

An Experimental Investigation of XML Compression Tools

Sherif Sakr
National ICT Australia
223 Anzac Parade, NSW 2052
Sydney, Australia
Sherif.Sakr@nicta.com.au

ABSTRACT

This paper presents an extensive experimental study of the state-of-the-art of XML compression tools. The study reports the behavior of nine XML compressors using a large corpus of XML documents which covers the different natures and scales of XML documents. In addition to assessing and comparing the performance characteristics of the evaluated XML compression tools, the study tries to assess the effectiveness and practicality of using these tools in the real world. Finally, we provide some guidelines and recommendations which are useful for helping developers and users for making an effective decision for selecting the most suitable XML compression tool for their needs.

1. INTRODUCTION

The eXtensible Markup Language (XML) has been acknowledged to be one of the most useful and important technologies that has emerged as a result of the immensed popularity of HTML and the World Wide Web. Due to the simplicity of its basic concepts and the theories behind, XML has been used in solving numerous problems such as providing neutral data representation between completely different architectures, bridging the gap between software systems with minimal effort and storing large volumes of semi-structured data. XML is often referred as *self-describing data* because it is designed in a way that the schema is repeated for each record in the document. On one hand, this self-describing feature grants the XML great flexibility and on the other hand, it introduces the main problem of *verbosity* of XML documents which results in huge document sizes. This huge size lead to the fact that the amount of information that has to be transmitted, processed, stored, and queried is often larger than that of other data formats. Since XML usage is continuing to grow and large repositories of XML documents are currently pervasive, a great demand for efficient XML compression tools has been exist. To tackle this problem, several research efforts have proposed the use of XML-conscious compressors which exploits the

well-known structure of XML documents to achieve compression ratios that are better than those of general text compressors. The usage of XML compressing tools has many advantages such as: reducing the network bandwidth required for data exchange, reducing the disk space required for storage and minimizing the main memory requirements of processing and querying XML documents.

Experimental evaluation and comparison of different techniques and algorithms which deals with the same problem is a crucial aspect especially in applied domains of computer science. This paper presents an extensive experimental study for evaluating the state-of-the-art of XML compression tools. We examine the performance characteristics of nine publicly available XML compression tools against a wide variety of data sets that consists of 57 XML documents. The web page of this study [1] provides access to the test files, examined XML compressors and the detailed results of this study.

The remainder of this paper is organized as follows. Section 2 briefly introduces the XML compression tools examined in our study and classifies them in different ways. Section 3 presents the data sets used to perform the experiments. Section 4 describes the test environments. Detailed and consolidated results of our experiments are presented in Section 5 before we draw our final conclusions in Section 6.

2. SURVEY OF XML COMPRESSION TOOLS

2.1 Features and Classifications

A very large number of XML compressors have been proposed in the literature of recent years. These XML compressors can be classified with respect to three main characteristics. The first classification is based on their awareness of the structure of the XML documents. According to this classification, compressors are divided into two main groups:

- **General Text Compressors:** Since XML data are stored as text files, the first logical approach for compressing XML documents was to use the traditional general purpose text compression tools. This group of XML compressors [6, 2, 23] is *XML-Blind*, treats XML documents as usual plain text documents and applies the traditional text compression techniques [33].
- **XML Conscious Compressors:** This group of compressors are designed to take the advantage of the awareness of the XML document *structure* to achieve better compression ratios over the general text compressors. This group of compressor can be further clas-

sified according to their dependence on the availability of the schema information of the XML documents as follows:

- *Schema dependent compressors* where both of the encoder and decoder must have access to the document schema information to achieve the compression process [3, 11, 25, 8].
- *Schema independent compressors* where the availability of the schema information is not required to achieve the encoding and decoding processes [29, 21, 9, 18].

Although schema dependent compressors may be, *theoretically*, able to achieve slightly higher compression ratios, they are not preferable or commonly used in practice because there is no guarantee that the schema information of the XML documents is always available.

The second classification of XML compressor is based on their ability of supporting queries.

- **Non-Queriable (Archival) XML Processor:** This group of the XML compressors does not allow any queries to be processed over the compressed format [29, 21, 11, 5, 9]. The main focus of this group is to achieve the highest compression ratio. By default, general purpose text compressors belong to the non-queriable group of compressors.
- **Queriable XML Processor:** This group of the XML compressors allow queries to be processed over their compressed formats [13, 34, 30]. The compression ratio of this group is usually worse than that of the archival XML compressors. However, the main focus of this group is to avoid full document decompression during query execution. In fact, the ability to perform direct queries on compressed XML formats is important for many applications which are hosted on resource-limited computing devices such as mobile devices and GPS systems. By default, all queriable compressors are XML conscious compressors as well.

The third classification considers whether the compression schemes operate in an online or offline manner.

- **Online Compressors:** are able to stream the compressed data to the decoder i.e the decode is able to begin the process of decompression before the encode has finished transmitting the compressed data.
- **Offline Compressors:** don't allow the decoder to begin the decompression process until the entire compressed file has been received.

The *online* feature of XML compression tools could be very important for the scenarios where the users are heavily exchanging compressed XML documents over networks. In these scenarios, online decompression processors can effectively decrease the network latency during the transmission process.

Table 1 lists the symbols that indicate the features of each XML compressor included in the list of Table 2.

Symbol	Description
G	General Text Compressor
S	Specific XML Compressor
D	Schema dependent Compressor
I	Schema Independent Compressor
A	Archival XML Compressor
Q	Queriable XML Compressor
O	Online XML Compressor
F	Offline XML Compressor

Table 1: Symbols list of XML compressors features

2.2 Examined Compressors

In our study we considered, to the best of our knowledge, *all* XML compression tools which are fulfilling the following conditions:

1. Is publicly and freely available either in the form of open source codes or binary versions.
2. Is a schema-independent. As previously mentioned, the set of compressors which is not fulfilling this condition is not commonly used in practice .
3. Be able to run under our Linux version of operating system.

Table 2 lists the surveyed XML compressors and their features where the **Bold** font is used to indicate the compressors which are fulfilling our conditions and included in our experimental investigation. The border line between the upper section and the lower section of Table 2 is used to differentiate between the non-queriable (upper section) and queriable (lower section) sets of the XML compressors. The three compressors (*DTDPPM*, *XAUST*, *rngzip*) have not been included in our study because they do not satisfy Condition 2. Although its source code is available and it can be successfully compiled, XGrid did not satisfy Condition 3. It always gives a fixed run-time error message during the execution process. The rest of the list (11 compressors) don't satisfy Condition 1. The status of a lack of source code/binaries for a large number of the XML compressors proposed in literature, to the best of our search efforts and contact with the authors, and especially from the queriable class [30, 32, 34] was a bit disappointing for us. This has limited a subset of our initially planned experiments especially those which targeted towards assessing the performance of evaluating the results of XML queries over the compressed representations. In the following we give a brief description of each examined compressor.

General Text Compressors numerous algorithms have been devised over the past decades to efficiently compress text data. In our evaluation study we selected three compressors which are considered to be the best representative implementations of the most popular and efficient text compression techniques. We selected gzip [6], bzip2 [2] and PPM [23] compressors to represent this group.

XMill in [29] Liefke and Suciuc have presented the first implementation of an XML conscious compressor. In XMill, both of the structural and data value parts of the source XML document are collected and compressed separately . In the structure part, XML tags and attributes are encoded in a dictionary-based fashion before passing it to a back-end general text compression scheme. Data values are grouped into

homogenous and semantically related containers according to their path and data type. Each container is then compressed separately using specialized compressor that is ideal for the data type of this container. In the latest versions of the XMill source distribution, the intermediate binaries of the compressed format can be passed to one of three alternative back-end general purpose compressor: gzip, bzip2 and PPM. In our experiments we evaluated the performance of the three alternative back-ends independently. Hence, in the rest of the paper we refer to the three alternative back-ends with the names *XMillGzip*, *XMillBzip* and *XMillPPM* respectively.

XMLPPM is considered as an adaptation of the general purpose *Prediction by Partial Matching* compression scheme (PPM) [23]. In [21], Cheney has presented XMLPPM as a streaming XML compressor which uses a *Multiplexed Hierarchical PPM Model* called (MHM). The main idea of this MHM model is to use four different PPM models for compressing the different XML symbols: element, attribute, character and miscellaneous data.

SCMPPM is described by Adiego et al. in [19] as a variant of the XMLPPM compressor. It combines a technique called *Structure Context Modelling* (SCM) with the PPM compression scheme. It uses a bigger set of PPM models than XMLPPM as it uses a separate model to compress the text content under each element symbol.

XWRT is presented by Skibinski et al. in [35]. It applies a dictionary-based compression technique called *XML Word Replacing Transform*. The idea of this technique is to replace the frequently appearing words with references to the dictionary which is obtained by a preliminary pass over the data. XWRT submits the encoded results of the preprocessing step to three alternative general purpose compression schemes: gzip, LZMA and PPM.

Axe chop is presented by Leighton et al. in [27]. It divides the source XML document into structural and data segments. The MPM compression algorithm is used to generate a context-free grammar for the structural segment which is then passed to an adaptive arithmetic coder. The data segment contents are organized into a series of containers (one container for each element) before applying the *Burrows-Wheeler Transformation* (BWT) compression [20] over each container.

Exalt in [36], Toman has presented an idea of applying a syntactical-oriented approach for compressing XML documents. It is similar to AXECHOP in that it utilized the fact that XML document could be represented using a context-free grammar. It uses the grammar-based codes encoding technique introduced by Kieffer and Yang in [26] to encode the generated context-free grammars.

3. OUR CORPUS

3.1 Corpus Characteristics

Determining the XML files that should be used for evaluating the set of XML compression tools is not a simple task. To provide an extensive set of experiments for assessing and evaluating the performance characteristics of the XML compression tools, we have collected and constructed a large corpus of XML documents. This corpus contains a wide variety of XML data sources and document sizes. Table

Compressor	Features	Code Available
GZIP (1.3.12) [6]	GAIF	Y
BZIP2 (1.0.4) [2]	GAIF	Y
PPM (j.1) [7]	GAIF	Y
XMill (0.7) [15]	SAIF	Y
XMLPPM (0.98.3) [16]	SAIO	Y
SCMPPM (0.93.3) [9]	SAI	Y
XWRT (3.2) [18]	SAI	Y
Exalt (0.1.0) [5]	SAIF	Y
AXECHOP [27]	SAIF	Y
DTDPPM [3]	SADO	Y
XAUST [11]	SAD	Y
rngzip [8]	SQD	Y
Millau [25]	SADO	N
XComp [28]	SAIF	N
XGrind [13]	SQIO	Y
XBzip [24]	SQI	N
XQueC [17]	SQI	N
XCQ [32]	SQIO	N
XPress [31]	SQIO	N
XQzip [22]	SQI	N
XSeq [30]	SQI	N
QXT [34]	SQI	N
ISX [37]	SQI	N

Table 2: XML Compressors List

3 describes the characteristics of our corpus. *Size* denotes the disk space of XML file in MBytes. *Tags* represents the number of distinct tag names in each XML document. *Nodes* represents the total number of nodes in each XML data set. *Depth* is the length of the longest path in the data set. *Data Ratio* represents the percentage of the size of data values with respect to the document size in each XML file. The documents are selected to cover a wide range of sizes where the smallest document is 0.5 MB and the biggest document is 1.3 GB. The documents of our corpus can be classified into four categories depending on their characteristics:

- **Structural documents** this group of documents has no data contents at all. 100 % of each document size is preserved to its structure information. This category of documents is used to assess the claim of XML conscious compressors on using the well known structure of XML documents for achieving higher compression ratios on the structural parts of XML documents. Initially, our corpus consisted of 30 XML documents. Three of these documents were generated by using our own implemented Java-based random XML generator. This generator produces completely random XML documents to a parameterized arbitrary depth with only *structural* information (no data values). In addition, we created a *structural* copy for each document of the other 27 *original* documents - with data values - of the corpus. Thus, each *structural* copy captures the structure information of the associated XML *original* copy and removes all data values. In the rest of this paper, we refer to the documents which include the data values as *original* documents and refer to the documents with no data values as *structural* documents. As a result, the final status of our corpus consisted of 57 documents, 27 *original* documents and 30 *structural* documents. The size of our own 3 randomly generated documents (R1,R2,R3) are indicated in Table 3 and the size of the *structural* copy of each *original* version of the document can be computed using the following equation:

$$size(structural) = (1 - DR) * size(Original)$$

where DR represents the data ratio of the document.

- **Textual documents:** this category of documents consists of simple structure and high ratio of its contents is preserved to the data values. The ratio of the data contents of these documents represent more than 70% of the document size.
- **Regular Documents** consists mainly of regular document structure and short data contents. This document category reflects the XML view of relational data. The data ratio of these documents is in the range of between 40 and 60 percent.
- **Irregular documents** consists of documents that have very deep, complex and irregular structure. Similar to purely structured documents, this document category is mainly focusing on evaluating the efficiency of compressing irregular structural information of XML documents.

3.2 Data Sets

Our data set consists of the following documents:

EXI-Group is a variant collection of XML documents included in the testing framework of the Efficient XML Interchange Working Group [4].

XMark-Group the XMark documents model an auction database with deeply-nested elements. The XML document instances of the XMark benchmark are produced by the *xmlegen* tool of the XML benchmark project [14]. For our experiments, we generated three XML documents using three increasing scaling factors.

XBench-Group presents a family of benchmarks that captures different XML application characteristics [12]. The databases it generates come with two main models: 1) Data-centric (DC) model contains data that are not originally stored in XML format such as e-commerce catalog data and transactional data 2) Text-centric (TC) model which represents text data that are more likely stored as XML. Each of these two models can be represented either in the form of a single document (SD) or multiple documents (MD). In short, these two levels of classifications are combined to generate four database instances: TCSD, DCSD, TCMD, DCMD. In addition, XBench can generate databases with 4 different sizes: small (11MB), normal (108MB) and large (1GB) and huge (10GB). In our experiments, we only use TCSD and DCSD instances of the small and normal sizes.

Wikipedia-Group Wikipedia offers free copies of all content to interested users [10]. For our corpus, we selected five samples of the XML dumps with different sizes and characteristics.

DBLP presents the famous database of bibliographic information of computer science journals and conference proceedings.

U.S House is a legislative document which provides information about the ongoing work of the U.S. House of Representatives.

SwissProt is a protein sequence database which describes the DNA sequences. It provides a high level of annotations and a minimal level of redundancy.

NASA is an astronomical database which is constructed by converting legacy flat-file formats into XML documents and

then making them available to the public.

Shakespeare represents the gathering of a collection of marked-up Shakespeare plays into a single XML file. It contains many long textual passages.

Lineitem is an XML representation of the transactional relational database benchmark (TPC-H).

Mondial provides the basic statistical information on countries of the world.

BaseBall provides the complete baseball statistics of all players of each team that participated in the 1998 Major League.

Treebank is a large collection of parsed English sentences from the Wall Street Journal. It has a very deep, non-regular and recursive structure.

Random-Group this group of documents has been generated using our own implementation of a Java-based random XML generator. This generator is designed in a way to produce *structural* documents with very random, irregular and deep structures according to its input parameters for the number of unique tag names, maximum tree level and document size. We used this XML generator for producing three documents with different size and characteristics. The main aim of this group is to challenge the examined compressors and assess the efficiency of compressing the structural parts of XML documents.

4. TESTING ENVIRONMENTS

To ensure the consistency of the performance behaviors of the evaluated XML compressors, we ran our experiments on two different environments. One environment with high computing resources and the other with considerably limited computing resources. Table 4 lists the setup details of our high resources environment and Table 5 lists the setup details of the limited one.

Operating System	Ubuntu 7.10 (Linux 2.6.22 Kernel)
CPU	Intel Core 2 Duo E6850 CPU 3.00 GHz, FSB 1333MHz 4MB L2 Cache
Hard Disk	Seagate ST3250820AS 250 GB
RAM	4 GB
Compilers	gcc/g++ 4.1

Table 4: Setup details of the powerful resources environment

Operating System	Ubuntu 7.10 (Linux 2.6.20 Kernel)
CPU	Intel Pentium 4 2.66GHz, FSB 533MHz 512KB L2 Cache
Hard Disk	Western Digital WD400BB 40 GB
RAM	512 MB
Compilers	gcc/g++ 4.1

Table 5: Setup details of the low resources environment

Data Set Name	Document Name	Size (MB)	Tags	Number of Nodes	Depth	Data Ratio
EXI [4]	EXI-Telecomp.xml	0.65	39	651398	7	0.48
	EXI-Weblog.xml	2.60	12	178419	3	0.31
	EXI-Invoice.xml	0.93	52	78377	7	0.57
	EXI-Array.xml	22.18	47	1168115	10	0.68
	EXI-Factbook.xml	4.12	199	104117	5	0.53
	EXI-Geographic Coordinates.xml	16.20	17	55	3	1
XMark [14]	XMark1.xml	11.40	74	520546	12	0.74
	XMark2.xml	113.80	74	5167121	12	0.74
	XMark3.xml	571.75	74	25900899	12	0.74
XBench [12]	DCSD-Small.xml	10.60	50	6190628	8	0.45
	DCSD-Normal.xml	105.60	50	6190628	8	0.45
	TCSD-Small.xml	10.95	24	831393	8	0.78
	TCSD-Normal.xml	106.25	24	8085816	8	0.78
Wikipedia [10]	EnWikiNews.xml	71.09	20	2013778	5	0.91
	EnWikiQuote.xml	127.25	20	2672870	5	0.97
	EnWikiSource.xml	1036.66	20	13423014	5	0.98
	EnWikiVersity.xml	83.35	20	3333622	5	0.91
	EnWikTionary.xml	570.00	20	28656178	5	0.77
DBLP	DBLP.xml	130.72	32	4718588	5	0.58
U.S House	USHouse.xml	0.52	43	16963	16	0.77
SwissProt	SwissProt.xml	112.13	85	13917441	5	0.60
NASA	NASA.xml	24.45	61	2278447	8	0.66
Shakespeare	Shakespeare.xml	7.47	22	574156	7	0.64
Lineitem	Lineitem.xml	31.48	18	2045953	3	0.19
Mondial	Mondial.xml	1.75	23	147207	5	0.77
BaseBall	BaseBall.xml	0.65	46	57812	6	0.11
Treebank	Treebank.xml	84.06	250	10795711	36	0.70
Random	Random-R1.xml	14.20	100	1249997	28	0
	Random-R2.xml	53.90	200	3750002	34	0
	Random-R3.xml	97.85	300	7500017	30	0

Table 3: Characteristics of XML data sets

5. EXPERIMENTS

We evaluated the performance characteristics of XML compressors by running them through an extensive set of experiments. The setup of our experimental framework was very challenging and complex. The details of this experimental framework is described as follows:

- We evaluated 11 XML compressors: 3 general purpose text compressors (gzip, bzip2, PPM) and 8 XML conscious compressors (XMillGzip, XMillBzip, XMillPPM, XMLPPM, SCMPMP, XWRT, Exalt, AXECHOP). For our main set of experiments, we evaluated the compressors under their default settings. The rationale behind this is that the default settings are considered to be the recommended settings from the developers of each compressors and thus can be assumed as the best behaviour. In addition to this main set of experiments, we run additional set of experiments with *tuned parameters* for the highest value of the level of compression parameter provided by some compressors (gzip, bzip2, PPM, XMillPPM, XWRT). That means in total we run 16 *variant* compressors. The experiments of the tuned version of XWRT could be only be performed on the high resource setup because they require at least 1 GB RAM.
- Our corpus consists of 57 documents: 27 *original documents*, 27 *structural copies* and 3 randomly generated *structural documents* (see Section 3.1).
- We run the experiments on two different platforms. One with limited computing resources and the other with high computing resources.
- For each combination of an XML test document and an XML compressor, we run two different operations (*compression* - *decompression*).

- To ensure accuracy, all reported numbers for our time metrics (*compression time* - *decompression time*) (see Section 5.2) are the average of five executions with the highest and the lowest values removed.

The above details lead to the conclusion that our number of runs was equal to 9120 on each experimental platform ($16 * 57 * 2 * 5$), i.e 18240 runs in total.

In addition to running this huge set of experiments, we needed to find the best way to collect, analyze and present this huge amount of experimental results. To tackle this challenge, we created our own mix of Unix shell and Perl scripts to run and collect the results of these huge number of runs. In this paper, we present an important part from results of our experiments. For full detailed results, we refer the reader to the web page of this experimental study [1].

5.1 Errors

During the run of our experiments, some tools failed to either compress or decompress some of the documents in our corpus. We consider the run as unsuccessful if the compressor fails to achieve either of the encoding and decoding processes of the test document. Thus, we had 57 runs for each compressor (one run per document). Figure 1 presents the percentage of unsuccessful runs of each compressor. For a detailed list of the errors generated during our experiments we refer to the web page of this study [1]. We have two main remarks about the results of Figure 1:

- The general purpose text compressors have shown complete stability. They were able to successfully perform the complete set of runs. They are *XML-Blind* thus require no knowledge of the document-structure. Hence, they can deal with any XML document even if it suffers from any syntax or well-formedness problems. However, XML conscious compressors are very sensitive to such problems. For example, some compressors which

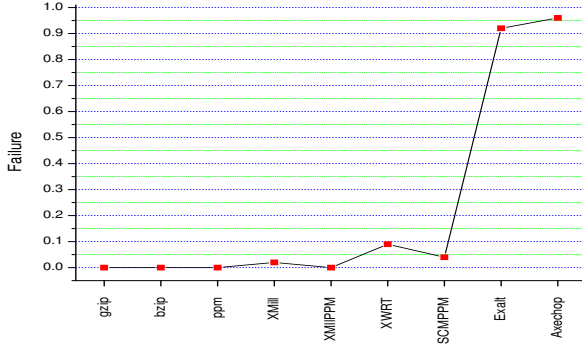


Figure 1: Percentage of unsuccessful runs of each compressor.

uses the Expat XML parser such as XMLPPM will fail to compress any XML document which uses external entity references if it does not have a *dummy* DTD declaration because the XML parser strictly applies the W3C specification and will consider this document as not well-formed.

- Except the latest version of XMLPPM (0.98.3), none of the XML conscious compressors was able to execute the whole set of runs successfully. Moreover, AXECHOP and Exalt compressors have shown very poor *stability*. They failed to run successful decoding parts of many runs. They were thus excluded from any consolidated results. Although an earlier version of XMLPPM (0.98.2) suffered from some problem in decompressing the Wikipedia data sets, the latest version of XMLPPM (0.98.3) released by Cheney during the time of doing the experiments of this work has fixed all earlier bugs and has shown to be the best XML conscious compressor from the *stability* point of view.

5.2 Performance Metrics

We measure and compare the performance of the XML compression tools using the following metrics:

Compression Ratio: represents the ratio between the sizes of compressed and uncompressed XML documents.

$$\text{Compression Ratio} = (\text{Compressed Size}) / (\text{Uncompressed Size})$$

Compression Time: represents the elapsed time during the compression process i.e the period of time between the start of program execution on a document until all the data are written to disk.

Decompression Time: represents the elapsed time during the decompression process i.e the period of time between the start of program execution on reading the decompressed format of the XML document until delivering the original document.

For all metrics: *the lower the metric value, the better the compressor.*

5.3 Experimental results

In this section we report the results obtained by running our exhaustive set of experiments. Figures 2 to 12 represents

an important part of the results of our experiments. Several remarks and guidelines can be observed from the results of our exhaustive set of experiments. Some key remarks are given as follows:

- The results of Figure 2 and Figure 3 show that the *tuned* run of XWRT with the highest level of compression ratio achieves the overall best average compression ratio with very expensive cost terms of compression and decompression times.
- Figure 10(a) shows that the three alternative back-ends of *XMill* compressor achieve the best compression ratio over the structural documents. Figure 4 shows that XMillPPM achieves the best compression ratio for all the datasets. The irregular structural documents (Treebank, R1, R2, R3) are very challenging to the set of the compressors. This explains why they all had the worst compression ratios.
- Figures 5 and 10 show that gzip-based compressors (gzip, XMLGzip) have the worst compression ratios. Excluding these two compressors, Figure 10 shows that the differences on the average compression ratios between the rest of compressors are very narrow. They are very close to each other, the difference between the best and the worst average compression ratios is less than 5%. Among all compressors, SCMPM achieves the best average compression ratio.
- Figures 6,7,8,9 show that the gzip-based compressors have the best performance in terms of compression time and decompression time metrics on both testing environments. The compression and decompression times of the PPM-Based compression scheme (XMillPPM, XMLPPM, SCMPM) are much slower than the other compressors. Among all compressors, SCMPM has the longest compression and decompression times.
- Figure 11 illustrates the overall performance of XML compressors on the high and limited resources setup where the values of the *performance metrics* are *normalized* respect to bzip2. The results of this figure illustrate the narrow differences between the XML compressors in terms of their compression ratios and the wide differences in terms of their compression and decompression times.

5.4 Ranking

Obviously, it is a nice idea to use the results of our experiments and our performance metrics to provide a global ranking of XML compression tools. This is however an especially hard task. In fact, the results of our experiments have not shown a *clear winner*. Hence, different ranking methods and different weights for the factors could be used for this task. Deciding the weight of each metric is mainly dependant on the scenarios and requirements of the applications where these compression tools could be used. In this paper we used three ranking functions which give different weights for our performance metrics. These three rankings function are defined as follows:

- $WF1 = (1/3 * CR) + (1/3 * CT) + (1/3 * DCT)$.
- $WF2 = (1/2 * CR) + (1/4 * CT) + (1/4 * DCT)$

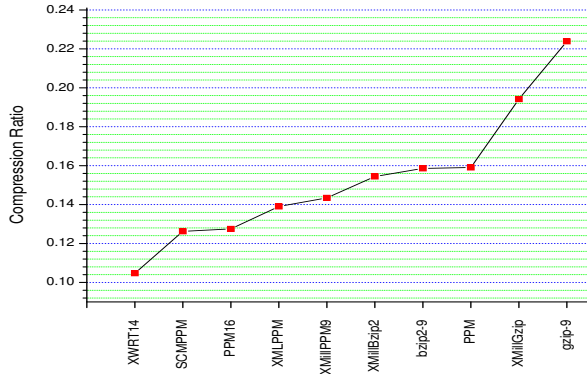


Figure 2: Average Compression Ratios (Tuned Parameters)

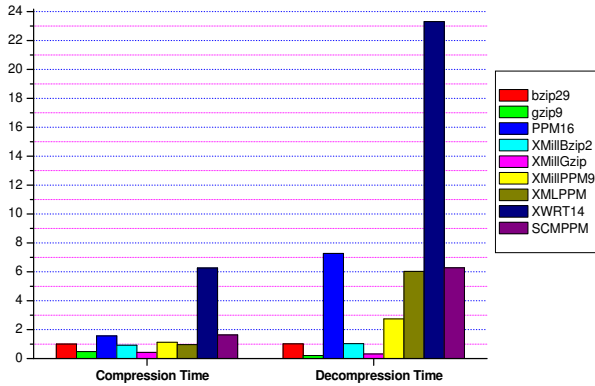


Figure 3: Average compression and decompression times over original documents on high resources setup (Tuned Parameters)

$$\bullet \text{ } WF3 = (3/5 * CR) + (1/5 * CT) + (1/5 * DCT)$$

where CR represents the compression ratio metric, CT represents the compression time metric and DCT represents the decompression time metric. In these ranking functions we used increasing weights for the compression ratio (CR) metric (33%, 50% and 60%) while CT and DCT were equally sharing the remaining weight percentage for each function. Figure 12 shows that *gzip* and *XMLGzip* are ranked as the best compressors using the three ranking functions and on both of the testing environments. In addition, Figure 12 illustrates that none of the XML compression tools has shown a significant or noticeable improvement with respect to the compression ratio metric. The increasing assignment for the weight of CR do not change the order of the global ranking between the three ranking functions.

6. CONCLUSION

We believe that this paper could be valuable for both the developers of new XML compression tools and interested users as well. For developers, they can use the results of this paper to effectively decide on the points which can be improved in order to make an effective contribution. For

this category of readers, we recommend tackling the area of developing stable efficient *queriable* XML compressors. Although there has been a lot of literature presented in this domain, our experience from this study lead us to the result that we are still missing efficient, scalable and stable implementations in this domain. For users, this study could be helpful for making an effective decision to select the suitable compressor for their requirements. For example, for users with highest compression ratio requirement, the results of Figure 2 recommends the usage of either the PPM compressor with the highest level of compression parameter (*ppmd e -o16 document.xml*) or the XWRT compressor with the highest level of compression parameter (*xwrt -l14 document.xml*)(if they have more than 1 GB RAM on their systems) while for the users with fastest compression time and moderate compression ratio requirements, *gzip* and *XMLGzip* are considered to be the best choice (Figure 12).

From the experience and the results of this experimental study, we can draw the following conclusions and recommendations:

- The primary innovation in the XML compression mechanisms was presented in the first implementation in this domain by XMill. It introduced the idea of separating the structural part of the XML document from the data part and then group the related data items into homogenous containers that can be compressed separately. This separation improves the further steps of compressing these homogenous containers using the general purpose compressors or any other compression mechanism because they can detect the redundant data easily. Most of the following XML compressors have simulated this idea in different ways.
- The dominant practice in most of the XML compressors is to utilize the well-known structure of XML documents for applying a pre-processing encoding step and then forwarding the results of this step to general purpose compressors. Consequently, the compression ratio of most XML conscious compressor is very dependent and related on the general purpose compressors such as: *gzip*, *bzip2* or *PPM*. Figure 10 shows that none of the XML conscious compressors has achieved an outstanding compression ratio over its back-end general purpose compressor. The improvements are always not significant with 5% being the best of cases. This fact could explain why XML conscious compressors are not widely used in practice.
- The compression time and decompression time metrics play a crucial role in the ranking of XML compressors.
- The authors of the XML compression tools should provide more attention to provide the source code of their implementations available. Many tools presented in the literature - specially the queriable ones - have no available source code which prevents the possibility of ensuring the repeatability of the reported numbers. It also hinders the possibility of performing fair and consistent comparisons between the different approaches. For example in [30], the authors compared the results of their implementation *Xseq* with *XBzip* using an inconsistent way. They used the reported query evaluation time of *XBzip* in [24] to compare with their times

although each of the implementation is running on a different environment.

- There are no publicly available solid implementations for grammar-based XML compression techniques and queriable XML compressors. These two areas provide many interesting avenues for further research and development.

As a future work, we are planning to continue maintaining and updating the web page of this study with further evaluations of any new evolving XML compressors. In addition, we will enable the visitor of our web page to perform their online experiments using the set of the available compressors and their own XML documents.

7. REFERENCES

- [1] Benchmark of XML compression tools. <http://xmlcompbench.sourceforge.net/>.
- [2] BZip2 Compressor. <http://www.bzip.org/>.
- [3] DTDPPM Compressor. <http://xmlppm.sourceforge.net/dtdppm/index.html>.
- [4] Efficient XML Interchange Benchmark Working Group. <http://www.w3.org/XML/EXI/>.
- [5] Exalt Compressor. <http://exalt.sourceforge.net/>.
- [6] GZip Compressor. <http://www.gzip.org/>.
- [7] PPM Compressor. <http://www.compression.ru/ds/>.
- [8] rngzip Compressor. <http://contrapunctus.net/league/haques/rngzip/>.
- [9] SCMPPM Compressor. <http://www.infor.uva.es/jadiego/files/scmppm-0.93.3.zip>.
- [10] Wikipedia Data Set. <http://download.wikipedia.org/backup-index.html>.
- [11] XAUST Compressor. <http://drona.csa.iisc.ernet.in/priti/xaust.tar.gz>.
- [12] XBench Benchmark. <http://softbase.uwaterloo.ca/ddbms/projects/xbench/>.
- [13] XGrind Compressor. <http://sourceforge.net/projects/xgrind/>.
- [14] XMark Benchmark. <http://monetdb.cwi.nl/xml/>.
- [15] XMill Compressor. <http://sourceforge.net/projects/xmill>.
- [16] XMLPPM Compressor. <http://xmlppm.sourceforge.net/>.
- [17] XQueC Compressor. <http://www.icar.cnr.it/angela/xquec/index.htm>.
- [18] XWRT Compressor. <http://sourceforge.net/projects/xwrt>.
- [19] J. Adiego, P. Fuente, and G. Navarro. Merging Prediction by Partial Matching with Structural Contexts Model. In *DCC '04: Proceedings of the Conference on Data Compression*, page 522, Washington, DC, USA, 2004. IEEE Computer Society.
- [20] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical Report 124, 1994.
- [21] J. Cheney. Compressing XML with Multiplexed Hierarchical PPM Models. In *DCC '01: Proceedings of the Data Compression Conference (DCC '01)*, page 163, Washington, DC, USA, 2001. IEEE Computer Society.
- [22] J. Cheng and W. Ng. XQzip: Querying Compressed XML Using Structural Indexing. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*, volume 2992 of *LNCS*, pages 219–236. Springer, 2004.
- [23] J. G. Cleary and I. H. Witten. Data compression using adaptive coding and partial string matching. *IEEE Transactions on Communications*, OM-32(4):396–402, April 1984.
- [24] P. Ferragina, F. Luccio, G. Manzini, and S. Muthukrishnan. Compressing and searching XML data via two zips. In *WWW '06: Proceedings of the 15th international conference on World Wide Web*, pages 751–760, New York, NY, USA, 2006. ACM.
- [25] M. Girardot and N. Sundaresan. Millau: an encoding format for efficient representation and exchange of XML over the Web. *Comput. Networks*, 33(1-6):747–765, 2000.
- [26] Kieffer and Yang. Grammar-Based Codes: A New Class of Universal Lossless Source Codes. *IEEE Transactions on Information Theory*, 46, 2000.
- [27] G. Leighton, J. Diamond, and T. Muldner. AXECHOP: A Grammar-based Compressor for XML. In *DCC '05: Proceedings of the Data Compression Conference*, pages 467–467, Washington, DC, USA, 2005. IEEE Computer Society.
- [28] W. Li. An XML compression tool. Master's thesis, University of Waterloo, 2003.
- [29] H. Liefke and D. Suciu. XMill: An efficient compressor for XML data. In W. Chen, J. F. Naughton, and P. A. Bernstein, editors, *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data, May 16-18, 2000, Dallas, Texas, USA*, pages 153–164. ACM, 2000.
- [30] Y. Lin, Y. Zhang, Q. Li, and J. Yang. Supporting efficient query processing on compressed XML files. In *SAC '05: Proceedings of the 2005 ACM symposium on Applied computing*, pages 660–665, New York, NY, USA, 2005. ACM.
- [31] J. Min, M. Park, and C. Chung. XPRESS: A queriable compression for XML data. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 122–133. ACM Press, 2003.
- [32] W. Ng, W. Lam, P. T. Wood, and M. Levene. XCQ: A queriable XML compression system. *Knowl. Inf. Syst.*, 10(4):421–452, 2006.
- [33] D. Salomon. *Data Compression: The Complete Reference*. pub-SV, 2004.
- [34] P. Skibinski and J. Swacha. Combining Efficient XML Compression with Query Processing. In *ADBIS*, pages 330–342, 2007.
- [35] P. Skibinski and J. Swacha. Fast Transform for Effective XML Compression. In *CADSM*, pages 323–326, 2007.
- [36] V. Toman. Compression of XML Data. Master's thesis, Charles University, Prague, 2004.
- [37] R. K. Wong, F. Lam, and W. M. Shui. Querying and maintaining a compact XML storage. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 1073–1082, New York, NY, USA, 2007. ACM.

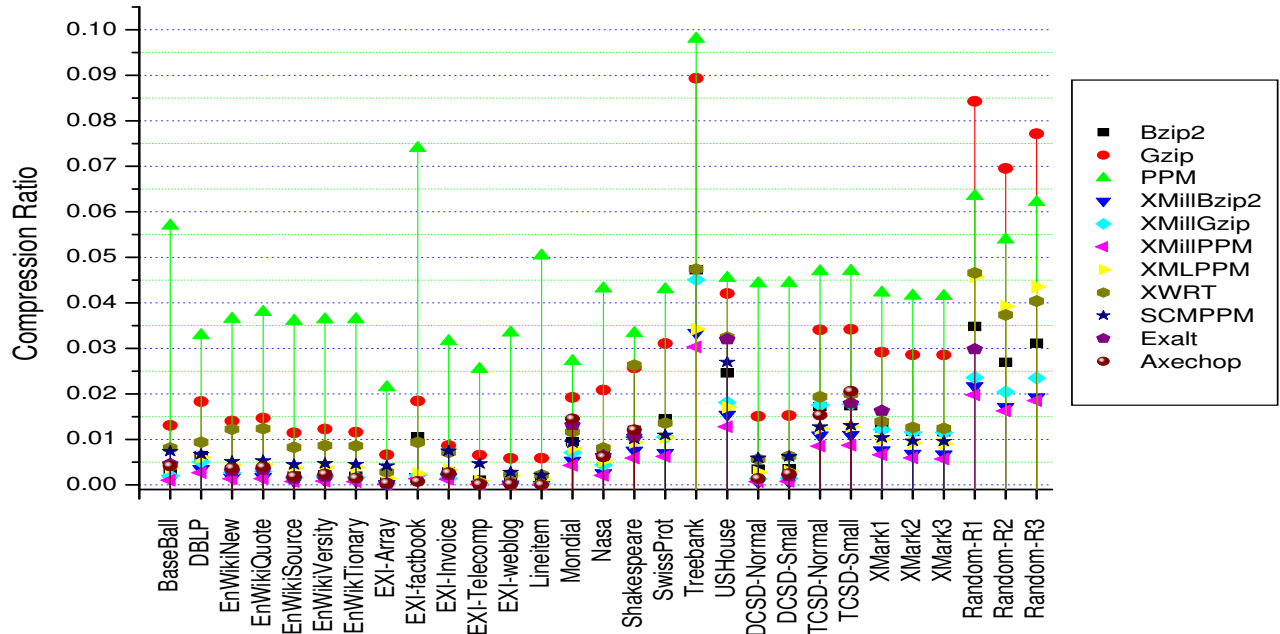


Figure 4: Detailed compression ratios of *structural* documents

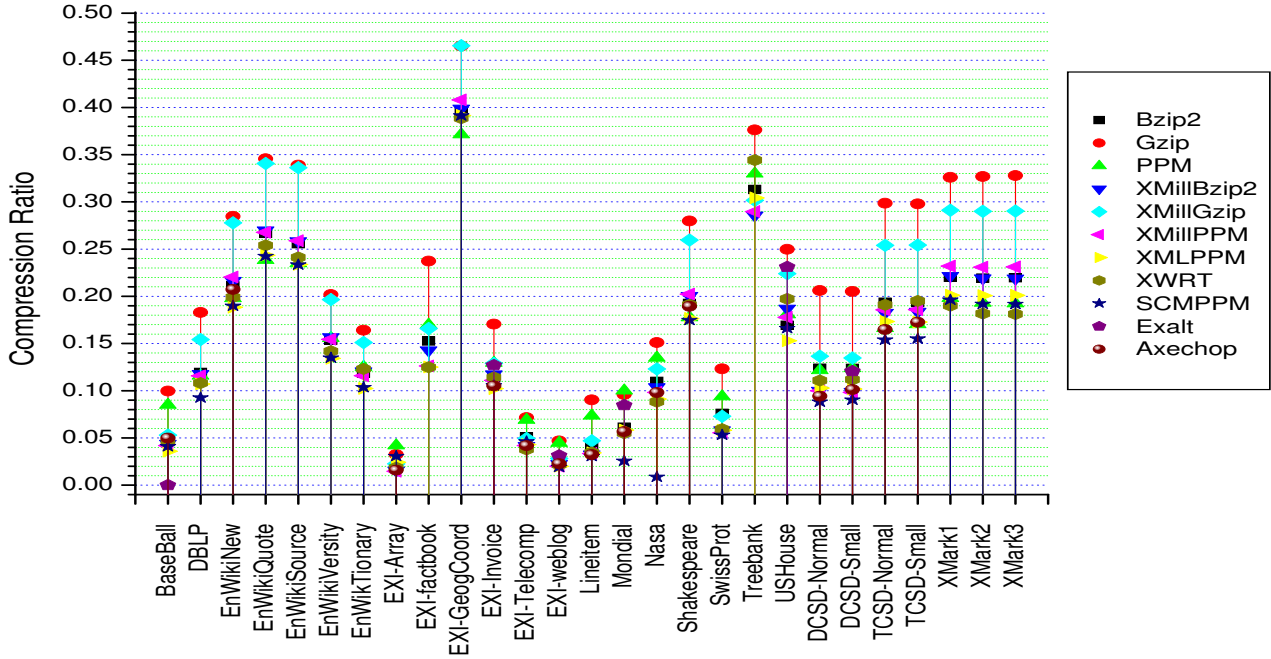


Figure 5: Detailed compression ratios of *original* documents

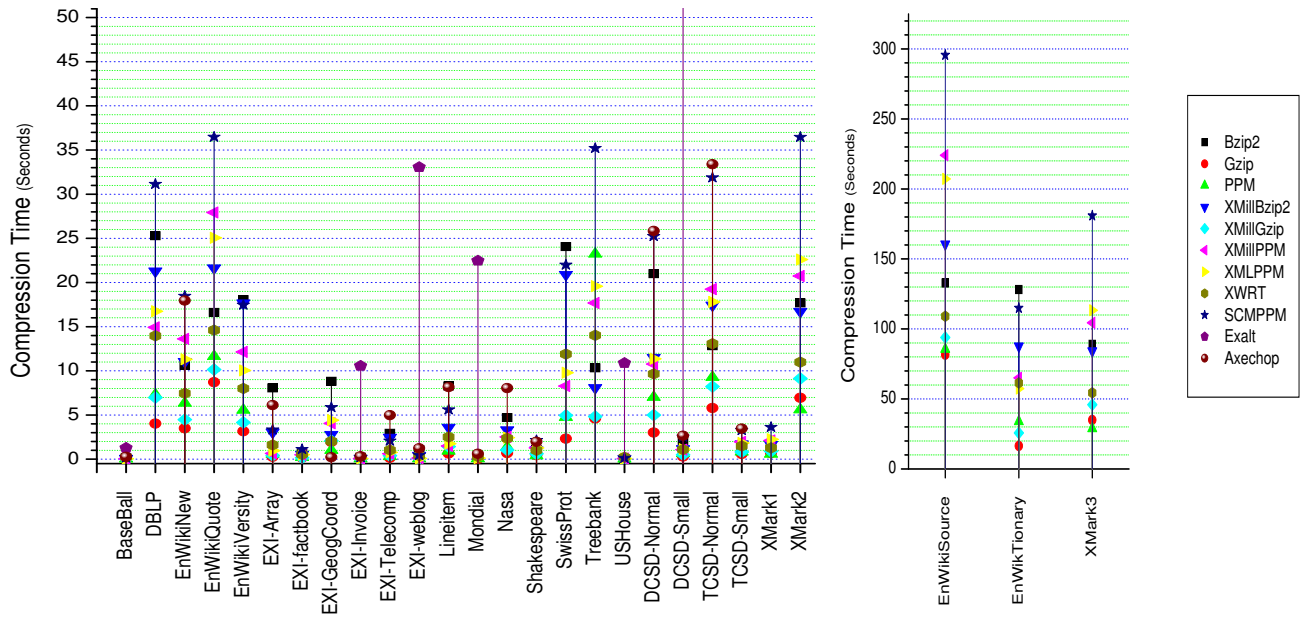


Figure 6: Detailed compression times on the high resources setup.

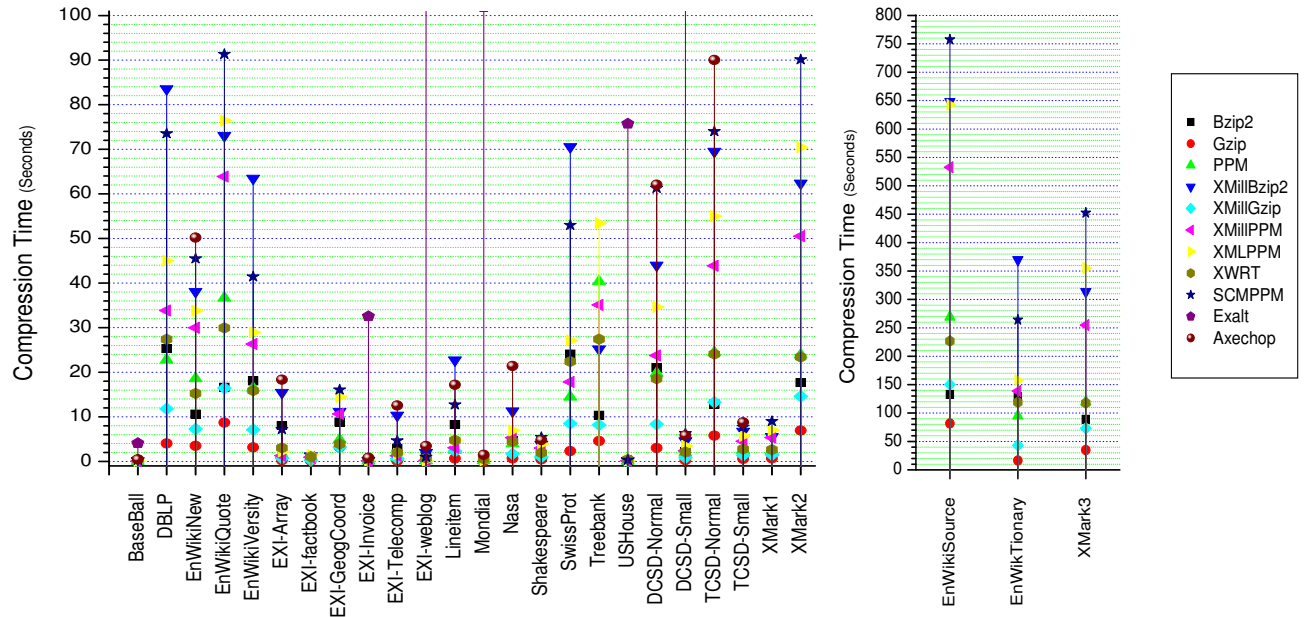


Figure 7: Detailed compression times on the limited resources setup.

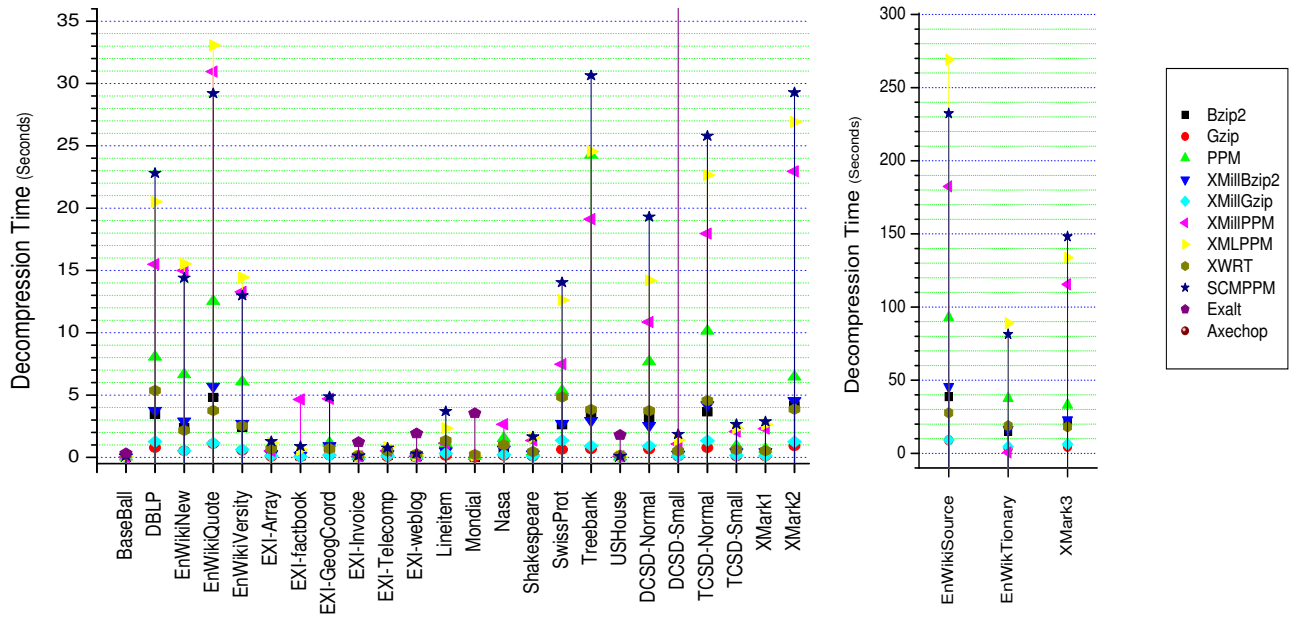


Figure 8: Detailed decomposition times on the high resources setup.

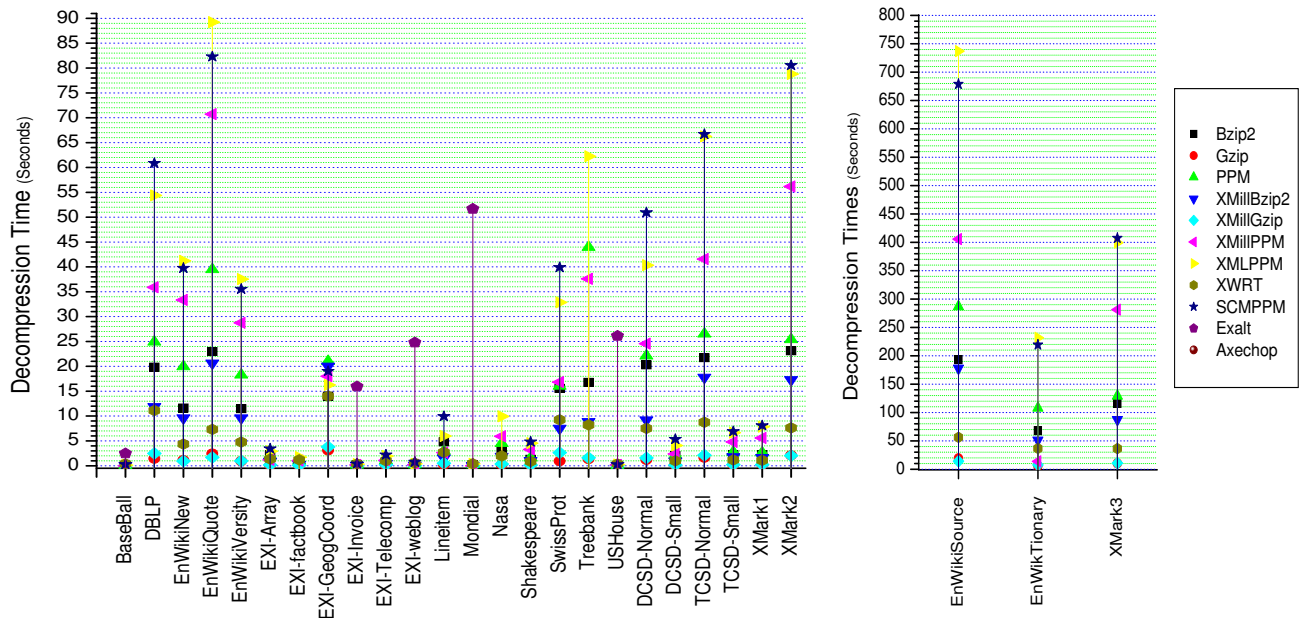
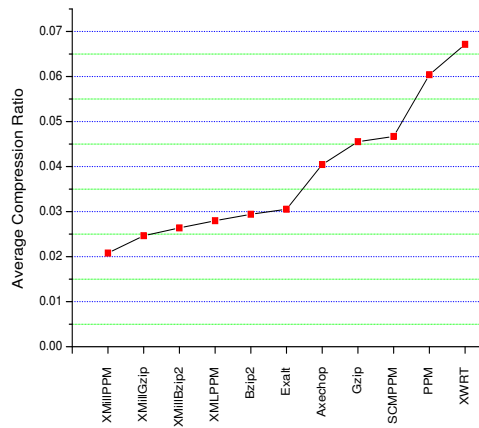
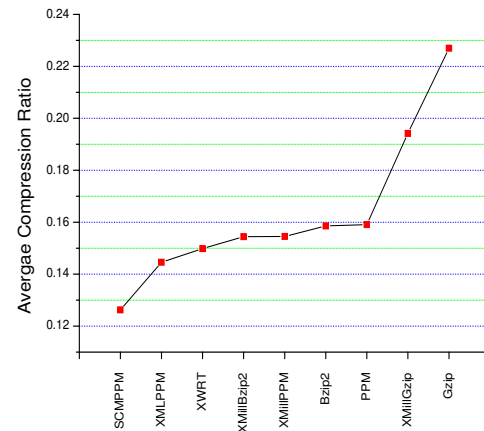


Figure 9: Detailed decomposition times on the limited resources setup.

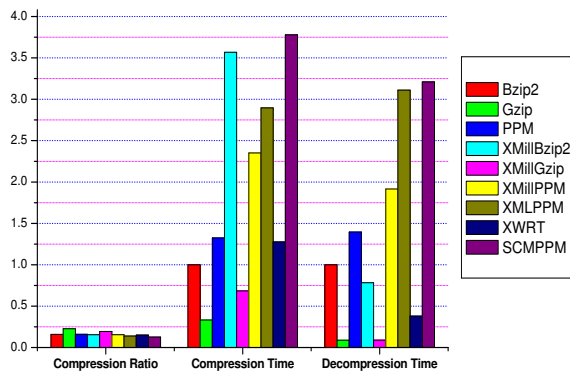


(a) Structural documents.

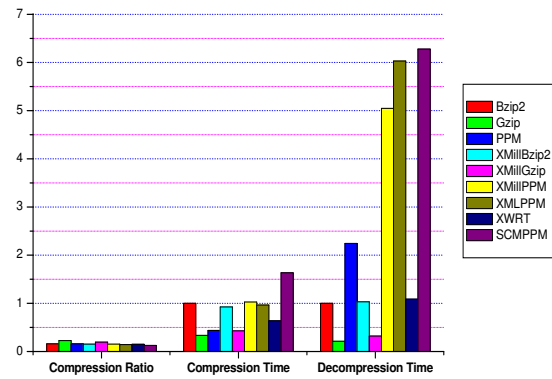


(b) Original documents.

Figure 10: Average compression ratios.

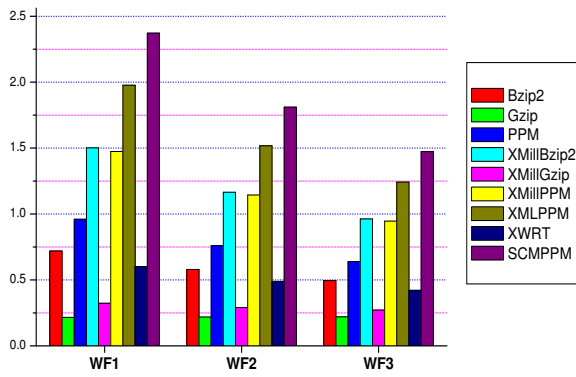


(a) Limited resources setup.

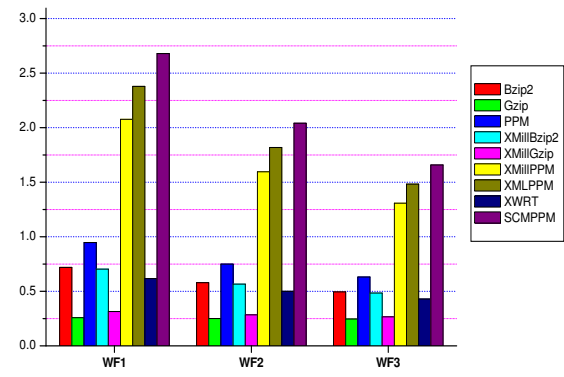


(b) High resources setup.

Figure 11: Overall performance of compressing *original* documents.



(a) Limited resources setup.



(b) High resources setup.

Figure 12: Ranking Functions of compressing *original* documents.