

UNIVERSIDAD DE CHILE FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN

TO INDEX OR NOT TO INDEX: TIME-SPACE TRADE-OFFS IN SEARCH ENGINES WITH POSITIONAL RANKING FUNCTIONS

TESIS PARA OPTAR AL GRADO DE MAGÍSTER EN CIENCIAS, MENCIÓN COMPUTACIÓN

SENEN ANDRÉS GONZÁLEZ CORNEJO

PROFESOR GUÍA: DIEGO ARROYUELO BILLIARDI GONZALO NAVARRO

MIEMBROS DE LA COMISIÓN: DIEGO ARROYUELO BILLIARDI GONZALO NAVARRO DIEGO SECO

Este trabajo ha sido parcialmente financiado por .

SANTIAGO DE CHILE DICIEMBRE 2013

Abstract

Positional ranking functions, widely used in web search engines, improve result quality by exploiting the positions of the query terms within documents. However, it is well known that positional indexes demand large amounts of extra space, typically about three times the space of a basic non-positional index. Textual data, on the other hand, is needed to produce text snippets. In this paper, we study time-space trade-offs for search engines with positional ranking functions and text snippet generation. We consider both index-based and non-index based alternatives for positional data. We aim to answer the question of whether one should index positional data or not.

We show that there is a wide range of practical time-space trade-offs. Moreover, we show that both position and textual data can be stored using about 66% of the space used by traditional positional indexes, with a minor increase in query time. This yields considerable space savings and outperforms, both in space and time, recent alternatives from the literature. We also propose several efficient compressed text representations for snippet generation, which are able to use about half of the space of current state-of-the-art alternatives with little impact in query processing time.

Dedicado a mi Esposa Carmen y a mi hija Jazmín, pues son la felicidad de mi vida.

Thanks

Quisiera partir agradeciendo a mi esposa y mi hija por darme todo su apoyo, y a mis padres por llevarme adelante. A Diego Arroyuelo quien estuvo con migo y me guió en el proceso del trabajo, un especial agradecimiento por su fuerte apoyo al final, durante la escritura, dado que tuvo que leer incontables errores (a estas alturas horrores) de ortografía, pero me ayudo mucho a darme cuenta lo difícil de escribir el trabajo realizado y lo importante de hacerlo correctamente. Al profesor Gonzalo Navarro, por aceptarme como alumno. Al laboratorio Yahoo! Reseach chile, por facilitarme las herramientas para poder realizar este trabajo. A mi segunda familia CROSS, por estar siempre con migo.

Contents

1.1	Basic Definitions
1.2	Text Representation
1.3	Inverted Indexes
1.4	Inverted Index Compression
	1.4.1 VByte
	1.4.2 Simple 9
	1.4.3 PforDelta
1.5	Snippet Generation
	1.5.1 Dictionary Compressors
	1.5.2 LZ77
	1.5.3 Lzma
	1.5.4 Snappy
	1.5.5 Lz4
1.6	Compressed Text Self-Indexes
1.7	Wavelet trees
	1.7.1 Supporting Operations
	1.7.2 Analysis of Space Usage
	1.7.3 Self-Indexes for IR Applications
Ind	exing for Positional Ranking: Classical Solutions
2.1	Basic Query Processing Steps for Positional Ranking and Snippet Ex-
	traction
2.2	Experimental Setup and Data
2.3	The Baseline: Positional Inverted Lists and Compressed Textual Data .
	2.3.1 Supporting the Query-Processing Step
	2.3.2 Supporting the Positional-Ranking Step
	2.3.3 Supporting the Snippet Generation Step
Co	omputing Term Positions from Textual Data
	Wandat Trace for Commuting Desitions and Fortugating Spingets
Δ	wavelet tree for Complifing Positions and Extracting Shippers

	4.4	Achieving Higher-Order Compression with the \mathtt{WT}	36
5 Co		omputing Term Positions and Snippets from the Compressed Text	
			38
	5.1	Using Standard Compressors	38
	5.2	Using Zero-Order Compressors with Fast Text Extraction	39
	5.3	Using Natural-Language Compression Boosters	40
	5.4	Further Comparison Between the Most Competitive Alternatives	43
		5.4.1 Space/Time Trade-Offs	43
		5.4.2 Average Position-Extraction Time as Function of the Query Lengths	44
		5.4.3 Average Position-Extraction Time as Function of k_1	47
	5.5	Conclusions	49
6	Disc	cussion and Further Experimental Results	50
	6.1	Scenario 1: Query Processing with Snippet Generation	50
	0	6.1.1 DAAT AND Queries	52
		6.1.2 BMW OR Queries	58
	6.2	Scenario 2: Query Processing without Snippet Generation	63
		6.2.1 DAAT AND Queries	64
		6.2.2 BMW OR Queries	67
	6.3	Discussion	70
7	Con	clusion and Future Work	71
\mathbf{A}	Adi	tional Experimental Results	76
	A.1	DAAT AND Queries	76
	A.2	BMW OR Queries	79

List of Tables

1.1	Number of bytes used to represent an integer number with <i>VByte.</i>	9
1.2	Number of nibbles used to represent an integer number with <i>VNibble</i> .	10
1.3	Meaning of the cases in the header of the S9 word	11
2.1	Experimental results for the initial query processing step (Step 1) for AND and OR queries. In both cases, BM25 ranking is used.	28
2.2	Experimental results for extracting term-position data (Step 2)	29
2.3	Experimental results for the snippet extraction phase (Step 3)	30
4.1	Experimental results for extracting term-position data using a \mathtt{WT}	36
5.1	Experimental results for extracting term-position data from text. \ldots	42
6.1	Glossary of the indexing schemes tested. All schemes include the inverted	F 1
6.2	Glossary of the indexing schemes for the figures. All schemes include the	51
	inverted index.	64

List of Figures

1.1	Example of an invertex index for a document collection, storing just docIDs.	6
1.2	Example of an inverted list for the term t_3	6
1.3	Byte layout of <i>VByte</i>	7
1.4	The 32-bit representation of 167.	8
1.5	The encoding of the integer 167 in VBYTE.	8
1.6	Decoding process for <i>VByte</i>	9
1.7	Decoding process for integer 167.	9
1.8	Example of <i>VNibble</i> bytes	10
1.9	Layout of an S9 word.	11
1.10	S9 word encoding the group of integers 98, 112, 117 and 121	12
1.11	Decodification process for an S9 word, which encodes the group of inte-	
	gers 98, 112, 117 and 121	12
1.12	Decoding process for an S9 word.	13
1.13	Layout of PforDelta block.	13
1.14	PforDelta block, encoding a group of integers, with exeptions 78, 110, 160	
	and 91	14
1.15	Decoding a PforDeltablock.	15
1.16	Decoding process for PforDelta block.	15
1.17	An example of user-query and the result snippets	16
1.18	Example of how the buffer of LZ77 is updated with the text $T = \{THIS_TR$	EXT_IS_TEXTE
1.19	Example of a WT.	22
2.1	Illustration of the query process with positional ranking and snippet	
	generation.	25
2.2	Diagram of the layer implementation.	26
2.3	A single block in the layered implementation of inverted lists, storing	
-	docIDs, frequencies and positions,	28
3.1	Illustration of the query process from Section 2.1, using the textual data	
	to obtain positions and for the snippet generation	33
11	Illustration of the position extraction process in a UT using exercise	
4.1	inustration of the position-extraction process in a wi , using operation	25
	select	99
5.1	The compression boosting scheme, with a block size of 200 KB	40
59	The compression boosting scheme, with a block size of 200 RD	10
0.2	Example of the decompression process for a document D_1 in the compression	-

5.3 5.4	Space usage for compression boosters, for different block sizes, and PILs. Position-extraction time for scheme VByte + 1z4, for block sizes 5 KB, 10 KB, 50 KB, 200 KB, 500 KB, and 1,000 KB. These block sizes correspond to the points in the plot from right to left. For comparison, we also show	44
5.5	the performance of PILs	45
5.6	Average position-extraction time for compression boosters and PIL, for	46
5.7	$k_1 = 50.$ Average position-extraction time for compression boosters and PIL, for $k_1 = 150$	47
5.8	Average position-extraction time for compression boosters and PIL, for	48
5.9	$\kappa_1 = 300$ Comparison between compression boosters and PILs, for different values	48
	of k_1	49
6.1	Time-space trade-offs for the overall query process for the GOV2 collec- tion. With $k_1 = 50$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step	
6.2	3	54
	tion. With $k_1 = 100$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3	55
6.3	Time-space trade-offs for the overall query process for the GOV2 collec- tion With $k_1 = 150$ and $k_2 \in \{10, 50\}$ including Step 1 Step 2 and	00
6 4	Step 3	56
0.4	tion. With $k_1 = 200$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater	
6.5	than 400 ms/q	57
	than 400 ms/q	58
6.6	Time-space trade-offs for the overall query process for the GOV2 collec- tion. With $k_1 = 50$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step	
6.7	3	59
6.8	Step 3	60
	than 400 ms/q	61

6.9	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 200$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 2. It is important to note that Scheme 1 has grown time greater	
	step 5. It is important to note that scheme 1 has query time greater than 400 ms/c	ດາ
6 10	Time space trade offs for the evenall guary process for the $COV2$ college	02
0.10	Time-space trade-ons for the overall query process for the $GOVZ$ conec- tion. With $h = 200$ and $h \in [10, 50]$ including Step 1. Step 2 and	
	tion. With $\kappa_1 = 500$ and $\kappa_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 2. It is important to note that Scheme 1 has grown time groaten	
	Step 5. It is important to note that Scheme 1 has query time greater	co
6 11	than 400 ms/q .	03
0.11	Time-space trade-ons for the overall query process for the GOV2 collec-	CF.
C 10	tion. With $\kappa_1 = 50$, including Step 1, Step 2	60
0.12	Time-space trade-ons for the overall query process for the GOV2 conec-	
	tion. With $k_1 = 100$, including Step 1, Step 2. It is important to note	<u>ر</u> ۲
0.10	that Scheme 1 has query time greater than 250 ms/q	65
6.13	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 150$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 250 ms/q	66
6.14	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 200$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 400 ms/q.	66
6.15	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 300$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 400 ms/q. $\ldots \ldots \ldots$	67
6.16	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 50$, including Step 1, Step 2	68
6.17	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 100$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 300 ms/q.	68
6.18	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 150$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 400 ms/q. $\ldots \ldots \ldots$	69
6.19	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 200$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 400 ms/q. \ldots	69
6.20	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 300$, including Step 1, Step 2. It is important to note	
	that Scheme 1 has query time greater than 400 ms/q. \ldots	70
A.1	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 50$ and $k_2 = 30$, including Step 1, Step 2 and Step	
	3	76
A.2	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 100$ and $k_2 = 30$, including Step 1, Step 2 and Step	_ .
	3	77
A.3	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 150$ and $k_2 = 30$, including Step 1, Step 2 and Step	
	3	77

A.4	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 200$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.	
	It is important to note that Scheme I has query time greater than 400	70
۸ E	ms/q	18
A.3	Time-space trade-ons for the overall query process for the GOV2 conec- tion. With $k = 200$ and $k = 20$, including Step 1. Step 2 and Step 2	
	tion. With $k_1 = 500$ and $k_2 = 50$, including Step 1, Step 2 and Step 5. It is important to note that Scheme 1 has query time greater than 400	
	ms/a	78
A.6	Time-space trade-offs for the overall query process for the GOV2 collec-	10
11.0	tion. With $k_1 = 50$ and $k_2 = 30$, including Step 1. Step 2 and Step	
	$3. \ldots \ldots$	79
A.7	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 100$ and $k_2 = 30$, including Step 1, Step 2 and Step	
	3	79
A.8	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 150$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.	
	It is important to note that Scheme 1 has query time greater than 400	
	ms/q	80
A.9	Time-space trade-offs for the overall query process for the GOV2 collec-	
	tion. With $k_1 = 200$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.	
	It is important to note that Scheme 1 has query time greater than 400 mg/g	<u>00</u>
A 10	$\operatorname{Time}_{\operatorname{Space}}$ trade offs for the overall every process for the COV2 collection	00
A.10	tion With $k_1 = 300$ and $k_2 = 30$ including Step 1 Step 2 and Step 3	
	It is important to note that Scheme 1 has query time greater than 400	
	ms/q	81
	·/ 1	

Introduction

Motivation

Web search has become an important part of day-to-day life, affecting even the way in which people think and remember things [41]. Indeed, web search engines are one of the most important tools that give access to the huge amount of information stored in the web. The success of a web search engine mostly depends on its efficiency and the quality of its ranking function. To achieve efficient processing of queries, search engines use highly optimized data structures, including inverted indexes [7, 11, 27]. State-of-the-art ranking functions, on the other hand, combine simple term-based ranking schemes such as BM25 [11], link-based methods such as Pagerank [8] or Hits [26], and up to several hundred other features derived from documents and search query logs.

Recent work has focused on *positional ranking functions* [36, 30, 42, 31, 38, 47, 11] that improve result quality by considering the positions of the query terms in the documents. Thus, documents where the query terms occur close to each other might be ranked higher, as this could indicate that the document is highly relevant for the query. To support such positional ranking, the search engine must have access to the position data. This is commonly done by building an index for all term positions within documents, called a *positional index*. The goal is to obtain an index that is efficient in terms of both index size and access time.

As shown in [36], positional ranking can be carried out in two phases. First, a simple term-based ranking scheme (such as BM25) defined over a Boolean filter is used to determine a set of high-scoring documents, say, the top 200 documents. In the second phase, the term positions are used to rerank these documents by refining their score values. (Additional third or fourth phases may be used to do further reranking according to hundreds of additional features [44], but this is orthogonal to our work.) Once the final set of top-scoring documents has been determined (say, the top 10), it is necessary to generate appropriate *text snippets*, typically text surrounding the term occurrences, to return as part of the result page. This requires access to the actual text in the indexed web pages. It is well known [47, 23] that storing position data requires a considerable amount of space, typically about 3 to 5 times the space of an inverted index storing only document identifiers and term frequencies. Furthermore, storing the documents for snippet generation requires significant additional space.

Positional Indexes

Compress positions is a problem where it has been difficult to make much progress. For instance, previous work [47] concludes that term positions in the documents do not follow simple distributions that could be used to improve compression (as is the case of, for instance, docIDs and frequencies). As a result, a positional index is about 3 to 5 times larger than a docID/frequency index, and becomes a bottleneck in index compression. Another important conclusion from [47] is that we may only have to access a limited amount of position data per query, and thus it might be preferable to use a method that compresses well even if its speed is slightly lower. Positions in inverted indexes are used mainly in two applications, phrase searching [27] and positional ranking schemes. In this work we focus in the use of positions to improve the ranking, where the positions of the query terms within the documents are used to enhance the performance of a standard ranking such as BM25. The rationale is that documents where the query terms appear close together could be more relevant for the query, so they should get a better score. Although in this work we focus only on the use of positions to improve ranking, the compression schemes used allow for phrase searching as well. This scenario is left for future work. A recent work on positional indexing is that of Shan et al [40]. They propose to use the *flat position indexes* [13, 16] as an alternative of positional inverted indexes. The result is that docIDs, term frequencies and position data can be stored in space close to that of positional inverted lists, yielding a reduction of space usage. However, this index does not store the text, which makes it less suitable in scenarios where text snippets must be generated.

Thesis contribution

In this thesis we study what are the best ways to organize and store the in-document positions and textual data, in order to efficiently support positional ranking and snippet generation in text search engines. One of our main conclusions is that some compressed representations of the textual data — which are needed to support snippet generation — can also be used to efficiently obtain the term positions needed by positional ranking methods. Therefore, *no positional index is needed in many cases*, thus saving considerable space at little cost in terms of running time.

Our main contributions can be summarized as follows:

1. A study of several trade-offs for compressing position data. Rather than storing a positional index, we propose to compute the term positions from a compressed representation of the text. This is shown in Chapters 4 and 5, following the description explained in Chapter 3. We explore and propose several compression alternatives. Our results significantly enhance current trade-offs between running time and memory space usage, enabling in this way more design choices for web search engines. One of our most interesting results is that both position and textual data can be stored in about 66% the space of current positional inverted indexes.

- 2. A study of several alternatives for compressing textual data, extending the alternatives studied in previous work [23]. In particular, we show that using the scheme in [43] (to compress text and support efficient snippet generation) before using a standard compressor yields good time-space trade-offs, extending the alternatives introduced in [20]. This study is explained in details in Chapter 5 It is important to note that variants of the scheme in [43] have been adopted by some commercial search engines, which makes our results of practical interest.
- 3. We propose several improvements over wavelet trees [25], in order to make them competitive for representing document collections. The details of the improvement are explained in Chapter 4 We show how to improve the compression ratio by compressing the sequence associated to every node of the tree with standard compressors. The result is a practical web-scale compressed self-index that is competitive with the best state-of-the-art compressors.

Chapter 1

Background and Related Work

1.1 Basic Definitions

Let $\mathcal{D} = \{D_1, \ldots, D_N\}$ be a document collection of size N, where each document is represented as a sequence $D_i[1..n_i]$ of n_i terms from a vocabulary $\Sigma = \{1, \ldots, V\}$. Notice that every term is represented by an integer, hence the documents are just arrays of integers.

Definition 1.1 Given a document D_i , we identified it with the unique document identifier (docID) i.

Definition 1.2 Given a term $t \in \Sigma$ and a document $D_i \in \mathcal{D}$, the *in-document positions* of t in D_i is the set $\{j | D_i[j] = t\}$.

1.2 Text Representation

Throughout this work, we assume that all term separators (like spaces, ', ', ';', '.', etc.) have been removed from the text. Also, we assume that all terms in the vocabulary have been represented in a case-insensitive way. This is in order to facilitate the search operations that we need to carry out over the documents in order to compute (on the fly) the positions of a given query term.

To be able to retrieve the original text (with separators and the original case) one can use the *presentation layer* introduced by Fariña et al. [19, Section 4]. This also supports removing stopwords and the use of stemming, among other vocabulary techniques. This extra layer requires extra space on top of that of the compressed text, as well as extra time to obtain the original text. However, this scheme must be used on all the alternatives that we consider in this work, and thus we disregard the overhead introduced by the presentation layer and focus only on the low-level details of compression (but keeping in mind that the original text can still be retrieved).

1.3 Inverted Indexes

The efficiency of query processing in search engines relies on *inverted indexes*. These data structures store a set of *inverted lists* I_1, \ldots, I_V , which are accessed through a *vocabulary table*.

Example 1.1 An example of such data structure can be seen in Figure 1.1, for a document collection of 4 documents $\{D_1, D_2, D_3, D_4\}$.

The list I_t maintains a *posting* for each document containing the term $t \in \Sigma$. Usually, a posting in an inverted list consists of a docID a term frequency, and the in-document positions of the term.

Example 1.2 An example can be seen in Figure 1.2, where the inverted list of a term t_3 is shown. Each posting on the list stores a docID the frequency, and the in-document positions.

In real systems, the docIDs, term frequencies and in-document positions are often stored separately. Indexes whose postings store in-document positions are called *positional inverted indexes*. The inverted lists of the query terms are used to produce the result for the query. Since the query results are usually large, the result set must be ranked by relevance.

1.4 Inverted Index Compression

For large document collections, the data stored in inverted indexes requires considerable amounts of space. Due to the large volume of data on the Web and its rapid growth, data compression has become important [46, 37]. In the case of inverted indexes for web search engines, this effect is reflected in the high space usage of the resulting inverted lists. Hence, the indexes must be compressed. To support efficient query processing (such as DAAT [11], WAND [10] or BMW OR [17]) and effective compression in the inverted lists, we sort them by increasing docID.

Also, to avoid decompressing whole lists at query time, we assume that an inverted list I_t is divided into chunks of 128 documents each. The particular choice of 128 documents per chunk is an implementation issue. Given a chunk of list I_t , the termposition data for all the documents in that chunk are stored in a *separate* chunk of variable size.

Let $d_t[1..|I_t|]$ denote the sorted list of docIDs for the inverted list I_t . Then, we replace $d_t[1]$ with $d_t[1] - 1$, and $d_t[i]$ with $d_t[i] - d_t[i - 1] - 1$ for $i = 2, ..., |I_t|$. In the





Figure 1.1: Example of an invertex index for a document collection, storing just docIDs.



Figure 1.2: Example of an inverted list for the term t_3 .

case of frequencies, every f_i is replaced with $f_i - 1$, since $f_i > 0$ always holds. For the positions, each $p_{i,j}$ is replaced with $p_{i,j} - p_{i,j-1} - 1$. Then these values are encoded with integer compression schemes that take advantage of the resulting smaller integers.

There has been a lot of progress on compressing docID and frequencies, with many compression methods available [48, 11]. Some of them achieve, in general, a very good compression ratio, but at the expense of a lower decompression speed [11], for example Elias γ and δ encodings [18], or Golomb/Rice parametric encodings [24], interpolative encoding [33]. Other methods achieve a (slightly) lower compression ratio, though with much higher decompression speed, for example VByte [45], S9 [4], and PforDelta [52], among others [11]. The best compression method depends on the scenario at hand.

To achieve compression, some of the following special features of posting lists are used:

- docID represent documents that are related to the term, so these are positive integers.
- There are no repeated elements in the lists.
- The lists are sorted by increasing docID, and hence the lists are represented differentially.

We review next the integer compression schemes that will be used along this thesis.

1.4.1 VByte

This method [39] represents an integer using a variable number of bytes, compressing each number separately. To do this, *VByte* use the most significant bit in a byte as a flag. This flag indicates whether the next byte corresponds to the next integer in the segment, or not. So for every byte, just 7 bits are used to represent the number. This is shown in Figure 1.3.



Figure 1.3: Byte layout of VByte.

This flag is called continuator. If the continuator is "0" it means that the current byte is not enough to encode the current integer and hence the encoding must use the next byte. If the continuator is "1" then the current byte is the last used by the encoding. The entire process of coding and decoding is shown below.



Example 1.3 Figure 1.4 shows the 32-bit representation of the integer 167.

Figure 1.4: The 32-bit representation of 167.

Notice that the most-significant bytes 4, 3 and 2 are not used, so that 167 could be represented with less than 4 bytes.

To encode a number in VByte, we try to accommodate its binary encoding in 7 bits. If so, we store the encoding in 7 bits and mark the continuator with a "1". Otherwise the least significant 7 bits are stored in a byte with continuator "0" and proceed with the remaining bits, until no extra bytes are needed to store the encoding.

Example 1.4 Figure 1.5 shows the encoding of 167 using VByte.



Figure 1.5: The encoding of the integer 167 in VBYTE.

The decoding process can be carried out very efficiently using bit-wise and arithmetic operations, as it can be seen in Figure 1.6.

Example 1.5 The decoding process for 167 is illustrated in Figure 1.7.

Table 1.1 shows the numbers of bytes used in VByte for different ranges of integers.

In the best case, for every number we use 1 byte, using a 1/4 of the original space (for 32-bit integers).

```
int decodeVByte (char * compressed) {
    int v = 0;
    int i = 0;
    while (compressed[i] & '128' == '0') {
        v = v || compressed[i];
        v = v<<7;
        i++;
    }
    v = v || compressed[i];
    return v;
}</pre>
```

Figure 1.6: Decoding process for VByte.

Byte 4 (bits from 25 to 32)	Byte 3 (bits from 17 to 24)	Byte 2 (bits from 9 to 16)	Byte 1 (bits from 1 to 8)
	<u></u>		
		Byte 2 VBYTE (bits drom 8 to 14)	Byte 1 VBYTE (bits from 1 to 7)
Byte 4 (bits from 25 to 32)	Byte 3 (bits from 17 to 24)	Byte 2 (bits from 9 to 16)	Byte 1 (bits from 1 to 8)
0 0 0 0 0 0 0	0 0 0 0 0 0 0	0 0 0 0 0 0 0	1 0 1 0 0 1 1 1

Figure 1.7: Decoding process for integer 167.

Table 1.1: Number of bytes used to represent an integer number with VB	Number of bytes used to represent an integer number wit	n <i>VByte</i>
--------------------------------------------------------------------------	---------------------------------------------------------	----------------

Integer range	Number of bytes used
from 0 to 127	1
from 128 to 16383	2
from 16384 to 2097151	3
from 2097152 to 268435455	4
from 268435456 to 42949672952	5

The main features of this method are:

- It is simple to encode and decode integers.
- Uses at most 1 byte more than the minimum amount of bytes to represent a number.
- Uses at least 1 byte per integer encoded.

When the integers to encode are small, using 1 byte is wasteful. An improvement to this is called *variable Nibble* (*VNibble*), which is similar to *VByte* but it works at the nibble level (4 bits). The most significant bit of each nibble is used as the continuator. This can be seen in the Figure 1.8.



Figure 1.8: Example of *VNibble* bytes.

The encoding and decoding processes are the same (or similar) than that of VByte. Table 1.2 shows the number of bytes used for VNibble for different range of integers. When the integers to encode are smaller or equal than 7, this method uses just 4 bits to represent it, which is better.

Integer range	Number of nibbles used
from 0 to 7	1
from $8 \text{ to} 63$	2
from 64 to 511	3
from 512 to 4095	4
from 4096 to 32767	5
from 32768 to 262143	6
from 262144 to 2097151	7
from 2097152 to 16777215	8
from 16777216 to 134217727	9
from 134217728 to 1073741823	10
from 1073741824 to 42949672952	11

Table 1.2: Number of nibbles used to represent an integer number with *VNibble*.

1.4.2 Simple 9

The VByte and VNibble schemes encode each integer individually. This means that in the best scenario VByte is able to encode 4 integers within a 32-bit machine word, and at best 8 integers for VNibble. However, if the integers to encode are snall (e.g. able to be encoded with 1 bit) we waste space.

The basic idea of S9 [4] is to take several consecutive integers in the sequence and try to fit them in a 32-bit machine word, encoding each integer with a fixed equal-size chunk. The size of the chunk is defined as the number of bits needed to encode the largest number in the group. To do so, S9 uses an S9 word of 32-bit. In each S9word, four bits are used as a header to define how many bits will be used to encode the integers, and how many integers will be stored in the S9 word. The remaining 28 bits are used to store the information of the integers encoded. The S9 word layout is shown in Figure 1.9.



Figure 1.9: Layout of an S9 word.

There are nine possible ways of dividing 28 bits into chunks of equal size (the number of cases is the reason behind the name of the method). The meaning of the cases is shown in Table 1.3. To encode integers in an S9 word, we try to fit the maximum

Case	Number of integers stored	Chunk size
0	28	1
1	14	2
2	8	3
3	7	4
4	5	5
5	4	7
6	3	8
7	2	14
8	1	28

Table 1.3: Meaning of the cases in the header of the S9 word.

amount of integers in the 28 available bits . First, we try to fit 28 integers. If that is not possible, then we try with 14, and so on, until we eventually get to the case were only one 28-bit integer can be stored. In the header of the S9 word we store the particular case used for it.

Example 1.6 If we encode the integers 98, 112, 117 and 121, each of these integers

can encoded in binary using 7 bits, which according to Table 1.3 correspond to case 5. The resulting $S9 \ word$ is shown in Figure 1.10.



Figure 1.10: S9 word encoding the group of integers 98, 112, 117 and 121.

To decode an $S9 \ word$, first we obtain the header by means of a bit mask. Then, the case is used on a switch statement where all 9 cases have been hard-coded (as it can be seen in Figure 1.12).

Example 1.7 Figure 1.11 illustrate the decoding process for an *S9 word* that encodes the integers 98, 112, 117 and 121.



Figure 1.11: Decodification process for an *S9 word*, which encodes the group of integers 98, 112, 117 and 121.

The main feature of this method that makes it highly efficient at decoding time is that all the integers in an *S9 word* are coded in the same amount of bits. Hence, the decoding process can be optimized by hard-coding all cases, as we already said.

An improvement of this technique is to use all the 16 possible cases that we can have with headers of 4 bits. This improvement is called S16 [49].

1.4.3 PforDelta

The PforDelta encoding [52] divides an inverted list into *blocks* of, usually, 128 integers each. To encode the integers within a given block, it gets rid of a given percentage—usually 10%—of the largest integers within the block, and stores them in a separate memory space. These are the *exceptions* of the block. Next, the method finds the largest

```
void decodeS9 (int s9word, int * result) {
   char case = getHeader(s9word);
   switch(case) {
    case 0 :
   result[0] = (s9word & 0x08000000) >> 27;
   result[1] = (s9word & 0x04000000) >> 26;
    . . .
   result[26] = (s9word & 0x0000002) >> 1;
   result[27] = (s9word & 0x0000001);
   break;
    case 1 :
    result[0] = (s9word & 0x0c000000) >> 26;
   result[1] = (s9word & 0x03000000) >> 24;
   result[12] = (s9word & 0x000000c) >> 1;
   result[13] = (s9word & 0x0000003);
   break;
    . . .
    case 8 :
   result[0] = (s9word & 0xffffff);
   break;
  }
```

Figure 1.12: Decoding process for an S9 word.

remaining integer in the block, let us say x, and represents each integer in the block in binary using $b = \lceil \log x \rceil$ bits. Though the exceptions are stored in a separate space, we still maintain the slots for them in their corresponding positions. This facilitates the decoding process. For each block we maintain a *header* that stores information about the compression used in the block, e.g., the value of b.

An ilustration of a PforDelta block is shown in Figure 1.13. To retrieve the positions



Figure 1.13: Layout of PforDelta block.

of the exceptions, we also store the position of the first exception in the header. In the slot of each exception, we store the offset to the next exception in the block. In the last exception, we store a '0'. This forms a linked list with the exception slots. In case that b is too small and we cannot accommodate the offset to the next exception (because the

offset to the next one cannot be accommodated within b bits), the algorithm forces to add extra exceptions between two original exceptions. This increases the space usage when the lists contain many small numbers.

Example 1.8 A PforDelta block that is encodes 128 integers is shown in Figure 1.14. The exceptions are 78, 110, 160 and 91.



Figure 1.14: PforDelta block, encoding a group of integers, with exeptions 78, 110, 160 and 91.

To decompress a block, we first take b from the header, and invoke a specialized function that obtains the b-bit integers. Each b has its own extracting function, so they can be hard-coded for high efficiency. Once we decode the integers, we traverse the list of exceptions of the block, storing the original integers in their corresponding positions. This step can be slower, yet it is carried out just for 10% of the block. The decoding process implementation is shown on Figure 1.16.

In typical implementations of PforDelta, the header is implemented in 32 bits, since we only need to store the values of b (in 6 bits, since $1 \le b \le 32$) and the position of the first exception (in 7 bits, since the block has 128 positions). PforDelta has shown to be among the most efficient compression schemes [48], achieving a high decompression speed.

Example 1.9 Figure 1.15 shows the size and the decoding process for the PforDelta block previously mentioned.

1.5 Snippet Generation

Besides providing a ranking of the most relevant documents for a query, search engines must show query snippets and support accessing the "in-cache" version of the documents. To achieve this, web search engines must store a copy of document collection.

Example 1.10 Figure 1.17 shows an example of user-query result, and the result snippets.

Each snippet shows a relevant portion of the result document, in order to help the user judge its likely relevance before accessing it. Turpin et al. [43] introduce a method to compress the text collection and support fast text extraction to generate snippets. However, to achieve fast extraction, they must use a compression scheme that uses more



Figure 1.15: Decoding a PforDeltablock.

```
/*Decode 128 integers of b bits from data and write them in results.*/
void getNumbers(void * data, char b, int * results);
/* Decode the pfordelta block*/
void decodePFD (void * block, int * results){
 int header = ((int *)block)[0];
 char b = (header & 0xfc000000) >> 26;
 int exceptionPos = (header & 0x7f) ;
 void * data = (void *)(((int *)(block)) +1);
 int exceptionStart = ((b * 128) + 32) / 32;
 int * exceptions = (int *)(((int *)(block))+(exceptionStart));
 getNumbers(data, b, results);
 int i = 0;
 while( results[exceptionPos] != '0'){
    int offset = results[exceptionPos];
   results[exceptionPos] = exceptions[i];
   exceptionPos += offset;
    i++;
 }
```

Figure 1.16: Decoding process for PforDelta block.



Figure 1.17: An example of user-query and the result snippets.

space than usual compressors. In a more recent work, Ferragina and Manzini [23] study how to store very large text collections in compressed form, such that the documents can be accessed when needed, and show how different compressors behave in such a scenario. One of their main concerns was how compressors can capture redundancies that arise very far apart in very long texts. Their results show that such large texts can often be compressed to just 5% of their original size, which was surprising since usual compression ratio is of about 20%.

1.5.1 Dictionary Compressors

The compressors based in dictionary, do not use codes of variable lengths like the methods previously mentioned. These use a dictionary to access the data previously read and in this way detect regularities in the data. For each element (a symbol or a chain of symbols) an entry is created, which indicates how interpret that element with the dictionary information. Dictionary quality and the entry features, allows detection and representation of large symbols chains , hence reaching high compression rates.

In this category the works of Abraham Lempel and Jacob Ziv [50, 51], known as LZ77 and LZ78, are the starting points for the compressors based in dictionaries. The ideas proposed in those works are the base for many new techniques, which try to improve the originals. In this thesis we will work with lzma, snappy and lz4 as LZ77 variants.

1.5.2 LZ77

LZ77 or sliding window, was the first method devised for Lempel and Ziv [50], which basically tries to detect redundant sequences of symbols in the data. To do this, the method keeps a buffer of (fixed size), divided in two sections, which are updated as the data is read (from here, the name of sliding window).

Example 1.11 The Figure 1.18 shows how the buffer and its two sections are updated while the data is read.



Figure 1.18: Example of how the buffer of LZ77 is updated with the text $T = \{THIS_TEXT_IS_TEXTED_EXT_THE\}.$

The right part of the buffer is used to read the new symbols, and is known as *lookahead buffer*. The left part of the buffer is called *search buffer* and correspond to the dictionary, due than that keeps the symbols that have been read previously.

The encoding process is computed as follow: we search for longest prefix of the *lookahead buffer* in the *search buffer*, then we store a lempelziv entry which has 3 values.

- Offset: Distance between the symbol and its copy in the *search buffer*.
- Prefix Length: Length of the copy.
- New symbol : Symbol in the *lookahead buffer* that difference the prefix from the copy.

Then, the buffer advances in the data one more symbols than Prefix Length. If the first symbol of the *lookahead buffer* cannot be found in the *search buffer*, which means that the length of the longest prefix is 0, the lempelziv entry is (0, 0, symbol) and is called literal.

The decoding process is very simple, we read sequentially the lempelziv entries, which contain (*offset*, *length*, *symbol*). If the entry is a literal, then we simply write the *symbol* in the output, else from the end of the output, we look back *offset* positions and from there copy the next consecutive *length* symbols at the end of the output, finally

write the *symbol*. We repeat the process while having lempelziv entries. In general the size of the *search buffer* is 2^{16} bytes (64 KB), while the *lookahead buffer* is 2^{8} bytes (256 B).

The following methods, lzma, snappy and lz4, basically share the coding process, which is find the longest prefix of the *lookahead buffer* in the dictionary (or *search buffer*), their differences are the dictionary size and the internal structure of the entries.

The base of the decoding process is very similar for all the alternatives mentioned, which is: read the entries, detect if is a copy or a literal, then write it. The minor differences are that some methods encode the entries not aligned to bytes, this is to use the exact space in the output.

1.5.3 Lzma

The method lzma, which is the principal algorithm of $7\text{-}zip^1$, is a variant of LZ77 designed to provide a high compression rate and a fast decompression speed.

The coding and decoding process are very similar to the LZ77 process (as explained in the previous section).

The main differences are the structure of the entries, and the size of the dictionary. first we explain the entry differences, lzma uses 2 types of entries, instead of 1, if the first bit in the entry is a '0' then the entry is a literal, else a copy, also the copy just store *offset* and *length*, the layout of the entries is in bits.

• Literal entry: stores the symbol, coded in range encoding [29].

- Layout : 0 + code of the symbol.

• Copy entry: stores the pair (*offset*, *length*), but depending of the values characteristic can use 6 types of representation.

Type 1: Simple copy, stores (*offset*, *length*), layout : 1 + 0 + length + offset.

- Type 2: Short copy, where *offset* is the last used and length is 1, layout: 1 + 1 + 0 + 0.
- Type 3: : Long copy [0], uses the last offset used, and can store a $length \in \{2, 273\}$, layout: 1+1+0+1+len.
- Type 4: : Long copy [2], uses the second last offset used, and can store a length $\in \{2, 273\}$, layout: 1+1+1+0 +len.
- Type 5: : Long copy [3], uses the third last offset used, and can store a $length \in \{2, 273\}$, layout: 1+1+1+1+0+len.
- Type 6: : Long copy [4], uses the forth last offset used, and can store a $length \in \{2, 273\}$, layout: 1+1+1+1+1+1.

And the layout of len depends of its value, as can be seen below.

¹http://www.7-zip.org/

- $len \in \{2, 9\}$: 0 + 3 bits.
- $len \in \{10, 17\}$: 1+0+3 bits.
- $len \in \{18, 273\}$: 1+1+8 bits.

One inconvenience, is that the entries are not aligned to bytes, so to read them implies several complex bit operations which increase the decompression time.

Second, the size of the dictionary is huge, dictionary size goes from 2^{23} (8 MB) to 2^{30} (1024 MB). Allowing the detection of very distant copies, yet the task of search them increases its complexity. To achieve a good compression time, 1zma uses hash tables that allow find, fast enough, the last appearances of the symbols in the dictionary, but due to the huge window size, the tables cannot store all the repetitions, producing that some symbol repetitions will not be encoded. Yet even if some regularities are not coded, the huge size of the dictionary allows larger chains of symbols to be detected and encoded, achieving higher compression ratio.

1.5.4 Snappy

The snappy [1], its other LZ77 variant, developed by $google^2$, with focus in compression speed and a reasonable compression ratio.

The coding and decoding process is the same as described above. The first difference is that instead of use a sliding window, **snappy** stores the uncompressed size of the data (max value of $2^{32} - 1$), coded in VByte, at the beginning of the output, then stores the compressed data as a stream of entries. This means that **snappy** compress a big block of information instead of sliding in the data.

The compressed stream, uses the same scheme than 1zma. So basically store 2 types of elements, literals and copies, both elements uses its first 2 bits to define its type (00 : literal, 01-10-00 : copies). Other difference is that in **snappy** the literals, can contain several symbols (from 1 to $2^{32} - 1$) instead of just 1, so after the firsts 2 bits, a literal stores the amount of symbols that contains (from 1 to $2^{32} - 1$) aligned to bytes, then stores the symbols. there are 2 types of literals:

- Short literals: Stores up to 60 symbols, using 6 bits to represent the number of literals (the values 60 -64, are reserved), layout : 0+0+[6 bits] + literals.
- Long literals: Stores up to 2^{32} symbols, the values from 60 to 64 of the 6 bits after the initial 00, defines how many bytes are used to represent the amount of literals, layout : 0+0+[6 bits]+[from 1 to 4 bytes]+[literals]

In the case of the copies, the entries just store the values *offset* and *length* of the copy, as lzma. Depending of the values for *offset* and *length*, its the type of entry copy that can be used, there are 3 types, which differs in size and utility :

• Type 1: Uses 2 consecutive bytes using

 $^{^{2}}$ http://www.google.com

- First 2 bits: (00).
- Next 3 bits: for the *length* of the copy (from 4 to 11).
- Next 11 bits : for the *offset* of the copy (from 0 to $2^{11} 1$).
- Type 1: Uses 3 consecutive bytes using
 - First 2 bits: (10).
 - Next 6 bits: for the *length* of the copy (from 1 to 64).
 - Next 16 bits : for the *offset* of the copy (from 0 to $2^{16} 1$).
- Type 1: Uses 4 consecutive bytes using
 - First 2 bits: (11).
 - Next 6 bits: for the *length* of the copy (from 1 to 64).
 - Next 32 bits : for the *offset* of the copy (from 0 to $2^{32} 1$).

Every element in the data stream is aligned to bytes, which allows fast read and write in the data stream.

The encode process is a little different from the previous, due that the literals can contain several symbols, when the copies contain less than 4 symbols then is reported as a literal.

1.5.5 Lz4

1z4 [2], its a method designed to achieve fast decompression speed and a high compression ratio, and be simple.

The difference with the two previous method is that 1z4 uses just 1 type of entry, as the original LZ77, so the coding and the decoding process remain simple. To encode, search for a prefix, code the result in the entry. To decode, read the entries, decode the information and write the result in the output.

The layout of the entry is as follow:

- Token: uses 1 byte, divided in 2 groups of 4 bits, t_1 and t_2 .
 - $-t_1$: Indicates the amount of literals in the entry, when $t_1 = 0$ means that the entry has no literals then read t_2 , if the value of $t_1 = 15$ means that other bytes are needed to decode the length of the literals, that byte is stored after the Token byte, and its information is added to t_1 , when the next byte is 255 then other byte is added to the codification, the process continue while bytes are needed (it has no limits in the amount of literals that can store).
 - t_2 : Indicates the value of the length of the copy, when $t_1 = 0$ means that the entry is not a copy (just a literal), if the value of $t_1 = 15$ means that other bytes are needed to decode the length of the copy, that byte is stored after the *offset*, and its information is added to t_2 , when the next byte is 255 then other byte is added to the codification, the process continue while

bytes are needed, notice that 1z4 just stores offset ≥ 4 , then it add 4 to the result length.

- Length of the literals: Can use from 0 to n bytes, those bytes stores the length of the literals when $t_1 \ge 15$.
- Literals: Can use from 0 to n bytes, depending of the amount of literals that are coded in the entry.
- Offset: Can use 2 bytes (if $t_2 = 0$, then this value do not exist), and represent values $\in \{0, 65535\}$. It represents the *offset*. Note that 0 is an invalid value, never used. 1 means "current position 1 byte". 65536 cannot be coded, so the maximum offset value is really 65535.
- Length of copy: Can use from 0 to n bytes, those bytes stores the length of the copy when $t_2 \ge 15$.
- Copy length: Can use from 0 to n bytes, depending of the value of the *length*.

This scheme provides high adaptability to the context, also its aligned to bytes, which means fast access to the information, additionally the entry structure allows a very fast decompression when the amount of literals and the copy length are ≤ 15 .

1.6 Compressed Text Self-Indexes

Succinct or compressed data structures use as little space as possible to support operations as efficiently as possible [35]. Thus, large data sets (like graphs, trees, and text collections) can be manipulated in main memory, avoiding the secondary storage. In particular, we are interested in compressed data structures for text sequences. A compressed self-index is a data structure that represents a text in compressed space, supports indexed search capabilities on the text, and is able to obtain any text substring efficiently [35]. They can be seen as compression tools with indexed search capabilities. The following operations have been the building block of many solutions in succinct data structures.

Definition 1.3 Given a sequence T[1..n] over an alphabet $\Sigma = \{1, \ldots, V\}$, we define operation $\operatorname{\mathsf{rank}}_c(T, \mathbf{i})$, for $c \in \Sigma$, as the number of occurrences of c in $T[1..\mathbf{i}]$. Operation $\operatorname{\mathsf{select}}_c(T, j)$ is defined as the position of the j-th occurrence of c in T.

1.7 Wavelet trees

A wavelet tree [25] (WT for short) is a succinct data structure that supports rank and select operations, among many virtues [21, 35, 34]

A WT representing a text T is a balanced binary search tree where each node v represents a contiguous interval $\Sigma^v = [i..j]$ of the sorted set Σ . The tree root represents the whole vocabulary. Σ^v is divided at node v into two subsets, such that the left child

 v_l of v represents $\Sigma^{v_l} = [i..\frac{i+j}{2}]$, and the right child v_r represents $\Sigma^{v_r} = [\frac{i+j}{2} + 1..j]$. Each tree leaf represents a single vocabulary term.

Hence, there are V leaves and the tree has height $O(\log V)$. For simplicity, in the following we assume that V is a power of two.

Let T^v be the subsequence of T formed by the symbols in Σ^v . Hence, $T^{root} = T$. Node v stores a bit sequence B^v such that $B^v[l] = 1$ if $T^v[l] \in \Sigma^{v_r}$, and $B^v[l] = 0$ otherwise. Given a WT node v of depth i, $B^v[j] = 1$ iff the i-th most-significant bit in the encoding of $T^v[j]$ is 1. In this way, given a term $c \in \Sigma$, the corresponding leaf in the tree can be found by using the binary encoding of c. Every node v stores B^v augmented with a data structure for rank/select over bit sequences [35]. The number of bits of the vectors B^v stored at each tree level sum up to n, and including the data structure every level requires n + o(n) bits. Thus, the overall space is $n \log V + o(n \log V)$ bits [25, 35].

Example 1.12 In Figure 1.19 we show an example WT for the text "CDEBFEGABBFCH-HCDEABG".



Figure 1.19: Example of a WT.

1.7.1 Supporting Operations

Since a WT replaces the text it represents, we must be able to retrieve T[i], for $1 \le i \le n$. The idea is to navigate the tree from the root to the leaf corresponding to the unknown T[i]. To do so, we start from the root, and check if $B^{root}[i] = 0$. If so, the leaf of T[i] is contained in the left subtree v_l of the root. Hence, we move to v_l looking for the symbol at position $\operatorname{rank}_0(B^{root}, i)$. Otherwise, we move to v_r looking for the symbol at position $\operatorname{rank}_1(B^{root}, i)$. This process is repeated recursively, until we reach the leaf of T[i], and runs in $O(\log V)$ time as we can implement the rank operation on bit vectors in constant time.

To compute $\operatorname{rank}_c(T, i)$, for any $c \in \Sigma$, we proceed mostly as before, using the binary encoding of c to find the corresponding tree leaf. On the other hand, to support $\operatorname{select}_c(T, j)$, for any $c \in \Sigma$, we must navigate the upward path from the leaf corresponding to term c. Both operations can be implemented in $O(\log V)$ time; see [35] for details.

1.7.2 Analysis of Space Usage

The space required by a WT is, in practice, about 1.1–1.2 times the space of the text [14]. In our application to IR, this would produce an index larger than the text itself, which is excessive. To achieve compression, we can generate the Huffman codes for the terms in Σ (this is a *word-oriented Huffman coding* [32]) and use these codes to determine the corresponding tree leaf for each term. Hence, the tree is not balanced anymore, but has a Huffman tree shape [14] such that frequent terms will be closer to the tree root than less frequent ones. This achieves a total space of $n(H_0(T) + 1) + o(n(H_0(T) + 1)))$ bits, where $H_0(T) \leq \log V$ is the zero-order empirical entropy of T [28]. In practice, the space is about 0.6 to 0.7 times the text size [14]. However, notice that we have no good worst-case bounds for the operations, as the length of the longest Huffman code assigned to a symbol could be O(V).

1.7.3 Self-Indexes for IR Applications

There have been some recent attempts to apply alternative indexing techniques, such as self-indexes, in large-scale IR systems. In particular, we mention the work by Brisaboa et al. [9] and Arroyuelo et al. [6, 5]. The former [9] concludes that WT are competitive when compared with an inverted index for finding all the occurrences of a given query term within a single text. The latter [6, 5] extends [9] by supporting more IR-like operations on a WT. The result is that a WT can represent a document collection using $n(H_0(T) + 1) + o(n(H_0(T) + 1))$ bits while supporting all the functionality of an inverted index. The experimental results in [6] compare with an inverted index storing just docIDs, which of course yields a smaller index. However, WTs also store extra information, such as the term frequencies and, most important for us here, the compressed text and thus the term-position data.

Recent work [23] also tried to use (among other alternatives) a compressed self-index to compress web-scale texts, in order to allow decompression of arbitrary documents. Their conclusion is that compressed self-indexes still need a lot of progress in order to be competitive with standard compression tools, both in compression ratio and decompression speed. A contribution of this thesis is a compressed self-index that is able to store web-scale texts and is competitive with the best state-of-the-art compressors. We think that this is a step forward in closing the gap between theory and practice in this area [22].

Chapter 2

Indexing for Positional Ranking: Classical Solutions

The ranking process in web search engines is a fundamental step of the search process. This is because it allows the user to find the documents that are more appropriate for their information needs. Although traditionally the tf-idf ranking model is the most used, nowadays search engines use more sophisticated ranking functions. One such example is that of positional ranking functions [11, 46], where the position of the query terms within the resulting documents are used to rank the documents: the closer they are the higher the ranking. The rationale behind positional ranking is that documents, where the query terms occur close to each other might be ranked higher, as this could indicate that the document is highly relevant to the query.

Another important aid used by nowadays web search engines is that of result snippets: a relevant portion of the document, which is used to help the user to decide about the relevance of a result.

In this chapter we review the state of the art for position indexing and snippet generation, and show experimental results that will be used to determine the improvements introduced by our contributions.

2.1 Basic Query Processing Steps for Positional Ranking and Snippet Extraction

From now on we assume a search engine where positional ranking is used to score documents, and where snippets must be generated for the top-scoring documents. Thus, solving a query involves the following steps:

1. Query Processing Step: Given a user query, use an inverted index to obtain the top- k_1 documents according to some standard query processing approach (e.g.,

DAAT) and ranking function (e.g., BM25).

- 2. Positional Ranking Step: Given the top- k_1 documents obtained on the previous step, obtain the positions of the query terms within these documents. Then re-rank the results using a positional ranking function [11, 47].
- 3. Snippet Generation Step: After the re-ranking of previous step, obtain snippets of length s for the top- k_2 documents, for a given $k_2 \leq k_1$.

An illustration of this process is shown in Figure 2.1, where term positions are obtained from a positional index (to be described next in this chapter) and text snippets are obtained from the compressed text database.



Figure 2.1: Illustration of the query process with positional ranking and snippet generation.

For instance, $k_1 = 200$ (as in [47]) and $k_2 = 10$ (as in most commercial search engines) are typical values for the query parameters. We assume s = 10 in this thesis. The different values for these parameters should be chosen according to the trade-off between query time and search effectiveness that we want to achieve. Step 2 is usually supported by a positional inverted index [27, 11, 47]. Step 3 is supported by compressing the document collection and supporting the extraction of arbitrary documents. We will describe these processes in detail in this chapter. The focus of this thesis is on alternative ways to implement the last two steps. Next, we evaluate experimentally the current solutions on the state of the art of positional indexing.

2.2 Experimental Setup and Data

Along this thesis, we use the following experimental setup. The machine where we ran our experiments is an HP ProLiant DL380 G7 (589152-001) server, with a Quadcore Intel(R) Xeon(R) CPU E5620 @ 2.40GHz processor, with 128KB of L1 cache, 1MB of L2
cache, 2MB of L3 cache, and 96GB of RAM, running version 2.6.34.8-68.fc13.i686.PAE of Linux kernel.

We use the TREC GOV2 document collection, with about 25.2 million documents and about 32.86 million terms in the vocabulary. We work just with the text content of the collection (that is, we ignore the html code from the documents). This requires about 130,048MB in ASCII format. When we represent the terms as integers, the resulting text uses 91,634 MB. We use a subset of 10,000 random queries from the TREC 2006 query log. All methods were implemented using C++, and compiled with g++ 4.4.5, with the full optimization flag -05.

To compare the different schemes propose in this thesis, we use two criteria. The first one is the compression ratio, which is the size of the compressed structure divided by the size of the uncompressed structure multiplied by 100. The second one is the time per query, which is the average time to answer a query with the defined scheme.

2.3 The Baseline: Positional Inverted Lists and Compressed Textual Data

This section describes and evaluates the baseline approaches to support term-position indexing and snippet extraction. These will be the starting points for the proposals of this thesis.

2.3.1 Supporting the Query-Processing Step

As we have said, query processing is supported by inverted indexes. To achieve competitive space and time with an inverted index, inverted lists are usually divided into chunks (of usually 128 docIDs each) and implemented using layers. The first layer stores the docIDs and the second one stores frequencies. See Figure 2.2 for an illustration.



Figure 2.2: Diagram of the layer implementation.

An advantage of this layered scheme is that every layer can be compressed differently, using the best alternative for each layer. In our experimental setting, these alternatives are PforDelta (Section 1.4.3) for the docID layer and S16 (Section 1.4.2) for the frequency layer, just as in [47]. Another advantage is that data is decompressed just when necessary at search time. A chunk in the second layer is decompressed only when a given docID in the first layer is a candidate to be in the top- k_1 . This process results in a fairly competitive query time to obtain the top- k_1 documents.

In our experiments the index for docIDs and frequencies for the GOV2 collection requires 47,854 MB of space in uncompressed form. Using PforDeltaompression for docIDs and S16 for frequencies, the space usage is reduced to 9,739 MB, achieving a compression ratio of about 20%. Notice also that the index uses 0.11 times the space of the imput text.

For query processing, we assume the well-known Document-at-a-Time (DAAT) approach, with the BM25 ranking [11] to obtain the top- k_1 most relevant documents in Step 1. BM25 is one of the most used and effective term-frequency scoring function. Given a query Q with terms $q_1, ..., q_n$ and a document \mathcal{D} , the score is computed as follows:

$$score_{BM25}(\mathcal{D}, Q) = \sum_{i=1}^{n} IDF(q_i) \cdot \frac{f(q_i, \mathcal{D}) \cdot (a+1)}{f(q_i, \mathcal{D}) + a \cdot (1 - b + b \cdot \frac{|\mathcal{D}|}{avgdl})}$$

where $f(q_i, D)$ is the frequency of term q_i in document \mathcal{D} , $|\mathcal{D}|$ is the number of words of document \mathcal{D} , and avgdl is the average document length in the text collection. Parameters a and b are for advanced optimization, typically $a \in [1.2, 2.0]$ and b = 0.75. Finally, $IDF(q_i)$ is the IDF (inverse document frequency) weight of the query term q_i . It is usually computed as:

$$IDF(q_i) = \log \frac{N - n(q_i) + 0.5}{n(q_i) + 0.5},$$

where N is the number of documents in the collection, and $n(q_i)$ is the number of documents that contain term q_i .

In Table 2.1 we show the average query time (in milliseconds per query) for the initial query processing Step 1, obtained from our experiments. We show results for two types of queries: traditional AND queries (using DAAT query processing) and the BMW OR approach from [17], which is one of the most efficient current solutions for disjunctive (OR) queries.

Notice that the query time for AND is almost constant (within two decimal digits) with respect to k_1 . The process to obtain the top- k_1 documents uses a heap (of size k_1). However, operating the heap takes negligible time, compared to the decompression of docIDs and the DAAT process. BMW OR, on the other hand, is an early-termination technique, and thus k_1 impacts the query time.

	, 0	
$\operatorname{top-}k_1$	DAAT AND (ms/q)	BMW OR [17] (ms/q)
50	14.75	35.70
100	14.77	43.39
150	14.80	47.90
200	14.81	51.74
300	14.81	58.19

Table 2.1: Experimental results for the initial query processing step (Step 1) for AND and OR queries. In both cases, BM25 ranking is used.

2.3.2 Supporting the Positional-Ranking Step

Positional inverted lists (PILs, for short) are the standard approach for indexing indocument position data in search engines [27, 11, 7].

PILs add another layer to the layered structure described in Section 2.3.1. See Figure 2.3 for a reference. As for docIDs and frequencies, we divide the positions in chunks, which now have variable size: if a docID chunk stores n docIDs, the corresponding positional chunk stores as many positions as as the sum of all the frequencies stored in the frequency chunk. Hence, because of the variable size of the positional chunks, each chunk of the frequency layer stores a pointer to the corresponding positional chunk.



Figure 2.3: A single block in the layered implementation of inverted lists, storing docIDs, frequencies and positions.

To obtain positional data at query time, after we obtain the top- $k_1 \text{ docIDs}$ for a given query, we identify the positional chunks containing the desired positional index entries. Then, these positional chunks are *fully decompressed*, and the corresponding positions are obtained. A drawback here is that we need to decompress the entire positional chunk, even if we only need a single entry in it. Thus, we might end up decompressing, in the worst case, k_1 positional chunks in each of the inverted lists involved in the query.

After obtaining the positions of all query terms within the top- k_1 , we proceed to re-rank them using a positional ranking score. Particular positional ranking functions are out of the scope of this thesis: any function than uses term positions could be used. In our experiments, we use the scoring model proposed by S. Büttcher and C. Clarke [12, 38], which is defined as follows. For each document we fetch the positions of all the query terms within the top- k_1 documents. Each term is associated an accumulator that contains the term proximity score, which is computed as follow: for every pair of consecutive terms T_i , T_j ($T_i \neq T_j$) in the query Q, we obtain all the positions of those terms (P_{T_i}, P_{T_j}) in the document \mathcal{D} and modify the accumulators for T_i and T_j as:

$$acc(T_{i}) \leftarrow acc(T_{i}) + w_{Tj} \cdot (\operatorname{dist}(P_{Ti} + P_{Tj}))^{-2}$$

 $acc(T_{j}) \leftarrow acc(T_{j}) + w_{Ti} \cdot (\operatorname{dist}(P_{Ti} + P_{Tj}))^{-2}$

where w_{T_x} is the IDF score of T_x (the same as in BM25). After computing all the accumulators, the score of the document \mathcal{D} is:

$$score_{BM25+Pos}(\mathcal{D}) = score_{BM25} + \sum_{T \in Q} \min(1, w_t) \cdot \frac{acc(T) \cdot (a+1)}{acc(T) + b},$$

where a and b are the same as in the BM25 formula.

It is important to remark that the space needed for in-document positions is usually high, since we must store the positions of all the occurrences of a term. For instance, in our experimental setting the total number of positions to be stored is 23,991,731,648. This is about twice the number of docIDs and frequencies that need to be stored in the inverted lists (about 12,512,013,184 integers). In uncompressed state, this means about 91,521 MB for positions, which is about twice the space used by the uncompressed docIDs and frequencies. Moreover, position data usually cannot be compressed as easily as in the case of docID and frequencies, which makes the situation even worse.

As we already said, the access pattern for position data is much sparser than that for docID and frequencies, since positions must be obtained only for the top- k_1 documents. Thus, just a few positions are decompressed from the PIL in each query. Given this sparse access pattern and the high space requirement of positions (as discussed above), it is better to use compression methods with a good compression ratio, like Golomb/Rice compression. These are slower to decompress, yet the fact that only a few positions are decompressed should not impact considerably in the overall query time. According to [47] and further personal communications with the main author of that work, the most effective compression techniques are Rice and S16. In Table 2.2 we show experimental results for obtaining positions with the baseline PILs, using the two compression schemes selected, Rice and S16. We also show query times for different values of k_1 , namely 50, 100, 150, 200 and 300 (the experiments in [47] only use $k_1 = 200$).

Compression Scheme	Space Usage	Average position extraction time (ms/q)				s/q)
	(MB)	$k_1 = 50$	$k_1 = 100$	$k_1 = 150$	$k_1 = 200$	$k_1 = 300$
PIL Rice	28,373	1.28	2.22	3.05	3.27	5.57
PIL S16	31,338	0.74	1.12	1.43	1.75	2.51

Table 2.2: Experimental results for extracting term-position data (Step 2).

As we can see, Rice requires only about 90% the space of S16, but takes twice as much time. Comparing the query times of Step 2 for Rice and S16 with the query times

of Step 1 in Table 2.1, we can see that position extraction is a small fraction of the overall time. Hence, we can use Rice to compress PILs and obtain a better space usage with only a minor increase in query time. For Rice, PILs use 2.91 times the space of the inverted index that stores docIDs and frequencies. For S16, this number is 3.22.

2.3.3 Supporting the Snippet Generation Step

In order to support snippet generation, web search engines must store a copy of the entire web in their servers. This requires, of course, considerable space. To compress the text collection and support decompressing arbitrary documents, a simple alternative that is used by several state-of-the-art search engines — for instance Lucene [15] — is to divide the whole collection into smaller text blocks, which are then compressed separately. The block size offers a time-space trade-off: larger blocks yield better compression, although decompression time is increased. Given the popularity [15, 23] and simplicity of this approach, we use it as the baseline for the compressed text.

Table 2.3 shows experimental results for the baseline for compressed textual data. Just as in [23], we divide the text into blocks of 0.2MB, 0.5MB or 1.0MB, and compress each block using different standard text compression tools. In particular, we show results for 1zma (which gives very good results in [23], so it serves as a comparison with the results in that work) and Google's snappy compressor [1], which is an LZ77 compressor that is optimized for speed rather than compression ratio. We also show results for 1z4 [2], which is another variant of LZ77 compression improved for speed. These three compressors offer the most interesting trade-offs among the alternatives we tried.

Compressor	Block size	Space usage	Compression Ratio	Average snippet extraction time (for different values of k_1		traction time (ms/q) t values of k_1
	(MB)	(MB)		$k_2 = 10$	$k_2 = 30$	$k_2 = 50$
lzma	0.2	14,987	16.35	29	84	136
	0.5	$13,\!489$	14.72	63	181	292
	1.0	$12,\!682$	13.84	117	335	540
snappy	0.2	$34,\!576$	37.73	2	6	9
	0.5	$34,\!426$	37.57	5	14	23
	1.0	34,390	37.53	10	28	46
lz4	0.2	$30,\!405$	33,18	1	4	7
	0.5	$30,\!070$	$32,\!81$	3	10	16
	1.0	29,953	32,68	7	19	31

Table 2.3: Experimental results for the snippet extraction phase (Step 3).

As it can be seen, 1zma achieves much better compression ratios than snappy and 1z4. Also, in all cases the compression ratio improves as we increase the block size.

This is because more text regularities can be detected. If we now consider the text as a whole (i.e., without the block structure, which is not useful for snippet generation, but for archival purposes) the compressed space achieved for the whole text is 8,133 MB for 1zma, 27,388 MB for snappy and 29,808 MB for 1z4. Notice that the compression ratio for 1zma is similar to that reported in [23]. Even when snappy uses less space than 1z4 compressing the whole text, separating the web in blocks shows that 1z4 has better performance in space and time. The differences in extraction time are also considerable, with 1z4 being faster than the other two alternatives, especially against 1zma. Note that [23] reports a decompression speed of about 35MB/sec for 1zma. However, to obtain a given document we must first decompress the entire block that contains it. Hence, most of the 35MB per second do not correspond to any useful data. In other words, this does not measure effective decompression speed for our scenario, and thus we report per-query times rather than MB/s for both methods.

Finally, notice that when we use lzma for the text, the space usage for the whole collection is much smaller than the space of positions. This is because compressors can take advantage of the text regularities, wherever they are. Positions, on the other hand, are stored separately for each term, so inter-term text regularities cannot be detected and compressed. A relevant question here is: How can one represent terms positions, in such a way that these text regularities are conserved? We answer this question as one of the most important contributions of this thesis.

Chapter 3

Computing Term Positions from Textual Data

The main conclusion from the previous chapter is that position and text data have high space requirements. The size of the two structures combined is 7 times the space used by the docID and frequency inverted index, which becomes a bottleneck [47].

This thesis focuses on alternative approaches to perform the aforementioned two-step document ranking process and the query snippet-generation phase, which are Step 2 and Step 3 in Section 2.1. The aim is to optimize both space and query processing time. One important feature of position data is that it only needs to be accessed for a limited number of promising documents, say a few dozens or hundreds of documents. This access pattern differs from that for document identifiers and term frequencies, which are accessed more frequently, making access speed much more important. For position data, on the other hand, we could consider somewhat slower but smarter alternative representations without losing too much efficiency at query time [47].

In this thesis, we push this idea further and consider not storing the position data (i.e, the positional index) at all. Instead, we propose to compute positions on the fly from a compressed representation of the text collection. We will study two alternative approaches to compressing the text collection:

- 1. Wavelet trees [25], which are succinct data structures from the combinatorial pattern matching community.
- 2. Compressed document representations, that support fast extraction of arbitrary documents.

It has been shown that, compared to positional indexes, web-scale texts can often be compressed in much less space [23]. More importantly, in our proposal the compressed text can be used for both positional re-ranking and snippet generation, saving additional space.

An example of the procedure that we propose is shown in Figure 3.1.



Figure 3.1: Illustration of the query process from Section 2.1, using the textual data to obtain positions and for the snippet generation.

Thus, solving a query with our proposal, includes the following steps:

- 1. Query Processing Step: Given a user query, use an inverted index to obtain the top- k_1 documents according to some standard query processing approach (e.g., DAAT) and ranking function (e.g., BM25).
- 2. Positional Ranking Step: Given the top- k_1 documents obtained on the previous step, obtain the in-document positions for the query terms from the documents themselves. Then re-rank the results using a positional ranking function [11, 47].
- 3. Snippet Generation Step: After the re-ranking of previous step, obtain snippets of length s for the top- k_2 documents, for a given $k_2 \leq k_1$.

One concern is how these alternatives impact query processing speed, as we must decompress documents and then search for the query terms within them. We will study the resulting trade-offs between running time and space requirement.

Thus, to index or not to index position data, that is the research question that we hope to answer in this thesis. To our knowledge, such alternative approaches for storing positional data have not been rigorously compared before. Our main result is that we can store all the information needed for query processing (i.e., document identifiers, term frequencies, position data, and text) using space close to that of state-of-the-art positional indexes (which only store position data and thus cannot be used for snippet creation), with only a minor increase in query processing time. Thus, we provide new alternatives for practical compression of position and text data, outperforming the recent approaches from [40].

Chapter 4

A Wavelet Tree for Computing Positions and Extracting Snippets

In this chapter we explore the alternative of representing the text collection using a Wavelet Tree (WT) data structure. We then use the WT functionality to obtain both, term-positions and snippets. This data structure replaces the positional inverted lists and the text collection, aiming at a better space usage.

Let $T = D_1 D_2 \cdots D_N$ be the text obtained from the concatenation (in arbitrary order) of the documents in the collection, where the documents are separated by a special symbol '\$'. We represent a text T with a WT to obtain term positions and text snippets. Given a position i in T, one can easily obtain both the docID of the document that contains T[i] as $rank_{\$}(T, i)$, and the starting position of a given document j by means of $select_{\$}(T, j)$ [6].

4.1 Byte-Oriented Huffman WT

Instead of a bit-oriented WT (as the one explained in Section 1.6), we use the byteoriented representation from [9], using a Huffman encoding of the words, which is the most efficient alternative reported in there. The idea is to first assign a Huffman code to each vocabulary term [32]. Then, we store the most significant *byte* of the encoding of each term in array B^{root} . That is, each WT node v stores an array of bytes B^v , instead of bit arrays as in Section 1.6. Next, each term in the text is assigned to one of the children of the root, depending on the first byte in the encodings. Notice that in this way the WT is 256-ary. See [9] for details.

To support rank and select, we use the simple approach from [9]. Given a WT node v, we divide the corresponding byte sequence B^v into super-blocks of s_b bytes each. For each super-block we store 256 super-block counters, one for each possible byte. These counters tell us how many occurrences of a given byte there are in the text up to the last position of the previous super-block. Also, each super-block is divided into blocks

of b bytes each. Every such block also stores 256 block counters, similarly as before. The difference is that the values of these counters are local to the super-block, hence less bits are used for them. To compute $\operatorname{rank}_c(T, i)$, we first compute the super-block j that contains i, and use the super-block counter for c to count how many c there are in T up to super-block j - 1. Then we compute the block i' that contains i and add (to the previous value) the block counter for c. Finally, we must count the number of c within block i'. This is done with a sequential scan over block i'. This block/super-block structure allows for time-space trade-offs. In our experiments we use $s_b = 2^{16}$. Hence, super-block counters can be stored in 16 bits each. We consider b = 1 KB, b = 3 KB and b = 7 KB. Operation select is implemented by binary searching the super-block/block counters; thus no extra information is stored for this [9].

4.2 Obtaining Term Positions from a WT

To obtain position data assume that, given docID i for a top- k_1 document and a query term t, we want to obtain the positions of t within D_i . A simple solution could be to extract document D_i from the WT, and then search for t within it. However, a more efficient way is to use operation select to find every occurrence of t within D_i , hence working in time proportional to the number of occurrences of the term (and not the document length). Let d be the starting position for document D_i in T. Hence, there are $r \leftarrow \operatorname{rank}_t(T, d)$ occurrences of t before document D_i , and the first occurrence of t within D_i is at position $j \leftarrow \operatorname{select}_t(T, r + 1)$, the second occurrence at position $j' \leftarrow \operatorname{select}_t(T, r + 2)$, and so on. Overall, if o is the number of occurrences of t within D_i , then we need 1 rank and o + 1 selects to find them. An illustration of this process can be seen in Figure 4.1.



Figure 4.1: Illustration of the position-extraction process in a $\mathtt{W}\mathtt{T}$, using operation select.

4.3 Experimental Results

We fully implemented the byte-oriented WT and the position extraction process. The aim was an implementation able to deal with large texts. Indeed, the text indexed in our experiments is the largest text indexed with a succinct/compressed data structure we are aware of in the literature.

In Table 4.1 we show the experimental trade-offs obtained for WT, for the different block sizes tested (see the rows "WT(7 KB)", "WT(3 KB)" and "WT(1 KB)"). For the sake of comparison, we also show the results for PILs from the previous chapter.

Note that WT (1 KB) obtains better times than PILs, yet requiring considerably more space. An advantage of the WT structure is that it can search for positions and get the snippets of the documents within the same space. WT (7 KB), on the other hand, is slower than PIL (Rice) and uses more space. Even though the WT, includes the textual data, its high space usage could leave it out of consideration on schemes where the text is not necessary. Next, we introduce extra improvements to make them competitive.

4.4 Achieving Higher-Order Compression with the WT

Basically, WTs are zero-order compressors, which explains their high space usage. To achieve higher-order compression, notice that B^{root} contains the most significant byte of the Huffman encodings of the original terms in the text. Thus, the original text structure is at least partially preserved in the structure of B^{root} , which might thus be as compressible as the original text. A similar behavior can be observed in internal nodes for the remaining bytes that form the encodings of the terms. Thus, we propose to compress the blocks of B^v in each WT node v by using standard compressors, based on LZ77 compressors (because decompression time is crucial for time efficiency in our case). Table 4.1 shows results for 1zma, snappy and 1z4, the best compressors we tried.

Compression Scheme	Space Usage	Avera	age Positi	on extrac	tion time	(ms/q)
	(MB)	$k_1 = 50$	$k_1 = 100$	$k_1 = 150$	$k_1 = 200$	$k_1 = 300$
PIL Rice PIL S16	28,373 31,338	$1.28 \\ 0.74$	2.22 1.12	$3.05 \\ 1.43$	$3.27 \\ 1.75$	$5.57 \\ 2.51$
WT(7 KB) WT(3 KB) WT(1 KB)	$40,534 \\ 42,577 \\ 56,917$	1.94 1.11 0.33	$3.68 \\ 2.12 \\ 0.62$	$5.30 \\ 3.07 \\ 1.04$	$6.85 \\ 3.97 \\ 1.15$	$9.80 \\ 5.70 \\ 1.75$
$ extsf{WT}(7 extsf{ KB}) + extsf{lzma} \\ extsf{WT}(1 extsf{ KB}) + extsf{lzma} \end{cases}$	$19,\!628 \\ 42,\!359$	$19.25 \\ 7.22$	$36.59 \\ 13.54$	$52.83 \\ 19.44$	$68.36 \\ 24.97$	$97.71 \\ 35.57$
${ t WT(7~{ m KB})+{ m snappy}} \ { t WT(1~{ m KB})+{ m snappy}} \ { t WT(1~{ m KB})+{ m snappy}}$	25,122 46,778	$14.35 \\ 2.07$	23.76 3.61	39.38 5.88	51.02 7.32	$74.56 \\ 10.47$
$ extsf{WT}(7 extsf{ KB}) + extsf{lz4} \ extsf{WT}(1 extsf{ KB}) + extsf{lz4} \ extsf{Lz4}$	24,911 46,600	14.55 2.08	23.77 3.60	39.38 5.90	51.05 7.31	74.60 10.48

Table 4.1: Experimental results for extracting term-position data using a WT.

Notice that WT (7 KB) + 1zma achieves 19,628 MB, almost half the space used by WT (7 KB). The time to obtain positions becomes, on the other hand, an order of magnitude slower. WT (7 KB) + snappy achieves slightly better times, using space that

is smaller to that of PILs. Overall, although the times are slower than that of PILs, this significant reduction in space could make WTs competitive in scenarios where storing the text is relevant.

Chapter 5

Computing Term Positions and Snippets from the Compressed Text

5.1 Using Standard Compressors

Using standard text compressors, our next approach is to obtain positions using the baseline for generating snippets from Section 2.3.3. At search time, the top- k_1 documents are obtained from their corresponding blocks. Hence, in the worst case we must decompress k_1 text blocks, compared with the (worst-case) k_1 positional chunks for each query term from the PILs. Then, the query terms are sought within these documents, obtaining their positions. Since we are looking for single terms, no sophisticated text search algorithm is used (like KMP, for instance). We just carry out a single scan on the document, and for each text position we check whether it is on the query terms or not. For the snippet extraction step, no further decompression is needed, because the documents for the top- k_2 are already decompressed.

In Table 5.1 (on page 42) we show experimental results for this approach, using lzma, snappy and lz4 compressors, and blocks of size 200 KB. We also compare with the times and space obtained using PILs (from Chapter 2) and the best alternatives for WTs (from Chapter 4). We can see that using lzma, we can store positions and text in about half the space of PIL (the latter just storing positions). Although this is promising, this approach is two orders of magnitude slower than the a positional index, which limits its use in real systems. Actually, the time is about 10 times slower than that of Step 1 (recall Table 2.1 on page 28). If we use snappy instead, we obtain an index that is 21.86% larger than PIL (Rice), or even using the results of 1z4 which is 7.16% larger than PIL (Rice). The times to obtain positions are about 6 times slower than using PILs, (which might be still acceptable in some cases since it corresponds to about 0.5 times the time of Step 1 in the query process).

Even though the space usage achieved with 1zma is competitive, the time needed to obtain the positions is too high. Hence the solution in not practical for web search engines. In what follows, we shall try several approaches to achieve (as much as possible)

the space usage of lzma, yet with a practical time for obtaining positions.

5.2 Using Zero-Order Compressors with Fast Text Extraction

An alternative to compressing the text that could support faster position lookups is the approach from Turpin et al. [43]. The idea is to first sort the vocabulary according to the term frequencies, and then assign term identifiers according to this order. In this way, the term identifier 0 corresponds to the most-frequent term in the collection, 1 to the second-most-frequent term, and so on. The document collection is then represented as a single sequence of identifiers, where each term identifier is encoded using VByte [3]. Note that the 128 most frequent terms in the collection are thus encoded in a single byte. Actually, the work in [43] uses a move-to-front strategy to store the encodings: the first time a term appears in a document, it is encoded with the original code assigned as before; the remaining appearances are represented as an offset to the previous occurrence of the term. We also use this approach in our experiments.

By using either an integer compression scheme, such as VByte or VNibble (a variant of VByte that represents any integer with a variable number of nibbles, i.e., half bytes, recall Section 1.4.1) for the text, or a word-based compression scheme like byte-oriented word Huffman [32], we are able to decompress *any text portion* very efficiently. No text blocks are needed this time, but just a small table indicating the starting position of each document. Table 5.1 shows the resulting trade-offs for the alternatives described until now, compared with the results obtained using PILs (from the previous chapter). We also whow results for VNibble and byte-oriented Huffman.

Notice that we improve the position extraction time significantly, making it competitive with PILs. This shows that being able to extract just the desired k_1 documents is an important fact (since we save time that is otherwise wasted when decompressing a whole block with standard compressors). The higher space usage, however, is a concern. Yet, note that we also represent the text within this space, not just the position data as in PILs. We also tried other compression schemes, such as PforDelta and S9, obtaining poorer compression ratios and similar decompression speed.

In our experimental results (see Table 5.1, rows "VByte", "VNibble" and "Byteoriented Huffman"), we obtain space savings of about 10% for VNibble, and about 2% for byte-oriented Huffman, in both cases compared to VByte's performance. Also, notice that we are now able to use space close to that of **snappy** (with blocks of 200 KB), yet with a better query time.

The faster position extraction time obtained is due to two facts. First, byte-oriented Huffman has a fairly fast decompression speed [32], and methods like VByte and VNibble are able to decompress hundred of million integers (which in our case correspond to terms) per second [48]. Second, these methods are able to decompress just the desired documents (as we already said), without negative impact on compressed size (to obtain

better times, standard compressors must use small blocks, hence achieving poor compression). However, this is basically zero-order compression (i.e., terms are encoded according to their frequency), and hence we are still far from the space usage of, for instance, lzma. The goal of the next approach is to maintain the position extraction times of VByte/VNibble, yet achieving higher-order compression.

5.3 Using Natural-Language Compression Boosters

To obtain higher-order compression, we propose to use a so-called natural-language compression booster [20]. The idea is to use first a (hopefully byte-oriented) zero-order compressor on the text (like byte-oriented word Huffman, or even VByte/VNibble on Turpin et al.'s approach). Then this compressed text is further compressed using some standard compression scheme, based on the LZ77 approach (e.g., lzma, snappy and lz4).

Since we will use a standard compressor again, we must organize the text collection in blocks of fixed size, as in Section 5.1. We compress, using a zero-order compressor, consecutive documents, until the accumulated size of the compressed documents surpasses the size defined for blocks. After reaching this condition, the block is recompressed with a high order compressor. To make the decompression process more efficient, we store a table which for each document, stores the block that contains it, as well as its position within the zero-order representation of the block. An illustration of this process can be seen in Figure 5.1.



Figure 5.1: The compression boosting scheme, with a block size of 200 KB.

It has been shown that this can yield better compression ratios than just using a

standard compression scheme [20] (especially for smaller block sizes). This is because the zero-order compressor virtually enlarges the LZ77 window (typically of 64 KB), hence more regularities can be detected and compressed. Also, block sizes are defined for the text compressed with the zero-order compressor. In other words, the block that is given as input to the higher-order compressor correspond to a bigger segment of the real text, In our case, we propose using Turpin et al.'s approach [43] as the booster (using VByte and VNibble as we explained above) on the sequence of term identifiers, rather than Word Huffman or End-Tagged as in [20]. Our experiments indicate that the former are faster and use only slightly more space than the latter.

To obtain the positions of a query term within a document D_i , as described in Section 5.1, we must first obtain the document and then search for the query terms. To do this, first we have to decompress the block containing D_i with the higher-order decompressor. Second, we have to decompress (using the zero-order decompressor) just the part of the block than contains the document needed. Then, the document is ready to be searched for the query terms. Figure 5.2 shows this process.



Figure 5.2: Example of the decompression process for a document D_1 in the compression-boosting scheme.

In Table 5.1 (on Page 42) We show results for the compression boosting approaches (see the rows for "VByte + 1zma", "VByte + snappy" and "VByte + 1z4", for different block sizes). We do not show in the table results for VNibble and byte-oriented Huffman as boosters. This is because our experiments show that using VByte as a booster produces the minimal space usage, compared to that achieved with VNibble and byte-oriented Huffman. That is, even though VNibble and byte-oriented Huffman on their own achieve better space usage than VByte, the combination of VByte plus a higher-order compressor uses from 9% to 16% less space than VNibble as a booster, and from 2% to 6% less space than byte-oriented Huffman as a booster, depending on the higher-order compressor used.

The better performance of VByte as booster compared to VNibble can be explained because VByte is byte aligned, hence higher-order compressors (which are also byte aligned) are able to detect and compress the regularities of the text in a better way. VNibble, on the other hand, is not byte aligned, then many regularities are not detected, and hence the compression ratio achieved is poorer. Finally, byte-oriented Huffman needs to store a data structure for the decoding process (typically, the Huffman tree). This makes Huffman use slightly more space than VByte. Regarding the coding process, VByte can be decoded very efficiently, using fairly simple code, whereas byte-oriented Huffman needs to traverse the Huffman tree to decode.

Compression	Space	Average position extraction time (ms/q)						
Scheme	Usage	for different values of k_1						
	(MB)	50	100	150	200	300		
PIL Rice	28,373	1.28	2.22	3.05	3.27	5.57		
PIL S16	$31,\!338$	0.74	1.12	1.43	1.75	2.51		
WT(7 KB)	$40,\!534$	1.94	3.68	5.30	6.85	9.80		
WT(1 KB)	56,917	0.33	0.62	1.04	1.15	1.75		
$\mathtt{WT}(7~\mathrm{KB}) + \mathtt{lzma}$	$19,\!628$	19.25	36.59	52.83	68.36	97.71		
$\mathtt{WT}(1~\mathrm{KB}) + \mathtt{lzma}$	$42,\!359$	7.22	13.54	19.44	24.97	35.57		
WT(7 KB) + lz4	24,911	14.55	23.77	39.38	51.05	74.60		
$\mathtt{WT}(1~\mathrm{KB}) + \mathtt{lz4}$	$46,\!600$	2.08	3.60	5.90	7.31	10.48		
lzma (200 KB)	$14,\!987$	137.60	260.36	375.21	482.09	684.94		
snappy (200 KB)	$34,\!576$	9.47	18.00	25.95	33.49	47.74		
lz4 (200 KB)	30,405	6.60	12.48	18.00	23.22	33.09		
VByte	38,339	0.95	1.91	2.86	3.81	5.72		
VNibble	$34,\!570$	1.86	3.71	5.57	6.75	8.10		
Byte-oriented Huffman	$38,\!070$	1.09	2.05	2.95	3.82	5.45		
$\mathrm{VByte}+\mathtt{lzma}\;(200\;\mathrm{KB})$	$12,\!486$	256.16	484.35	716.41	906.54	1,284.87		
$\mathrm{VByte}+\mathtt{lzma}~(50~\mathrm{KB})$	$13,\!981$	70.32	133.18	192.09	246.94	351.71		
$\mathrm{VByte}+\mathtt{lzma}\;(10\;\mathrm{KB})$	16,762	19.26	36.51	52.67	68.00	97.04		
VByte + lzma (1 KB)	$22,\!340$	6.11	11.60	16.80	21.72	31.10		
VByte + snappy (200 KB)	$20,\!158$	9.71	18.86	26.41	34.01	48.69		
$\mathrm{VByte}+\mathtt{snappy}~(50~\mathrm{KB})$	20,366	2.36	4.48	6.47	8.36	11.95		
$\mathrm{VByte}+\mathtt{snappy}\;(10\;\mathrm{KB})$	$22,\!086$	0.82	1.56	2.25	2.91	4.17		
VByte + snappy (1 KB)	$27,\!919$	0.45	0.86	1.24	1.60	2.30		
$\mathrm{VByte} + \mathtt{1z4} \; (200 \; \mathrm{KB})$	$18,\!361$	6.14	11.64	16.77	21.64	30.81		
VByte + 1z4 (50 KB)	$18,\!845$	1.82	3.45	4.98	6.44	9.20		
VByte + 1z4 (10 KB)	$21,\!665$	0.68	1.29	1.87	2.42	3.46		
VByte + lz4 (1 KB)	$27,\!680$	0.41	0.76	1.11	1.43	2.05		

Table 5.1: Experimental results for extracting term-position data from text.

Comparing now the performance of VByte as booster when used with the different

higher-order compressors tested, we can see that, overall, the reduction in space usage (compared to the original VByte approach) is considerable. Comparing VByte + 1zma (200 KB) with 1zma (200 KB), the result is a reduction in space usage of 16.68% (12,486 MB vs 14,987 MB), but at the cost of twice the running time of the original 1zma. For VByte + snappy (200 KB), on the other hand, we obtain a reduction of 41.69% in space for blocks of size 200 KB against snappy (200 KB), with a minor increase in average position-extraction time.

When we use VByte as booster of 1z4, the size of the structure is reduced compared with 1z4 (for all block sizes), without affecting much the average position-extraction time. Notice also that for all compressors, when using smaller blocks, the time to obtain positions rapidly improves, while the size does not increase considerably in some cases. For example, using a block size of 50 KB with VByte + snappy, the average positionextraction time decreases to 2.36 milliseconds per query, which is competitive with the time to obtain positions from PIL (Rice). We can conclude that the best alternative is the use of VByte as booster of 1z4, obtaining (depending on the block size) better average position-extraction time and smaller size than PILs, making both techniques competitive in both space and time. However, VByte + 1z4 also contains the text within this space, allowing its use during snippet generation.

5.4 Further Comparison Between the Most Competitive Alternatives

The most important result from previous sections is that we have found an alternative that is competitive with PILs. Our alternative does not use any index data structure for positions, but just the compressed text. Our results indicate that having an index for positions is not always the best strategy. In this section we will do a more detailed analysis about the behavior of our alternatives. The idea is to find in which cases the alternative of not indexing positions is effective. In particular, we will study the performance of the proposed schemes from different perspectives:

- Space/time trade-offs provided.
- Average position-extraction time as a function of the query length.
- Average position-extraction time as a function of k_1 .

Our analysis will be based on alternative will be based on alternative VByte + 1z4, since it was the best performer in previous section. We will use block sizes of 5 KB, 10 KB, 50 KB, 200 KB, 500 KB and 1,000 KB.

5.4.1 Space/Time Trade-Offs

One drawback of PILs is that they have very few parameters that can be tunned to provide space/time trade-offs. Basically, the parameters are the type of compression

used (in our case we use Rice) and the chunk size for the docID layer. The latter is usually 128, which has shown to be the most effective value. In our setting the minimum space achieved by PILs is 28,373 MB. Hence, if the space allowed for the positional structure is less than that, PILs cannot be used. In the case of compression boosters, the use of larger blocks allows us to achieve a smaller space usage. Indeed, the minimum space achieved is (obviously) when no block structure is used, but the text is compressed as a whole. For VByte + 1zma, this is 8,394 MB, whereas for VByte + 1z4 it is 17,413 MB. Also, our experiments indicate that using blocks of size bigger than 200 KB only yields little improvements in space usage, of up to 8% for VByte + 1zma and up to 2% for VByte + 1z4. This can be seen in Figure 5.3, where we show the experimental space usage as a function of the block size. Moreover, using a block size bigger than 200 KB increases dramatically the average position-extraction time. This is illustrated in Figure 5.4, where it can be seen how (for VByte + 1z4) the time increases with almost no further improvement in space for blocks of size 500 KB and 1,000 KB



Figure 5.3: Space usage for compression boosters, for different block sizes, and PILs.

5.4.2 Average Position-Extraction Time as Function of the Query Lengths

In this section we study how the schemes behave for queries of different length (in number of terms). The time needed to extract positions using compression boosters, can be computed as follows. For each document D_i in the top- k_1 for a query Q, let $T(B_{D_i})$ be the time needed to obtain the positions of the document D_i (which is stored in block B_{D_i}). To obtain the positions for the query terms from Q in document D_i , we



Figure 5.4: Position-extraction time for scheme VByte + 1z4, for block sizes 5 KB, 10 KB, 50 KB, 200 KB, 500 KB, and 1,000 KB. These block sizes correspond to the points in the plot from right to left. For comparison, we also show the performance of PILs.

follow the steps described in Section 5.3 (and Figure 5.2). In practice, even when the number of comparisons needed to extract the positions for a decompressed document D_i are $|Q| \cdot |D_i|$, where |Q| is the number of query terms of Q and $|D_i|$ is the number of words in document D_i , the time $T(B_{D_i})$ is basically dominated by the time needed to decompress the block B_{D_i} , which can be seen as a constant with respect to |Q| (since it just depends on the size of the original block for compression boosters). This can be observed in Figure 5.5, where the average position-extraction remains almost constant as we increase the number of query terms. The worst case of position-extraction time is when every document of the top- k_1 are in a different block. We can conclude that the average position-extraction time depends only on the block size and on k_1 .



Figure 5.5: Average position-extraction times (for $k_1 \in \{50, 300\}$) per number of query terms, for block sizes of 5 KB, 10 KB, 50 KB, 200 KB, 500 KB, 1,000 KB.

Due to the high average position-extraction times, for the compression boosters schemes with blocks larger than 50 KB, the following comparison with PILs is done using blocks of size 5 KB, 10 KB and 50 KB.

The time needed to extract positions from PILs can be computed as follow: given a query Q and its top- k_1 results, for each document D_i in the top- k_1 , we must obtain the positions of the |Q| query terms within D_i . This means that in the worst case $k_1 \cdot |Q|$ positional chunks must be decompressed. Hence and unlike the case of compression boosters, the average position-extraction time for PILs depends on |Q|. This effect can be observed in Figures 5.6, 5.7 and 5.8, for different values of k_1 . As in can be seen, the average position-extraction time for PILs grows lineally with |Q|.



Figure 5.6: Average position-extraction time for compression boosters and PIL, for $k_1 = 50$.

To conclude, the increase of query length affects PILs considerably, whereas, compression boosters are not affected at all by this fact.

5.4.3 Average Position-Extraction Time as Function of k_1

In this section we show the behavior of the schemes as the value of k_1 varies. Figure 5.9 shows that the average position-extraction time of both schemes is lineal in k_1 . Notice that the growth for PILs as k_1 increases, is the slowest. This is because as k_1 grows, many documents in the top- k_1 are within the same chunk in the inverted index, hence the cost of decompressing these chunks grows slowly. Text blocks, on the other hand, store much less documents (tens of them), so the probability of decompressing k_1 text blocks is high (which we also observed on our experiments). Moreover, as we increase the block size, more unuseful documents are decompressed, which explains why the cost for blocks of 50 KB grows faster.



Figure 5.7: Average position-extraction time for compression boosters and PIL, for $k_1 = 150$.



Figure 5.8: Average position-extraction time for compression boosters and PIL, for $k_1 = 300$.





Figure 5.9: Comparison between compression boosters and PILs, for different values of k_1 .

5.5 Conclusions

From our experiments we can conclude that our alternative VByte + 1z4 is a very good choice to support positional ranking and snippet generation. If we consider the query lengths, our experiments indicate that PILs are more adequate for short queries. For long queries, on the other hand, the number of inverted lists involved makes PILs less competitive. In such a case, VByte + 1z4 would be preferred. Long queries appear in many relevant applications, such as search engines for targeted advertising (where queries contain many terms involving the characterization of a user, for instance). Finally, we conclude that approach VByte + 1z4 is more sensitive to the value of k_1 than PILs, in the sense that the position-extraction time for the former degrades faster than that of PILs when the value of k_1 is increased.

Chapter 6

Discussion and Further Experimental Results

In this chapter we consider the best indexing alternatives from previous chapters to build schemes that allow us to carry out the search process as described in Section 2.1 (on page 24). The idea is to show the available trade-offs for the complete search process.

6.1 Scenario 1: Query Processing with Snippet Generation

The first scenario we test implements the three steps from Section 2.1 (i.e., query processing, positional ranking step, and snippet generation). In this scenario we must store the document collection, and must be able to compute position data. Latter in this chapter we will study an alternative scenario where snippets are not needed, hence the text is not needed. Table 6.1 defines the schemes used in our experiments. All schemes include the inverted index to carry Step 1 in the query process. The inverted index includes docID s (compressed with PforDelta) and frequencies (compressed with S16). The total space usage for the GOV2 collection is 9,739 MB. For the query process, we set $k_1 \in \{50, 100, 150, 200, 300\}$ and $k_2 \in \{10, 30, 50\}$. As on the experiments of previous chapters, we use the TREC 2006 query log.

The space/time trade-offs for the schemes tested are determined as follows:

lzma: We used text blocks of size 200 KB, 500 KB and 1000 KB.

lz4: We used text blocks of size 200 KB, 500 KB and 1000 KB.

WT: For WT, we used text blocks of size 1 KB, 2 KB and 7 KB.

WT + lzma: For WT, we used text blocks of size 1 KB, 2 KB and 7 KB.

- WT + lz4: For WT, we used text blocks of size 1 KB, 2 KB and 7 KB.
- **VByte/VNibble:** The two points in this case are obtained by using VNibble compression for the text (the least space) and VByte compression.
- VByte + lzma: We used text blocks of size 5 KB, 10 KB, 50 KB and 200 KB.
- **VByte** + **lz4:** We used text blocks of size 5 KB, 10 KB, 50 KB and 200 KB.
- PIL + lzma: We used text blocks of 200 KB (for lzma). The time points in the trade-off are obtained using Rice compression for the PILs (the least space) and S16.
- **PIL** + lz4: We used text blocks of 200 KB (for lz4). The time points in the trade-off are obtained using Rice compression for the PIL s (the least space) and S16.
- **PIL** + **VNibble:** The time points in the trade-off are obtained using Rice compression for the PILs (the least space) and S16.
- PIL + VByte + lzma: The time points in the trade-off are obtained using Rice compression for the PILs (the least space) and S16. For the compressor booster we used text blocks of size 50 KB.
- PIL + VByte + lz4: The time points in the trade-off are obtained using Rice compression for the PILs (the least space) and S16. For the compressor booster we used text blocks of size 50 KB.

index, which for the GOV2 collection represent 9,739 MB.	Table	6.1:	Gloss	ary	of the	indexing	schemes	tested.	All	schemes	include	the	inverted
	index,	whic	h for	the	GOV2	collectio	n represe	nt 9,739) MI	3.			

. ...

. .

-

- ·

Indexing Scheme	Description
lzma	Text compressed with 1zma for positions and text.
	Text compressed with 124 for positions and text.
WT WT+ lzma WT+ lz4	WT for positions and text. WT compressed with lzma for positions and text. WT compressed with lz4 for positions and text.
VByte/VNibble	Text compressed with VByte/VNibble for positions and text.
VByte + 1zma VByte + 1z4	VByte compression booster on lzma for positions and text. VByte compression booster on lz4 for positions and text.
PIL+ lzma PIL+ lz4 PIL+ VNibble PIL+ VByte + lzma PIL+ VByte + lz4	 PILs for positions, text compressed with lzma. PILs for positions, text compressed with lz4. PILs for positions, text compressed with VNibble. PILs for positions, text compressed with VByte + lzma. PILs for positions, text compressed with VByte + lz4.

6.1.1 DAAT AND Queries

We test first DAAT AND queries for Step 1. This kind of queries are usually supported very efficiently, and correspond just to a small fraction of the overall query process. This means that the time needed to obtain position data should be also efficient, and small differences among alternatives could influence the total processing time.

For $k_1 = 50$ and $k_2 = 10$ (Figure 6.1, above), it can be seen that schemes that use PILs to index position data have a high space usage, as they need to store the text to obtain snippets. Some of these schemes offer a competitive query time, yet their space usage makes them unfordable. PIL + VByte + 1z4 is the most space-efficient alternative using PILs, yet it cannot compete at query time. This shows in practice what we have mentioned along this thesis: state-of-the-art solutions for positional data have a high space usage. This is because they need to store both positions and text separately, and in particular positions are not as compressible as other index components (even the text is more compressible than position).

Scheme WT offers also highly competitive query time, yet using less space than schemes using PIL s. Yet the space usage is still high. Alternatives WT+ 1zma and WT + 1z4 yield a significant reduction of about 42% in space usage. However, the resulting query time degrades quickly, and becomes less competitive.

Scheme VByte/VNibble (which implements Turpin et al.'s idea [43] to compress the text) yields a good trade-off, with small space usage than WT and a highly competitive query time. Up to this point we have achieved a reduction of about 21% in space usage compared to the most time-efficient alternative that uses PILs (i.e., the state of the art up to now), with almost the same query time. However, the space usage is still high compared to the alternatives we study next.

Scheme 1z4 reduces the space usage even more, providing a competitive query time. Scheme 1zma, on the other hand, reduces the space usage significantly, yet at the price of a much slower (unpractical) query time.

When we add the compressor boosters to lzma and lz4 (i.e., schemes VByte + lzma and VByte + lz4), we obtain the most important trade-offs among all alternatives. For instance, if we compare with PIL + lz4 (which is the fastest alternative in the state of the art, and is among the most space-efficient alternatives that use PIL), VByte + lz4 (with text block of size 50 KB, which is the third point in the curve from the right) yields a significant reduction in space usage of about 49.81% (56,958 MB versus 28,584 MB), with a query time that is just 1.03 times slower (from 16.95 msecs/query to 17.49 msecs/query). This shows the effectiveness of our proposal. Moreover, recall that the positional index PIL (using Rice compression) requires 28,373 MB, wheres VByte + lz4 requires, as we already said, 28,584 MB. That is, in 1.007 times the space of just PILs we can store the inverted index (docIDs and frequencies), the positional data, and the text collection. In other words, using only slightly more space than PIL (Rice), scheme VByte + lz4 includes everything needed for query processing. This is one of the most important results and conclusions in this thesis: "not to index" positional

data can be a highly-competitive alternative in practical scenarios.

If we compare now PIL+ 1z4 with VByte + 1zma (for blocks of size 10 KB, which corresponds to the second point in the curve, from the right), we obtain a reduction in space usage of about 53.47% (56,958 MB versus 26,501 MB), while the query time is 1.99 times slower (from 17.47 msecs/query to 34.93 msecs/query). In other words, we are able to accommodate the complete index in 0.94 times the space of PILs, and still are able to offer a competitive query time.

Finally, the smallest space alternatives we tested (which are not shown in the figures) are the ones that use the inverted index for query processing and 1zma compression for positions and snippets. This achieves about 22,225 MB of space. This scheme includes everything needed for query processing, and uses only 78% the space of PIL. However, query processing time increases significantly, to more than 400 ms per query. This scheme could be useful in some cases where the available memory space is very restricted, such that a larger index would mean going to disk.

A recent alternative [40] proposes to use *flat positional indexes* [13, 16] to support phrase querying; this index could also be used for positional ranking. This is basically a positional index from which docID and frequency information can also be obtained. The results reported for the GOV2 collection in [40] give an index of size 30,310 MB that includes docIDs and frequencies, but not the text needed for snippet generation, making this approach uncompetitive for our scenario.

If we increase k_2 to 50, the result is almost the same (see Figure 6.1, below). The only schemes that are affected (i.e., the query time is increased) are those using PILs (because more documents than before must be decompressed per query, in order to obtain their snippets) and those using WTs (as we need to show snippets for more documents, this means that we need to obtain more words from the text, which means traversing the WT more intensively). Schemes that use the compressed text for snippets and positions must decompress k1 documents to obtain position data, to then extract snippets for the already decompressed top-k2 documents (using negligible extra space). Figures 6.2 to 6.5 show experimental results for the remaining values of k_1 and $k_2 \in 10, 50$. Experiments for $k_2 = 30$ are shown in Appendix A. As it can be seen, the same conclusions can be drawn.



Figure 6.1: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3.



Figure 6.2: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3.



Figure 6.3: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3.





Figure 6.4: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.





Figure 6.5: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.

6.1.2 BMW OR Queries

Figures 6.6 to 6.10 show experimental results for BMW OR queries. The results now are slightly different, since now Step 1 of query processing is more expensive than for AND queries. Hence, the differences in query time are smaller. For instance, for $k_1 = 50$ and $k_2 = 10$ (Figure 6.6, above) we can conclude that scheme VByte + 1zma with text blocks of size 10 KB uses 0.94 times the PIL size (which just stores positional data).

The query time is 1.43 times slower than that of scheme PILs + 1z4.

For scheme VByte + 1z4, the results are even more impressive. Using just 1.007 times the space of just PILs we can achieve a query time that is 1.005 times slower than scheme PIL+ 1z4.

As in previous section, we conclude that VByte + 1z4 and VByte + 1zma (using text blocks of size 50 KB) offer very relevant trade-offs, being able to replace PILs in many cases.



Figure 6.6: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3.



Figure 6.7: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3.



Figure 6.8: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.




Figure 6.9: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure 6.10: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$ and $k_2 \in \{10, 50\}$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.

6.2 Scenario 2: Query Processing without Snippet Generation

We now study the practical performance of all proposed schemes in a scenario where text snippets are not needed. That is, Step 3 of Section 2.1 is not carried out. Notice that the text is not needed in the state-of-the-art schemes based on PILs. Hence, we now have just one scheme using PILs (in previous section the various schemes with PILs only differed in how they compressed the text). See Table 6.2 for a summary of the alternatives we tested. In our experiments we set $k_1 \in \{50, 100, 150, 200, 300\}$. All schemes include the inverted index, as in previous section.

Table 6.2: Glossary of the indexing schemes for the figures. All schemes include the inverted index.

Indexing Scheme	Description
lzma	Text compressed with 1zma for positions.
lz4	Text compressed with 1z4 for positions.
WT	WT for positions.
WT+ lzma	WT compressed with 1zma for positions.
WT+ lz4	WT compressed with 1z4 for positions.
VByte/VNibble	Text compressed with VByte/VNibble for positions.
VByte + 1zma	VByte compression booster on lzma for positions.
VByte + 1z4	VByte compression booster on lz4 for positions.
PIL	PILs for positions.

6.2.1 DAAT AND Queries

Figures 6.11 up to 6.15 show results for AND queries, for all values of k_1 we tested. The most important result to highlight is that the state-of-the-art scheme based on PILs (which this time does not store the text, saving considerable space) cannot compete with schemes VByte + 1z4 and VByte + 1zma, neither in space usage nor query time. That is, "not to index" positional data is also effective in scenarios where snippets are not needed.



Figure 6.11: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$, including Step 1, Step 2.



Figure 6.12: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 250 ms/q.



Figure 6.13: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 250 ms/q.



Figure 6.14: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure 6.15: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 400 ms/q.

6.2.2 BMW OR Queries

Figures 6.16 up to 6.20 show results for BMW OR queries, for all values of k_1 Notice that the differences are even smaller than in Section 6.1.2, because in this scenario the text is not needed, then PILs are not affected by the snippet extraction time. Using the same scenario as in Section 6.1.2, where the scheme VByte +1z4 is 1.007 times the space of just PILs we can achieve a query time just 1.01 times slower than the PIL+1z4 alternative. Concluding that "not to index" positional data is effective in all the scenarios studied.



Figure 6.16: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$, including Step 1, Step 2.



Figure 6.17: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 300 ms/q.



Figure 6.18: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure 6.19: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure 6.20: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$, including Step 1, Step 2. It is important to note that Scheme 1 has query time greater than 400 ms/q.

6.3 Discussion

After all the scenarios described in this chapter, we can conclude that for web search engines that follow the steps of Secction 2.1. Using PILs to extract the positions, is not an efficient solution.

This is due the fact that to answer the Step 3, another structure is needed and its size is near to the PIL size. In this scenario all the alternatives studied (WT and compression boosters) have a better space/time trade-off than PILs.

However, in the case when the Step 3, is not needed, just the scheme of compression boosters can be competitive. Yet in the configuration of Vbyte + 1z4, using the same space that just PILs, can store the inverted index, the positions and the text information, showing a better performance.

Chapter 7

Conclusion and Future Work

From our study we can conclude that there exists a wide range of practical timespace trade-offs, other than just the classical positional inverted indexes. We studied several alternatives, trying to answer the question of whether it is necessary to index position data or not. As one of the most relevant points in the trade-off, we propose a compressed document representation based on the approach in [43] combined with 1z4 compression [2]. This allows us to compute position and snippet data using less space than a standard positional inverted index that only stores position data. Even if we include the space used for document identifiers and term frequencies, this approach uses just 1.007 times the space of a positional inverted index, with the basically the same query time.

This means that in many practical cases, "not to index" position data may be the most efficient approach. This provides new practical alternatives for positional index compression, a problem that has been considered difficult to address in previous work [47, 23]. Finally, we also showed that compressed self-indexes such as wavelet trees [25] can be competitive with the best solutions in some scenarios.

In our experiments we found that the time to extract positions is proportional to the number of blocks decompressed, this is valid for all the alternatives. Yet in the special case of compression boosters to decompress k_1 documents to extract positions, the structure had to decompress k_1 blocks. The problem to solve in future works, is how to reorder the Web, keeping in mind the most relevant top- k_1 documents, and store them together, achieving that for a query just have to decompress $f < k_1$ blocks. decreasing the position-extraction time.

Another interesting future work is the study of phrases searching, this is because in the solution proposed, we have to decompress the whole document to search for positions, in this scenario we could also search for the complete phrase in the query, which can be done in the same time due than we traverse the whole document to search for positions.

Bibliography

- [1] http://code.google.com/p/snappy/.
- [2] https://code.google.com/p/lz4/.
- [3] V. N. Anh and A. Moffat. Compressed inverted files with reduced decoding overheads. In Proc. of 21st Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pages 290–297, 1998.
- [4] V. N. Anh and A. Moffat. Inverted index compression using word-aligned binary codes. Inf. Retr., 8(1):151–166, 2005.
- [5] D. Arroyuelo, V. Gil-Costa, S. González, M. Marin, and M. Oyarzún. Distributed search based on self-indexed compressed text. *Information Processing and Man*agement, 2012. To appear.
- [6] D. Arroyuelo, S. González, and M. Oyarzún. Compressed self-indices supporting conjunctive queries on document collections. In *SPIRE*, LNCS 6393, pages 43–54, 2010.
- [7] R. Baeza-Yates and B. Ribeiro-Neto. Modern Information Retrieval the Concepts and Technology Behind Search, Second Edition. Pearson Education Ltd., Harlow, England, 2011.
- [8] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. J. of Computer Networks, 30(1-7):107-117, 1998.
- [9] N. Brisaboa, A. Fariña, S. Ladra, and G. Navarro. Implicit indexing of natural language text by reorganizing bytecodes. *Information Retrieval*, 2012. To appear.
- [10] A. Z. Broder, D. Carmel, M. Herscovici, A. Soffer, and J. Zien. Efficient query evaluation using a two-level retrieval process. In *Proc. of 12th International Conference* on Information and Knowledge Management, pages 426–434. ACM, 2003.
- [11] S. Büttcher, C. Clarke, and G. Cormack. Information Retrieval: Implementing and Evaluating Search Engines. MIT Press, 2010.
- [12] Stefan Büttcher, Charles L. A. Clarke, and Brad Lushman. Term proximity scoring for ad-hoc retrieval on very large text collections. In *Proceedings of the 29th Annual*

International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '06, pages 621–622, New York, NY, USA, 2006. ACM.

- [13] C. Clarke, G. Cormack, and F. Burkowski. An algebra for structured text search and a framework for its implementation. *Computer Journal*, 38(1):43–56, 1995.
- [14] F. Claude and G. Navarro. Practical rank/select queries over arbitrary sequences. In SPIRE, LNCS 5280, pages 176–187. Springer, 2008.
- [15] D. Cutting. Apache Lucene. http://lucene.apache.org/.
- [16] J. Dean. Challenges in building large-scale information retrieval systems: invited talk. In WSDM, page 1, 2009.
- [17] S. Ding and T. Suel. Faster top-k document retrieval using block-max indexes. In Proc. of 34th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pages 993–1002, 2011.
- [18] P. Elias. Universal codeword sets and representations of the integers. IEEE Transactions on Information Theory, 21(2):194–203, 1975.
- [19] A. Fariña, N. Brisaboa, G. Navarro, F. Claude, A. Places, and E. Rodríguez. Wordbased self-indexes for natural language text. ACM Transactions on Information Systems (TOIS), 30(1):article 1, 2012.
- [20] A. Fariña, G. Navarro, and J. Paramá. Boosting text compression with word-based statistical encoding. *Computer Journal*, 55(1):111–131, 2012.
- [21] P. Ferragina, R. Giancarlo, and G. Manzini. The myriad virtues of wavelet trees. Information and Computation, 207(8):849–866, 2009.
- [22] P. Ferragina, R. González, G. Navarro, and R. Venturini. Compressed text indexes: From theory to practice. ACM Journal of Experimental Algorithmics, 13, 2008.
- [23] P. Ferragina and G. Manzini. On compressing the textual web. In WSDM, pages 391–400, 2010.
- [24] S. Golomb. Run-length encoding. IEEE Transactions on Information Theory, 12(3):399–401, 1966.
- [25] R. Grossi, A. Gupta, and J. S. Vitter. High-order entropy-compressed text indexes. In SODA, pages 841–850, 2003.
- [26] J. M. Kleinberg. Authoritative sources in a hyperlinked environment. J. of ACM, 46(5):604–632, 1999.
- [27] C. Manning, P. Raghavan, and H. Schütze. Introduction to Information Retrieval. Cambridge University Press, 2008.
- [28] G. Manzini. An analysis of the Burrows-Wheeler transform. J. ACM, 48(3):407–

430, 2001.

- [29] G Nigel N Martin. Range encoding: an algorithm for removing redundancy from a digitised message. In Proc. Institution of Electronic and Radio Engineers International Conference on Video and Data Recording, 1979.
- [30] D. Metzler and W. B. Croft. A markov random field model for term dependencies. In Proc. of 28th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, 2005.
- [31] G. Mishne and M. Rijke. Boosting web retrieval through query operations. In *Proc. of 27th European Conference on IR Research*, 2005.
- [32] A. Moffat. Word-based text compression. Software, Practice, and Experience, 19(2):185–198, 1989.
- [33] A. Moffat and L. Stuiver. Binary interpolative coding for effective index compression. Inf. Retr., 3(1):25–47, 2000.
- [34] G. Navarro. Wavelet trees for all. Journal of Discrete Algorithms, 2013. To appear.
- [35] G. Navarro and V. Mäkinen. Compressed full-text indexes. ACM Computing Surveys, 39(1), 2007.
- [36] Y. Rasolofo and J. Savoy. Term proximity scoring for keyword-based retrieval systems. In Proc. of 25th European Conference on IR Research, 2003.
- [37] David Salomon. *Data Compression: The Complete Reference*. 2007. With contributions by Giovanni Motta and David Bryant.
- [38] R. Schenkel, A. Broschart, S. Hwang, M. Theobald, and G. Weikum. Efficient text proximity search. In 14th String Processing and Information Retrieval Symposium, 2007.
- [39] Falk Scholer, Hugh E. Williams, John Yiannis, and Justin Zobel. Compression of inverted indexes for fast query evaluation. In SIGIR, pages 222–229. ACM, 2002.
- [40] D. Shan, W. X. Zhao, J. He, R. Yan, H. Yan, and X. Li. Efficient phrase querying with flat position index. In *CIKM*, pages 2001–2004, 2011.
- [41] B. Sparrow, J. Liu, and M. Wegner. Google effects on memory: Cognitive consequences of having information at our fingerprints. *Science*, 333(6043):776–778, 2011.
- [42] T. Tao and C. Zhai. An exploration of proximity measures in information retrieval. In Proc. of 30th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, 2007.
- [43] A. Turpin, Y. Tsegay, D. Hawking, and H. Williams. Fast generation of result

snippets in web search. In Proc. of 30th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pages 127–134, 2007.

- [44] L. Wang, J. J. Lin, and D. Metzler. A cascade ranking model for efficient ranked retrieval. In Proc. of 34th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pages 105–114, 2011.
- [45] H. Williams and J. Zobel. Compressing integers for fast file access. Computer Journal, 42(3):193–201, 1999.
- [46] I. H. Witten, A. Moffat, and T. C. Bell. Managing Gigabytes: Compressing and Indexing Documents and Images (2nd ed.). Morgan Kaufmann Publishing, 1999.
- [47] H. Yan, S. Ding, and T. Suel. Compressing term positions in web indexes. In Proc. of 32nd Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pages 147–154, 2009.
- [48] H. Yan, S. Ding, and T. Suel. Inverted index compression and query processing with optimized document ordering. In WWW, pages 401–410, 2009.
- [49] Jiangong Zhang, Xiaohui Long, and Torsten Suel. Performance of compressed inverted list caching in search engines. In *Proceedings of the 17th International Conference on World Wide Web*, WWW '08, pages 387–396, New York, NY, USA, 2008. ACM.
- [50] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *IEEE TRANSACTIONS ON INFORMATION THEORY*, 23(3):337–343, 1977.
- [51] Jacob Ziv and Abraham Lempel. Compression of individual sequences via variablerate coding, 1978.
- [52] M. Zukowski, S. Héman, N. Nes, and P. Boncz. Super-scalar ram-cpu cache compression. In *ICDE*, page 59, 2006.

Appendix A

Aditional Experimental Results

In this section we shows the trade-offs for the alternatives described in Table 6.1, with $k_3 = 30$.

A.1 DAAT AND Queries



Figure A.1: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.



Figure A.2: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.



Figure A.3: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.



Figure A.4: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$ and $k_2 = 30$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure A.5: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$ and $k_2 = 30$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.

A.2 BMW OR Queries



Figure A.6: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 50$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.



Figure A.7: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 100$ and $k_2 = 30$, including Step 1, Step 2 and Step 3.



Figure A.8: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 150$ and $k_2 = 30$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure A.9: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 200$ and $k_2 = 30$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.



Figure A.10: Time-space trade-offs for the overall query process for the GOV2 collection. With $k_1 = 300$ and $k_2 = 30$, including Step 1, Step 2 and Step 3. It is important to note that Scheme 1 has query time greater than 400 ms/q.