

INVESTIGATING BINARY STRING ENCODING FOR COMPACT REPRESENTATION OF XML DOCUMENTS

Ramez Alkhatib¹

Department of Computer Technology, Hama University, Hama, Syria

ABSTRACT

Since Extensible Markup Language abbreviated as XML, became an official World Wide Web Consortium recommendation in 1998, XML has emerged as the predominant mechanism for data storage and exchange, in particular over the World Web. Due to the flexibility and the easy use of XML, it is nowadays widely used in a vast number of application areas and new information is increasingly being encoded as XML documents. Because of the widespread use of XML and the large amounts of data that are represented in XML, it is therefore important to provide a repository for XML documents, which supports efficient management and storage of XML data. Since the logical structure of an XML document is an ordered tree consisting of tree nodes, establishing a relationship between nodes is essential for processing the structural part of the queries. Therefore, tree navigation is essential to answer XML queries. For this purpose, many proposals have been made, the most common ones are node labeling schemes. On the other hand, XML repeatedly uses tags to describe the data itself. This self-describing nature of XML makes it verbose with the result that the storage requirements of XML are often expanded and can be excessive. In addition, the increased size leads to increased costs for data manipulation. Therefore, it also seems natural to use compression techniques to increase the efficiency of storing and querying XML data. In our previous works, we aimed at combining the advantages of both areas (labeling and compaction technologies), Specially, we took advantage of XML structural peculiarities for attempting to reduce storage space requirements and to improve the efficiency of XML query processing using labeling schemes. In this paper, we continue our investigations on variations of binary string encoding forms to decrease the label size. Also We report the experimental results to examine the impact of binary string encoding on reducing the storage size needed to store the compacted XML documents.

KEYWORDS

XML Compaction, XML Labeling, XML Storage, Binary encoding

1. INTRODUCTION

The ability to efficiently manage XML data is essential because the potential benefits of using XML as a representation method for any kind of data. There have been many proposals to manage XML documents. However, XML Labeling and compaction techniques are considered as two major approaches able to provide robust XML document storage and manipulation.

¹ Part of this work was done while the author was member of the Database and Information Systems Research Group, University of Konstanz

Since the logical structure of an XML document is an ordered tree consisting of tree nodes that represent elements, attributes and text data, establishing a relationship between nodes is essential for processing the structural part of the queries. Therefore, tree navigation is essential to answer XML queries. However standard tree navigations (such as depth- first or breadth-first traversals) are not sufficient for efficient evaluation of XML queries, especially the evaluation of ancestor and descendant axes. For this purpose, many node labeling schemes have been made. The use of labeling schemes to encode XML nodes is a common and most beneficial technique to accelerate the processing of XML queries and in general to facilitate XML processing when XML data is stored in databases [15].

The power of XML comes from the fact that it provides self-describing capabilities. XML repeatedly uses tags to describe the data itself. At the same time this self-describing nature of XML makes it verbose with the result that the storage requirements of XML are often expanded and can be excessive. In addition, the increased size leads to increased costs for data manipulation. The inherent verbosity of XML causes doubts about its efficiency as a standard data format for data exchange over the internet. Therefore, compression of XML documents has become an increasingly important research issue and it also seems natural to use compression techniques to increase the efficiency of storing and querying XML data [3, 4, 6, 8]. In our works, we focused on combining the strengths of both labeling and compaction technologies and bridging the gap between them to exploit their benefits and avoid their drawbacks to produce a level of performance that is better than using labeling and compression independently.

In this paper, we continue our investigations on variations of binary encoding forms that would provide for opportunities to further minimize the storage costs of the labels. The rest of the paper is structured as follows: Section 2 and 3 review The CXQU and CXDLS compaction approaches respectively. In Section 4, we present variations of binary encoding schemes can be used to minimize the storage costs of the labels. Experimental results to study the impact of prefix free encoding schemes on reducing the storage size are presented in Section 5. Finally, we conclude and outline future work in Section 6.

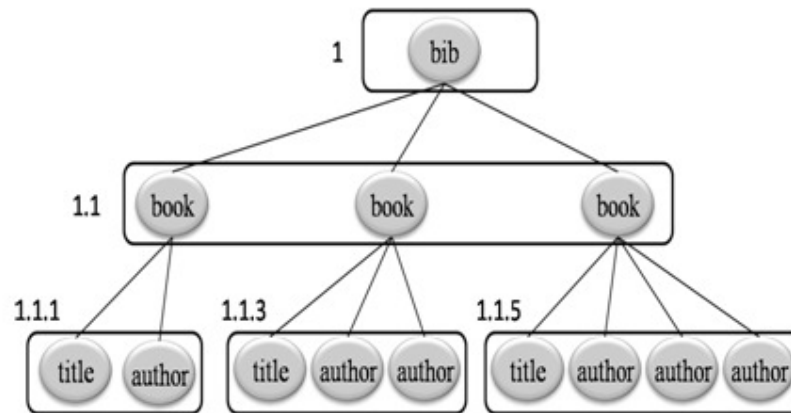


Figure 1. Simple XML document with cluster labels

2. THE CXQU COMPACTION APPROACH

CXQU is our proposed approach [1] to represent XML documents. It not only supports queries and updates but also compacts the structure of an XML document based on the exploitation of repetitive consecutive tags in the structure of the XML documents by using our proposed labeling scheme called Cluster Labeling Scheme (CLS) [1]. CLS assigns a unique identifier to each group

of elements which have the same parent (i.e. sibling element nodes). CLS preserves the hierarchal structure of XML documents after the compaction and supports the managing compacted XML documents efficiently. It allows insertion of nodes anywhere in the XML tree without the need for the subsequent relabeling of existing nodes. To compact an XML document with CXQU, first, it separates its structural information from the content to improve query processing performance by avoiding scans of irrelevant data values. CXQU then compacts the structure using our algorithm, which basically exploits the repetition of similar sibling nodes of XML structure, where “similar” means: elements with the same tag name. CXQU stores the compacted XML structure and the data separately in a robust compact storage that includes a set of access support structures to guarantee fast query performance and efficient Updates. Figure 1 displays the cluster labels and Figure 2 displays the compacted structure of a simple XML document, where the crossed-out nodes will not be stored.

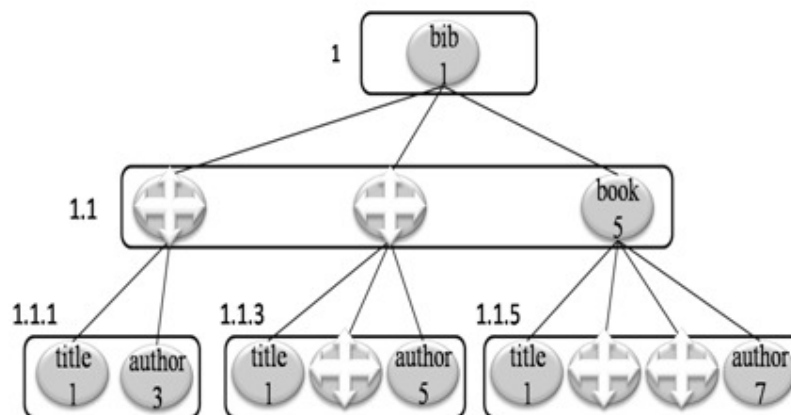


Figure 2. The compacted structure using CXQU

3. THE CXDLS COMPACTION APPROACH

We also proposed an improved technique called CXDLS [2] combining the strengths of both labeling and compaction techniques. CXDLS bridges the gaps between numbering schemes and compaction technology to provide a solution for the management of XML documents that produces better performance than using labeling and compaction independently. CXDLS compacts the regular structure of XML efficiently. At the same time, it works well when applied to less regular or irregular structures. While this technique has the potential for compact storage, it also supports efficient querying and update processing of the compacted XML documents by taking advantage of the ORDPATH labeling scheme. ORDPATH [14] is a particular variant of a hierarchical labeling scheme, which is used in Microsoft SQL Server's XML support. It aims to enable efficient insertion at any position of an XML tree, and also supports extremely high performance query plans for native XML queries.

CXDLS helps to remove the redundant, duplicate subtrees and tags in an XML document. It takes advantage of the principle of separately compacting structure from data and it also uses the ORDPATH labeling scheme for improving the query and update processing performance on compacted XML structures.

In CXDLS, the XML structure is compacted based on the basic principle of exploiting the repetitions of similar nodes in the XML structure, where two nodes N and N' of XML structure are said to be „similar“ if they are consecutive elements, i.e. sibling nodes, in the structure and have exactly the same tag name. Another principle is to exploit the repetitions of identical

subtrees, where two subtrees S and S' of XML structure are said to be „identical“ if they are consecutive and have exactly the same structure. Figure 3 shows the ORDPATH labels and Figure 4 displays the compacted structure using CXDLS.

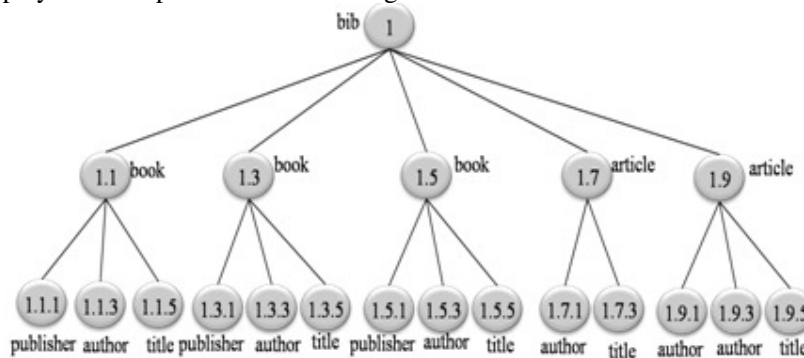


Figure 3. Simple XML document with ORDPATH labels

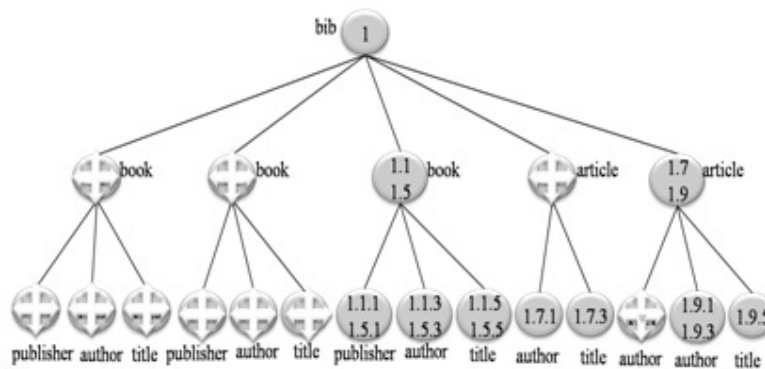


Figure 4. The compacted structure using CXDLS

4. BYTE REPRESENTATION OF THE LABELS

To achieve low storage consumption for XML documents, we have to reduce the size of node labels. Therefore, both ORDPATH and Cluster labeling schemes used Unicode-like compact representation that consists of a compressed binary representation and a prefix free encoding. It uses successive variable length L_i/O_i bitstrings and is generated to maintain document order and allow cheap and easy node comparisons. One L_i/O_i bitstring pair represents a component of a label. L_i bitstring specifies the number of bits of the succeeding O_i bitstring. The L_i bitstrings are represented using a prefix free encoding that can be constructed using a Huffman tree, an example for a prefix free encoding shown in figure 5(a). The binary encoding of a label is produced by locating each component value in the O_i value ranges and appending the corresponding L_i bitstring followed by the corresponding number of bits specifying the offset for the component value from the minimum O_i value within that range.

Example: Let us consider the bitstring pairs translation for the label (1.3.22). Note that the first component '1' is located in the O_i value range of [0, 7]. So that the corresponding L_0 bitstring is 01 and the length $L_0 = 3$, indicating a 3-bit O_0 bitstring. We therefore encode the component "1" with $L_0 = 01$ and $O_0 = 001$. Similar to that the binary encoding of the component "3" is the bitstring pair $L_1 = 01$, $O_1 = 011$. The component 22 is located in the O_i value range of [8,23] and its corresponding L_2 bitstring 100 and the length $L_2 = 4$. Thus the O_2 bitstring is 1111 that is the

offset of 15 from 8 specified in 4 bits. As final result the bitstring 01001010111001111 is the binary encoding of the cluster label (1.3.22).

Variations of prefix free encoding schemes can be created using the idea of Huffman trees, Figure 5 show different forms of prefix free encoding schemes.

Because the labels are binary encoded and stored in a byte array, in the case of use the codes in Figure 5(a) or the codes in Figure 5(b), the last byte may be incomplete. Therefore, it is padded on the right with zeros to end on an 8-bit boundary. This padding can lead to an increase in the storage requirements. For example, by using codes in Figure 5(a), the binary encoding of 1.9 is 010011000001 but its total length in bytes is 2 bytes and will be stored as the following bitstring 0100110000010000. Also by using codes in Figure 5(a), the label 1.9, for example, results in the bit sequence 0111100001, but it is padded by zeros to store it in 2 byte arrays as 0111100001000000 bitstring. In order to avoid padding with zeros, prefix free encoding scheme, shown in figure 5(c), was designed in a way that each division already observes byte boundaries.

Bitstring	Li	Oi value range
01	3	[0, 7]
100	14	[8, 23]
101	6	[24, 87]
1100	8	[88, 343]
1101	12	[344, 4439]
11100	16	[4440, 69975]
11101	32	[69976, 4.3×10 ⁹]
11110	48	[4.3×10 ⁹ , 2.8×10 ¹⁴]

(a)

Bitstring	Li	Oi value range
01	0	[1, 1]
10	1	[2, 3]
110	2	[4, 7]
1110	4	[8, 23]
11110	8	[24, 279]
111110	12	[280, 4375]
1111110	16	[4376, 69911]
11111110	20	[69912, 1118487]

(b)

Bitstring	Li	Oi value range
0	7	[1, 127]
10	14	[128, 16511]
110	21	[16512, 2113663]
1110	28	[2113664, 270549119]
1111	36	[270549120, ~ 237]

(c)

Figure 5. Variations of prefix free encoding schemes

5. THE IMPACTS OF PREFIX FREE ENCODING SCHEMES

In order to examine the impact of prefix free encoding schemes, mentioned in the previous section, on reducing the storage size needed to store the XML documents that are compacted using our approaches CXQU and CXDLS. We did experiment to measure and compare the storage requirements of our approaches and our cluster labeling scheme with other labeling schemes, such as OrdPath and Dewey [7,14]. In the experiment, each approach is suffixed with a number that refers to a prefix free encoding scheme, where number 1 refers to Figure 5(a) and so on respectively. We conducted our experiment using a variety of both synthetic and real datasets that covered a variety of sizes [5, 9, 10, 11, 12, 13, 16], application domains, and document characteristics. Table 1 displays the different structural properties of the used datasets.

Table 1. XML datasets used in the experiments

Datasets	File name	Topics	Size	No. of elements	Max depth
D1	Mondial	Geographical database	1,77MB	22423	6
D2	OT	Religion	3,32 MB	25317	6
D3	NT	Religion	0,99 MB	8577	6
D4	BOM	Religion	1,47 MB	7656	6
D5	XMark	XML benchmark	113 MB	1666315	12
D6	NCBI	Biological data	427,47 MB	2085385	5
D7	SwissPort	DB of protein sequences	112 MB	2977031	6
D8	Medline02n0378	Bibliography medicine science	120 MB	2790422	8
D9	medline02n0001	Bibliography medicine science	58,13 MB	1895193	8
D10	Part	TPC-H benchmark	6,02 MB	200001	4
D11	Lineitem	TPC-H benchmark	30,7 MB	1022976	4
D12	Customer	TPC-H benchmark	5,14 MB	135001	4
D13	Orders	TPC-H benchmark	5,12 MB	150001	4
D14	TOL	Organisms on Earth	5,36MB	80057	243

It is clearly visible from the results of the experiment in Figures 6, 7 and 8, that the use of third prefix free encoding scheme in our approaches made them more efficient in term of storage requirements for various XML data sets, when compared to other prefix free encoding schemes. These results confirm that the success rate of the use of our approaches (SCQX, CXDLS and cluster labeling scheme) is very high and they can dramatically reduce the storage requirements for almost all the datasets.

From result in Figure 8, it can be observed that the storage requirements, by using the approach CXDLS, are very small for the documents such as PART, Lineitem, Order and Customer because they have a regular structure and CXDLS focuses on compacting regular XML structures. At the same time the storage requirements are still relatively small for other documents that have either an irregular structure or less regular structure.

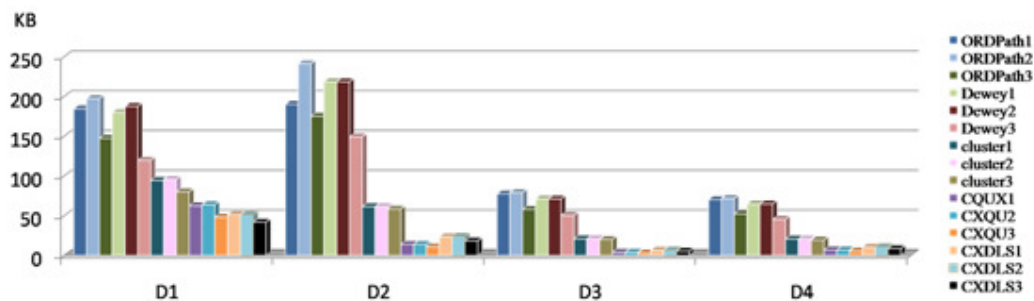


Figure 6. The storage requirements

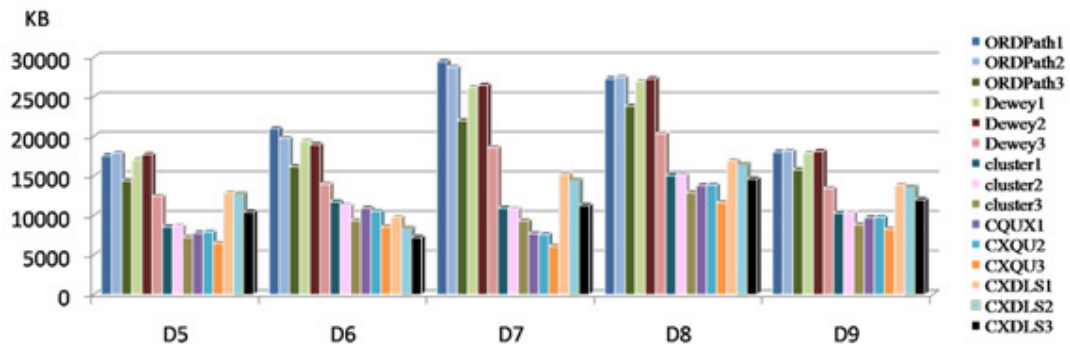


Figure 7. The storage requirements

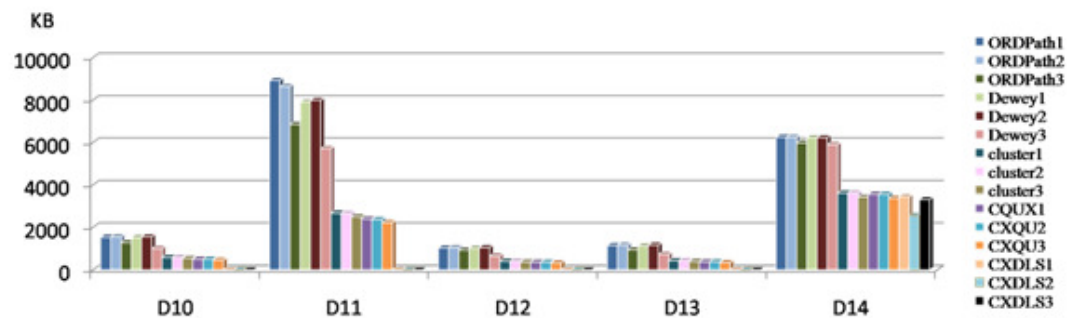


Figure 8. The storage requirements

6. CONCLUSIONS

This work investigated prefix free encoding technique created using the idea of Huffman trees. Our experimental results indicate that it is possible to provide significant benefits in terms of the storage requirements by using prefix free encoding, our compaction and labeling scheme techniques. An interesting future research direction is to explore more encoding formats and study how our compaction techniques could be extended to these formats. Since minimizing the storage costs can further improve query and update performance, one other possible future direction is to test the influence of prefix free encoding schemes on the query and update performance.

REFERENCES

- [1] R. Alkhatib and M. H. Scholl. Cxqu: A compact xml storage for efficient query and update processing. In P. Pichappan and A. Abraham, editors, ICDIM, pages 605–612. IEEE, 2008.
- [2] R. Alkhatib and M. H. Scholl. Compacting xml structures using a dynamic labeling scheme. In A. P. Sexton, editor, BNCOD, volume 5588 of Lecture Notes in Computer Science, pages 158–170. Springer, 2009.
- [3] M. Ali, M. A. Khan: Efficient parallel compression and decompression for large XML files. Int. Arab J. Inf. Technol. 13(4):403-408, 2016
- [4] H. AlZadjali, Siobhán North: XML Labels Compression using Prefix-encodings. WEBIST, pages 69-75, 2016
- [5] J. Bosak. Xml-tagged religion. Oct 1998. <http://xml.coverpages.org>.

- [6] S. Böttcher, R. Hartel, C. Krislin: CluX - Clustering XML Sub-trees. ICEIS, pages 142, 150, 2010
- [7] T. Härder, M. P. Haustein, C. Mathis, and M. W. 0002. Node labeling schemes for dynamic xml documents reconsidered. *Data Knowl. Eng.*, 60(1):126–149, 2007.
- [8] M. Lohrey, S. Maneth, R. Mennicke: XML tree structure compression using RePair. *Inf. Syst.* 38(8): 1150-1167, 2013
- [9] D. R. MADDISON, K.-S. SCHULZ, and W. P. MADDISON. The tree of life web project. *ZOOTAXA*, pages 19–40, 20 Nov. 2007.
- [10] W. May. Information extraction and integration with FLORID: The MONDIAL case study. Technical Report 131, Universität Freiburg, Institut für Informatik, 1999. Available from <http://dbis.informatik.uni-goettingen.de/Mondial>.
- [11] G. Miklau. Xml repository. <http://www.cs.washington.edu/research/xmldatasets>.
- [12] NCBI. National center for biotechnology information(ncbi) xml data format. <http://www.ncbi.nlm.nih.gov/index.html>.
- [13] NLM. National library of medicine (nlm) xml data format. <http://xml.coverpages.org>
- [14] P. E. O’Neil, E. J. O’Neil, S. Pal, I. Cseri, G. Schaller, and N. Westbury. Ordpaths: Insert-friendly xml node labels. In G. Weikum, A. C. König, and S. Deßloch, editors, *SIGMOD Conference*, pages 903–908. ACM, 2004.
- [15] Tatarinov, S. Viglas, K. S. Beyer, J. Shanmugasundaram, E. J. Shekita, and C. Zhang. Storing and querying ordered xml using a relational database system. In M. J. Franklin, B. Moon, and A. Ailamaki, editors, *SIGMOD Conference*, pages 204–215. ACM, 2002.
- [16] T. web project. the tol tree structure. 1998. <http://tolweb.org/tree>