Any symbol  $a \in \Sigma$  is represented with 8 bits by using the ascii code for symbols.
4. The flag information is not explicitly encoded because it is completely predictable. Indeed, after a dictionary pointer there is a symbol and after a symbol there is a dictionary pointer.
5. The parsing impose a strictly alternation between dictionary pointers and symbols. At any time  $i$ , if  $i$  is the position in the text at which the already chosen parsing ends up, then the match between the longest prefix of the *look-ahead buffer*  $T[i + 1 : i + Q]$  and a dictionary phrase  $w \in D_i$  is chosen to be outputted followed by the mismatch symbol. For instance, if  $w$  is the longest match between the dictionary and the look-ahead buffer, with  $w$  represented by the couple  $(p, q)$ , then  $F_d p q F_s T_{i+|w|}$  are concatenated to the parsing. Otherwise, the already chosen parsing overpass position  $i$  in the text and nothing has to be added to the current parsing.

This new formalization allows to describe dictionary algorithms in a more natural way. Moreover, it allows to easily describe those variants where just a single point of the algorithm is different. For instance, let us focus on LZSS, the LZ77-based algorithm due to Storer and Szymanski of the '82 (see [36]). The main idea of this algorithm is to relax the parsing constrain of LZ77 about dictionary pointers and symbols alternation, allowing to use dictionary pointers every time that it is possible and to use symbols just when it is needed.

### **LZSS characterization**

Given a text  $T \in \Sigma^*$  and processing it left to right, at time  $i$  the text up to the  $i$ th character has been read.

1. Let be  $D_i = \{w, \text{ such that } w \in \text{Fact}(T[i - P : i]) \text{ and } |w| < Q\}$ , where  $P$  is the maximum offset for text factors,  $Q$  is the maximum length for dictionary phrases. Let  $T[i - P : i]$  be called the *search buffer* and  $T[i + 1 : i + Q]$  be called the *look-ahead buffer*.
2. The dictionary phrase  $w = T[i - p : i - p + q]$  is represented by the vector  $(p, q)$  where  $p$  is the backward offset in the search buffer at which the phrase  $w$  appears. The vector  $(p, q)$ , also called dictionary pointer, is coded by using two fixed length sequences of bits, where  $p$  has length  $\log_2(P)$  and  $q$  has length  $\log_2(Q)$ .
3. Any symbol  $a \in \Sigma$  is represented with 8 bits by using the ascii code for symbols.
4. The flag information is explicitly encoded by using 1 bit with conventional meaning. For instance,  $F_d = 0$  and  $F_s = 1$ .
5. At any time  $i$ , if  $i$  is the position in the text at which the already chosen parsing ends up, the match between the longest prefix of the *look-ahead buffer*  $T[i + 1 : i + Q]$  and a dictionary phrase  $w \in D_i$  is chosen to be outputted. For instance, if  $w$  is the longest match between the dictionary and the look-ahead buffer, with  $w$  represented by the couple  $(p, q)$ , then  $\{F_d \ p \ q\}$  are concatenated to the parsing. If there is no match between dictionary and look-ahead buffer, then a single symbol is emitted, i.e.  $F_s \ T[i + 1]$ . Otherwise, the above chosen parsing overpass position  $i$  in the text and nothing has to be added to the current parsing.

## Dictionary-Symbolwise Scheme

Let now focus on the parsing point. For any dictionary compression algorithm we can build a class of variants taking fixed the first four points and changing the parsing strategy or, if it is needed, arranging a bit the first four points to allow using a different parsing. Usually the variants of the same class maintain decoder compatibility, i.e. their outputs can be decoded by the same decoding algorithm. Algorithms of the same class can

be compared looking for the best parsing, i.e. the parsing that minimizes compression ratio. We call *scheme* such class of algorithms.

**Definition 2.3.** Let a dictionary-symbolwise *scheme* be a nonempty set of dictionary-symbolwise algorithms having in common the same first four specifics, i.e. they differ from each other by the parsing methods only.

A scheme does not need to contain *all* the algorithms having the same first four specifics. Let us notice that any of the above points from 1 to 5 can depend on all the others, i.e. they can be mutually interdependent. The word *scheme* has been used by other authors with other meaning, e.g. *scheme* is sometimes used as synonymous of algorithm or method. In this thesis, *scheme* always refers to the above Definition 2.3.

*Remark 2.1.* For any dictionary-symbolwise scheme  $S$  and for any parsing method  $P$ , a dictionary-symbolwise compression algorithm  $\mathcal{A}_{S,P}$  is completely described by the first four specifics of any of the algorithms belonging to  $S$  together with the description of the parsing method  $P$ .

Let us here briefly analyze some LZ-like compression algorithms. The LZ78 algorithm is, following the above definitions, a dictionary-symbolwise algorithm. It is easy to naturally arrange its original description to a dictionary-symbolwise complaint definition. Indeed, its dictionary building description, its dictionary pointer encoding, its symbolwise encoding, its parsing strategy and the null encoding of the flag information are, all together, a complete dictionary-symbolwise algorithm definition. The flag information in this case is not necessary, because there is not ambiguity about the nature of the encoding to use for any of the parse segments of the text as the parsing strategy imposes a rigid alternation between dictionary pointers and symbols. Similar arguments apply for LZ77. Later on we refer to these or similar dictionary-symbolwise algorithms that have null flag information as “pure” dictionary algorithms and to scheme having only “pure” dictionary algorithms in it as “pure” scheme.

LZW (see [30, Section 3.12]) naturally fits Definition 2.1 of dictionary algorithms and, conversely, LZSS naturally fits Definition 2.2 on dictionary-symbolwise algorithms as well as the LZMA algorithm (see [30, Section 3.24]).





empty word  $\varepsilon$  is associated with vertex 0, that is also called the *origin* of the graph. The set of directed edges is

$$E = \{(p, q) \subset (V \times V) \mid p < q \text{ and } \exists w_{p,q} = T[p+1 : q] \in D_p\}$$

where  $T[p+1 : q] = a_{p+1}a_{p+2}\cdots a_q$  and  $D_p$  is the dictionary relative to the  $p$ -th processing step, i.e. the step in which the algorithm either has processed the input text up to character  $a_p$ , for  $p > 0$ , or it has begun, for  $p = 0$ . For each edge  $(p, q)$  in  $E$ , we say that  $(p, q)$  is *associated with* the dictionary phrase  $w_{p,q} = T[p+1 : q] \in D_p$ . In the case of a static dictionary,  $D_i$  is constant along all the algorithm steps, i.e.  $D_i = D_j, \forall i, j = 0 \cdots n$ . Let  $L$  be the set of edge labels  $L_{p,q}$  for every edge  $(p, q) \in E$ , where  $L_{p,q}$  is defined as the cost (weight) of the edge  $(p, q)$  when the dictionary  $D_p$  is in use, i.e.  $L_{p,q} = C((p, q))$ .

Let us consider for instance the case where the cost function  $C$  associates the length in bit of the encoded dictionary pointer of the dictionary phrase  $w_{p,q}$  to the edge  $(p, q)$ , i.e.  $C((p, q)) = \text{length}(\text{encode}(\text{pointer}(w_{p,q})))$ , with  $w_{p,q} \in D_p$ . In this case the weight of a path  $\mathcal{P}$  from the origin to the node  $n = |T|$  on the graph  $G_{\mathcal{A},T}$  corresponds to the size of the output obtained by using the parsing induced by  $\mathcal{P}$ . The path of minimal weight on such graph corresponds to the parsing that achieves the best compression. The relation between path and parsing will be investigated in Section 2.4.

If the cost function is a total function, then  $L_{p,q}$  is defined for each edge of the graph.

**Definition 2.4.** Let us say that  $G_{\mathcal{A},T}$  is *well defined* iff  $L_{p,q}$  is defined for each edge  $(p, q)$  of the graph  $G_{\mathcal{A},T}$ .

For instance, the use of common variable-length codes for dictionary pointers, as Elias or Fibonacci codes or a static Huffman code, leads to a well defined graph. Sometimes the cost function is a partial function, i.e.  $L_{p,q}$  is not defined for some  $p$  and  $q$ , and  $G_{\mathcal{A},T}$  in such cases is not well defined. For instance, encoding the dictionary pointers via statistical codes, like a Huffman code or an arithmetic code, leads to partial cost functions. Indeed the encoding of pointers and, accordingly, the length of the encoded dictionary pointers may depend on how many times a code is used. For

instance, when a dynamic code like the *online* adaptive Huffman code is used, the codeword lengths depend on how frequently the codewords have been used in the past. In the case of a semi static Huffman code, since it is an *offline* encoding, the codeword lengths depend on how frequently the codewords are used in the compression process of the whole text. In these cases, the cost function obviously depends on the parsing and the labels of the graph edges are not defined until a parsing is fixed. Moreover, the cost function may be undefined for edges that represent phrases never used by the parsing. The latter case is still an open problem, i.e. it is not known how to find an optimal parsing strategy when the encoding costs depend on the parsing itself.

*Remark 2.2.* We call  $G_{\mathcal{A},T}$  the “Schuegraf’s graph” in honour of the first author of [33] where a simpler version was considered in the case of static-dictionary compression method.

We can naturally extend the definition of the graph associated with an algorithm to the dictionary-symbolwise case. Given a text  $T = a_1 \dots a_n$ , a dictionary-symbolwise algorithm  $\mathcal{A}$ , and a cost function  $C$  defined on edges, the graph  $G_{\mathcal{A},T} = (V, E, L)$  is defined as follows. The vertices set is  $V = \{0 \dots n\}$ , with  $n = |T|$ . The set of directed edges  $E = E_d \cup E_s$ , where

$$E_d = \{(p, q) \subset (V \times V) \mid p < q - 1, \text{ and } \exists w = T[p + 1 : q] \in D_p\}$$

is the set of dictionary edges and

$$E_s = \{(q - 1, q) \mid 0 < q \leq n\}$$

is the set of symbolwise edges.  $L$  is the set of edge labels  $L_{p,q}$  for every edge  $(p, q) \in E$ , where the label  $L_{p,q} = C((p, q))$ . Let us notice that the cost function  $C$  hereby used has to include the cost of the flag information to each edge, i.e.  $C((p, q))$  is equal to the cost of the encoding of  $F_d$  ( $F_s$ , resp.) plus the cost of the encoded dictionary phrase  $w \in D_p$  (symbolwise  $a_q$ , resp.) associated with the edge  $(p, q)$  where  $(p, q) \in E_d$  ( $E_s$ , resp.). Moreover, since  $E_d$  does not contain edges of length one by definition,  $G_{\mathcal{A},T} = (V, E, L)$  is not a multigraph. Since this graph approach can be extended to multigraph, with an overhead of formalism, one can relax the  $p < q - 1$  constrain in the

definition of  $E_d$  to  $p \leq q - 1$ . All the results we will state in this thesis, naturally extend to the multigraph case.

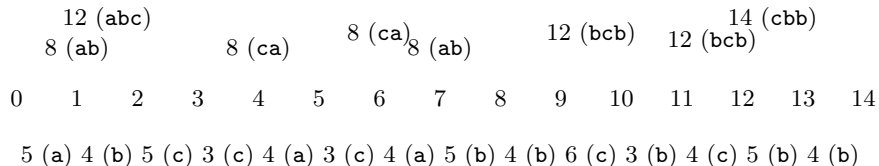


Figure 2.2: Graph  $G_{\mathcal{A},T}$  for the text  $T = \text{abccacabbcbcb}$ , for the dictionary-symbolwise algorithm  $\mathcal{A}$  with static dictionary  $D = \{\text{ab}, \text{abc}, \text{bcb}, \text{ca}, \text{cbb}\}$  and cost function  $C$  as defined in the graph. The dictionary phrase or the symbol associated with an edge is reported near the edge label within parenthesis.

## 2.4 On Parsing Optimality

In this section we assume that the reader is well acquainted with LZ-like dictionary encoding and with some simple statistical encodings such as the Huffman code.

**Definition 2.5.** Fixed a dictionary description, a cost function  $C$  and a text  $T$ , a *dictionary (dictionary-symbolwise) algorithm* is *optimal* within a set of algorithms if the cost of the encoded text is minimal with respect to all others algorithms in the same set. The parsing of an *optimal* algorithm is called *optimal* within the same set.

When the bit length of the encoded dictionary pointers is used as a cost function, the previous definition of optimality is equivalent to the classical well known definition of bit-optimality for dictionary algorithm. Notice that the above definition of optimality strictly depends on the text  $T$  and on a set of algorithms. A parsing can be optimal for a certain text but not for another one. Obviously, we are mainly interested on parsings that are optimal either for *all* texts over an alphabet or for classes of texts. Whenever it is not explicitly written, from now on when we talk about optimal parsing we

mean optimal parsing for *all* texts. About the set of algorithm it makes sense to find sets as large as possible.

Classically, there is a bijective correspondence between parsings and paths in  $G_{\mathcal{A},T}$  from vertex 0 to vertex  $n$ , where optimal parses correspond to minimal paths and vice-versa. We say that a parse (path, resp.) *induces* a path (parse, resp.) to denote this correspondence. This correspondence was firstly stated in [33] only in the case of sets of algorithms sharing the same *static* dictionary and where the encoding of pointers has constant cost.

For example the path along vertices (0, 3, 4, 5, 6, 8, 11, 12, 13, 14) is the shortest path for the graph in Fig. 2.2. Authors of [12] were the first to formally extend the Shortest Path approach to dynamically changing dictionaries and variable costs.

**Definition 2.6.** A scheme  $\mathcal{S}$  has the *Schuegraf property* if, for any text  $T$  and for any pair of algorithms  $\mathcal{A}, \mathcal{A}' \in \mathcal{S}$ , the graph  $G_{\mathcal{A},T} = G_{\mathcal{A}',T}$  with  $G_{\mathcal{A},T}$  well defined.

This property of schemes is called *property of Schuegraf* in honor of the first of the authors in [33]. In this case we define  $G_{\mathcal{S},T} = G_{\mathcal{A},T}$  as the graph of (any algorithm of) the scheme. The proof of the following proposition is straightforward.

**Proposition 2.4.1.** *There is a bijective correspondence between optimal parsings and shortest paths in  $G_{\mathcal{S},T}$  from vertex 0 to vertex  $n$ .*

**Definition 2.7.** Let us consider an algorithm  $\mathcal{A}$  and a text  $T$  and suppose that graph  $G_{\mathcal{A},T}$  is well defined. We say that  $\mathcal{A}$  is *graph optimal* (with respect to  $T$ ) if its parsing induces a shortest path in  $G_{\mathcal{A},T}$  from the origin (i.e. vertex 0) to vertex  $n$ , with  $n = |T|$ . In this case we say that its parsing is *graph optimal*.

Let  $\mathcal{A}$  be an algorithm such that for any text  $T$  the graph  $G_{\mathcal{A},T}$  is well defined. We want to associate a scheme  $\mathcal{SC}_{\mathcal{A}}$  with it in the following way. Let  $S$  be the set of all algorithms  $\mathcal{A}$  such that, for any text  $T$ ,  $G_{\mathcal{A},T}$  exists (i.e. it is well defined). Let  $\mathcal{B}$  and  $\mathcal{C}$  be two algorithms in  $S$ . We say that  $\mathcal{B}$  and  $\mathcal{C}$  are equivalent or  $\mathcal{B} \equiv \mathcal{C}$  if, for any text  $T$ ,  $G_{\mathcal{B},T} = G_{\mathcal{C},T}$ .

We define the scheme  $\mathcal{SC}_A$  to be the equivalence class that has  $\mathcal{A}$  as a representative. It is easy to prove that  $\mathcal{SC}_A$  has the Schuegraf property.

We can connect the definition of *graph* optimal parsing with the previous definition of  $\mathcal{SC}_A$  to obtain the next proposition, which proof is an easy consequence of the Proposition 2.4.1 and of the Schuegraf property of  $\mathcal{SC}_A$ . Roughly speaking, the graph optimality within the scheme  $\mathcal{SC}_A$  implies scheme (or global) optimality.

**Proposition 2.4.2.** *Let us consider an algorithm  $\mathcal{A}$  such that for any text  $T$  the graph  $G_{\mathcal{A},T}$  is well defined. Suppose further that for a text  $T$  the parsing of  $\mathcal{A}$  is graph optimal. Then the parsing of  $\mathcal{A}$  of the text  $T$  is (globally) optimal within the scheme  $\mathcal{SC}_A$ .*

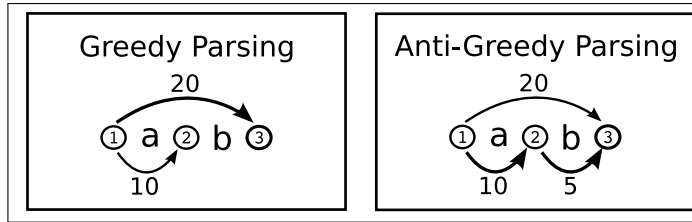


Figure 2.3: Locally but not globally optimal parsing

We have simple examples (see Figure 2.3), where a parsing of a text is graph optimal and the corresponding algorithm belongs to a scheme that has not the Schuegraf property and it is not optimal within the same scheme.

For instance, let us now consider the text  $T = ab$  and two algorithms  $\mathcal{A}$  and  $\mathcal{A}'$  in the same scheme, where  $\mathcal{A}$  is the algorithm that uses the greedy parsing and  $\mathcal{A}'$  uses the anti-greedy parsing, as in Figure 2.3. The parsing of  $\mathcal{A}$  is the greedy parsing, that at any reached position chooses the longest match between text and dictionary. The graph  $G_{\mathcal{A},T}$  for this greedy algorithm has three nodes, 0, 1, 2, and only two edges, both outgoing 0, one to node 1 that costs 10 and another to node 2 that costs 20. The greedy parsing reaches the end of the text with this second edge which has global cost 20 and then it is graph optimal. The parsing of  $\mathcal{A}'$  is the anti-greedy that at any reached position chooses the shortest nonempty match between text and dictionary. The graph  $G_{\mathcal{A}',T}$  for this anti-greedy algorithm has three nodes, 0, 1, 2, and three edges, two outgoing 0, one to node 1 that costs 10

and another to node 2 that costs 20 and a third outgoing 1 to node 2 that costs 5. The parsing of the anti-greedy algorithm is  $(a)(b)$  and it costs 15. The dictionary is composed by  $\langle a, ab \rangle$  if the parsing of the processed text has reached an even position (starting from position 0) with a cost of 10 and 20 respectively. The dictionary is  $\langle a, b \rangle$  if the parsing of the processed text has reached an odd position with a cost of 5 each. Notice that now the dictionary phrase “a” has a different cost than before. The dictionary and the edge costs are changing as a function of the reached position, depending if this position is even or odd, and, in turn, it depends on the parsing. Therefore, we have that  $G_{\mathcal{A},T}$  and  $G_{\mathcal{A}',T}$  are well defined but they are different from each other and the scheme has not the Schuegraf property. Therefore both the greedy and the anti-greedy parsing are graph optimal but none of them is (globally) optimal.

## 2.5 Dictionary-Symbolwise Can Be Better

So, why should we use dictionary-symbolwise compressors?

From a practical point of view, coupling a fast symbolwise compressor to a dictionary compressor gives one more degrees of freedom to the parsing, increasing compression ratio without slowing up the entire compression process. Or, on the contrary, a dictionary compressor coupled with a powerful symbolwise compressor can speed up the decompression process without decreasing the compression ratio. This approach that mixes together dictionary compression and symbolwise compression methods is already widely used in practical compression software solutions, even if its scientific basis were not clearly defined and it was treated just as a practical trick to enhance compression ratio and to take under control and improve the decompression speed. Several viable algorithms and most of the commercial data compression programs, such as *gzip*, *zip* or *cabarc*, are, following our definition, dictionary-symbolwise. Still from a practical point of view, some experimental results are showed and discussed in next section.

In this section instead we study some theoretical reasons for using dictionary-symbolwise compression algorithms.

First of all, it is not difficult to give some “artificial” and trivial example

where coupling a dictionary and a symbolwise compressor give rise to a better optimal solution. Indeed let us consider the static dictionary  $D = \{a, b, ba, bb, abb\}$  together a cost function  $C$  that could represents the number of bits of a possible code:  $\{C(a) = 8, C(b) = 12, C(ba) = 16, C(bb) = 16, C(abb) = 4\}$ .

A greedy parsing of the text  $babb$  is  $(ba)(bb)$  and the cost of this parsing is 32. An optimal parsing for this dictionary is  $(b)(abb)$  that has cost 16. This example shows, as also the one of Figure 2.3, that a greedy parsing is not always an optimal parsing in dictionary compressors.

Let us consider further the following static symbolwise compressor that associates with the letter  $a$  a code of cost 8 and that associates with the letter  $b$  a code of cost 4 that could represent the number of bits of this code. The cost of coding  $babb$  following this symbolwise compressor is 20.

If we connect them in a dictionary-symbolwise compressor then an optimal parsing is  $S(b)D(abb)$  where the flag information is represented by the letter  $S$  for symbolwise and by  $D$  for dictionary. The cost of the trivially encoded flag information is one bit for each letter or phrase. Therefore the cost of this parsing is 10.

In this subsection, however, we will prove something more profound than artificial examples such the one above. Indeed, from a theoretical point of view, Ferragina et al. (cf. [17]) proved that the compression ratio of the classic greedy-parsing of a LZ77 pure dictionary compressor may be far from the bit-optimal pure dictionary compressor by a multiplicative factor  $\Omega(\log n / \log \log n)$ , which is indeed unbounded asymptotically. The family of strings that is used in [17] to prove this result, is a variation of a family that was used in [24].

In next two subsections we show a similar result between the bit optimal dictionary compressor and a dictionary-symbolwise compressor. Therefore a bit optimal dictionary-symbolwise compressor can use, in some pathological situation, the symbolwise compressor to avoid them and be provably better than a simple bit optimal dictionary compressor.

## LZ77 Case

Let us define these two compressors. The first is a LZ77-based compressor with unbounded window size as dictionary and with a Huffman code on lengths and an optimal parser. It allows overlaps between the search buffer, i.e. the dictionary, and the look ahead buffer. The encoding of pointers can be any of the classical encoding for integer. We just impose a Huffman code on the lengths.

We further denote by  $\text{OPT-LZH}(s)$  the bit length of the output of this compressor on the string  $s$ .

The same LZ77 is used as dictionary compressor in the dictionary-symbolwise compressor. Clearly we do not include the parser in the dictionary-symbolwise compressor, but, analogously as above, we suppose to have an optimal parser for the dictionary-symbolwise compressor, no matter about the description. The flag information  $\{D, S\}$  is coded by a run-length encoder. The cost of a run is subdivided over all symbolwise arcs of the run, i.e. if there is a sequence of  $n$  consecutive symbolwise arcs in the optimal parsing then the cost of these  $n$  flag information  $S$  (for Symbolwise) will be in total  $O(\log n)$  and the cost of each single flag information in this run will be  $O(\frac{\log n}{n})$ .

It remains to define a symbolwise compression method.

In the next result we could have used a PPM\* compressor but, for simplicity, we use a longest match symbolwise. That is, the symbolwise at position  $k$  of the text searches for the closest longest block of consecutive letters in the text up to position  $k - 1$  that is equal to a suffix ending in position  $k$ . This compressor predicts the  $k + 1$ -th character of the text to be the character that follows the block. It writes a symbol 'y' (that is supposed not to be in the text) if this is the case. Otherwise it uses an escape symbol 'n' (that is supposed not to be in the text) and then writes down the correct character plainly. A temporary output alphabet has therefore two characters more than the characters in the text. This temporary output will be subsequently encoded by a run-length encoder (see [15]).

This is not a very smart symbolwise compressor but it fits our purposes, and it is simple to analyze.



We further denote by  $\text{OPT-DS}(s)$  the bit length of the output of this dictionary-symbolwise compressor on the string  $s$ .

**Theorem 2.5.1.** *There exists a constant  $c > 0$  such that for every  $n' > 1$  there exists a string  $s$  of length  $|s| \geq n'$  satisfying*

$$\text{OPT-LZH}(s) \geq c \frac{\log |s|}{\log \log |s|} \text{OPT-DS}(s).$$

*Proof.* For every  $n'$  let us pick a binary word  $w$  of length  $3n, n \geq n', w = a_1 a_2 \cdots a_{3n}$  that has the following properties.

1. For any  $i, i = 1, 2 \cdots n$ , the compressor  $\text{OPT-LZH}(s)$  cannot compress the word  $a_i a_{i+1} \cdots a_{2i+n-1}$  of length  $n + i$  with a compression ratio greater than 2.
2. every factor (i.e. every block of consecutive letters) of  $w$  having length  $3 \log 3n$  of  $w$  is unique, i.e. it appears in at most one position inside  $w$ .

Even if it could be hard to explicitly show such a word, it is relatively easy to show that such a word exists. Indeed, following the very beginning of the Kolmogorov's theory, the vast majority of words are not compressible. A simple analogous counting argument can be used to prove that Property 1 is satisfied by the vast majority of strings of length  $3n$ , where, for vast majority we mean that the percentage of strings not satisfying Property 1 decreases exponentially in  $n$ . Here, to play it safe, we allowed a compression "two to one" for all the  $n$  considered factors.

A less known result (see [2, 7, 13, 14, 18, 37]) says that for random strings and for any  $\epsilon > 0$  the percentage of strings of length  $n$  having each factor of length  $2 \log n + \epsilon$  *unique* grows exponentially to 1 (i.e. the percentage of strings not having this property decreases exponentially). Here we took as  $\epsilon$  the number 1. Therefore such a string  $a_1 \cdots a_{3n}$  having both properties surely exists for some  $n \geq n'$ .

Let us now define the word  $s$  over the alphabet  $\{0, 1, c\}$  in the following way.

$$s = a_1 a_2 \cdots a_{n+1} c^{2^n} \cdots a_i a_{i+1} \cdots a_{2i+n-1} c^{2^n} \cdots a_n a_{n+1} \cdots a_{3n-1} c^{2^n} a_{n+1} a_{n+2}$$

Let us now evaluate  $\text{OPT-LZH}(s)$ . By Property 1 each binary word that is to the left or to the right of a block of  $2^n$   $c$ 's cannot be compressed in less than  $\frac{1}{2}n$  bits in a “stand-alone” manner. If one such a string is compressed by a pointer to a previous string then the offset of this pointer will be greater than  $2^n$  and, so, its cost in bit is  $O(n)$ . We defined the string  $s$  in such a manner that all “meaningful” offsets are different, so that even a Huffman code on offsets (that we do not use, because we use a Huffman code only for lengths) cannot help. Therefore there exists a constant  $c'$  such that  $\text{OPT-LZH}(s) \geq c'n^2$ .

Let us now evaluate  $\text{OPT-DS}(s)$ . We plan to show a parse that will give a string of cost  $\text{P-DS}(s) \leq \hat{c}n \log n$  as output. Since  $\text{OPT-DS}(s) \leq \text{P-DS}(s)$  then also  $\text{OPT-DS}(s) \leq \hat{c}n \log n$ .

The blocks of  $2^n$   $c$ 's have all the same length. We parse them with the dictionary compressor as  $(c)(c^{2^n} - 1)$ . The dictionary compressor is not used in other positions in the parse  $P$  of the string  $s$ . The Huffman code on lengths of the dictionary compressor would pay  $n$  bits for the table and a constant number of bits for each occurrence of a block of  $2^n$   $c$ 's. Hence the overall cost in the parse  $P$  of all blocks of letters  $c$  is  $O(n)$ . And this includes the flag information that consists into two bits  $n$  times.

Parse  $P$  uses the symbolwise compressor to parse all the binary strings. The first one  $a_1a_2 \cdots a_{n+1}$  costs  $O(n)$  bits. Starting from the second  $a_2a_3 \cdots \cdots a_{n+3}$  till the last one, the symbolwise will pay  $O(\log n)$  bits for the first  $3 \log 3n$  letters and then, by Property 2, there is a long run of ‘y’ that will cover the whole string up to the last two letters. This run will be coded by the run-length code of the symbolwise. The overall cost is  $O(\log n)$  and this includes the flag information that is a long run of  $S$  coded by the run-length of the flag information. The cost of the symbolwise compressor including the flag information over the whole string is then  $O(n \log n)$ , that dominates the cost of the dictionary-symbolwise parse  $P$ .

The length of the string  $s$  is  $O(n2^n + n^2)$  and therefore  $\log |s| = n + o(n)$  and the thesis follows.  $\square$

*Remark 2.3.* In the theorem above it is possible to improve the constants in the statement. This can be done simply using for instance a word  $a_1 \cdots a_{n^2}$

instead of  $a_1 \cdots a_{3n}$ . It is possible to optimize this value, even if, from a conceptual point of view, it is not important.

We want to underline that the Huffman code on the lengths is essential in this statement. At the moment we are not able to find a sequence of strings  $s$  where the dictionary-symbolwise compressor is provably better than the optimal dictionary version without using a Huffman code. It is an open question whether this is possible.

We finally notice that if the dictionary is coupled with a ROLZ technique then the optimal solution of the pure dictionary compressor reaches the same level of the dictionary-symbolwise compressor. This is not surprising because the ROLZ technique is sensible to context and do not “pay” for changing the source of the text.

### **LZ78 Case**

Matias and Sahinalp in [28] already shown that Flexible Parsing is optimal with respect to all the prefix-closed dictionary algorithms, including LZ78, where optimality stands for phrase optimality. Flexible Parsing is also optimal in the suffix-close dictionary algorithm class. Phrase optimality is equal to bit optimality under the fixed length codeword assumption, so we say just optimality. From now on we assume *FP* or its extension as optimal parsing and the bit length of the compressed text as coding cost function.

In this subsection we prove that there exists a family of strings such that the ratio between the compressed version of the strings obtained by using an optimal LZ78 parsing (with constant cost encoding of pointers) and the compressed version of the strings obtained by using an optimal dictionary-symbolwise parsing is unbounded. The dictionary, in the dictionary-symbolwise compressor, is still the LZ78 dictionary, while the symbolwise is a simple Last Longest Match Predictor that will be described later. We want to notice here that similar results were proved in [28] between flexible parsing and the classical LZ78 and in [17] between a compressor that uses optimal parsing over a LZ77 dictionary and the standard LZ77 compressor (see also [24]). Last but not least we notice that in this example, analogously as done in [28], we use an unbounded alphabet just to make the example more clear.

An analogous result can be obtained with a binary alphabet with a more complex example.

Let us define a Dictionary-Symbolwise compressor that uses LZ78 as dictionary method, the Last Longest Match Predictor as symbolwise method, Run Length Encoder to represent the flag information and one optimal parsing method. Let us call it OptDS-LZ78. We could have used a PPM\* as symbolwise compressor but Last Longest Match Predictor (LLM) fits our purposes and it is simple to analyze. LLM Predictor is just a simple symbolwise compression method that uses the last longest seen match to predict next character.

The symbolwise searches, for any position  $k$  of the text, the closest longest block of consecutive letters up to position  $k - 1$  that is equal to a suffix ending in position  $k$ . This compressor predicts the  $(k + 1)$ -th character of the text to be the character that follows the block. It writes a symbol 'y' (that is supposed not to be in the text) if this is the case. Otherwise it uses an escape character 'n' (that is supposed not to be in the text) and then writes down the correct character plainly. A temporary output alphabet has therefore two characters more than the characters in the text. This temporary output will be subsequently encoded by a run-length encoder. This method is like the Yes?No version of Symbol Ranking by P. Fenwick (see [15]).

It costs  $\log n$  to represent a substring of  $n$  characters that appear after the match. For each position  $i$  in the uncompressed text if  $m_i$  is the length of the longest match in the already seen text it produces  $n$  that costs  $O(\log n)$  bits as output, i.e.  $C(T[i + 1 : i + n]) = n$  and  $Cost(n) = O(\log n)$  where

$$\forall m, j \quad m_i = \max_m (T[i - m : i] = T[j - m : j] \text{ with } j < i \text{ and } T[i - m - 1] \neq T[j - m - 1])$$

Let us consider a string  $S$

$$S = \sum_{z=1}^k 1 + \dots + z = [1 + 12 + 123 + \dots + 1..z + \dots + 1..k]$$

that is the concatenation of all the prefixes of  $1..k$  in increasing order. Let us consider the string  $T'$  that is the concatenation of the first  $\sqrt{k}$  suffixes of  $2..k$ , i.e.  $T' = 2..k \cdot 3..k \cdot \dots \cdot \sqrt{k}..k$  and a string  $T = S \cdot T'$ . We

use  $S$  to build a dictionary formed by just the string  $1..k$  and its prefixes and no more. We assume that both the dictionary and the symbolwise methods work the same up to the end of the string  $S$ , so they produce an output that is very similar in terms of space. It is not difficult to prove that an optimal LZ78 compressor would produce on  $T$  a parse having cost at least  $O(k + k \log k) = O(k \log k)$  while the optimal dictionary-symbolwise compressor (under the constant cost assumption on encoding pointers) has a cost that is  $O(k + \sqrt{k} \log k) = O(k)$ .

*Proof.* (Sketch) An optimal constant-cost LZ78 compressor must use  $k$  phrases to code  $S$ . Then each phrase used to code the subword  $2 \dots k$  of  $T'$  has length at most 2 and therefore the number of phrases that it must use to code  $2 \dots k$  is at least  $(k-1)/2 \geq \frac{1}{2}k/2$ . Analogously, each phrase used to code the subword  $3 \dots k$  of  $T'$  has length at most 3 and therefore the number of phrases that it must use to code  $3 \dots k$  is at least  $(k-2)/3 \geq \frac{1}{3}k/2$ . We keep on going up to conclude that number of phrases that it must use to code  $\sqrt{k} \dots k$  is at least  $(k - \sqrt{k} + 1)/\sqrt{k} \geq \frac{1}{\sqrt{(k)}}k/2$ . Adding all these numbers we get that the total number of phrases is smaller than or equal to  $O(k + \log \sqrt{k} \times k/2) = O(k \log k)$ .

Let us now prove that an optimal dictionary-symbolwise compressor has a cost that is  $O(k)$  by showing that there exists at least one parse that has cost  $O(k)$ .

The parse that we analyze parses  $S$  with the LZ78 dictionary and spend for this part of the string  $O(k)$ . Then it uses the LLM Predictor to compress the subword  $2 \dots k$  of  $T'$ . Firstly it outputs a symbol 'n' followed by the symbol 2 because it is unable to predict the symbol 2 and then it outputs  $k-2$  symbols 'y' that, in turn, are coded by the run length encoder with a cost that is  $O(\log k)$ . The whole cost of subword  $2 \dots k$  is then  $O(\log k)$ . Then the LLM Predictor compresses sequentially the subword  $i \dots k$  of  $T'$ , with  $3 \leq i \leq \sqrt{k}$  and any time it spends at most  $O(\log k)$ . The total cost of this parse is then  $O(k + \sqrt{k} \log k) = O(k)$ .  $\square$

## Chapter 3

# About the Parsing Problem

In this chapter we survey some of the milestone results about the parsing problem, starting from those concerning static dictionaries through the dynamic case. In the last section we present a small new contribution that complete the picture of fixed costs case. It is the generalization of the greedy parsing of Cohn for static dictionary to the dynamic dictionary case.

### 3.1 The Optimal Parsing Problem

Optimal with respect to what? Obviously in data compression we are mainly interested to achieve the best compression ratio, that corresponds to minimizing the size of the compressed data. This notion of optimality is sometimes called bit-optimality. But our question has a deeper sense. When can we say that a parsing algorithm is an optimal parsing algorithm? We could say that a parsing is optimal with respect to the input data or with respect to the compression algorithms in which it is involved.

Following Proposition 2.4.2, we say that, fixed a compression scheme, e.g. the LZ77-like algorithms, an optimal parsing is the parsing algorithm that applied to the scheme gives the compression algorithm (see the Remark 2.1) with the best compression ratio for any input text with respect to all the compression algorithms in the same scheme.

Therefore, given a dictionary-based compression algorithm, e.g. LZW, we take the largest compression scheme that contains this algorithm. This

scheme is the set of all the algorithms that use the same dictionary description and the same encodings (e.g. the LZW-like algorithms or LZW scheme). All the algorithms in this scheme differ each other in just the parsing method. If the scheme has the Schuegraf property (see Definition 2.6), then an optimal parsing will be any parsing that minimize the compression ratio within the algorithms of the scheme, i.e., from Proposition 2.4.1, any parsing that induces a shortest path on the graph associated to the scheme. Notice that if the scheme has not the Schuegraf property, then it is not clear how to obtain an optimal parsing that is optimal for the entire scheme.

### 3.2 An Overview about Theory and Practice

In this section we do not want to give an exhaustive dissertation about the parsing problem as there is a vast literature concerning this topic and there are many practical solutions. We want just to recall some of the milestones in the data compression field that are strongly relevant with the scope of this thesis.

In '73, the Wagner's paper (see [38]) shows a  $O(n^2)$  dynamic programming solution, where a text  $T$  of length  $|T| = n$  is provided at ones.

In '74 Schuegraf (see [33]) showed that the parsing problem is equal to the shortest path problem on the graph associated to both a text and a static dictionary (see Section 2.3 for the graph-based model of the parsing problem). Since that the full graph for a text of length  $n$  can have  $O(n^2)$  edges in the worst case and the minimal path algorithm has  $O(V + E)$  complexity, we have another solution for the parsing problem of  $O(n^2)$  complexity.

In '77 and '78 years the foundational dynamic dictionary-based compression methods LZ77 and LZ78 have been introduced. They both use an online greedy parsing that is simple and fast. Those compression methods use both an uniform (constant) cost model for the dictionary pointers. The online greedy approach is realized by choosing the longest dictionary phrase that matches with the forwarding text to extend the parsing, moving on the text left to right, until the whole text is covered.

In '82, the LZSS compression algorithm, based on the LZ77 one, was presented (see [36]). It improves the compression rate and the execution

time without changing the original parsing approach. The main difference is that a symbol is used only when there is no match between dictionary and text. But when there is any match, the longest one is always chosen. It uses a flag bit to distinguish symbols from dictionary pointers. In the same paper Storer proved the optimality of the greedy parsing for the original LZ77 dictionary with unbounded size (see Theorem 10 in [36] with  $p = 1$ ).

In '84, LZW variant of LZ78 was introduced by Welch (see [39]). This is, to our best knowledge, the first theoretical compression method that uses a dynamic dictionary and variable costs of pointers.

In '85, Hartman and Rodeh proved in [20] the optimality of the *one-step-lookahead parsing* for prefix-closed static dictionary and uniform pointer cost. The main point of this approach is to choose the phrase that is the first phrase of the longest match between two dictionary phrases and the text, i.e. if the current parsing covers the text up to the  $i$ th character, then we choose the phrase  $w$  such that  $ww'$  is the longest match with the text after position  $i$ , with  $w, w'$  belong to the dictionary.

In '89 and later in '92, the *deflate* algorithm was presented and used in the *PKZIP* and *gzip* compressors. It uses a LZ77-like dictionary, the LZSS flag bit and, for some options, a non-greedy parsing. Both dictionary pointers and symbols are encoded by using a Huffman code.

In '95, Horspool investigated in [21] about the effect of non-greedy parsing in LZ-based compression. He showed that using the above *one-step-lookahead parsing* in the case of dynamic dictionaries leads to better compression. Horspool tested this parsing on the LZW algorithm and on a new LZW variant that he presented in the same paper.

In '96 the greedy parsing was ultimately proved by Cohn et al. (see [5]) to be optimal for static suffix-closed dictionary under the uniform cost model. They also proved that the right to left greedy parsing is optimal for prefix-closed dictionaries. Notice that the greedy approach is linear and online. These results will be reported and extended in the next section.

In '99, Matias and Sahinalp (see [28]) gave a linear-time optimal parsing algorithm in the case of prefix-closed dynamic dictionary and uniform cost of dictionary pointer, i.e. the codeword of all the pointers have equal length. They extended the results given in [20], [21] and [23] to the dynamic case.



Matias and Sahinalp called their parsing algorithm *Flexible Parsing*. It is also known as semi-greedy parsing.

In '09, Ferragina et al. (see [17]) introduced an optimal parsing algorithm for LZ77-like dictionary and variable-length code, where the code length is assumed to be the cost of a dictionary pointer.

In '10, Crochemore et al. (see [8]) introduced an optimal parsing for prefix-closed dictionaries and variable pointer costs. It was called *dictionary-symbolwise flexible parsing* and it is extensively treated and extended in this thesis. It works for the original LZ77 and LZ78 algorithms and for almost all of their known variants.

### 3.3 A Generalization of the Cohn's Theorem

A static dictionary  $D$  is prefix-closed (suffix-closed) if and only if for any phrase  $w \in D$  in the dictionary all the prefixes (suffixes) of  $w$  belong to the dictionary, i.e.  $\text{suff}(w) \subset D$  ( $\text{pref}(w) \subset D$ ). For instance, the dictionary  $\{bba, aba, ba, a\}$  is suffix-closed.

The classic Cohn' and Khazan's result of '96 (see [5]) states that if  $D$  is a static suffix-closed dictionary, then the greedy parsing is optimal under the uniform cost assumption. Symmetrically, the reverse of the greedy parsing on the reversed text, i.e. the right to left greedy parsing, is optimal for static prefix-closed dictionary. Roughly speaking, the original proof concerns with suffix-closed dictionaries and shows that choosing the longest dictionary phrases guarantees to cover the text with the minimum number of phrases. The main idea is that, fixed a point  $a$  on the text, the greedy choice is the longest match with the dictionary and text on the right of  $a$ , e.g the phrase that cover the text from  $a$  to  $b$ . Now,  $b$  is the rightmost point that can be reached from  $a$  or from any point on the left of  $a$ . For absurd, given  $\alpha$  on the left of  $a$ , if in the dictionary there is a phrase equal to the text from  $\alpha$  to  $\beta$  on the right of  $b$ , then, for the suffix-closed property, the phrase  $a - \beta$  is in the dictionary too. This contradicts that  $a - b$  is the longest match.

Notice that the optimal greedy solution for suffix-closed static dictionary can be computed *online*.

Let us compare these results to the LZ family of algorithms. We notice

that above results concern with static dictionaries and cannot be directly applied in these cases. Moreover, the greedy parsing of the LZ-78 algorithm is not the one that was proved to be optimal for the prefix-closed case.

For the LZ77 algorithm, which dictionary is both prefix-closed and suffix-closed, we need to generalize the Cohn' and Khazan's result for suffix-closed dictionary to the dynamic case. Within this thesis we refer to the Cohn' and Khazan's result on suffix-closed dictionary simply as the Cohn's theorem in honour of its first author.

Let us start focusing on the prefix and suffix-closed definition of dynamic dictionaries. The classic definition of suffix-closed dynamic dictionary just require that the dictionary is suffix-closed at any time. Therefore, given a suffix-closed dynamic dictionary  $D$ , if  $D_i$  is the dictionary  $D$  at the time  $i$ , then  $D_i$  is suffix-closed, for any  $i$ . The prefix-closed case is analogously defined. Notice that this definition does not make any assumption on the relationship between dictionaries in two different moments. We call above properties *weak* suffix-closed property and *weak* prefix-closed property respectively. We introduce the *strong* suffix-closed property and the *strong* prefix-closed property of dynamic dictionary as follows.

**Definition 3.1.** A dynamic dictionary  $D$  has the *strong* prefix-closed property *iff*, for any dictionary phrase  $w \in D_i$ , for any  $k$  in  $[0..|w|]$ ,  $w_k$  is in  $D_i$  and in  $D_{i+k}$ , where  $w_k$  is the prefix of  $w$  of length  $k$ .

**Definition 3.2.** A dynamic dictionary  $D$  has the *strong* suffix-closed property *iff*, for any dictionary phrase  $w \in D_i$ , for any  $k$  in  $[0..|w|]$ ,  $w^k$  is in  $D_i$  and in  $D_{i+k}$ , where  $w^k$  is the suffix of  $w$  of length  $|w| - k$ .

Notice that the strong prefix-closed property implies the weak prefix-closed property. Analogously, the strong suffix-closed property implies the weak property.

We say that a dictionary is *non-decreasing* when  $D_i \subset D_j$  for any  $i, j$  points in time, with  $0 \leq i \leq j$ . A static dictionary is obviously non-decreasing. Practically speaking, a dynamic dictionary is non-decreasing when it can only grow along the time, as for instance the original LZ78 where, at each algorithm step, at most one phrase is inserted in the dynamic dictionary.

*Remark 3.1.* For non-decreasing dictionaries, the weak suffix-closed property and the weak prefix-closed property are equal to the strong suffix-closed property and to the strong prefix-closed property, respectively.

Notice that while the original LZ78 dictionary is non-decreasing, many of its practical implementation and its variants are not. This is because the size of the dictionary is bounded in practice for space saving purpose. Since the original LZ77 dictionary is defined as the set of the substrings of the search buffer, i.e. the backward text up to a certain distance, LZ77 has not the non-decreasing property at all.

**Proposition 3.3.1.** *The LZ77 dictionary is both weak prefix-closed and weak suffix-closed. The LZ77 unbounded dictionary is non-decreasing. The original LZ77 dictionary is strong suffix-closed.*

*Proof.* We refer to the LZ77 dictionary definition given in Section 2.2. The dictionary at time  $i$  is defined as  $D_i = \{w, \text{ such that } w \in \text{Fact}(T[i - P : i]) \text{ and } |w| < Q\}$ , where  $P$  is the maximum offset for text factors,  $Q$  is the maximum length for dictionary phrases.

Since any set of factors is prefix and suffix closed and since for any phrase  $w$  each of its suffixes or prefixes has length less than or equal to  $|w|$ , then LZ77 dictionary is weak prefix-closed and weak suffix-closed by definition.

The LZ77 dictionary is unbounded when  $P \geq |T|$ . In this case the dictionary is equal to the set  $\{w, \text{ such that } w \in \text{Fact}(T[0 : i]) \text{ and } |w| < Q\}$  that is obviously non-decreasing.

Let us focus on the general set  $\text{Fact}(T[i - P : i])$ . This is a sliding window of size  $P$  over the text  $T$ . For any value  $i$ , let be  $T[i - P : i] = au$  and  $T[i - P + 1 : i + 1] = ub$  with  $a, b$  in  $\Sigma$  and  $u$  in  $\Sigma^*$ . Since all the proper suffixes of  $au$  are also suffixes of  $u$ , then for any  $w \in \text{Fact}(T[i - P : i])$  the proper suffixes  $w^k$  of length  $|w| - k$  where  $1 < k \leq |w|$  are also in  $\text{Fact}(T[i - P + 1 : i + 1])$ . Therefore this property also holds for the dictionaries  $D_i$  and  $D_{i+1}$  and it easy to see that this property is equivalent to the strong suffix-closed property defined in Definition 3.2.  $\square$

**Corollary 3.3.2.** *The LZ77 dictionary has the strong suffix-closed property in both the bounded and the unbounded variants.*

The proof of the above proposition comes straightforward from the Proposition 3.3.1, the Remark 3.1 and the dictionary definition.

Let us now show the effect of the suffix-closed and the prefix-closed properties on the Schuegraf graph to visualize those concepts. Given a text  $T$  and an algorithm  $A$  that uses a dictionary  $D$ , we have that if  $D$  has the strong suffix-closed property, then for any edge  $(i, j)$  of the graph  $G_{A,T}$  associated with the phrase  $w \in D_i$ , with  $|w| = j - i$  and  $w = T[i + 1 : j]$ , then all the edges  $(k, j)$ ,  $i < k < j$  are into  $G_{A,T}$ . In the case of prefix-closed dictionaries, as prefix edges start from the same node, the prefix of a dictionary phrase are all represented in the graph if the dictionary has just the weak prefix-closed property.

We want now to extend the elegant proof of Cohn et al. (see [5]) to the case of strong suffix-closed dynamic dictionaries.

Given a text  $T$  of length  $n$  and a dynamic dictionary  $D$  where, at the moment  $i$ -th with  $0 \leq i < n$ , the text  $T_i$  has been processed and  $D_i$  is the dictionary at time  $i$ . Recall that we are under the uniform cost assumption.

**Theorem 3.3.3.** *The greedy parsing of  $T$  is optimal for strong suffix-closed dynamic dictionaries.*

*Proof.* The prove is by induction. We want to prove that for any  $n$  smaller than or equal to the number of phrases of an optimal parsing, there exists an optimal parsing where the first  $n$  phrases are greedy phrases. The inductive hypothesis is that there exists an optimal parsing where the first  $n - 1$  phrases are greedy phrases. We will prove that there is an optimal parsing where the first  $n$  phrases are greedy and, therefore, any greedy parsing is optimal.

Fixed a text  $T$  and a strong suffix-closed dynamic dictionary  $D$ , let  $\mathcal{O} = o_1 o_2 \cdots o_p = T$  be an optimal parsing and let  $\mathcal{G} = g_1 g_2 \cdots g_q = T$  be the greedy parsing, where, obviously,  $p \leq q$ .

The base of the induction with  $n = 0$  is obviously true. Let us prove the inductive step.

By inductive hypothesis,  $\forall i < n$  we have that  $o_i = g_i$ . Since  $g_n$  is greedy, then the  $n$ -th phrase of the greedy parsing is longer than or equal

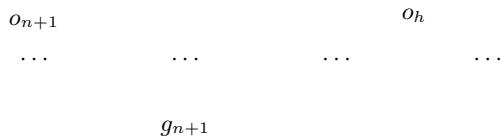


Figure 3.1: Detail of the differences between parsing  $\mathcal{O}$  and parsing  $\mathcal{G}$  over a text  $T$  between positions  $|o_1 \cdots o_n|$  and  $|o_1 \cdots o_h|$ . Nodes and dots represent the text and edges represent parsing phrases as reported on edge labels.

to the  $n$ -th phrase of the optimal parsing, i.e.  $|g_n| \geq |o_n|$  and therefore  $|o_1 \cdots o_n| \leq |g_1 \cdots g_n|$ .

If  $|g_n| = |o_n|$ , then the thesis follows. Otherwise,  $|g_n| > |o_n|$  and  $o_n$  is the first phrase in the optimal parsing that is not equal the  $n$ -th greedy phrase.

Let  $h$  be the minimum number of optimal parsing phrases that overpass  $g_n$  over the text, i.e.  $h = \min\{i \mid |o_1 \cdots o_i| \geq |g_1 \cdots g_n|\}$ . Since  $|g_n| > |o_n|$ , then  $h > n$ . If  $|o_1 \cdots o_h| = |g_1 \cdots g_n|$ , then the parsing  $g_1 \cdots g_n o_{h+1} \cdots o_p$  uses a number of phrases strictly smaller than the number of phrases used by the optimal parsing that is a contradiction. Therefore  $|o_1 \cdots o_h| > |g_1 \cdots g_n|$ . The reader can see this case reported in Figure 3.1.

Let  $|o_1 \cdots o_{h-1}| = T_j$  be the text up to the  $j$ -th symbol. Then  $o_h \in D_j$ , where  $D_j$  is the dynamic dictionary at the time  $j$ . Let  $o_h^k$  the  $k$ -th suffix of  $o_h$  with  $k = |o_1 \cdots o_h| - |g_1 \cdots g_n|$ . For the Property 3.2 of  $D$ ,  $o_h^k \in D_{j+k}$  and then there exists a parsing  $o_1 \cdots o_{n-1} g_n o_h^k o_{h+1} \cdots o_p$ , where  $g_n o_h^k = o_n \cdots o_h$ .

From the optimality of  $\mathcal{O}$ , it follows that  $h = n + 1$ , otherwise there exists a parsing with less phrases than an optimal one. See Figure 3.2.

Therefore  $o_1 \cdots o_{n-1} g_n o_{n+1}^k o_{n+2} \cdots o_p$  is also an optimal parsing. Since  $o_1 \cdots o_{n-1}$  is equal to  $g_1 \cdots g_{n-1}$ , the thesis follows. □

**Corollary 3.3.4.** *The greedy parsing is an optimal parsing for any version of the LZ77 dictionary.*

The proof of the above corollary comes straightforward from the Theorem 3.3.3 and the Corollary 3.3.2.

To our best knowledge, this is the first proof of optimality of the greedy

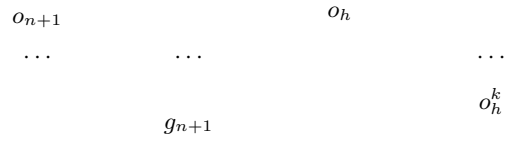


Figure 3.2: Detail of the differences between parsing  $\mathcal{O}$  and parsing  $\mathcal{G}$  over a text  $T$  between positions  $|o_1 \cdots o_n|$  and  $|o_1 \cdots o_h|$ . Nodes and dots represent the text and edges represent parsing phrases as reported on edge labels. The dashed edge  $o_h^k$  represents a suffix of  $o_h$ .

parsing that cover the original LZ77 dictionary case and almost all of the practical LZ77 dictionary implementations where the search buffer is a sliding windows on the text.



## Chapter 4

# Dictionary-Symbolwise Flexible Parsing

In this chapter we consider the case of dictionary-symbolwise algorithms where the parsing is a free mixture of dictionary phrases and symbols. We present the *dictionary-symbolwise flexible parsing* that is a dictionary-symbolwise optimal parsing algorithm for prefix-closed dictionaries and *variable* costs. This is a generalization of the Matias' and Sahinalp's *flexible parsing* algorithm (see [28]) to variable costs.

The algorithm is quite different from the original Flexible Parsing but it has some analogies with it. Indeed, in the case of LZ78-like dictionaries, it makes use of one of the main data structures used for the original flexible parsing in order to be implemented in linear time.

In next sections we will show some properties of the graph  $G_{\mathcal{A},T}$  when the dictionary of the algorithm  $\mathcal{A}$  is prefix-closed and the encoding of the dictionary pointers leads to a *nondecreasing* cost function. We will call *c-supermaximal* some significant edges of  $G_{\mathcal{A},T}$  and we will use those edges to build the graph  $G'_{\mathcal{A},T}$  that is a subgraph of  $G_{\mathcal{A},T}$ . Then, we will show that any minimal path from the origin of  $G'_{\mathcal{A},T}$  is a minimal path in  $G_{\mathcal{A},T}$ . We will introduce the *dictionary-symbolwise flexible parsing* algorithm that build the graph  $G'_{\mathcal{A},T}$  and then find a minimal weight path on it in order to parse the text. We will prove that this parsing is optimal within any scheme having the Schuegraf Property. We will show that the *dictionary-*



*symbolwise flexible parsing* has linear time complexity w.r.t. the text size in the LZ78-like dictionary cases. In the case of LZ77-like dictionaries, the same algorithm has  $O(n \log n)$  time complexities and uses  $O(n)$  space, where  $n$  is the size of the text.

## 4.1 The *c-supermaximal* Edges

We suppose that a text  $T$  of length  $n$  and a dictionary-symbolwise algorithm  $\mathcal{A}$  are given. We assume here that the dictionary is prefix closed at any moment.

Concerning the costs of the dictionary pointer encodings, we recall that costs are variable, that costs assume positive values and that they must include the cost of the flag information. Concerning the symbolwise encodings, the costs of symbols must be positive, including the flag information cost. They can vary depending on the position of the characters in the text and on the symbol itself. Furthermore, we assume that the graph  $G_{\mathcal{A},T}$  is well defined following our Definition 2.4, i.e. we assume that all edges have a well defined cost.

We denote by  $d$  the function that represents the distance of the vertices of  $G_{\mathcal{A},T}$  from the origin of the graph. Such a distance  $d(i)$  is classically defined as the minimal cost of all possible weighted paths from the origin to the vertex  $i$ , with  $d(0) = 0$ . This distance obviously depends on the cost function. We say that cost function  $C$  is *prefix-nondecreasing* at any moment if for any  $u, v \in D_p$  phrases associated with edges  $(p, i), (p, q)$ , with  $p < i < q$  (that implies that  $u$  is prefix of  $v$ ), one has that  $C((p, i)) \leq C((p, q))$ .

**Lemma 4.1.1.** *Let  $\mathcal{A}$  be a dictionary-symbolwise algorithm such that for any text  $T$  the graph  $G_{\mathcal{A},T}$  is well defined. If the dictionary is always prefix-closed and if the cost function is always prefix-nondecreasing then the function  $d$  is nondecreasing monotone.*

*Proof.* It is sufficient to prove that for any  $i$ ,  $0 \leq i < n$  one has that  $d(i) \leq d(i + 1)$ . Let  $j \leq i$  be a vertex such that  $(j, i + 1)$  is an edge of the graph and  $d(i + 1) = d(j) + C((j, i + 1))$ . If  $j$  is equal to  $i$  then  $d(i + 1) = d(i) + C((i, i + 1))$  and the thesis follows. If  $j$  is smaller than

$i$  then, since the dictionary  $D_j$  is prefix closed,  $(j, i)$  is still an edge in  $D_j$  and  $d(i) \leq d(j) + C((j, i)) \leq d(j) + C((j, i + 1)) = d(i + 1)$  and the thesis follows. The last inequality in the previous equation comes from the prefix-nondecreasing property of the cost function.  $\square$

Let us call vertex  $j$  a *predecessor* of vertex  $i \iff \exists(j, i) \in E$  such that  $d(i) = d(j) + C((j, i))$ . Let us define  $pre(i)$  being the smallest of the predecessors of vertex  $i, 0 < i \leq n$ , that is  $pre(i) = \min\{j \mid d(i) = d(j) + C((j, i))\}$ . In other words  $pre(i)$  is the smallest vertex  $j$  that contributes to the definition of  $d(i)$ . Clearly  $pre(i)$  has distance smaller than  $d(i)$ . We notice that a vertex can be a predecessor either via a dictionary edge or via a symbol edge. It is also possible to extend previous definition to pointers having a cost smaller than or equal to a fixed  $c$  as follows.

**Definition 4.1.** For any cost  $c$  we define  $pre_c(i) = \min\{j \mid d(i) = d(j) + C((j, i)) \text{ and } C((j, i)) \leq c\}$ . If none of the predecessor  $j$  of  $i$  is such that  $C((j, i)) \leq c$  then  $pre_c(i)$  is undefined.

If all the costs of pointers are smaller than or equal to  $c$  then for any  $i$  one has that  $pre_c(i)$  is equal to  $pre(i)$ .

Analogously to the notation of [27], we want to define two boolean operations *Weighted-Extend* and *Weighted-Exist*.

**Definition 4.2. (Weighted-Extend)** Given an edge  $(i, j)$  in  $G_{\mathcal{A}, T}$  and a cost value  $c$ , the operation *Weighted-Extend* $((i, j), c)$  finds out whether the edge  $(i, j + 1)$  is in  $G_{\mathcal{A}, T}$  having cost smaller than or equal to  $c$ .

More formally, let  $(i, j)$  in  $G_{\mathcal{A}, T}$  be such that  $w = T[i + 1 : j] \in D_i$ . Operation *Weighted-Extend* $((i, j), c) = \text{“yes”} \iff wa_{j+1} = T[i+1 : j+1] \in D_i$  with  $j < n$  such that  $(i, j + 1)$  is in  $G_{\mathcal{A}, T}$  and  $C((i, j + 1)) \leq c$ , where  $C$  is the cost function associated with the algorithm  $\mathcal{A}$ . Otherwise *Weighted-Extend* $((i, j), c) = \text{“no”}$ . Let us notice that *Weighted-Extend* always fails to extend any edge ending at node  $n$ .

**Definition 4.3. (Weighted-Exist)** Given  $i, j$  with  $0 \leq i < j \leq n$  and a cost value  $c$ , the operation *Weighted-Exist* $(i, j, c)$  finds out whether or not the phrase  $w = T[i + 1 : j]$  is in  $D_i$  and the cost of the corresponding edge  $(i, j)$  in  $G_{\mathcal{A}, T}$  is smaller than or equal to  $c$ .

Let us notice that doing successfully the operation *Weighted-Extend* on  $((i, j), c)$  means that  $wa_{j+1} \in D_i$  is the weighted extension of  $w$  and the encoding of  $(i, j+1)$  has cost less or equal to  $c$ . Similarly, doing successfully a *Weighted-Exist* operation on  $(i, j, c)$  means that an edge  $(i, j)$  exists in  $G_{\mathcal{A}, T}$  having cost less than or equal to  $c$ .

**Definition 4.4. (c-supermaximal)** Let  $E_c$  be the subset of all the edges of the graph having cost smaller than or equal to  $c$ . For any cost  $c$ , let the set  $M_c \subseteq E_c$  be the set of *c-supermaximal* edges, where  $(i, j) \in M_c \iff (i, j) \in E_c$  and  $\forall p, q \in V$ , with  $p < i$  and  $j < q$ , the arcs  $(p, j), (i, q)$  are not in  $E_c$ . For any  $(i, j) \in M_c$  let us call  $i$  a *c-starting point* and  $j$  a *c-ending point*.

**Proposition 4.1.2.** *Suppose that  $(i, j)$  and  $(i', j')$  are in  $M_c$ . One has that  $i < i'$  if and only if  $j < j'$ .*

*Proof.* Suppose that  $i < i'$  and that  $j \geq j'$ . Since the dictionary  $D_i$  is prefix closed we have that  $(i, j')$  is still in  $D_i$  and therefore it is an edge of  $G_{\mathcal{A}, T}$ . By the prefix-nondecreasing property of function  $C$  we have that  $C((i, j')) \leq C((i, j)) = c$ , i.e.  $(i, j') \in E_c$ . This contradicts the fact that  $(i', j')$  is in  $M_c$  and this proves that if  $i < i'$  then  $j < j'$ . Conversely suppose that  $j < j'$  and that  $i \geq i'$ . If  $i > i'$  by previous part of the proof we must have that  $j > j'$ , that is a contradiction. Therefore  $i = i'$ . Hence  $(i, j)$  and  $(i, j')$  both belongs to  $M_c$  and they have both cost smaller than or equal to  $c$ . This contradicts the fact that  $(i, j)$  is in  $M_c$  and this proves that if  $j < j'$  then  $i < i'$ .  $\square$

By previous proposition, if  $(i, j) \in M_c$  we can think  $j$  as a function of  $i$  and conversely. Therefore it is possible to represent  $M_c$  by using an array  $M_c[ ]$  such that if  $(i, j)$  is in  $M_c$ , then  $M_c[j] = i$  otherwise  $M_c[j] = Nil$ . Moreover the non-*Nil* values of this array are strictly increasing. The positions  $j$  having value different from *Nil* are the ending positions.

We want to describe a simple algorithm that outputs all *c-supermaximal* edges scanning the text left-to-right. We call it *Find Supermaximal(c)*. It uses the operations *Weighted-Extend* and *Weighted-Exist*. The algorithm starts with  $i = 0, j = 0$  and  $w = \epsilon$ , i.e. the empty word. The word  $w$  is

indeed implicitly defined by the arc  $(i, j)$  when  $i < j$  or it is the empty word when  $i = j$ . Therefore  $w$  will not appear explicitly in the algorithm. Since the values of  $i$  and  $j$  are only increased by one and  $i$  is always less than or equal to  $j$ , the word  $w$  can be seen as a sliding window of variable size that scans the text left-to-right.  $w$  is moved along the text either by extensions or by contractions to its suffixes.

At each step of the algorithm,  $j$  is firstly increased by one. This extends  $w$  concatenating it to  $T[j]$ . The algorithm executes then a series of *Weighted-Exist* increasing  $i$  by one, i.e. it contracts many times  $w$ . This series of *Weighted-Exist* ends when  $w$  is the empty word or an edge  $(i, j) \in E_c$  is found such that  $(i, j)$  is not contained in any already found *c-supermaximal* edge (see 4.1.4). Indeed, since the increment on  $j$  at line 3, if such edge  $(i, j)$  exists, then we have that  $\forall (p, q) \in M_c$  with  $p < i$ ,  $(i, j) \in E_c$ . Moreover, if such edge  $(i, j)$  exists,  $i$  is a *c-starting* point and a series of *Weighted-Extend* is executed looking for the corresponding *c-ending* point. After each *Weighted-Extend* positive answer,  $j$  is incremented by one. Once that *Weighted-Extend* outputs “no”, i.e. once that  $(i, j)$  cannot be weighted-extended any more,  $(i, j)$  is a *c-supermaximal* and it is inserted into  $M_c$  to be outputted later. The step of the algorithm ends when a *c-supermaximal* is found or when  $w$  is equal to the empty word. The algorithm runs as long as there are unseen characters, i.e. until  $j$  reaches  $n$ .

The algorithm is stated more formally in Table 4.1.

**Proposition 4.1.3.** *Given a cost value  $c$ , the Find Supermaximal algorithm correctly computes  $M_c$ .*

*Proof.* First of all let us prove that if  $(\hat{i}, \hat{j})$  is inserted by the algorithm in  $M_c$  then  $(\hat{i}, \hat{j})$  is *c-supermaximal*.

If  $(\hat{i}, \hat{j})$  is inserted into  $M_c$  at line 11, then an edge  $(\hat{i}, j')$  at line 4 was previously proved to exist and to have cost  $C((\hat{i}, j')) \leq c$ . It caused the termination of the loop at lines 4 – 6. For the line 7 we know that  $\hat{i} < j'$  and by the loop 8 – 10 we know that all the edges  $(\hat{i}, q)$  with  $j' \leq q \leq \hat{j}$  exist and they all are such that  $C((\hat{i}, q)) \leq c$ . Therefore  $(\hat{i}, \hat{j})$  costs at most  $c$  and then the first part of the definition is verified. Since the *Weighted-Extend* $((\hat{i}, \hat{j}), c) = \text{“no”}$  at line 8, that was the exit condition of that loop,

<i>Find Supermaximal</i> ( $c$ )	
01.	$i \leftarrow 0, j \leftarrow 0, M_c \leftarrow \emptyset$
02.	WHILE $j < n$ DO
03.	$j \leftarrow j + 1$
04.	WHILE $i < j$ AND <i>Weighted-Exist</i> ( $i, j, c$ ) = “no” DO
05.	$i \leftarrow i + 1$
06.	ENDWHILE
07.	IF $i < j$ THEN
08.	WHILE <i>Weighted-Extend</i> ( $(i, j), c$ ) = “yes” DO
09.	$j \leftarrow j + 1$
10.	ENDWHILE
11.	INSERT ( $(i, j), M_c$ )
12.	ENDIF
13.	ENDWHILE
14.	RETURN $M_C$

Table 4.1: The pseudocode of the *Find Supermaximal* algorithm. The function INSERT simply insert the edge  $(i, j)$  in the dynamical set  $M_c$ .

then  $(\hat{i}, \hat{j} + 1) \notin E_c$ . Since  $D_i$  is prefix closed and the cost function  $C$  is prefix-nondecreasing  $\forall q \in V$  with  $\hat{j} < q$  the arc  $(\hat{i}, q)$  is not in  $E_c$ , because otherwise  $(\hat{i}, \hat{j} + 1)$  would be in  $E_c$ .

It remains to prove that  $\forall p \in V$  with  $p < \hat{i}$  the arc  $(p, \hat{j})$  is not in  $E_c$ .

Suppose by contradiction that such arc  $(p, \hat{j})$  exists in  $E_c$ . Since the variables  $i, j$  never decrease along algorithm steps, the variable  $i$  reaches the value  $p$  before that  $(\hat{i}, \hat{j})$  is inserted in  $M_c$ . Let  $j_p$  be the value of  $j$  when  $i$  reached the value  $p$ . Since the variable  $i$  is increased only inside the loop at lines 4 – 6, we have that  $p \leq j_p$ . If  $p = j_p$  the algorithm terminates the current step by the conditions at lines 4 and 7 and it enters the next step with  $j = j_p + 1$  due to line 3. Therefore  $j$  will reach the value  $j_p + 1$  for  $p = j_p$  otherwise  $j$  will be equal to  $j_p$ . In both cases, since  $i < j$ , the condition at line 7 is satisfied and the loop 8 – 10 is reached. Since  $D_p$  is prefix closed and the cost function is prefix-nondecreasing then  $\forall q$  such that  $j \leq q < \hat{j}$ , *Weighted-Extend*( $(p, q)$ ) = “yes”. Then, the loop 8 – 10 increases the  $j$  up

to at least the value  $\hat{j}$ , i.e. the algorithm reaches the line 11 with  $\hat{j} \leq j$ . At this point, an edge  $(p, j)$  is inserted in  $M_c$  and the algorithm moves on the next step. Since the increment of the variable  $j$  at line 3, we have that in the rest of the algorithm only edges where  $j$  is greater than  $\hat{j}$  may be considered and then  $(\hat{i}, \hat{j})$  will not be inserted. That is a contradiction. Therefore, if  $(\hat{i}, \hat{j})$  is inserted by the algorithm in  $M_c$  then  $(\hat{i}, \hat{j})$  is *c-supermaximal*.

We have now to prove that if  $(\hat{i}, \hat{j})$  is *c-supermaximal* then it is inserted by the algorithm in  $M_c$ .

Suppose that variable  $i$  never assumes the value  $\hat{i}$ . The algorithm ends when variable  $j$  is equal to  $n$ . Let  $i_n$  be the value of variable  $i$  when  $j$  becomes  $n$ , then we have that  $i_n < \hat{i} < \hat{j} < n = j$ . If the variable  $j$  reaches the value  $n$  inside the loop 8 – 10 then the operation *Weighted-Extend* $((i_n, n - 1), c)$  has outputted “yes” just before. At line 11 the edge  $(i_n, n)$  is inserted into  $M_c$  and then  $(i_n, n)$  is *c-supermaximal*. This contradicts that  $(\hat{i}, \hat{j})$  is *c-supermaximal*. Otherwise, if the variable  $j$  reaches the value  $n$  at line 3, then we have two cases. In the first one, *Weighted-Exist* $(i_n, n, c)$  outputs “yes”, i.e. the edge  $(i_n, n)$  is in  $E_c$ . Since  $i = i_n < n = j$  line 7 condition is satisfied, *Weighted-Extend* $((i_n, n), c)$  outputs “no” by definition and then  $(i_n, n)$  is in  $M_c$ , i.e. it is a *c-supermaximal*. That is a contradiction again. In the second case, *Weighted-Exist* $(i_n, n, c)$  outputs “no” one or multiple times while  $i$  grows up to a value  $i'_n < \hat{i}$  by hypothesis. Using the same argumentation as before,  $(i'_n, n)$  in  $M_c$  leads to a contradiction.

Therefore at a certain moment variable  $i$  assumes the value  $\hat{i}$ . Let  $j_{\hat{i}}$  be the value of variable  $j$  in that moment. We suppose that  $j_{\hat{i}} \leq \hat{j}$ . Since the dictionary  $D_{\hat{i}}$  is prefix closed and the cost function is prefix nondecreasing, *Weighted-Exist* $(\hat{i}, j_{\hat{i}}, c)$  outputs “yes” causing the exit from the loop at lines 4 – 6. At this point, inside the loop 8 – 10, the variable  $j$  reaches the value  $\hat{j}$  since *Weighted-Extend* $((\hat{i}, j), c)$  outputs “yes” for any  $j$  less than  $\hat{j}$ , while *Weighted-Extend* $((\hat{i}, \hat{j}), c)$  outputs “no”. Finally,  $(\hat{i}, \hat{j})$  is inserted into  $M_c$  at line 11. Suppose by contradiction that  $j_{\hat{i}} > \hat{j}$  when  $i$  assumes the value  $\hat{i}$  at line 5. This may happen only if the edge  $(\hat{i} - 1, j_{\hat{i}})$  has been inserted in  $M_c$  in the previous step of the algorithm. Since  $\hat{i} - 1 < \hat{i} < \hat{j} < j_{\hat{i}}$  this contradicts the hypothesis that  $(\hat{i}, \hat{j})$  is *c-supermaximal*.

□

**Proposition 4.1.4.** *For any edge  $(i, j) \in E_c$  there exists a  $c$ -supermaximal edge  $(\hat{i}, \hat{j})$  containing it, i.e. such that  $\hat{i} \leq i$  and  $j \leq \hat{j}$ .*

*Proof.* We build  $(\hat{i}, \hat{j})$  in algorithmic fashion. The algorithm is described in what follows in an informal but rigorous way. If edge  $(i, j)$  is not  $c$ -supermaximal then we proceed with a round of *Weighted-Extend* $((i, j), c)$  analogously as described in algorithm *Find Supermaximal* and we increase  $j$  of one unit until *Weighted-Extend* outputs “no”. Let  $j'$  be the value of  $j$  for which *Weighted-Extend* output “no”. Clearly  $(i, j') \in E_c$  and  $(i, j' + 1)$  is not. If  $(i, j')$  is not  $c$ -supermaximal the only possibility is that there exists at least one  $i' < i$  such that  $(i', j') \in E_c$ . At this point we keep iterating previous two steps starting from  $(i - 1, j')$  instead of  $(i, j)$  and we stop whenever we get a  $c$ -supermaximal edge, that we call  $(\hat{i}, \hat{j})$ .  $\square$

By previous proposition, for any node  $v \in G_{\mathcal{A}, T}$  if there exists a node  $i < v$  such that  $C((i, v)) = c$  and  $d(v) = d(i) + c$  then there exists a  $c$ -supermaximal edge  $(\hat{i}, \hat{j})$  containing  $(i, v)$  and such that  $\hat{j}$  is the *closest* arrival point greater than  $v$ . Let us call this  $c$ -supermaximal edge  $(\hat{i}_v, \hat{j}_v)$ . We use  $\hat{i}_v$  in next proposition.

**Proposition 4.1.5.** *Suppose that  $v \in G_{\mathcal{A}, T}$  is such that there exists a previous node  $i$  such that  $C((i, v)) = c$  and  $d(v) = d(i) + c$ . Then  $\hat{i}_v$  is a predecessor of  $v$ , i.e.  $d(v) = d(\hat{i}_v) + C((\hat{i}_v, v))$  and, moreover,  $d(\hat{i}_v) = d(i)$  and  $C((\hat{i}_v, v)) = c$ .*

*Proof.* Since  $(\hat{i}_v, \hat{j}_v)$  contains  $(i, v)$  and the dictionary at position  $\hat{i}_v$  is prefix closed then  $(\hat{i}_v, v)$  is an edge of  $G_{\mathcal{A}, T}$ . Since  $(\hat{i}_v, \hat{j}_v)$  has cost smaller than or equal to  $c$  then, by the suffix-nonincreasing property, also  $(\hat{i}_v, v)$  has cost smaller than or equal to  $c$ . Since the distance  $d$  is nondecreasing we know that  $d(\hat{i}_v) \leq d(i)$ . By very definition of the distance  $d$  we know that  $d(v) \leq d(\hat{i}_v) + C((\hat{i}_v, v))$ .

Putting all together we have that

$$d(v) \leq d(\hat{i}_v) + C((\hat{i}_v, v)) \leq d(i) + c = d(v).$$

Hence the inequalities in previous equation must be equalities and, furthermore,  $d(\hat{i}_v) = d(i)$  and  $C((\hat{i}_v, v)) = c$ .  $\square$

**Corollary 4.1.6.** *For any vertex  $v$ , the edge  $(\hat{i}_v, v)$  is the last edge of a path of minimal cost from the origin to vertex  $v$ .*

*Proof.* Any edge  $x$  in  $G_{\mathcal{A},T}$  such that  $d(v) = d(x) + C((x, v))$  is the last edge of a path of minimal cost from the origin to vertex  $v$ .  $\square$

*Remark 4.1.* Let us notice that the variable  $i$  is increased only at line 05 along the *Find Supermaximal* algorithm.

## 4.2 The Subgraph $G'_{\mathcal{A},T}$

In the following we describe a graph  $G'_{\mathcal{A},T}$  that is a subgraph of  $G_{\mathcal{A},T}$  and that is such that for any node  $v \in G_{\mathcal{A},T}$  there exists a minimal path from the origin to  $v$  in  $G'_{\mathcal{A},T}$  that is also a minimal path from the origin to  $v$  in  $G_{\mathcal{A},T}$ . The proof of this property, that will be stated in the subsequent proposition, is a consequence of Proposition 4.1.5 and Corollary 4.1.6.

We describe the building of  $G'_{\mathcal{A},T}$  in an algorithmic way.

The set of nodes of  $G'_{\mathcal{A},T}$  is the same of  $G_{\mathcal{A},T}$ . First of all we insert all the symbolwise edges of  $G_{\mathcal{A},T}$  in  $G'_{\mathcal{A},T}$ . Let now  $\mathcal{C}$  be the set of all possible costs that any dictionary edge has. This set can be built starting from  $G_{\mathcal{A},T}$ , but, in many meaningful cases, the set  $\mathcal{C}$  is usually well known and can be ordered and stored in an array in a time that is linear in the size of the text.

For any  $c \in \mathcal{C}$  we use algorithm *Find Supermaximal* to obtain the set  $M_c$ . Then, for any  $(i, j) \in M_c$ , we insert in  $G_{\mathcal{A},T}$  all the prefix of  $(i, j)$  except those which are contained in another  $c$ -supermaximal edge  $(i', j') \in M_c$ . In detail, for any  $c$ -supermaximal edge  $(i, j) \in M_c$ , let  $(i', j') \in M_c$  be the previous  $c$ -supermaximal edge overlapping  $(i, j)$ , i.e.  $j' = \max_h \{(s, h) \in M_c | i < h < j\}$ . Notice that this  $j'$  could not exist but, if it exists then by Proposition 4.1.2 there exists a unique  $i'$  such that  $(i', j') \in M_c$ . If  $(i', j')$  exists, then we add in  $G'_{\mathcal{A},T}$  all the edges of the form  $(i, x)$ , where  $j' < x \leq j$ , with label  $L_{(i,x)} = c$ . If  $(i', j')$  does not exist, then we add in  $G'_{\mathcal{A},T}$  all the edges of the form  $(i, x)$ , where  $i < x \leq j$ , with label  $L_{(i,x)} = c$ . In both cases, If such an edge  $(i, x)$  is already in  $G'_{\mathcal{A},T}$ , we just set the label  $L_{(i,x)}$  to  $\min\{L_{(i,x)}, c\}$ . This concludes the construction of  $G'_{\mathcal{A},T}$ .

The algorithm BUILD  $G'_{\mathcal{A},T}$  is formally stated in the Tabel 4.2.



BUILD $G'_{\mathcal{A},T}$
01. CREATE node 0
02. FOR $v = 1$ TO $ T $
03.     CREATE node $v$
04.     CREATE symbolwise edge $(v - 1, v)$
05. $L(v - 1, v) \leftarrow C((v - 1, v))$
06. ENDFOR
07. FOR ANY increasing $c \in \mathcal{C}$
08. $M_c \leftarrow \text{Find Supermaximal}(c)$
09. $j' \leftarrow 0$
10.     FOR ANY $(i, j) \in M_c$ left-to-right
11.         FOR ANY $x \mid \max\{j', i\} < x \leq j$
12.             IF $(i, x) \notin G'_{\mathcal{A},T}$ THEN
13.                 CREATE edge $(i, x)$
14. $L(i, x) \leftarrow c$
15.             ELSE
16. $L(i, x) \leftarrow \min\{L(i, x), c\}$
17.             ENDIF
18.         ENDFOR
19. $j' \leftarrow j$
20.     ENDFOR
21. ENDFOR

Table 4.2: The pseudocode of the Build  $G'_{\mathcal{A},T}$  algorithm.

*Remark 4.2.* Notice that for any cost  $c$  the above algorithm add in  $G'_{\mathcal{A},T}$  at most a linear number of edges. For any node in  $G'_{\mathcal{A},T}$  there is an incoming symbolwise edge and there also can be at most one incoming dictionary edge for any cost  $c$ .

Let us notice that the graph  $G'_{\mathcal{A},T}$  is a subgraph of  $G_{\mathcal{A},T}$ . Nodes and symbolwise edges are the same in both graphs by definition of  $G'_{\mathcal{A},T}$ . The edges  $(i, x)$  we add to  $G'_{\mathcal{A},T}$ , are the prefix of a  $c$ -supermaximal edge  $(i, j)$  of  $G_{\mathcal{A},T}$ . Since that the dictionary  $D_i$  is prefix closed, then all the edges  $(i, x)$  are also edges of  $G_{\mathcal{A},T}$ .

**Proposition 4.2.1.** *For any node  $v \in G_{\mathcal{A},T}$ , any minimal path from the origin to  $v$  in  $G'_{\mathcal{A},T}$  is also a minimal path from the origin to  $v$  in  $G_{\mathcal{A},T}$ .*

*Proof.* The proof is by induction on  $v$ . If  $v$  is the origin there is nothing to prove. Suppose now that  $v$  is greater than the origin and let  $(i, v)$  be the last edge of a minimal path in  $G_{\mathcal{A},T}$  from the origin to  $v$ . By inductive hypothesis there exists a minimal path  $\mathcal{P}$  from the origin to  $i$  in  $G'_{\mathcal{A},T}$  that is also a minimal path from the origin to  $i$  in  $G_{\mathcal{A},T}$ . If  $(i, v)$  is a symbolwise arc then it is also in  $G'_{\mathcal{A},T}$  and the concatenation of above minimal path  $\mathcal{P}$  with  $(i, v)$  is a minimal path from the origin to  $v$  in  $G'_{\mathcal{A},T}$  that is also a minimal path from the origin to  $v$  in  $G_{\mathcal{A},T}$ .

Suppose now that  $(i, v)$  is a dictionary arc and that its cost is  $c$ . Since it is the last edge of a minimal path we have that  $d(v) = d(i) + c$ . By Proposition 4.1.5  $d(v) = d(\hat{i}_v) + C((\hat{i}_v, v))$  and, moreover,  $d(\hat{i}_v) = d(i)$  and  $C((\hat{i}_v, v)) = c$ . By Corollary 4.1.6, the edge  $(\hat{i}_v, v)$  is the last edge of a path of minimal cost from the origin to vertex  $v$ . By inductive hypothesis there exists a minimal path  $\mathcal{P}$  from the origin to  $\hat{i}_v$  in  $G'_{\mathcal{A},T}$  that is also a minimal path from the origin to  $\hat{i}_v$  in  $G_{\mathcal{A},T}$ . Since  $(\hat{i}_v, v)$  has been added by construction in  $G'_{\mathcal{A},T}$ , the concatenation of above minimal path  $\mathcal{P}$  with  $(\hat{i}_v, v)$  is a minimal path from the origin to  $v$  in  $G'_{\mathcal{A},T}$  that is also a minimal path from the origin to  $v$  in  $G_{\mathcal{A},T}$ .  $\square$

Let us notice that it is possible to create the dictionary edges of  $G'_{\mathcal{A},T}$  without an explicit representation in memory of all the  $M_c$  arrays. This is just an implementation detail that enhances speed and the usage of memory of the Build  $G'_{\mathcal{A},T}$  algorithm in practice, without changing its order of complexity. The point is that we can insert the  $c$ -*supermaximal* edges and their prefix directly in the graph as soon as they are found along a *Find Supermaximal* execution. The correctness of this approach is a direct consequence of the following remark.

*Remark 4.3.* Given a cost  $c$ , the edges  $(i, x)$  used by the Build  $G'_{\mathcal{A},T}$  algorithm inside the block at lines 10 – 20 are those which the *Weighted-Extend* and the *Weighted-Exist* operations of the *Find Supermaximal(c)* algorithm report a positive answer for.

### 4.3 The Dictionary-Symbolwise Flexible Parsing

We can now finally describe the *Dictionary-symbolwise flexible parsing*.

The *Dictionary-symbolwise flexible parsing* algorithm firstly uses the algorithm BUILD  $G'_{\mathcal{A},T}$  and then uses the classical SINGLE SOURCE SHORTEST PATH (SSSP) algorithm (see [6, Ch. 24.2]) to recover a minimal path from the origin to the end of graph  $G_{\mathcal{A},T}$ . The correctness of the above algorithm is stated in the following theorem and it follows from the above description and from Proposition 4.2.1.

**Theorem 4.3.1.** *Dictionary-symbolwise flexible parsing is graph optimal.*

Notice that graphs  $G_{\mathcal{A},T}$  and  $G'_{\mathcal{A},T}$  are directed acyclic graphs (DAG) and their nodes from 1 to  $n$ , where 1 is the origin or the unique source of the graph and  $n = |T|$  is the last node, are topologically ordered and linked by symbolwise edges. Recall that, given a node  $v$  in a weighted DAG, the classic solution to the SSSP is composed by two steps. The first one computes the distance and a predecessor of any node in the graph. It is accomplished by performing a visit on all the nodes in topological order and making a *relax* on any outgoing edge. Therefore, for any node  $v$  from 1 to  $n$  and for any edge  $(v, v')$  in the graph, the *relax* of  $(v, v')$  sets the distance and the predecessor of  $v'$  to  $v$  if  $d(v) + C((v, v')) < d(v')$ . The classic algorithm uses two arrays,  $\pi[]$  and  $p[]$ , to store distance and predecessor of nodes. The second step recovers the shortest path by following backward the predecessors chain from the last node to the origin of the graph and reverting it. From this simple analysis it follows that if we know all the outgoing edges of any node in topological order then we can do directly the *relax* operation on them without having an explicit representation of the graph.

Let us suppose to have an *online* version of the BUILD  $G'_{\mathcal{A},T}$  algorithm, where for any  $i$  from 1 to  $|T|$ , only edges  $(i, j)$  are created on the graph. We want now to merge the *online* BUILD  $G'_{\mathcal{A},T}$  algorithm to the *relax* step of the SSSP algorithm. We maintain the two arrays  $\pi[]$  and  $p[]$  of linear size w.r.t. the text size, containing the distance and the predecessor of any node and we replace any edge creation or label updating with the *relax* operation.

About the *online* version of the BUILD  $G'_{\mathcal{A},T}$ , we can use the Remark 4.1 to make a kind of parallel run of the *Find Supermaximal* algorithm

for any cost  $c$ , maintaining the variables  $i$  synchronized on the same value. Moreover, we use the Remark 4.3 to directly handle the edge creation as soon as they are found. We address all of these variations to the BUILD  $G'_{\mathcal{A},T}$ , the *Find Supermaximal* as well as the merge with the SSSP algorithm in order to obtain the DICTIONARY-SYMBOLWISE FLEXIBLE PARSING algorithm. The pseudocode of the DICTIONARY-SYMBOLWISE FLEXIBLE PARSING algorithm is reported in Table 4.3.

Let us notice that above algorithm uses only one dictionary at one time and it never needs to use previous version of the dynamic dictionary. Recall that the dictionary is used by the *Weighted-Exist* and the *Weighted-Extend* operations. This is a direct consequence of the fact that any edge  $(i, j)$  refers to the dictionary  $D_i$  and that after edge  $(i, j)$  creation, only edge  $(p, q)$  with  $p \geq i$  can be created.

**Proposition 4.3.2.** *Any practical implementations of the Dictionary-symbolwise flexible parsing does not require to explicitly represent the graph  $G'_{\mathcal{A},T}$  regardless of its size. Since  $G'_{\mathcal{A},T}$  nodes are visited in topological order by classic SSSP solutions, the algorithm needs to maintain just two linear size arrays, i.e. the array of node distances and the array of node predecessors, in order to correctly compute an optimal parsing.*

Let us summarize the Dictionary-Symbolwise Flexible Parsing algorithm requirements. Given a text  $T$  of size  $n$  the Dictionary-Symbolwise Flexible Parsing algorithm uses

- $O(n)$  space for the  $\pi[]$  and  $p[]$  arrays, regardless of the graph  $G'_{\mathcal{A},T}$  size that is not really built, plus the dictionary structure.
- $O(|E|)$  time to analyze all the edges of the graph  $G'_{\mathcal{A},T}$ .
- it is not *online* because the backward recovering of the parsing from the  $p[]$  array.

With respect to the original Flexible Parsing algorithm we gain the fact that it can work with variable costs of pointers and that it is extended to the dictionary-symbolwise case. This covers for instance the LZW-like and

<i>Dictionary-Symbolwise Flexible Parsing</i>	
01.	FOR $i$ FROM 0 TO $ T  - 1$
02.	<i>Relax</i> ( $i, i + 1, C((i, i + 1))$ )
03.	FOR ANY $c \in \mathcal{C}$
04.	IF $i = j_c$ THEN
05.	$j_c \leftarrow 1 + j_c$
06.	ENDIF
07.	IF $j_c \leq  T $ AND <i>Weighted-Exist</i> ( $i, j_c, c$ ) = “yes” THEN
08.	<i>Relax</i> ( $i, j_c, C((i, j_c))$ )
09.	WHILE <i>Weighted-Extend</i> ( $(i, j), c$ ) = “yes” DO
10.	$j_c \leftarrow 1 + j_c$
11.	<i>Relax</i> ( $i, j_c, C((i, j_c))$ )
12.	ENDWHILE
13.	$j_c \leftarrow 1 + j_c$
14.	ENDIF
15.	ENDFOR
16.	ENDFOR
17.	RETURN <i>Reverse</i> ( $v$ )
 <i>Relax</i> ( $u, v, c$ )	
01.	IF $\pi[u] + c < \pi[v]$ THEN
02.	$\pi[v] \leftarrow \pi[u] + c$
03.	$p[v] \leftarrow u$
04.	ENDIF
 <i>Reverse</i> ( $v$ )	
01.	IF $v > 0$ THEN
02.	<i>Reverse</i> ( $p[v]$ )
03.	ENDIF
04.	RETURN $v$

Table 4.3: The pseudocode of *Dictionary-Symbolwise Flexible Parsing* algorithm, the *Relax* and the *Reverse* procedures. The distance array  $\pi[]$  and the predecessor array  $p[]$  are initialized to 0. Notice that the algorithm uses a different  $j_c$  variable for any  $c$  value.

the LZ77-like cases. But we lose the fact that the original one was “on-line”. A minimal path has to be recovered, starting from the end of the graph backward. But this is an intrinsic problem that cannot be eliminated. Even if the dictionary edges have just one possible cost, in the dictionary-symbolwise case it is possible that any minimal path for a text  $T$  is totally different from any minimal path for the text  $Ta$ , that is the previous text  $T$  concatenated to the symbol  $a$ . The same can happen when we have a (pure) dictionary case with variable costs of dictionary pointers. In both cases, for this reason, it is unlikely that there exists an “on-line” optimal parsing algorithm, and, indeed, the original flexible parsing fails being optimal in the dictionary case when costs are variable.

On the other hand our algorithm is suitable when the text is divided in several contiguous blocks and, therefore, in practice there is not the need to process the whole text but it suffices to end the current block in order to have the optimal parsing (relating to that block).

## 4.4 Time and Space Analyses

In this section we analyze the *Dictionary-symbolwise flexible parsing* in both LZ78 and LZ77-like algorithm versions.

### LZ78 Case

Concerning LZ78-like algorithms, the dictionary is prefix closed and it is implemented by using the LZW variant. We do not enter into the details of this technique. We just recall that the cost of pointers increases by one unit whenever the dictionary size is “close” to a power of 2. The moment when the cost of pointers increases is clear to both encoder and decoder. In our dictionary-symbolwise setting, we suppose that the flag information has a constant cost. We assume therefore that it takes  $O(1)$  time to determine the cost of a dictionary edge.

The maximal cost that a pointer can assume is smaller than  $\log_2(n)$  where  $n$  is the text size. Therefore the set  $\mathcal{C}$  of all possible costs of dictionary edges has a logarithmic size and it is cheap to calculate.

In [27] the operations *Extend* and *Contract* are presented. It is also presented a linear size data structure called trie-reverse-trie-pair that allows to execute both those operations in  $O(1)$  time. The operation  $Extend(w, a)$  says whether the phrase  $wa$  is in the currently used dictionary. The operation  $Contract(w)$  says whether the phrase  $w[2 : |w|]$  is in the current dictionary.

Since at any position we can calculate in  $O(1)$  time the cost of an edge, we can use the same data structure to perform our operations of *Weighted-Extend* and of *Weighted-Exist* in constant time as follows. In order to perform a  $Weighted-Extend((i, j), c)$  we simply execute the operation  $Extend(w, a_{j+1})$  with  $w = T[i + 1 : j]$ , i.e. the phrase associated to the edge  $(i, j)$ , and then, if the answer is “yes”, we perform a further check in  $O(1)$  time on the cost of the found edge  $(i, j + 1)$ . Therefore,  $Weighted-Extend((i, j), c)$  is equal to  $Extend(T[i + 1 : j], a_{j+1}) \text{ AND } C((i, j + 1)) \leq c$ .

In order to perform a  $Weighted-Exist((i, j), c)$  we simply use the *contract* on the phrase  $a_i w$ , where  $w = T[i + 1 : j]$ , and, if the answer is “yes” we perform a further check in  $O(1)$  time on the cost of the found edge  $(i, j)$ . Therefore,  $Weighted-Exist(i, j, c)$  is equal to  $Contract(a_i T[i + 1 : j]) \text{ AND } C((i, j)) \leq c$ .

At first glance, the algorithm  $BUILD\ G'_{\mathcal{A}, T}$  seems to take  $O(n \log n)$  time. But, since there is only one active cost at any position in any LZW-like algorithms, then if  $c < c'$  then  $M_c \subseteq M_{c'}$ , as stated in the following proposition.

**Definition 4.5.** We say that a cost function  $C$  is *LZW-like* if for any  $i$  the cost of all dictionary pointers in  $D_i$  is a constant  $c_i$  and that for any  $i$ ,  $0 \leq i < n$  one has that  $c_i \leq c_{i+1}$ .

**Proposition 4.4.1.** *If the cost function  $C$  is LZW-like, one has that if  $c < c'$  then  $M_c \subseteq M_{c'}$ .*

*Proof.* We have to prove that for any  $(i, j) \in M_c$  then  $(i, j) \in M_{c'}$ . Clearly if  $(i, j) \in M_c$  then its cost is smaller than or equal to  $c < c'$ . It remains to prove that  $(i, j)$  is  $c'$ -supermaximal, e.g. that  $\forall p, q \in V$ , with  $p < i$  and  $j < q$ , the arcs  $(p, j), (i, q)$  are not in  $E_{c'}$ . Since  $(i, j) \in M_c$  and since the cost of  $(i, j)$  is by hypothesis equal to  $c_i$ , we have that  $c_i \leq c$ . If arc  $(p, j)$  is in

$E_{c'}$  then its cost is  $c_p \leq c_i \leq c$  and therefore it is also in  $E_c$  contradicting the  $c$ -supermaximality of  $(i, j)$ . If arc  $(i, q)$  is in  $E_{c'}$  then its cost is  $c_i \leq c$  and therefore it is also in  $E_c$  contradicting the  $c$ -supermaximality of  $(i, j)$ .  $\square$

At this point, in order to build  $G'_{\mathcal{A},T}$  we proceed in an incremental way. We build  $M_c$  for the smallest cost. Then, we start from the last built  $M_c$  to build  $M_{c'}$ , where  $c'$  is the smallest cost greater than  $c$ . And so on until all the costs are examined. We insert any edge  $(i, j)$  only in the set  $M_c$  where  $c$  is the real cost of the  $(i, j)$  edge. In this way, we avoid to insert the same edge  $(i, j)$  in more than one  $M_c$  since that the algorithm will insert eventually the edge  $(i, j)$  from the set  $M_c$  with the minimal cost  $c = C((i, j))$ .

A direct consequence of above approach, we have that only a linear number of edges are inserted in the graph  $G'_{\mathcal{A},T}$ .

The overall time for building  $G'_{\mathcal{A},T}$  is therefore linear, as well as its size. The SINGLE SOURCE SHORTEST PATH over  $G'_{\mathcal{A},T}$ , that is a DAG topologically ordered, takes linear time (see [6, Ch. 24.2]).

In conclusion we state the following proposition.

**Proposition 4.4.2.** *Suppose that we have a dictionary-symbolwise scheme, where the dictionary is LZ78-like and the cost function is LZW-like. The symbolwise compressor is supposed to be, as usual, linear time. Using the trie-reverse-trie-pair data structure, Dictionary-Symbolwise flexible parsing is linear.*

## LZ77 Case

Concerning LZ77, since the dictionary has the weak prefix-closed property, i.e. the dictionary is always prefix closed, we have that the Dictionary-Symbolwise Flexible Parsing is an optimal parsing. We exploit the discreteness of the cost function  $C$  when it is associated to the length of the codewords of a variable-length code, like Elias codes or a Huffman code, to bound the cardinality of the set  $\mathcal{C}$  to  $O(\log n)$ . Indeed let us call  $\hat{c}$  the maximum cost of any dictionary pointer, i.e. the length in bits of the pointer with the maximum possible offset and the maximum possible length.

Even if the cost actually depends on the text  $T$ , it usually has an upper



bound that depends on the encoding and on the dictionary constrains and we can assume it being  $\hat{c} = O(\log n)$ , with  $|T| = n$ .

Operations *Weighted-Exist* and *Weighted-Extend* can be implemented in linear space and constant time by using classical suffix tree or other solutions when the dictionary is a LZ77-like one. For instance, in [10] it is shown how to compute the Longest Previous Factor (LPF) array in linear time. Recall that  $T[i : LPF[i]]$  is the longest factor already seen in the text at some position  $i' < i$ . It is easy to see that following relations hold. The operation *Weighted-Exist*  $(i, j, c)$  outputs “yes”  $\iff j \leq i + LPF[i]$  AND  $C((i, j)) \leq c$  and the operation *Weighted-Extend*  $((i, j), c)$  outputs “yes”  $\iff j < i + LPF[i]$  AND  $C((i, j + 1)) \leq c$ . We recall that we are also assuming that it is possible to compute the cost of a given edge in constant time. Therefore, we use linear time and space to build the LPF array and then any operation *Weighted-Exist* or *Weighted-Extend* take just constant time.

Suppose to have a dictionary-symbolwise scheme, where the dictionary is LZ77-like and the dictionary pointer encoding, the symbolwise encoding and the flag information encoding are any variable-length encoding one. The use of the codeword length as cost function leads to a function that assumes integer values. Given  $\hat{c}$  the maximum cost of any dictionary pointer with  $\hat{c} \leq \log(n)$ , the Dictionary-Symbolwise Flexible Parsing runs in  $O(n \log n)$  time and space.

Let us notice that in most of the common LZ77 dictionary implementation, as it is in the *deflate* compression tool, our assumption about the computation of edge cost in  $O(1)$  time is not trivial to obtain.

Obviously, we are interested, for compression purpose, to the smallest cost between all the possible encoding of a phrase. For instance, the use of the length-distance pair as dictionary pointer leads to multiple representation of the same (dictionary) phrase since this phrase can occur more than once in the (already seen) text.

Since the closest occurrence uses the smallest distance to be represented, the cost of encoding the phrase using this distance is usually the smallest one, accordingly to the used encoding method.

A practical approach that looks for the above smallest distance makes use of hash tables, built on fixed length phrases.

A new data structure able to answer to the edge cost query in constant time and able to support the *Weighted-Exist* and the *Weighted-Extend* operation in the case of LZ-77 dictionaries will be introduced in next chapter.



## Chapter 5

# The Multilayer Suffix Tree

We introduce here an online full-text index data structure that is able to find the rightmost occurrence of any factor or an occurrence which bit representation has equal length (Query 1). It has linear space complexity and it is built in  $O(n \log n)$  amortized time, where  $n$  is the size of the text. It is able to answer to the Query 1, given a pattern  $w$ , in  $O(|w| \log \log n)$ . Furthermore, we will show how to use this structure to support the Weighted-Exist and the Weighted-Extend operations used by the *Dictionary-Symbolwise Flexible Parsing* algorithm in  $O(1)$  amortized time.

### 5.1 Preliminary Definitions

Let  $Pos(w) \subset \mathbb{N}$  the set of all the occurrences of  $w \in Fact(T)$  in the text  $T \in \Sigma^*$ , where  $Fact(T)$  is the set of the factors of  $T$ . Let  $Offset(w) \subset \mathbb{N}$  the set of all the occurrence offsets of  $w \in Fact(T)$  in the text  $T$ , i.e.  $x \in Offset(w)$  iff  $x$  is the distance between the position of an occurrence of  $w$  and the end of the text  $T$ . For instance, given the text  $T = \text{babcbabbabbb}$  of length  $|T| = 12$  and the factor  $w = \text{abb}$  of length  $|w| = 3$ , the set of positions of  $w$  over  $T$  is  $Pos(w) = \{4, 9\}$ . The set of the offsets of  $w$  over  $T$  is  $Offset(w) = \{7, 2\}$ . Notice that  $x \in Offset(w)$  iff exists  $y \in Pos(w)$  such that  $x = |T| - y - 1$ . Since the offsets are function of occurrence positions, there is a bijection between  $Pos(w)$  and  $Offset(w)$ , for any factor  $w$ .

Given a number encoding method, let  $Bitlen : \mathbb{N} \rightarrow \mathbb{N}$  a function that

associates to a number  $x$  the length in bit of the encoding of  $x$ . Let us consider the equivalence relation *having equal codeword bit-length* on the set  $Offset(w)$ . The numbers  $x, y \in Offset(w)$  are bit-length equivalent *iff*  $Bitlen(x) = Bitlen(y)$ . Let us notice that the *having equal codeword bit-length* relation induces a partition on  $Offset(w)$ .

**Definition 5.1.** The *rightmost* occurrence of  $w$  over  $T$  is the offset of the occurrence of  $w$  that appears closest to the end of the text, if  $w$  appears at least once over  $T$ , otherwise it is not defined.

Notice that for any factor  $w \in Fact(T)$ , the rightmost offset of  $w$  is defined as follows.

$$rightmost(w) = \begin{cases} \min\{x \mid x \in Offset(w)\} & \text{if } Offset(w) \neq \emptyset \\ \text{not defined} & \text{if } Offset(w) = \emptyset \end{cases}$$

Let us notice that referring to the rightmost occurrence of a factor in an *online* algorithmic fashion, where the input text is processed left to right, corresponds to referring to the rightmost occurrence over the text *already processed*. Indeed, if at a certain algorithm step we have processed the first  $i$  symbols of the text, the rightmost occurrence of  $w$  is the occurrence of  $w$  closest to the position  $i$  of the text reading left to right.

**Definition 5.2.** Let  $rightmost_i(w)$  be the rightmost occurrence of  $w$  over  $T_i$ , where  $T_i$  is the prefix of the text  $T$  ending at the position  $i$  in  $T$ . Obviously,  $rightmost_n(w) = rightmost(w)$  for  $|T| = n$ .

In many practical algorithms, like in the data compression field, the text we are able to refer to is just a portion of the whole text. Let  $T[j : i]$  be the factor of the text  $T$  starting from the position  $j$  and ending to the position  $i$ . We generalize the definition of  $rightmost(w)$  over a factor  $T[j : i]$  of  $T$  as follows.

**Definition 5.3.** Let  $rightmost_{j,i}(w)$  be the rightmost occurrence of  $w$  over  $T[j : i]$ , where  $T[j : i]$  is the factor of the text  $T$  starting at the position  $j$  and ending at the position  $i$  of length  $i - j + 1$ . Obviously,  $rightmost_{1,n}(w) = rightmost(w)$  for  $|T| = n$ .

The online full-text index we are going to introduce is able to answer to the *rightmost equivalent length* query in constant time, also referred hereby as *Query 1*. The Query 1 is more formally stated as below.

**Definition 5.4. (Query 1)** Let us suppose that we have a text  $T \in \Sigma^*$ , a pattern  $w \in \Sigma^*$  and a point  $i$  in time. Assume that the prefix  $T_i$  of the text has been read. If  $w$  appears at least once in  $T_i$ , then the *rightmost equivalent length* query provides an occurrence of  $w$  which offset is in  $[rightmost_i(w)]$ , where  $[rightmost_i(w)]$  is the equivalence class induced by the relation *having equal codeword bit-length* containing  $rightmost_i(w)$ . Otherwise the *rightmost equivalent length* query provides the value *nil*.

## 5.2 The Data Structure

There are many classic full-text index able to represent the set  $Fact(T[1 : i])$ , like the suffix tree, the suffix array, the suffix automaton and others. Many of them can easily be preprocessed in the *offline* fashion to make them able to find efficiently the rightmost occurrence of any factor over the whole text, but none of them can directly answer to the above Query 1.

The main idea of this new data structure, is based on two observations. The first one is that the equivalence relation *having equal codeword bit-length* that induces a partition on  $Offset(w)$ , for any  $w$ , also induces a partition on the set of all the possible offsets over a text  $T$  independently of a specific factor, i.e. on the set  $[1..|T|]$ . The second observation is that for any encoding method for which the *Bitlen* function is a monotonic function, each equivalence class in  $[1..|T|]$  is composed by contiguous points in  $[1..|T|]$ . Indeed, given a point  $p \in [1..|T|]$ , the equivalence class  $[p]$  is equal to the set  $[j..i]$ , with  $j \leq p \leq i$ ,  $j = \min\{x \in [p]\}$  and  $i = \max\{x \in [p]\}$ .

Putting these observations all together suggests that the Query 1 can be addressed by a set of classic full text indexes, each one devoted to one or more classes of equivalence of the relation *having equal code bit-length*.

Fixed an encoding method for numbers and a text  $T$ , we assume that *Bitlen* is a monotonic function and that we know the set  $B = \{b_1, b_2, \dots, b_s\}$ , with  $b_1 < b_2 < \dots < b_s$ , that is the set of the *Bitlen* values of all the

possible offsets over  $T$ . We define the set  $SW = \{sw_1, sw_2, \dots, sw_s\}$ , with  $sw_1 < sw_2 < \dots < sw_s$ , where  $sw_i$  is the greatest integer (smaller than or equal to the length of the text  $T$ ) such that  $Bitlen(j)$  is less than or equal to  $b_i$ . More formally,  $sw_i = \max\{j \leq |T| \mid Bitlen(j) \leq b_i\}$ . Notice that  $sw_s = |T|$ .

All the known standard non-unary representation of numbers satisfy the following property.

**Property 5.1.** *There exists a constant  $k > 1$  and an integer  $\hat{k}$  such that for any  $\hat{k} \leq i < s$  one has  $sw_i \geq k sw_{i-1}$*

*Remark 5.1.* A similar result can be stated about the lengths. Let  $B' = \{b'_1, b'_2, \dots, b'_j\}$ , with  $b'_1 < b'_2 < \dots < b'_j$ , that is the set of the *Bitlen* values of all the possible lengths over  $T$ . Notice that if we use the same encoding for the lengths and the offsets then  $B' = B$  and  $j = s$ . Analogously we can define the set  $SW'$ . In what follows within this thesis we will use the same encoding for the lengths and the offsets and therefore the set  $SW'$  is equal to  $SW$  and previous property holds for the lengths too with the same constants  $\hat{k}$  and  $k$ .

Let us consider an example.

The Table 5.1 reports the Elias  $\gamma$  codes and the *Bitlen* values for integers from 1 to 18. Suppose, for instance, that our cost function is associated to this Elias codes and that we have a text  $T$  of length 18. The set  $B$  therefore is  $B = \{b_1 = 1, b_2 = 3, b_3 = 5, b_4 = 7, b_5 = 9\}$  and we have that  $sw_1 = 1$ ,  $sw_2 = 3$ ,  $sw_3 = 7$ ,  $sw_4 = 15$  and  $sw_5 = 18$ . Notice that, indeed, Property 5.1 is satisfied for  $\hat{k} = 2$  and  $k = 2$ .

Let us now introduce the Multilayer Suffix Tree data structure.

We suppose that a text  $T$  is provided *online* and at time  $i$  the first  $i$  characters of  $T$  have been read, i.e. at time  $i$  the prefix  $T_i$  of length  $i$  of the text  $T$  has been read.

**Definition 5.5.** The Multilayer Suffix Tree is a data structure composed by the set  $S = \{S_{sw_1}, S_{sw_2}, \dots, S_{sw_s}\}$  of suffix trees where, for any  $\alpha \in SW$  and at any moment  $i$ ,  $S_\alpha$  is the suffix tree for sliding window of  $T_i$  with sliding

$i$	$\gamma(i)$	$Bitlen(\gamma(i))$
1	1	1
2	010	3
3	011	3
4	00100	5
5	00101	5
6	00110	5
7	00111	5
8	0001000	7
9	0001001	7
10	0001010	7
11	0001011	7
12	0001100	7
13	0001101	7
14	0001110	7
15	0001111	7
16	000010000	9
17	000010001	9
18	000010010	9

Table 5.1: Elias  $\gamma$  code for integers from 1 to 18 and their *Bitlen* value.

window of size  $\alpha$  such that  $S_\alpha$  represents all the factors of  $T[i - \alpha : i]$ . We call  $S_\alpha$  simply the layer  $\alpha$  or the layer of size  $\alpha$ .

From now on we will refer to suffix trees or layers indifferently.

We use a *online* suffix tree for sliding window introduced by Larson in [26] and later refined by Senft in [34], in order to represent each layer of our multilayer suffix tree. Therefore, for any  $S_\alpha \in S = \{S_{sw_1}, S_{sw_2}, \dots, S_{sw_s}\}$ ,  $S_\alpha$  is a full-text index for  $T[i - \alpha : i]$ , where  $T$  is a given text and  $i$  is a point in time.

We think that it is possible to adapt our data structure to work with other classic indexes for sliding window (see for instance [22, 29, 35]).

Let us recall that in [26] an *online* and linear time construction for the



suffix tree is reported. The suffix tree uses linear space w.r.t. the text size, i.e.  $S_\alpha$  uses  $O(\alpha)$  space. Moreover, any common full-text index allows to find an occurrence of a given pattern in constant time, but cannot be directly used to answer to the Query 1.

**Proposition 5.2.1.** 1. *If a pattern  $w$  is in layer  $\alpha$  with  $\alpha \in SW$ , then  $w$  is also in layer  $\beta$  for any  $\beta \in SW$  with  $\alpha \leq \beta$ .* 2. *If a pattern  $w$  is not in a layer  $\alpha$ ,  $\alpha \in SW$ , then  $w$  is not in layer  $\beta$  with  $\beta \leq \alpha$*

*Proof.* The proof of the property at point 1 comes immediately from suffix trees properties. Since layer  $\alpha$  is a full text index for  $T[i - \alpha : i]$  and layer  $\beta$  is a full-text index for  $T[i - \beta : i]$ , for any  $\alpha, \beta \in SW$  with  $\alpha \leq \beta$  and for any  $i$ ,  $T[i - \alpha : i]$  is a suffix of  $T[i - \beta : i]$ . The property at point 2 can be deduced by point 1.  $\square$

**Proposition 5.2.2.** *Fixed a text  $T$  of size  $|T| = n$ , at any moment  $i$  with  $0 \leq i \leq n$  and for any standard variable-length code, the multilayer suffix tree uses  $O(i)$  space.*

*Proof.* Since at time  $i$  the maximum offset of all the occurrences in the text  $T_i$  is  $O(i)$ , for any standard variable-length code the maximum value of the set  $SW$  is  $O(i)$ . Since Property 5.1 holds for the set  $SW$  and since the multilayer suffix tree space is equal to the sum of the space of its layers, as an immediate consequence we have that space used by the multilayer suffix tree is  $\sum_{\alpha \in SW} \alpha = O(i)$ .  $\square$

From the linear time of the *online* construction of the suffix tree for sliding window and since the number of layers is  $|SW| = O(\log |T|)$ , we can immediately state the following proposition.

**Proposition 5.2.3.** *Given a text  $T$  of length  $|T| = n$ , for any standard variable-length code, it is possible to build online the multilayer suffix tree in  $O(n \log n)$  amortized time.*

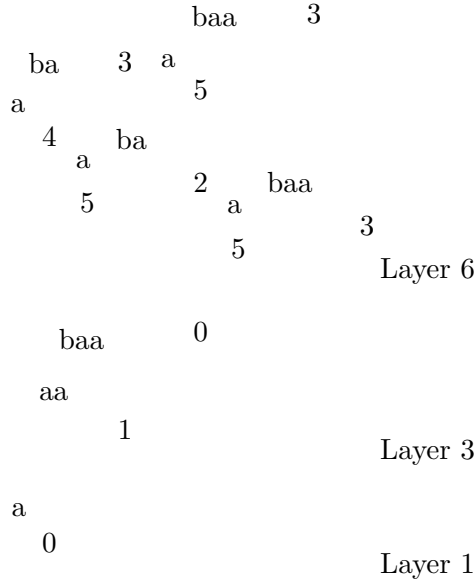


Figure 5.1: The Multilayer Suffix Tree for the text  $T = \text{ababaa}$  and for the Elias  $\gamma$ -code, where  $B = \{1, 3, 5\}$ ,  $SW = \{1, 3, 6\}$ . The solid edges are regular ones and the dashed links are the suffix-links of internal nodes. For convenience, we added edge labels with the substring of  $T$  associated to edges. The node value is the position over the text of the incoming edge. We omitted the string depth of nodes. Let consider, for instance, the phrase  $w = \text{“ba”}$  with  $Pos(w) = \{1, 3\}$ ,  $Offset(w) = \{4, 2\}$ ,  $rightmost(w) = 2$  and  $\gamma(2) = 010$ . Since  $w$  is in layer  $sw_2 = 3$  and is not in layer  $sw_1 = 1$ , we have that  $Bitlen(rightmost(w)) = 3$  is equal to  $Bitlen(sw_2 = 3) = b_2 = 3$ .

### 5.3 On the Query 1

We want now to show how to answer to the Query 1 for a given pattern in  $O(|pattern| \log \log n)$ , with  $n$  equal to the length of the text.

**Proposition 5.3.1.** *If a pattern  $w$  is in layer  $\beta$  and is not in layer  $\alpha$ , where  $\alpha$  is the maximum of the values in  $SW$  smaller than  $\beta$ , then any occurrence of  $w$  in the layer  $\beta$  correctly answer to the **Query 1**.*

*Proof.* If a pattern  $w$  is in layer  $\beta$  and is not in layer  $\alpha$ , where  $\alpha$  is the

maximum of the values in  $SW$  smaller than  $\beta$ , then from Prop. 5.2.1, it follows that  $\beta$  is the smallest layer where it appears  $w$ . Therefore  $w$  has at least one occurrence in  $T[i - \beta : i - \alpha - 1]$ , i.e.  $rightmost_i(w) \in (\alpha..\beta]$ .

Since  $(\alpha..\beta]$  is the equivalence class  $[\beta]$  of the *having equal codeword bit-length* relation, we have that any occurrence of  $w$  in the layer  $\beta$  correctly answer to the **Query 1**.  $\square$

*Remark 5.2.* Let us notice that if  $S_{sw_x}$  is the smallest layer where  $rightmost_i(w)$  appears, then the *Bitlen* value of the offset of  $rightmost_i(w)$  is equal to  $b_x$ .

Using above proposition, we are able to answer to the Query 1 once we find the smallest layer containing the rightmost occurrence of the pattern, if any, otherwise we just report *nil*.

What follows is the trivial search of the smallest layer that contains an occurrence of a given pattern.

Given a pattern  $w$  at time  $i$ , we look for  $w$  in the layer  $sw_1$ , i.e. the smallest layer. If  $w$  is in  $S_{sw_1}$ , then all the occurrences of  $w$  in  $T[i - sw_1 : i]$  belong to the class of the rightmost occurrence of  $w$  over  $T$ . If  $w$  is not in  $S_{sw_1}$ , then we look for any occurrence of  $w$  in  $S_{sw_2}$ , the next layer in increasing order. If  $w$  is in  $S_{sw_2}$ , since it is not in  $S_{sw_1}$ , for the Prop. 5.3.1, any occurrence of  $w$  in this layer belong to the rightmost occurrence of  $w$  over  $T_i$ . Continuing in this way, as soon as we found an occurrence of  $w$  in a layer, this occurrence correctly answer to the Query 1.

**Proposition 5.3.2.** *Using the trivial search of a given pattern in the layers on the multilayer suffix tree we are able to answer to the Query 1 in time proportional to the pattern size times the cardinality of the set  $SW$ .*

Since many of the classic variable-length codes for integers, like the Elias's  $\gamma$ -codes, produce codewords of length proportional to the logarithm of the represented value, we can assume that the cardinality of  $SW$  is  $O(\log |T|)$ . Since that  $|T_i| = i$ , in the *online* fashion, we have that above proposition becomes as follows.

**Proposition 5.3.3.** *Using any classic variable-length code method, the above data structure is able to answer to the Query 1 in  $O(|\text{pattern}| \log i)$  time.*

A similar result can be obtained by using a variant of the Amir et al. algorithm presented [1], but it does not support *Weighted-Exist* and *Weighted-Extend* operations in constant time.

Since Prop. 5.2.1 holds for the layers of our data structure, we can use the binary search to find the smallest layer containing a given pattern. Since for any classic variable-length code  $|SW| = O(\log i)$ , the number of layers in our structure is  $O(\log i)$  and the proof of following proposition comes straightforward.

**Proposition 5.3.4.** *Using any classic variable-length code, at any time  $i$  the multilayer suffix tree is able to answer to the Query 1 for a given pattern in  $O(|\text{pattern}| \log \log i)$  time.*

## 5.4 On *Weighted-Exist* and *Weighted-Extend*

Let us now focus on the *Weighted-Exist* and on the *Weighted-Extend* operations used by the Dictionary-Symbolwise Flexible Parsing algorithm. We want to show how to use the multilayer suffix tree to support the *Weighted-Exist* and the *Weighted-Extend* operations in amortized constant time when the dictionary is LZ77-like. For simplicity in the following we focus on the case of LZ77 dictionary with unbounded size. The case of LZ77 dictionary with bounded size, like in the original LZ77 case, is a simple generalization of the unbounded case that comes from the sliding window capability of the layers of the multilayer suffix tree. Indeed, the unbounded case is the special case of the bounded dictionary where the search buffer size is greater than or equal to the text size.

We assume that the Dictionary-Symbolwise Flexible Parsing has to parse the text  $T$  and that at any time  $i$  with  $0 \leq i \leq |T|$ , the dynamic dictionary  $D_i$  is represented by the multilayer suffix tree of  $T_i$ .

Let us consider a run of the *Find Supermaximal* algorithm described in Section 4.1 and reported for convenience in Table 5.4, where  $c$  is a fixed value.

Recall that the cost of any dictionary phrase  $w \in D_i$  associated to the edge  $(i, j) \in G_{\mathcal{A}, T}$  is given by the cost of the flag information  $F_d$  for the dic-

<pre> <i>Find Supermaximal</i>(<i>c</i>) 01.  <math>i \leftarrow 0, j \leftarrow 0, M_c \leftarrow \emptyset</math> 02.  WHILE <math>j &lt; n</math> DO 03.    <math>j \leftarrow j + 1</math> 04.    WHILE <math>i &lt; j</math> AND <i>Weighted-Exist</i>(<math>i, j, c</math>) = “no” DO 05.      <math>i \leftarrow i + 1</math> 06.    ENDWHILE 07.    IF <math>i &lt; j</math> THEN 08.      WHILE <i>Weighted-Extend</i> (<math>(i, j), c</math>) = “yes” DO 09.        <math>j \leftarrow j + 1</math> 10.      ENDWHILE 11.      INSERT (<math>(i, j), M_c</math>) 12.    ENDIF 13.  ENDWHILE 14.  RETURN <math>M_C</math> </pre>
--

Table 5.2: The pseudocode of the *Find Supermaximal* algorithm as in the Table 4.1 of Page 52.

tionary phrases plus the cost of the encoding of the dictionary pointer. Let us consider  $F_d$  constant within this section. Moreover, since any dictionary pointer in LZ77-like algorithms is encoded by the couple (*length,offset*) as described in Section 2.2, the cost of the encoding of a dictionary pointer is equal to the cost of the encoding of the length of the dictionary phrase plus the cost of the encoding of the offset of the rightmost occurrence of the phrase inside the dictionary.

Let us recall that for any cost value  $c$  and for any phrase  $w \in D_i$  associated to the edge  $(i, j)$  in  $G_{\mathcal{A},T}$  with  $w = T[i + 1 : j]$ , the operation *Weighted-Extend*( $(i, j), c$ ) finds out whether the phrase  $wa = T[i + 1 : j + 1]$  is in  $D_i$  and the cost of the edge  $(i, j + 1)$  in  $G_{\mathcal{A},T}$  is smaller than or equal to  $c$ . Let us also recall that given  $i, j$  with  $0 \leq i < j \leq n$  and a cost value  $c$ , the operation *Weighted-Exist*( $i, j, c$ ) finds out whether or not the phrase  $w = T[i + 1 : j]$  is in  $D_i$  and the cost of the corresponding edge  $(i, j)$  in  $G_{\mathcal{A},T}$  is smaller than or equal to  $c$ .

Therefore, we can summarize as it follows.

$$C((i, j)) = C(F_d) + \text{Bitlen}(|w|) + \text{Bitlen}(\text{rightmost}_i(w))$$

$$\text{Weighted-Exist}(i, j, c) = \text{“yes”} \iff T[i + 1 : j] \text{ is in } D_i \text{ and } C((i, j)) \leq c$$

$$\text{Weighted-Extend}((i, j), c) = \text{“yes”} \iff T[i + 1 : j + 1] \text{ is in } D_i \text{ and } C((i, j + 1)) \leq c$$

Notice that previous statements depend on  $\text{rightmost}_i(w)$  that cannot be found in constant time. Notice also that we could use *Query 1* of  $w$  instead of  $\text{rightmost}_i(w)$ , since they have equal *Bitlen* values. Unfortunately, also the *Query 1* cannot be done in constant time. Indeed, using the *Query 1* we can perform any of the needed operations in  $O(|w| \log \log i)$  time.

The idea to perform these operations in constant time is based on the observation that we don't actually need to know the cost of the given edge, but we can just check if a phrase is in a specific layer of our data structure that is the greatest layer that satisfy the  $C((i, j)) \leq c$  inequality.

Since for any  $sw_k \in SW$  there is a  $b_k \in B$  associated to it by definition of  $SW$ , we can restate the operations as it follows.

$$\text{Weighted-Exist}(i, j, c) = \text{“yes”} \iff \exists k \text{ such that } T[i + 1 : j] = w \text{ is in the layer } sw_k \text{ and } C(F_d) + \text{Bitlen}(|w|) + b_k \leq c$$

$$\text{Weighted-Extend}((i, j), c) = \text{“yes”} \iff \exists k \text{ such that } T[i + 1 : j + 1] = wa \text{ is in the layer } sw_k \text{ and } C(F_d) + \text{Bitlen}(|wa|) + b_k \leq c$$

**Definition 5.6.** Let be  $L(w, c) = \max\{k \mid C(F_d) + \text{Bitlen}(|w|) + b_k \leq c\}$ .

Let us notice that an edge has cost less than or equal to  $c$  if and only if the phrase associated to this edge is in the layer  $sw_{L(w,c)}$ . More formally,

**Proposition 5.4.1.** *Let be  $(i, j) \in G_{\mathcal{A},T}$ . We have that for any  $c$ ,  $C((i, j)) \leq c \iff w = T[i + 1 : j]$  is in the layer  $sw_{L(w,c)}$ .*

*Proof.* The phrase  $w$  is the dictionary phrase associated to the edge  $(i, j)$ . If  $(i, j)$  is in  $G_{\mathcal{A},T}$ , then  $w$  is in  $D_i$ . Let us call  $sw_p \in SW$  the smallest

layer where  $w$  appears. For the Prop. 5.2.1 we have that  $w$  is also in all the layers  $sw_x$  with  $sw_x \geq sw_p$ . If  $C((i, j)) \leq c$ , then for some value  $b_k$  we have that  $C(F_d) + \text{Bitlen}(|w|) + b_k \leq c$ . From the definition of  $L(w, c)$  we have that  $b_{L(w, c)}$  is the greatest cost that satisfies the cost inequality and  $sw_{L(w, c)}$  is the layer associated to it. We have that  $sw_{L(w, c)} \geq sw_p$ . Indeed, by contradiction, if  $sw_{L(w, c)} < sw_p$ , then  $b_{L(w, c)} < b_p$  and  $b_p$  does not satisfy the cost inequality, i.e.  $C(F_d) + \text{Bitlen}(|w|) + b_p > c$ . This contradicts  $C((i, j)) \leq c$ . This completes the “if” part of the proof.

If  $w$  is in layer  $sw_{L(w, c)}$ , then  $w$  is in  $D_i$  and by definition of  $L(w, c)$ , we have that  $C(F_d) + \text{Bitlen}(|w|) + sw_{L(w, c)} \leq c$ . Therefore, we have that  $(i, j)$  is in  $G_{\mathcal{A}, T}$  by definition of  $G_{\mathcal{A}, T}$  and since  $b_{L(w, c)}$  is the maximum cost that satisfies the cost inequality, we have that  $C((i, j)) \leq C(F_d) + \text{Bitlen}(|w|) + sw_{L(w, c)} \leq c$ . This concludes the proof.  $\square$

**Corollary 5.4.2.**

Weighted-Exist( $(i, j, c)$ ) = “yes”  $\iff T[i + 1 : j] = w$  is in the layer  $sw_{L(w, c)}$

Weighted-Extend( $(i, j), c$ ) = “yes”  $\iff T[i + 1 : j + 1] = wa$  is in the layer  $sw_{L(wa, c)}$

Let us show how to associate a point in our data structure to a pattern.

Since that edges in suffix trees may have labels longer than one character, the ending point of a pattern may be either a node or a point in the middle of an edge. A classic notation to refer to a general *point* in a suffix tree is a triple composed by a *locus*, a symbol and an offset. The locus is the node closest to the *point* on the path from the root, the symbol is the discriminant character between the outgoing edges of the locus, and the offset tells how many characters in the edge precede the *point*.

More formally, given a pattern  $w \in \text{Fact}(T)$ , let  $S_T$  be the suffix tree for the text  $T$ . Let be  $w = uv$ , with  $u, v \in \text{Fact}(T)$ , where  $u$  is the longest prefix of  $w$  ending at a node  $p$  in  $S_T$  and  $v$  is the prefix of the edge where the pattern  $w$  ends out. Let call  $p$  the *locus* of  $w$ ,  $|v|$  the offset and  $v[1]$  the discriminant character between the outgoing edges of  $p$ . If  $w$  ends to a node

in  $S_T$ , then  $w = u\epsilon$ , this node is the locus  $p$ , the length is equal to 0 and the symbol is any  $c \in \Sigma$ .

In the multilayer suffix tree we denote by  $P(\text{pattern}, x)$  the *point* of the given *pattern* in the layer  $sw_x$ .

### Constant Case

Let us notice that if, for any  $w \in \Sigma^*$ ,  $C(F_d)$  and  $\text{Bitlen}(|w|)$  are constant along a run of the *Find Supermaximal*( $c$ ) algorithm, then also  $L(w, c) = L$  is a constant. It would be easy to obtain an amortized constant time for our operations. Let us take a closer look to this case and let us show how to perform the *Weighted-Exist* and the *Weighted-Extend* operations.

Let suppose to have a pointer to  $P(w, L)$  in our data structure for the text  $T_i$ , with  $w = T[i + 1 : j]$ . Let us define the procedure *Find* as in Table 5.3. In order to perform a *Weighted-Extend*(( $i, j$ ),  $c$ ) we use the procedure *Find*( $j + 1$ ). In order to perform a *Weighted-Exist*( $i, j, c$ ) we use the procedure *Find*( $j$ ). The correctness of this approach immediately follows from Corollary 5.4.2.

Essentially, we handle a point  $p$  in the layer  $sw_L$ . A *Weighted-Extend* corresponds to read a character in the layer  $sw_L$ . A *Weighted-Exist* corresponds to follow a suffix-link from the locus of the point  $p$  together with the skip-count trick (see [19, Sec. 6.1.3]) and to read one character. The time spent to do the two operations *Weighted-Extend* and *Weighted-Exist* in this way turns out to be amortized constant.

We are here using standard techniques concerning suffix trees.

For the purists, we report in the Appendix B all the details concerning how to use the procedure *Find* inside the algorithm *Find Supermaximal* together with the adapted pseudocode and the proof of correctness and time analysis.

From the time analysis in Appendix B next proposition directly follows.

**Proposition 5.4.3.** *Let  $c$  be a cost value and  $L$  a constant. The *Weighted-Exist* and the *Weighted-Extend* operations of the *Find Supermaximal*( $c$ ) algorithm are performed in amortized constant time by using the multilayer suffix tree data structure.*



<pre> <i>Find</i>(<math>k</math>) 1. <math>a \leftarrow T[k]</math> 2. IF <math>Succ(p, a) \neq nil</math> THEN 3.   <math>p \leftarrow Succ(p, a)</math> 4.   RETURN “yes” 5. ELSE 6.   <math>p \leftarrow Suff(p)</math> 7.   RETURN “no” 8. ENDIF  <i>Succ</i>(<math>P(w, k), a</math>) 1. IF <math>wa</math> is in layer <math>sw_k</math> THEN 2.   RETURN <math>P(wa, k)</math> 3. ELSE 4.   RETURN <math>nil</math> 5. ENDIF  <i>Suff</i>(<math>P(w, k)</math>) 1. RETURN <math>P(w, k)</math> </pre>
--

Table 5.3: The pseudocode of the *Find* procedure and its subroutines. The *Succ*( $p, a$ ) routine reads the character  $a$  starting from the point  $p$  in constant time. The *Suff*( $p$ ) routine finds the point of the longest proper suffix of  $p$  by using one suffix-link and the skip-count trick in  $O(1)$  amortized time.

## General Case

In previous subsection we showed how to perform the two operations *Weighted-Extend* and *Weighted-Exist* used inside the algorithm *Find Supermaximal* in amortized constant time by using the multilayer data structure under one hypothesis. The hypothesis is that  $Bitlen(|w|)$  (and therefore  $L(w, c)$ ) is constant.

In this subsection we want to extend this result to the general case. Here we consider just one cost  $c$ . The algorithm that we will describe in the

following is immediately parallelizable to all costs to support the *Dictionary-Symbolwise Flexible Parsing* algorithm.

In the general case, along a run of the *Find Supermaximal* algorithm, since  $L(w, c)$  may change following  $\text{Bitlen}(|w|)$  changes, we may have to change layer while we move the point  $p$ . The main goal of what we are going to describe is that we can handle the layer changes in amortized constant time by using just two points instead of one. The following description is informal but rigorous.

Let  $p, q$  be two points in our data structure. Fixed a cost  $c$ , the *Find Supermaximal* algorithm set  $i = 0$  and  $j = 0$  (that means that the text is just starting). The points  $p$  and  $q$  are set to be equal to the root of the layer  $sw_{L(\epsilon, c)}$ , where  $\epsilon$  is the empty word. Let be  $L(\epsilon, c) = L$ . We assume that we can calculate  $L(w, c)$ ,  $\text{Bitlen}(|w|)$  and therefore  $b_{L(w, c)}$  in constant time, that is an absolutely natural hypothesis that is verified for standard Variable-Length Codes. The algorithm starts to perform the operations *Weighted-Extend* $((i, j), c)$  and *Weighted-Exist* $(i, j + 1, c)$  in the order described by the algorithm *Find Supermaximal*. As described in previous subsection, we support these operations and we maintain  $p = q = P(w, L)$  in the same layer  $sw_L$  until  $L(w, c)$  changes, i.e. until when  $L(w, c) = L'$  with  $L \neq L'$ .

Whenever  $L(w, c)$  changes from  $L$  to  $L'$ , we maintain  $q$  on the layer  $L$  and we move  $p$  to the layer  $L'$ . Roughly speaking,  $sw_{L'}$  is the current layer and  $sw_L$  is the previous layer. This simple dynamic programming trick allows to save time when  $L(w, c)$  turns back to its previous value. In this case we just swap  $p$  and  $q$ . Otherwise, when  $sw_{L'}$  is different from the layer of  $q$ , we firstly save  $p$  by using  $q$ , i.e. we set  $q$  equal to  $p$ , and then we set  $p$  to be equal to  $P(w, L')$ . In order to reach the point  $P(w, L')$ , we simply start reading  $w$  from the root of the layer  $sw_{L'}$  in  $O(|w|)$  time. We will see in the time analysis that this quantity can be amortized to a constant in the overall run of the algorithm.

Notice that for  $L' \leq L$ , i.e. when the new layer  $sw_{L'}$  is smaller than the layer  $sw_L$ , we have that if  $w$  is in  $sw_{L'}$ , then  $w$  is also in  $sw_L$  for the Prop. 5.2.1. In this case, we can maintain  $p = P(w, L')$  and  $q = P(w, L)$  as stated before. Otherwise, for  $L' > L$  and  $w$  is in  $sw_{L'}$ , we have not guarantees that

$w$  is in  $sw_L$ . In this case we let  $q$  points to the longest prefix of  $w$  in  $sw_L$ . Therefore, if  $v$  is the longest prefix of  $w$  in  $sw_L$ , then we set  $q = P(v, L)$ . As soon as a new *Weighted-Extend* or a new *Weighted-Exist* operation is accomplished, we maintain  $q$  to the longest prefix of  $w$ .

Since  $v$  is a prefix of  $w$ , the number of steps used to reach  $P(v, L)$  in the layer  $sw_L$  is less than or equal to the number of steps needed to reach  $P(w, L)$ . Furthermore, the number of steps used later to keep  $q$  to the longest prefix of  $w$  is less or equal to the number of steps saved before.

The redefined *Find* procedure in Table 5.4 is in charge of handling the  $p, q$  points as described. The same table reports also some subroutines.

We present the time analysis of the overall pointer handling.

From the above description and from the results of previous subsection, we have that the overall time to move two pointers  $p$  and  $q$  inside the multilayer suffix tree inside two fixed layers is linear on the size of the text. Let  $L_p$  be the cost of the phrase pointed by  $p$  and  $L_q$  the cost of the phrase pointed by  $q$ . Let us now focus on the steps used to move  $p$  and  $q$  from a layer to another. Since  $q$  simply follows point  $p$ , we focus just on the pointer  $p$ . Two cases are possible. If  $L(w, c) = L_q$ , then  $p$  is set to  $q$  in constant time. Otherwise,  $L(w, c) \neq L_q$  and we have to move  $p$  from  $P(\epsilon, L(w, c))$  to  $P(w, L(w, c))$  in  $O(|w|)$  time. We want to show that these  $O(|w|)$  steps are amortizable over the number of algorithm steps between two consecutive change of  $L(w, c)$ .

Suppose that a change of  $L(w_t, c)$  happens at time  $i_t$  with  $j = j_t$ ,  $w_t = T[i_t + 1 : j_t]$  of length  $h = |w_t|$ . Suppose that the first changes after time  $i_t$  happens at time  $i_{t'}$  with  $j = j_{t'}$ ,  $w_{t'} = T[i_{t'} + 1 : j_{t'}]$  of length  $h' = |w_{t'}|$ .

If  $h' > h$ , then, from the Property 5.1 and the Remark 5.1, we have that  $h' \geq kh$  unless  $h' \leq \hat{k}$ , where  $\hat{k}$  is a constant and there is no need of amortizing. Notice that, from time  $i_t$  to time  $i_{t'}$ , the algorithm do  $(i_{t'} - i_t + j_{t'} - j_t) \geq (h' - h)$  steps. We amortize the  $O(|w_{t'}| = h')$  steps of the pointer  $p$  over the  $(h' - h)$  steps of the *Find Supermaximal* algorithm. For each algorithm step we make at most

$$\frac{h'}{h' - h} \leq \frac{k}{k - 1}$$

steps on the multilayer suffix tree. Indeed, from  $h' \geq k h$  follows

$$\frac{h}{h'} \leq \frac{h}{k h} = \frac{1}{k} \Rightarrow 1 - \frac{h}{h'} \geq 1 - \frac{1}{k} \Rightarrow \frac{1}{1 - \frac{h}{h'}} \leq \frac{1}{1 - \frac{1}{k}} \Rightarrow \frac{h'}{h' - h} \leq \frac{k}{k - 1}$$

where  $\frac{k}{k-1}$  is the amortization constant over the  $h' - h$  last algorithm steps.

If  $h' < h$ , then  $h' \leq \frac{1}{k} h$ . Therefore, for each algorithm step we make at most

$$\frac{h'}{h - h'} \leq \frac{1}{k - 1}$$

steps on the multilayer suffix tree. This complete the time analysis. The following theorem is an immediate consequence of this time analysis.

**Proposition 5.4.4.** *The total number of steps on the multilayer suffix tree used to perform all the Weighted-Exist and the Weighted-Extend operations of a run of the Find Supermaximal(c) algorithm is linear on the size of the text. Therefore, the multilayer suffix tree uses amortized constant time for any operation.*

Analogously to what done in Subsection 5.4 concerning cost parallelization, we extend above proposition to the *Dictionary-Symbolwise Flexible Parsing* as it follows.

**Corollary 5.4.5.** *The multilayer suffix tree supports the Weighted-Exist and the Weighted-Extend operations inside the Dictionary-Symbolwise Flexible Parsing algorithm in amortized constant time.*

<pre> <i>Find</i>(<i>k</i>) 01.  <math>L \leftarrow L(T[i + 1 : k], c)</math> 02.  IF <math>L_p = L</math> THEN 03.    RETURN <i>Sub-Find</i>(<i>k</i>) 04.  ELSE 05.    IF <math>L_q \neq L</math> THEN 06.      <math>gap \leftarrow k - 1 - i + 1</math> 07.      <math>L_q \leftarrow L</math> 08.      <math>q \leftarrow P(\epsilon, L)</math> 09.      <i>Fillgap</i>(<math>k - 1</math>) 10.    ENDIF 11.    IF <math>gap = 0</math> THEN 12.      <i>Swap</i>() 13.      RETURN <i>Sub-Find</i>(<i>k</i>) 14.    ELSE 15.      <math>p \leftarrow Suff(p)</math> 16.      <math>q \leftarrow Suff(q)</math> 17.      <i>Fillgap</i>(<math>k - 1</math>) 18.      RETURN “no” 19.    ENDIF 19.  ENDIF  <i>Succ</i>(<math>P(w, k), a</math>) 1.  IF <math>wa</math> is in layer <math>sw_k</math> THEN 2.    RETURN <math>P(wa, k)</math> 3.  ELSE 4.    RETURN <i>nil</i> 5.  ENDIF </pre>	<pre> <i>Sub-Find</i>(<i>k</i>) 01.  <math>a \leftarrow T[k]</math> 02.  IF <math>Succ(p, a) \neq nil</math> THEN 03.    <math>p \leftarrow Succ(p, a)</math> 04.    <math>gap \leftarrow gap + 1</math> 05.    <i>Fillgap</i>(<i>k</i>) 06.    RETURN “yes” 07.  ELSE 08.    <math>p \leftarrow Suff(p)</math> 09.    <math>q \leftarrow Suff(q)</math> 10.    <i>Fillgap</i>(<math>k - 1</math>) 11.    RETURN “no” 12.  ENDIF  <i>Fillgap</i>(<i>k</i>) 1.  <math>t \leftarrow Succ(q, T[k - gap + 1])</math> 2.  WHILE <math>gap &gt; 0</math> &amp; <math>t \neq nil</math> DO 3.    <math>q \leftarrow t</math> 4.    <math>gap \leftarrow gap - 1</math> 5.    <math>t \leftarrow Succ(q, T[k - gap + 1])</math> 6.  ENDWHILE  <i>Swap</i>() 1.  Swap <math>p</math> and <math>q</math> 2.  Swap <math>L_p</math> and <math>L_q</math>  <i>Suff</i>(<math>P(aw, k)</math>) 1.  RETURN <math>P(w, k)</math> </pre>
--	--

Table 5.4: The pseudocode of the redefined *Find* procedure and its subroutines. Variables  $L_p, L_q, p, q, gap$  and  $i$  are variables of the *Find Supermaximal* algorithm.  $L_p$  and  $L_q$  are initialized to  $L(\epsilon, c)$ ,  $p$  and  $q$  are initialized to  $P(\epsilon, L_p)$  and the *gap* counter is initialized to 0. Notice that *Find* begins always with  $p = P(T[i - 1 : k - 1], L_p)$  and  $q = P(T[i - 1 : k - 1 - gap], L_p)$ .

## Chapter 6

# Conclusion

In this thesis we present some advancement on dictionary-symbolwise theory. We describe the *Dictionary-Symbolwise Flexible Parsing*, a parsing algorithm that extends the *Flexible Parsing* (see [28]) to *variable* costs and to the dictionary-symbolwise domain. We prove its optimality for prefix-closed dynamic dictionaries under some reasonable assumption. *Dictionary-Symbolwise Flexible Parsing* is *linear* for LZ78-like dictionaries. In the case of LZ77-like dictionary, we obtain the  $O(n \log n)$  complexity as authors of [17] recently did by using a completely different subgraph, where  $n$  is the text size.

We introduce the Multilayer Suffix Tree data structure that is able to represent the LZ77 dictionary in  $O(n)$  space and  $O(n \log n)$  building time. This data structure is also able to perform the cost-based queries of the *Dictionary-Symbolwise Flexible Parsing* algorithm in amortized constant time.

We also prove that dictionary-symbolwise compressors can be asymptotically better than optimal pure dictionary compressors in terms of compression ratio.

Furthermore, we ultimately prove the optimality of the greedy parsing for LZ77-like dictionary under the uniform cost assumption about the encoding of dictionary-phrases.

Last but not least, our *Dictionary-Symbolwise Flexible Parsing* algorithm allows to couple classical LZ-like compressors with several symbolwise meth-

ods to obtain new dictionary-symbolwise algorithms with proof of parsing optimality.

We conclude this thesis with two open problems.

1. Theoretically, LZ78 is better on memoryless sources than LZ77. Experimental results say that when optimal parsing is in use it happens the opposite. To prove this fact both in pure dictionary case and in dictionary-symbolwise case.

2. Common symbolwise compressors are based on the arithmetic coding approach. When these compressors are used, the costs in the graph are almost surely noninteger and, moreover, the graph is usually not well defined. The standard workaround is to use an approximation strategy. A big goal should be finding an optimal solution for these important cases.

# Bibliography

- [1] Amihoud Amir, Gad M. Landau, and Esko Ukkonen. Online time-stamped text indexing. *Information Processing Letters*, 82(5):253 – 259, 2002.
- [2] Richard Arratia and Michael S. Waterman. The Erdős-Rényi strong law for pattern matching with a given porportion of mismatches. *Annals of Probability*, 17:1152–1169, 1989.
- [3] Timothy C. Bell, John G. Cleary, and Ian H. Witten. *Text compression*. Prentice Hall, 1990.
- [4] Timothy C. Bell and Ian H. Witten. The relationship between greedy parsing and symbolwise text compression. *Journal of the ACM*, 41(4):708–724, 1994.
- [5] Martin Cohn and Roger Khazan. Parsing with prefix and suffix dictionaries. In James A. Storer and Martin Cohn, editors, *Data Compression Conference*, pages 180 – 189. IEEE Computer Society, 1996.
- [6] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT Press, Cambridge, MA, second edition, 2001.
- [7] Maxime Crochemore, Chiara Epifanio, Alessandra Gabriele, and Filippo Mignosi. On the suffix automaton with mismatches. In Jan Holub and Jan Zdárek, editors, *CIAA*, volume 4783 of *Lecture Notes in Computer Science*, pages 144–156. Springer, 2007.
- [8] Maxime Crochemore, Laura Giambruno, Alessio Langiu, Filippo Mignosi, and Antonio Restivo. Dictionary-symbolwise flexible parsing.



- In Costas S. Iliopoulos and William F. Smyth, editors, *Combinatorial Algorithms: IWOCA '2010 Revised Selected Papers*, volume 6460 of *Lecture Notes in Computer Science*, pages 390–403. Springer, 2010.
- [9] Maxime Crochemore, Laura Giambruno, Alessio Langiu, Filippo Mignosi, and Antonio Restivo. Dictionary-symbolwise flexible parsing. *Journal of Discrete Algorithms - IWOCA '10 Special Issue*, 2011.
- [10] Maxime Crochemore and Lucian Ilie. Computing longest previous factor in linear time and applications. *Inf. Process. Lett.*, 106(2):75–80, 2008.
- [11] Maxime Crochemore and Thierry Lecroq. Pattern-matching and text-compression algorithms. *ACM Comput. Surv.*, 28(1):39–41, 1996.
- [12] Giuseppe Della Penna, Alessio Langiu, Filippo Mignosi, and Andrea Ulisse. Optimal parsing in dictionary-symbolwise data compression schemes. <http://math.unipa.it/~alangiu/OptimalParsing.pdf>, unpublished manuscript, 2006.
- [13] Chiara Epifanio, Alessandra Gabriele, and Filippo Mignosi. Languages with mismatches and an application to approximate indexing. In *Proceedings of the 9th International Conference Developments in Language Theory (DLT05)*, LNCS 3572, pages 224–235, 2005.
- [14] Chiara Epifanio, Alessandra Gabriele, Filippo Mignosi, Antonio Restivo, and Marinella Sciortino. Languages with mismatches. *Theor. Comput. Sci.*, 385(1-3):152–166, 2007.
- [15] Peter M. Fenwick. Symbol ranking text compression with shannon recodings. *Journal of Universal Computer Science*, 3(2):70–85, 1997.
- [16] Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *Journal of the ACM*, 52(4):552–581, 2005.
- [17] Paolo Ferragina, Igor Nitto, and Rossano Venturini. On the bit-complexity of Lempel-Ziv compression. In *SODA '09: Proceedings of*

*the Nineteenth Annual ACM - SIAM Symposium on Discrete Algorithms*, pages 768–777, Philadelphia, PA, USA, 2009. Society for Industrial and Applied Mathematics.

- [18] Alessandra Gabriele, Filippo Mignosi, Antonio Restivo, and Marinella Sciortino. Indexing structure for approximate string matching. In *Proc. of CIAC'03*, volume 2653 of *LNCS*, pages 140–151, 2003.
- [19] Dan Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computational Biology*. Cambridge University Press, 1997.
- [20] Alan Hartman and Michael Rodeh. Optimal parsing of strings. In *Combinatorial Algorithms on Words*, pages 155–167. Springer – Verlag, 1985.
- [21] R. Nigel Horspool. The effect of non-greedy parsing in Ziv-Lempel compression methods. In *Data Compression Conference*, pages 302–311, 1995.
- [22] Shunsuke Inenaga, Ayumi Shinohara, Masayuki Takeda, and Setsuo Arikawa. Compact directed acyclic word graphs for a sliding window. In *SPIRE*, pages 310–324, 2002.
- [23] Tae Young Kim and Taejeong Kim. On-line optimal parsing in dictionary-based coding adaptive. *Electronic Letters*, 34(11):1071–1072, 1998.
- [24] S. Rao Kosaraju and Giovanni Manzini. Compression of low entropy strings with Lempel–Ziv algorithms. *SIAM J. Comput.*, 29(3):893–911, 2000.
- [25] Alessio Langiu. Optimal parsing in dictionary-symbolwise compression algorithms. Master’s thesis, University of Palermo, 2008. [http://math.unipa.it/~alangiui/Tesi\\_Alessio\\_Langiui\\_MAc.pdf](http://math.unipa.it/~alangiui/Tesi_Alessio_Langiui_MAc.pdf).
- [26] N. Jesper Larsson. Extended application of suffix trees to data compression. In *Data Compression Conference*, pages 190–199, 1996.

- [27] Yossi Matias, Nasir Rajpoot, and Süleyman Cenk Sahinalp. The effect of flexible parsing for dynamic dictionary-based data compression. *ACM Journal of Experimental Algorithms*, 6:10, 2001.
- [28] Yossi Matias and Süleyman Cenk Sahinalp. On the optimality of parsing in dynamic dictionary based data compression. In *SODA*, pages 943–944, 1999.
- [29] Joong Chae Na, Alberto Apostolico, Costas S. Iliopoulos, and Kunsoo Park. Truncated suffix trees and their application to data compression. *Theor. Comput. Sci.*, 304:87–101, July 2003.
- [30] David Salomon. *Data compression - The Complete Reference, 4th Edition*. Springer, 2007.
- [31] David Salomon. *Variable-length Codes for Data Compression*. Springer-Verlag, 2007.
- [32] Khalid Sayood. *Introduction to Data Compression*. Morgan Kaufmann, 1996.
- [33] Ernst J. Schuegraf and H. S. Heaps. A comparison of algorithms for data base compression by use of fragments as language elements. *Information Storage and Retrieval*, 10(9-10):309–319, 1974.
- [34] Martin Senft. Suffix tree for a sliding window: An overview. In *Proceedings of WDS'05, Part 1*, pages 41–46, 2005.
- [35] Martin Senft and Tomáš Dvořák. Sliding cdawg perfection. In *Proceedings of the 15th International Symposium on String Processing and Information Retrieval, SPIRE '08*, pages 109–120, Berlin, Heidelberg, 2009. Springer-Verlag.
- [36] James A. Storer and Thomas G. Szymanski. Data compression via textural substitution. *J. ACM*, 29(4):928–951, 1982.
- [37] Wojciech Szpankowski. *Average Case Analysis of Algorithms on Sequences*. Wiley, 2001.

- [38] Robert A. Wagner. Common phrases and minimum-space text storage. *Commun. ACM*, 16(3):148–152, 1973.
- [39] Terry A. Welch. A technique for high-performance data compression. *IEEE Computer*, january:8–19, 1984.



## Appendix A

# Experiments

We now discuss about some experiments. Readers must keep into account that the results of this thesis are mainly theoretical and that they apply to a very large class of compression algorithms. Due to this, the use of different methods of encoding for dictionary pointers as well as for symbolwise encoding and for the flag information encoding together with the dictionary constrains leads to different performances. Performances about time and space are strongly dependent on the programming language in use and on the programmers abilities. Therefore we decided to focus only on compression ratio.

Here we discuss two particular cases that allow to compare our results with some well know commercial compressors. The first one is related to LZ78-like dictionary and a Huffman code. The second one concerns LZ77-like dictionaries with several window sizes and a Huffman code. We compare the obtained compression ratio with the `gzip`, `zip` and `cabarc` compression tools. The encoding method in use is a semi static Huffman code.

In the first experiment, using a simple semi static Huffman code as symbolwise compressor, we improved the compression ratio of the *Flexible Parsing* with LZW-dictionary by 3 to 5 percent on texts such as the `bible.txt` file or the prefixes of English Wikipedia data base (see Table A.1). We obtain that the smaller is the file the greater is the gain.

We have experimental evidence that many of the most relevant LZ77-like commercial compressors are, following our definition, dictionary-symbolwise

File (size)	bible.txt (4047392 Byte)	enwik (100 MB)
gzip -9	29.07%	36.45%
lzwp	30.09%	35.06%
lzwlds	25.84%	31.79%

Table A.1: Compression ratio comparison of some LZW-like compressors and the gzip tool. (*gzip -9* is the gzip compression tool with the -9 parameter for maximum compression. *lzwp* is the Flexible Parsing algorithm of Matias-Rajpoot-Sahinalp with a LZW-like dictionary. *lzwlds* is our Dictionary-Symbolwise Flexible Parsing algorithm with LZW-like dictionary and a Huffman code.)

File (size)	bible.txt (4047392 Byte)	enwik (100 MB)
gzip -9	29.07%	36.45%
gzip by 7zip	27.44%	35.06%
zip by 7zip	25.99%	33.72%
cabarc	22.13%	28.46%
lzhds-32KB	27.47%	35.02%
lzhds-64KB	26.20%	33.77%
lzhds-2MB	22.59%	28.82%
lzhds-16MB	22.51%	26.59%

Table A.2: Compression ratio comparison of some LZ77-like compressors. (*gzip -9* is the gzip compression tool with the -9 parameter for maximum compression. *gzip by 7zip* is the gzip compression tool implemented in the 7-Zip compression suite. *zip by 7zip* is the 7-Zip implementation of the zip compression tool. *cabarc* is the MsZip cabinet archiving tool also known as cabarc (version 5.1.26 with -m lzx:21 option used). *lzhds-x* is our Dictionary-Symbolwise Flexible Parsing with LZ77-like dictionary of different dictionary sizes, as stated in the suffix of the name, and a Huffman code.)

File (size)	bible.txt (4047392 Byte)	enwik (100 MB)
gzip -9 / lzhd-32KB	105.82%	104.08%
gzip by 7zip / lzhd-32KB	99.89%	100.11%
zip by 7zip / lzhd-64KB	99.19%	99.85%
cabarc / lzhd-2MB	97.96%	98.75%

Table A.3: Ratio between the compression ratio of different LZ77-like compressors. All the involved compressors, except for the gzip one, seem to have an optimal parsing strategy. (See Table A.2 caption for compressor descriptions.) Notice that on each row there are compressors having the same windows size.

algorithms and they use an optimal parsing (see Table A.2 and Table A.3). In Table A.3 is shown the ratio between compression performances of compressors with similar constrains and encoding. Indeed, *gzip* and *lzhd-32KB* use a LZ77-like dictionary of 32KB, *zip* and *lzhd-64KB* have dictionary size of 64KB. *cabarc* and *lzhd-2MB* use 2MB as dictionary size. They all use a Huffman code. We notice that a difference of about 5 percent is due to parsing optimality while small differences of about 2 percent are due to implementation details like different codeword space and different text block handling. We think that *gzip* and *zip* implementations in the 7-Zip compression suite and *cabarc* have an optimal parsing, even if this fact is not clearly stated or proved.





## Appendix B

# Time Analysis of the Constant Case

We present the time analysis of the *Find Supermaximal* algorithm variant for the constant  $Bitlen(|w|)$  case, defined as in Table B.1, where the *Weighted-Exist* and the *Weighted-Extend* operations are done by using the *Find* procedure of the Multilayer Suffix Tree of  $T_i$ .

For a fixed text  $T$ , let us define the set  $J = \{j_0, j_1, \dots, j_q\} \subset [0..n]$  to be the set of the  $j$  values after each of the positive answer of the *Weighted-Exist*( $i, j, c$ ) operations in the *Find Supermaximal* algorithm for a fixed integer ( $c$ ). Notice that, since at each step of the *Find Supermaximal* algorithm either  $i$  or  $j$  is increased and they are never decreased, we have that for each  $j_p \in J$  there exists one and only one value  $i$  such that *Weighted-Exist*( $i, j, c$ ) gives a positive answer. For any  $j_p \in J$ , let us call  $i_p$  such corresponding value.

We want to show how to perform the *Weighted-Exist* and the *Weighted-Extend* operations from the generic step  $j_p$  to the step  $j_{p+1}$  with  $0 \leq p < q$ . In the meantime, we want to proof following proposition.

**Proposition B.0.6.** *If at time  $i_p$  we have a pointer  $P(w, L)$  such that  $w = T[i_p + 1 : j_p]$ , then at time  $i_{p+1}$  we will have correctly computed the pointer  $P(w, L)$  with  $w = T[i_{p+1} + 1 : j_{p+1}]$  using at most  $i_{p+1} - i_p + j_{p+1} - j_p$  steps on our data structure.*

*Proof.* The proof is by induction on  $p$ . For  $p = 0$ , let  $j_0$  be the first element in  $J$ .

When  $i = j$ ,  $T[i + 1 : j] = \epsilon$  by definition.  $P(\epsilon, L)$  is the pointer to the root of the layer  $sw_L$ . From a simple analysis of the pseudocode of *Find Supermaximal* in Table 4.1 we have that  $i_0 = j_0 - 1$ . Therefore, *Weighted-Exist*( $i, i + 1, c$ ) returns a positive answer as soon as there is a character  $T[i + 1] \in D_i$ . In order to perform this check, we maintain the pointer  $P(\epsilon, L)$  fixed to the root of the layer  $L$  and we just check for an outgoing edge starting with the character  $a = T[i + 1]$  in constant time for  $0 \leq i \leq i_0 = j_0 - 1$ . Once that  $i = i_0$ , we move in constant time  $P(\epsilon, L)$  to  $P(a, L)$ . So, for  $p = 0$  and consequently at time  $i_0$ , we have the pointer  $P(a, L)$  such that  $a = T[i_0 + 1 : j_0]$  and *Weighted-Exist*( $i_0, j_0, c$ ) gives a positive answer. This is the base of the induction.

The inductive hypothesis is that for a generic  $p < q$ , i.e. at time  $i_p$ , we have the pointer  $P(w, L)$  to  $w$  in the layer  $sw_L$  with  $w = T[i_p + 1 : j_p]$  and  $L = L(w, c)$  for any  $w$ . At time  $i_p$ , the *Find Supermaximal* algorithm has successfully done *Weighted-Exist*( $i_p, j_p, c$ ), i.e. the edge  $(i, j)$  is in  $G_{\mathcal{A}, T}$  and  $C(i, j) \leq c$  with  $i = i_p$  and  $j = j_p$ .

We want to proof for  $p + 1$  that we can find  $P(T[i_{p+1} + 1 : j_{p+1}], L)$ .

We consider the steps from time  $i_p$  to the time  $i_{p+1}$ . At time  $i_p$  we have  $P(w, L)$  with  $w = T[i_p + 1 : j_p]$  by inductive hypothesis. Since that the algorithm has yet successfully done *Weighted-Exist*( $i, j, c$ ) with  $i = i_p$  and  $j = j_p$ , it now performs a series of *Weighted-Extend*( $(i, j), c$ ) operations for increasing values of  $j$ , until it found a  $j$  such that  $(i, j)$  is in  $G_{\mathcal{A}, T}$  with  $C(i, j) \leq c$  and  $(i, j + 1)$  is not. In order to perform such series of operations we have just to check, for any  $j$ , if  $wa$  is in the layer  $L(wa, c)$ , where  $a = T[j + 1]$  is the  $j + 1$ -th character of  $T$ . Since we are assuming  $L(w, c)$  constant, we have that  $L(wa, c) = L(w, c) = L$  and we have just to check if an  $a$  can be read starting from the point  $P(w, L)$ . If there is an  $a$  after  $P(w, L)$ , then  $(i, j + 1)$  is in  $D_i$  with  $C(i, j + 1) \leq c$  and the *Weighted-Extend*( $(i, j), c$ ) returns "yes". We also let  $P(w, L)$  become  $P(wa, L)$  and the algorithm increase the value of  $j$  of an unit. If there is not an  $a$  after  $P(w, L)$ , then *Weighted-Extend*( $(i, j), c$ ) returns "no" and the series of *Weighted-Extend* terminate.

At this point, the algorithm starts a series of *Weighted-Exist*( $i+1, j+1, c$ ) increasing the value of  $i$ , until *Weighted-Exist*( $i+1, j+1, c$ ) returns "yes".

Notice that the edge  $(i+1, j)$  represents the longest proper suffix of  $w$  and since the dictionary has the strong suffix-closed property (see Corollary 3.3.2),  $(i+1, j)$  is in  $D_i$ . Let  $u$  be such suffix of  $w$ . Since the cost function  $C$  is *prefix-nondecreasing*, we know that  $C(i+1, j) \leq C(i, j)$  and then  $u$  is in the layer  $sw_{L(w,c)}$ . We move  $P(w, L)$  to  $P(u, L)$  in amortized constant time by using standard property of suffix trees.

Since we have  $P(u, L)$ , in order to perform the operation *Weighted-Exist*( $i+1, j+1, c$ ), we have just to check if there is an  $a$  after  $P(u, L)$ , where  $a$  is equal to  $T[j+1]$ . As we have already seen, we can check this in constant time. If there is not an  $a$  after  $P(u, L)$  the *Weighted-Exist*( $i+1, j+1, c$ ) returns "no". The algorithm increases the value of  $i$  of an unit and we move again  $P(u, L)$  to its longest proper suffix. Otherwise, If there is an  $a$  after  $P(u, L)$  the *Weighted-Exist*( $i+1, j+1, c$ ) returns "yes".

The algorithm now increases  $i$  and  $j$  by one. Since last *Weighted-Exist* operation returned "yes", we know that  $ua \in D_i$ , where  $ua$  is the dictionary phrase associated to the edge  $(i, j)$ . Therefore we can move  $P(u, L)$  to  $P(ua, L)$ . Since this is the first positive *Weighted-Exist* answer after  $i_p$ , we know that current value of  $i, j$  are  $i = i_{p+1}$  and  $j = j_{p+1}$ . Since  $ua$  is equal to  $T[i+1 : j]$ , this concludes the proof.

□

<pre> <i>Find Supermaximal</i>(<math>c, L</math>) 01. <math>i \leftarrow 0, j \leftarrow 0, M_c \leftarrow \emptyset, p \leftarrow P(\epsilon, L)</math> 02. WHILE <math>j &lt; n</math> DO 03.   <math>j \leftarrow j + 1</math> 04.   WHILE <math>i &lt; j</math> AND <math>Find(j) = \text{"no"}</math> DO 05.     <math>i \leftarrow i + 1</math> 06.   ENDWHILE 07.   IF <math>i &lt; j</math> THEN 08.     WHILE <math>Find(j + 1) = \text{"yes"}</math> DO 09.       <math>j \leftarrow j + 1</math> 10.     ENDWHILE 11.     INSERT <math>((i, j), M_c)</math> 12.   ENDIF 13. ENDWHILE 14. RETURN <math>M_c</math>    <i>Find</i>(<math>k</math>) 1. <math>a \leftarrow T[k]</math> 1. IF <math>Succ(p, a) \neq nil</math> THEN 2.   <math>p \leftarrow Succ(p, a)</math> 3.   RETURN "yes" 4. ELSE 5.   <math>p \leftarrow Suff(p)</math> 6.   RETURN "no" 7. ENDIF    <i>Succ</i>(<math>P(w, k), a</math>) 1. IF <math>wa</math> is in layer <math>sw_k</math> THEN 2.   RETURN <math>P(wa, k)</math> 3. ELSE 4.   RETURN <math>nil</math> 5. ENDIF    <i>Suff</i>(<math>P(aw, k)</math>) 1. RETURN <math>P(w, k)</math> </pre>
--

Table B.1: The pseudocode of the *Find Supermaximal* algorithm for the constant *Bitlen*( $|w|$ ) case, together with the *Find* procedure and its subroutines. The *Succ*( $p, a$ ) routine reads the character  $a$  starting from the point  $p$  in constant time. The *Suff*( $p$ ) routine find the point of the longest proper suffix of  $p$  by using one suffix-link and the skip-count trick in amortized constant time.