# The Imperative Study on XML Parsing and Parsing Techniques

**Swati Gupta[1]* Dr. Ashish Chourasia[2]**

[1] Research Scholar, University of Technology, Jaipur, Rajasthan

[2] University of Technology, Jaipur, Rajasthan

*Abstract – XML plays a very important role as a global data exchange interface. This requires XML documents to be accessed by users. By using an XML parser, This compliance with the pattern preserve to verified. The parser creates it simple to access the data and moreover guarantee so it is reliable. Parsers can be contrasted by checking their compliance with the WWW Consortium's XML guidelines. Thus, extracting data from XML documents and creating XML documents become important topics of discussion. There are many APIs (Application Program Interfaces) available which can perform these operations. In this paper we are analysis of different XML parsers is useful in assessing the consistency and weaknesses of the objects. Different tests were conducted that contrast with gauges, speed, memory use, etc. DOM is a bad memory because it needs to keep the entire record tree in memory, rendering it unable to manage extremely large records. Objects in the DOM tree can addressed& manipulated by utilizing approach on the objects.*

*Keywords – XML, Parsers, DOM, API, SAX, Parsing Techniques*

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - x - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

## 1.1    INTRODUCTION

Mega knowledge is exchanging and distributing in the present world, the XML is a very important function as a global information exchange interface. 21th century is known as an era of information. Most of the businesses are shifted on the internet. These business transactions create huge amount of data. Thus today, designing of data management system is a major challenge for computer system which can exchange the information between two different applications or two different transactions over internet. This operation is supported by an XML. As an accepted standard for web-based data exchange, an XML is growing. This helps clients to exchange documents from XML. XML is competent to extract information from an XML document that does not contain any facts or information on the content. To achieve this clarity, XML documents will adhere to XML specifications[1].

This specific compliance can be verified by using an XML parser. The parser makes the information easy to access and also ensures that it is valid. Within today's countless XML parsers, coded in a truth of dialects, these parsers may not offer the comparative output of parsing speeds, accuracy, and preconditions for ability. They has two crucial aim to parse XML are follows: tree and stream. Parsing of the tree form requires stacking the entire XML

document in memory. The layout of the tree record takes into account arbitrary access to the components of the document and to modify XML. Tree-type parsing instances include the DOM & Simple XML. The DOM is the most recognizable parser based on the tree. The simplest tree-based code parser is Simple XML. Stream-based parsers are so called in the light of the fact that they parse the XML in a stream with a lot of the same basis as spilling speech, operating with a common hub, and, when finished with that hub, completely overlooking its truth. XML Reader is a draw parser and you likewise code the result table in a cursor for a database query. It makes working with new or excentric XML documents easier. This is eventually marked out that a parser must be selected to match all specific prerequisites for the execution of time-reserve funds, correct parsing, and power needs. In this investigation, XML document will be checked the assorted operating systems such as WIN7, WIN8, UBUNTU, Red Hat, and so on. with three DOM APIs, such as PHP, JAVA and Microsoft. The main purpose of this inquiry is to check whether all the operating system affects the parsing speed and, in the event of the perfect combination of parser and operating framework, this will be found at a certain stage. Due to segment we depict regarding various XML Parsers[2].

www.ignited.in

## 1.2    XML INDEXED STRUCTURE

Files pre-create on XML information can encourage XML inquiries processing by finding data rapidly. Index is built on traversing paths with the element set or on the value of elements and attributes. Broadly, XML indexes are classified into two types: 'Value Indexes' and 'Structural Indexes'. Value indexes are built on XML data values i.e. from value node. The structural indexes are built on structure of XML documents. The 'Value Indexes' answer content based query effectively while 'Structural Indexes' answer all types queries i.e. Contentment relationship, Order query etc. The XML indexes are discussed below:-

1)    **Structural Summary**: - A structural summary is smaller version of an XML tree where all paths from root node to any leaf node in the actual XML tree are preserved. The main purpose of a structural summary is to eliminate redundant structural information without losing structural constraints; that is the structural relationship between XML element such as P-C & A-D relationships. The drawback of these proposals is it has very high index size. It grows exponentially with XML file, sometime it is equal to the XML file itself. These proposals can solve total matching queries effectively but cannot solve partial matching and twig queries.

2)    **Selectivity Indexes**: -These index techniques are implemented from the RDBMS literature. In RDBMS literature, the indexes are created on certain fields which are selected by Database Administrator (DBA). The field of indexing is selected on different criteria like on frequently triggered XML queries, the fields which satisfy certain SQL queries.

3)    **Structural Join Indexes**: - The basic framework of any structural index is to encode some structural relationships between XML nodes (elements attribute etc) so that query processor is able to predict results by simply approaching a corresponding index without accessing the actual data file.

   A.    **Node Index approach: -** The major purpose of the node-labeling approach is: a) assigning a unique code for all node in the XML tree. b) Preserving the nested hierarchical relationship (structural) during XML updates, c) minimizing the re-labelling cost (including the processing time and I/O accessibility) in the case of data updates, and d) reducing the storage space for store generated code.

   B.    **Path Encoding Approach:** - This type of index has idea of creating a path summary of XML data to speed up the processes of query evaluation. This index stores all path information starting from root node to any arbitrary XML node and indexes are created on these paths.

   C.    **Sequence based Approach**:- Sequence based indexes convert both the XML query and XML documents into sequences and use the well-established sequences matching techniques to obtain query answer. In this approach, XML documents and branching queries are represented as sequence, and subsequent matching provides the query answers.

## 1.3    DIFFERENT PARSER AND PARSING TECHNIQUE

Parsing levels are variable. Through Lexical Review to Precise Parsing. Fuzzy Sorting, Island Grammars, Skeleton Grammars & Error Repair are among those two barriers.
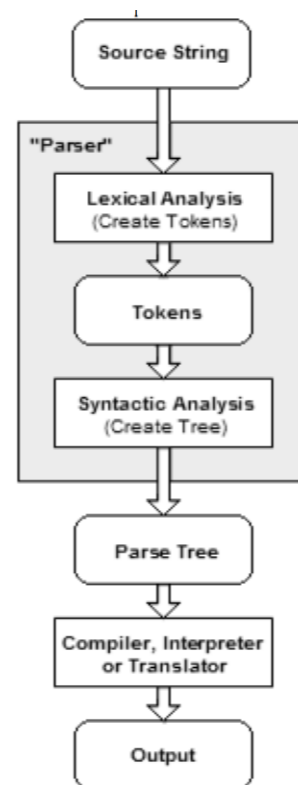


**Fig 1.1: parsing process**

a.    **Fuzzy Parsing**

To endorse multiple programming languages or different dialects of the same programming language, often reengineering systems use a

**Swati Gupta[1]\* Dr. Ashish Chourasia[2]**

method of fuzzy parsing. The purpose of a Fuzzy Parser is to obtain a partial source code model based on a syntactical. Fuzzy parsing's key idea is that some anchor terminals exist. The parser skips all feedback until an anchor terminal is identified and then attempts context-free analysis with an output beginning with the anchor terminal found

### b. Island Grammars

We get accommodating parsers with grammar on the island. It is a grammar consisting of detailed productions defining some interesting constructs (islands) and liberal productions capturing the rest (water). We can trade off precision, completeness and speed of creation by varying the quantity and information in outputs for value constructs. In addition to the one we just described[ MOON 01], there are some different versions of island grammars known. Leon Moonen is thinking about the following:

- Lake grammar: We get a lake grammar when we begin with a complete grammar of a language and expand it with a number of liberal (water) outputs. These grammar is useful in enabling arbitrary embedded code in the software that we want to process.

- Islands with lakes: This is a combination of island and sea productions. We may designate nested buildings as lake-shaped islands.

- Island lakes: this is another mix of island and water production.

## 1.4 MARKUP LANGUAGES USED IN PARSING

### A. Standard Generalized Markup Language (SGML)

It interacts with the online documentation technical markings. The standard SGML paper consists of a declaration of DTD or Paper Form, one of the components of the top level (or else defined as marks or markups), paragraphs and text. Based on the type of document you are writing, the top-level element should & we are going to use it for our records. This is an example of a good document in SGML.

```
<!DOCTYPE article PUBLIC "-//OASIS//DTD DocBook
V3.1//EN">
<article>
  <sect1      id="introduction"><title>Hello      world
introduction</title>
    <para>
    Hello world!
    </para>
  </sect1>
</article>
```

Remember on the paper how we reported on the authorization that used. This is essential; you would get all kinds of errors if you run the file through the SGML parser, if you forget that. Once you build it, that information will not be viewable. The reason it's not visible is that the parser thinks it's just a comment (and it's!) so it just drops out of the final document that has been parsed.

### B. Hyper Text Markup Language (HTML)

HTML is not a language for programming but rather a language for marking up. If you recognize XML already, HTML is a snap to remember. We advise you not to seek and blast in one sitting through this tutorial. Alternatively, we suggest that you learn HTML for 15 minutes to an hour a day and then take a break and let the details sink in. We really don't go anywhere! For several years, HTML hasn't been available.

### C. Extensible Markup Language (XML)

XML is a markup language which distinguished several rules for encoding documents in a smart and machine-discernible structure. XML's program objectives emphasize effortlessness, all-inclusive comment and various human languages, which is fictitious knowledge system by solid maintain via Unicode.

Despite the fact that the XML design revolves around documents, it is typically used to represent self-assertive information systems, for example in web administration Several APIs have been created to assist software developers process XML content, and there are a few blueprint constructs to assist with XML-based language sense. XML assertion XML documents get started by proclaim such data regarding themselves, as follows example: XML is utilized for data organization. The Structured Data includes items like spreadsheets, address books, parameters of design, exchanges related to money, and detailed drawings. XML is a lot of rules for planning content organizations (you can also call them rules or shows) that allow you to organize your data.

## 1.5 XML PARSERS

Oracle includes Java, C, C++, & PL / SQL XML parsers. This chapter mainly addresses Java's parser. All of these parsers is a standalone XML feature that parses an XML document (& probably a standalone content type specification (DTD) or XML schema) that could interpret by your code. The examples set out in this section are in Java

**Swati Gupta[1]\* Dr. Ashish Chourasia[2]**

**Fig 1.2 XML Parser for Java**

## 1.6 DOCUMENT OBJECT MODEL [DOM]

The DOM is a cross-platform & language-independent display in HTML, XHTML and XML documents to communicate to and interact with objects. By using strategies on the products, articles in the DOM tree may be tendered and managed. In its application programming interface (API), the transparent interface of a DOM is defined. Historical experience of the Document DOM is interwoven with the historical background of the late 1990's "tech wars" between Netscape Navigator & Microsoft Internet Explorer, and that of JavaScript & JScript, the very first scripting languages to be extensively updated in internet browser programming engines.



**Fig 1.3 DOM Model**

## 1.7 XML AND THE JAVA™

XML & Java Framework are paradise-build businesses. XML defines a cross-platform architecture of knowledge, and Java provides a traditional cross-platform programming interface. At the same period, XML and Java code assign programmers for one-time registration, operate everywhere TM main for data collection, and generate documentation for both Java-based & non-Java-based systems.

## 1.8 PARSING

Parsing may often be viewed as a descriptive concept, for example when explaining how words are broken up into sentences on the garden road. Parsing, syntax analysis, or syntactic analysis is the method of interpreting a set of symbols that conforms to the principles of formal grammar, as in human language, programming languages or data structures. Parsing in this context relates to how humans view a sentence or expression rather than machines (in spoken language or text) "in form of syntax representatives, recognizing sections of words, syntactic associations, etc." It is especially important when addressing what linguistic signals help speakers decode sentences on the garden road. Parsers range from very simple functions like scanf, to complex programs like a C++ compilers frontend, or a web browser's HTML parser. Regular expressions are already being used in certain ways prior to parsing, as the lexing phase whose performance is then utilized by the parser. In programming, a parser is one of the components in an interpreter or compiler that tests for appropriate syntax & constructs a data structure (also some sort of parse tree, abstract syntax tree or some hierarchical structure) implied in the input tokens.



**Fig 1.4:- parsing model**

## 1.8 SIMPLE API FOR XML PARSING (SAX)

The SAX is an open domain API which is considerably urbanized by the mailing catalog partner XML-DEV. It gave the procedure of parsing an XML document an event-driven interface. An event driven interface presents a method for a "callback" warning to the application's cipher as the basic parser distinguishes the document's XML syntactic structure.

David Megginson, Head of Megginson Technologies, led the development of the XML (SAX) Simple API, a widely used specification that depicted how XML parsers can efficiently pass information from XML documents to application software. SAX was first introduced in Java and so now nearly all big programming languages support it.

The SAX API utilizes a "push parsing" approach where the XML document is interpreted by a SAX

**Swati Gupta[1]\* Dr. Ashish Chourasia[2]**

parser & invokes methods on a delegate (a Content Handler) to handle which ever event is found in the XML document. Frequently, one never writes a parser, but one does provide a handler to collect all the information needed from the XML document.

The SAX interface conquers the drawbacks of the DOM interface by keeping only the lowest possible data necessary at the level of the parser (example; namespaces contexts, validation state), so only information that is kept in memory by the Content Handler-for which you, the developer, are responsible. The tradeoff is that with an strategy, there is no way to "go back in time / the XML document": whereas DOM allows a node to go back to its parent, because there's no such possibility in SAX.

The StAX API adopts a similar approach to XML processing as the SAX API (i.e., event driven), the only very significant difference being that StAX is a pull parser (where SAX was a push parser). For SAX the parser is in charge, and the Content Handler uses callbacks. In Stax, you call the parser and check when / if you want to get the next "event" XML. Figure shows the basic outline of the SAX parsing APIs. To start the process, a SAX Parser Factory class instance is utilized to generate a parser instance. A SAX Reader object is wrapped in by the parser. When invoking the parser's parse) (technique, the reader invokes several of the callback methods that are implemented within the application. Those techniques are described by the Content Handler, Error Handler, DTD Handler, & Entity Resolver interfaces.
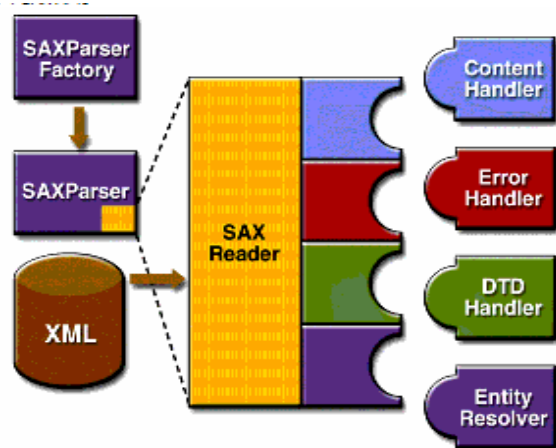


**Fig 1.5 SAX APIs**

## 1.9    PULL PARSING

Pull parsing observes the record as development of object read in grouping exploit the iterated pattern proposes. Its regard as make of recursive plummet parsers in that the configuration of the cipher live out the parsing reflect the construction of the XML being parsed, and center of the highway parsed outcomes get utilize and found nearby aspect within the

capacities playing out the parsing, or went down (as capacity parameters) into lower-level capacities, or returned (as capacity return esteems) to more significant level functions.[21] Instances of pull parsers include Data::Edit::Xml in Perl, StAX in Java, XML Pull Parser in Smalltalk, XML Reader in PHP, Element Tree. iterparse in Python, System. Xml. Xml Reader in the. NET Framework, and DOM traversal API (NodeIterator & Tree Walker).

In an XML text, a pull parser allows an iterator that visits the different elements, features, & details successively. Code that utilizes this iterator will test the present thing (for example, to tell whether it's a starting tag or end tag, or message) & evaluate its attributes (nearby name, namespace, XML property estimates, content estimates, & so on), and can also move the iterator to the following. The code would thus be able to eliminate information from the report such as navigates. The recursive plunge advance will in general fit keeping data as composed neighborhood feature in the cryptogram doing the parsing, whereas SAX, for instance, commonly involve a parser to actually remain middle of the road data within a mound of section that are parent element of the factor being parsed. Pull-parsing cryptogram gets further straight to realize and continue than SAX parsing cipher.

## CONCLUSION

In this paper analysis of different XML parsers is useful in assessing the consistency and weaknesses of the objects. Different tests were conducted that contrast with gauges, speed, memory use, etc. DOM is a bad memory because it needs to keep the entire record tree in memory, rendering it unable to manage extremely large records. Objects in the DOM tree can be tackled & exploited using entity methods. DOM API has taken less time as compared to DOM4J and JDOM for files of size 5KB. But as the size of the XML file increases, every tree based API (DOM, DOM4J and JDOM) takes more or less the same amount of time. To conclude, we can say that streaming type APIs (SAX and StAX) have worked faster that tree type APIs, in both reading & script scenarios. The query language for XML, XPath, is easy to code but is the slowest of all the APIs considered.

## REFERENCES

[1].    Extensible Markup Language, http://www.w3.org/TR/REC-xml.

[2].    OASIS, http://www.oasis-open.org.

[3].    Juancarlo Anez (1999). "Java XML Parsers-A Comparative Evaluation of 7 Free Tools," Java Report Online.

[4]. Yi Chen, Susan B. Davidson, Yifeng Zheng (2010). "A bi-labeling based XPath processing system" Information Systems 35, 2010, doi: 10.1016/j.is.2009.05. ACM

[5]. W3C Website. Extensible Markup Language (XML) 1.0 (Fourth Edition). Online:http://www.w3.org/TR/REC-xml/, Accessed on: 30/10/2006

[6]. Barbara Catania and Anna Maddalena, Athena Vakali (2005). "XML Document Indexes: A Classification " SEPTEMBER OCTOBER 2005 Published by the IEEE Computer Society 1089-7801/05/$20.00 © 2005 IEEE IEEE INTERNET COMPUTING

[7]. Mikael Fernandus Simalango (2008). "XML Query Processing and Query Languages: A Survey" Property of Amikelive.com – Technical Paper Series 25/10/2008

[8]. Byron Choi, Mary Fernandez, Jerome SimeonThe XQuery Formal Semantics: A Foundation for Implementation and Optimization May 31, 2002

[9]. W3C Consortium, http://www.w3.org, 2006

[10]. W3C Consortium, XML Path Language (XPath) 2.0, http://www.w3.org/TR/xpath20/, 2006.

[11]. W3C Consortium, XQuery 1.0: An XML Query Language, http://www.w3.org/TR/xquery/, 2006

[12]. D.D. Chamberlin (2002). "XQuery: An XML Query Language," IBM Systems J., Vol. 41, No. 4.

[13]. H. Jagadish, S. Al-Khalifa, A. Chapman, L. Lakshmanan, A. Nierman, S. Paparizos, J. Patel, D. Srivastava, N. Wiwatwattana, Y. Wu, and C. Yu. (2002). TIMBER: A Native XML Database. The VLDB Journal- Volume 11, pages 274-291.

[14]. Byron Choi, Mary Fernandez, Jerome Simeon (2002). The XQuery Formal Semantics: A Foundation for Implementation and Optimization.

[15]. Z. Chen, H.V. Jagadish, L.V.S. Lakshmanan, and S. Paparizos (2003). "From Tree Patterns to Generalized Tree Patterns: On Efficient Evaluation of XQuery," Proc. 29th Int'l Conf. Very Large Data Bases (VLDB '03).

[16]. S. Paparizos, Y. Wu, L.V.S. Lakshmanan, and H.V. Jagadish (2004). "Tree Logical Classes for Efficient Evaluation of XQuery," Proc. 23rd ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '04).

[17]. Mohammed Al-Badawi, Dr. Siobhán North, Dr. Barry Eaglestone (2007). Research memorandum Indexing XML Databases: Classifications, Problems Identification and a New Approach, 15th November, 2007 in The University of Sheffield Department of Computer Science.

[18]. J.-K. Min, C.-H. Lee, and C.-W. Chung (2008). "XTRON: An XML data management system using relational database," Information and Software Technology, vol. 50, pp. 462-479.

[19]. Jong P. Yoon, Vijay Raghavan, Venu Chakilam, Larry Kerschberg BitCube (2001). A Three-Dimensional Bitmap Indexing for XML Documents Jong Yoon, et. al., BitCube, Journal of Intelligent Information Systems, Vol. 17.

[20]. S. Al-Khalifa, H.V. Jagadish, N. Koudas, J.M. Patel, D. Srivastava, and Y. Wu (2002). "Structural Joins: A Primitive for Efficient XML Query Pattern Matching," Proc. 18th IEEE Int'l Conf. Data Eng. (ICDE '02).

[21]. N. Bruno, N. Koudas, and D. Srivastava (2002). "Holistic Twig Joins: Optimal XML Pattern Matching," Proc. 21st ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '02), 2002 Conf. Management of Data (SIGMOD '02).

**Corresponding Author**

**Swati Gupta***

Research Scholar, University of Technology, Jaipur, Rajasthan

**Swati Gupta[1]* Dr. Ashish Chourasia[2]**