

XML: Data and Document Processing

Some Basic XML Concepts

David Edmond
School of Information Systems
Queensland University of Technology

February 20, 2007

Introduction

This week's topics

- ▶ Why XML?
- ▶ Elements and attributes
- ▶ A range of examples of XML documents
- ▶ The original design goals for XML
- ▶ The XML family
- ▶ Some of the rules for *producing* an XML document

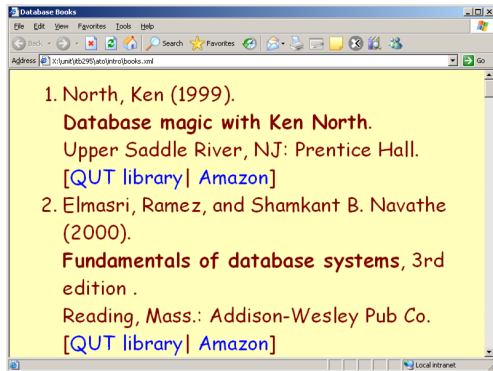
✓ <http://www.w3.org/TR/REC-xml>

What does this have to do with XML?

genuit genuit: re autem dñi filiorum numquam.
cum proprietate et significatione cognoscimus.
non enim ait filiorum meorum genuit et exaltavit
sed in octavo tantum filiorum genuit et exaltavit.
nisi forte in eo quod ait filius primogenitus
meus israel, quinquam hoc quod primogenitus
meus dixit ad deum trahendum filio proprietate,

HTML = Data + Markup

Look at the HTML document shown below.



The links

Each book has two distinct hypertext links. For the first book, these are:

`http://libcat.qut.edu.au/search/i?SEARCH=0136471994`>

`http://www.amazon.com/exec/obidos/ASIN/0136471994/`>

- ▶ One links to the QUT library catalogue.
- ▶ The other links to Amazon.

The underlying HTML

```
1. <HTML>
2. <TITLE>Database Books</TITLE>
3. <BODY STYLE="margin-left: 40px;
4.     margin-right: 20px;
5.     background-color: #FFFFBB;
6.     color:maroon; font-family:'Comic Sans MS';
7.     font-size:12pt;
8.     padding:0px 6px">
9. <OL>
10. <LI>North, Ken (1999).<BR/>
11.     <B>Database magic with Ken North</B>.<BR/>
12.     Upper Saddle River, NJ: Prentice Hall.<BR/>
13.     [<A HREF="....search/i?SEARCH=0136471994">QUT library</A>|
14.     <A HREF="..../ASIN/0136471994/">Amazon</A>]
15. <BR/></LI>
16. <LI> Elmasri, Ramez, and Shamkant B. Navathe (2000 ).<BR/>
17.     <B>Fundamentals of database systems</B>, 3rd edition.<BR/>
18.     Reading, Mass.: Addison-Wesley Pub Co.<BR/>
19.     [<A HREF="....search/i?SEARCH=0201542633">QUT library</A>|
20.     <A HREF="..../ASIN/0201542633/">Amazon</A>]
21. <BR/></LI>
22. .
23. .
24. .
25. </OL>
26. </BODY>
27. </HTML>
```

Only books

```
1. <?xml version="1.0"?>
2. <?xml-stylesheet href="BIBstyle3.xsl" type="text/xsl"?>
3. <bibliography>
4.   <book>
5.     <author>North, Ken</author>
6.     <title>Database magic with Ken North</title>
7.     <address>Upper Saddle River, NJ</address>
8.     <publisher>Prentice Hall</publisher>
9.     <year>1999</year>
10.    <isbn>0136471994</isbn>
11.  </book>
12.
13.  <book>
14.    <author> Elmasri, Ramez, and Shamkant B. Navathe</author>
15.    <title>Fundamentals of database systems</title>
16.    <address>Reading, Mass.</address>
17.    <publisher>Addison-Wesley Pub Co</publisher>
18.    <year>2000</year>
19.    <edition>3rd</edition>
20.    <isbn>0201542633</isbn>
21.  </book>
22.
23. <!-- Lots more <book/> elements here. -->
24.
25. </bibliography>
```


Elements and attributes

Example 1: Tags, tags, tags, ...

Consider the following snippet of information from the QUT phone book:

1.	Edgar	Miss Pam	Optometry	KG B501	35695
2.	Edmond	Dr David	Information Systems	GP S842	32240
3.	Edmonds	Dr Ian	Physical Sciences	GP M206	32584

We could write this as an XML document:

```
1. <Phonebook>
2.   <Entry>
3.     <LastName>Edgar</LastName>
4.     <Title>Miss</Title>
5.     <FirstName>Pam</FirstName>
6.     <School>Optometry</School>
7.     <Campus>GP</Campus>
8.     <Room>B501</Room>
9.     <Extension>35695</Extension>
10.  </Entry>
11. <!-- Lots more <Entry/> elements here. -->
12. </Phonebook>
```

XML elements

- ▶ An element is *not* the same as a tag!
- ▶ XML is case-sensitive. (`<LastName>` \neq `<lastname>`)
- ▶ The names XML, xML ... xml are reserved.
- ▶ A name cannot start with a digit:
 - ▶ illegal: `<911/>`
 - ▶ legal: `<_911/>`
- ▶ A name may only contain letters, digits, `'-'`, `'_'`, `'.'`
 - ▶ illegal: `<Jack's Diary/>`
 - ▶ legal: `<Jacks_Diary/>`
- ▶ Elements cannot overlap:
 - ▶ illegal: `<OutTag><InTag><OutTag/><InTag>`
 - ▶ legal: `<OutTag><InTag><InTag/><OutTag>`

Empty elements

- ▶ An *empty* element is one that has no “content”. There are two ways of writing them:
 - ▶ `<AnEmptyElement />`.
 - ▶ `<AnEmptyElement></AnEmptyElement>`.
- ▶ An empty element can be useful. For example, the line-break element in XHTML may be written `
` or `
</br>`.
- ▶ Empty elements may still carry information by means of attributes. For example, the image element in XHTML may be written ``.

Example 2: Attributes

We could make a slight variation to each phonebook entry in the following way:

```
1. <Entry>
2.   <LastName Title="Miss">Edgar</LastName>
3.   <FirstName>Pam</FirstName>
4.   <School>Optometry</School>
5.   <Campus>GP</Campus>
6.   <Room>B501</Room>
7.   <Extension>35695</Extension>
8. </Entry>
```

- ▶ The `<Title/>` element has been turned into an attribute of the `<LastName/>` element.
- ▶ Why might we make such a change?

Example 3: Yet more attributes

We could make even more extensive use of attributes, by expressing each entry in the following way:

```
1. <Entry Extension="35695">
2.   <Name Title="Miss"
3.     LastName="Edgar"
4.     FirstName="Pam"/>
5.   <School>Optometry</School>
6.   <Campus>GP</Campus>
7.   <Room Building="B">501</Room>
8. </Entry>
```

XML attributes

- ▶ The names of XML attributes follow the same conventions as elements.
- ▶ Their values must be demarcated by quotes (") or apostrophes (').
 - ▶ illegal: `<LastName Title=Mr>`
 - ▶ legal: `<LastName Title="Mr">`
 - ▶ illegal: `<if test="Age<25">`
 - ▶ illegal: `<Album Category='Rock'n'Roll'>`
 - ▶ legal: `<Album Category="Rock'n'Roll">`
 - ▶ legal: `<Album Category='Rock"n"Roll'>`
- ▶ What if you need to use both a quotation mark *and* an apostrophe in an attribute?
- ▶ When should you use attributes to convey information?

Some other XML concepts

XML comments

- ▶ Comments begin `<!--` and end `-->`.
- ▶ For example `<!-- Created \today\ -->`
- ▶ Two consecutive dashes (--) should *never* appear within a comment.

The XML prolog

At the very beginning of any XML document there may be a “prolog”, that is, some general preliminary information about the document. This could include:

- ▶ The way the document is physically encoded.
- ▶ A pointer to a stylesheet to be used to transform the document.
- ▶ A definition of how the document should be validated before any processing occurs.
- ▶ Comments on the document.

Some exercises

Example 4: Duncan

Suppose we want to send a letter to Duncan:

1. Duncan Hunter
2. 45 Somerset Street
3. Clayfield

Encode the above as XML.

1. <envelope>
2. <name>
3. <given>.....</given>
4. <family>.....</family>
5. </name>
6. <street>
7. <number>.....</number>
8. <name>.....</name>
9. </street>
10. <suburb>.....</suburb>
11. </envelope>

Example 4: Duncan

Suppose we want to send a letter to Duncan:

1. Duncan Hunter
2. 45 Somerset Street
3. Clayfield

Encode the above as XML.

1. <envelope>
2. <name>
3. <given>Duncan</given>
4. <family>Hunter</family>
5. </name>
6. <street>
7. <number>45</number>
8. <name>Somerset</name>
9. </street>
10. <suburb>Clayfield</suburb>
11. </envelope>

Example 5: It's simply not cricket!

Sports results are suitable targets for XML. Consider this cricket innings:

1. Team: Australia
2. Haydon b Jones 58
3. Sanger c Kent, b Jones 87
4. Wonting not out 55
5. Lemon not out 13

Suppose we encoded Haydon's performance as follows:

1. <bat>
2. <name>Haydon</name>
3. <runs>58</runs>
4. <out style="bowled">
5. <bowler>Jones</bowler>
6. </out>
7. </bat>

How might we encode Sanger?

1. <bat>
2. <name>Sanger</name>
3. <runs>87</runs>
4. <out style=".....">
5. <bowler>Jones</bowler>
6. <.....>Kent</.....>
7. </out>
8. </bat>

Example 5: It's simply not cricket!

Sports results are suitable targets for XML. Consider this cricket innings:

1. Team: Australia
2. Haydon b Jones 58
3. Sanger c Kent, b Jones 87
4. Wonting not out 55
5. Lemon not out 13

Suppose we encoded Haydon's performance as follows:

1. <bat>
2. <name>Haydon</name>
3. <runs>58</runs>
4. <out style="bowled">
5. <bowler>Jones</bowler>
6. </out>
7. </bat>

How might we encode Sanger?

1. <bat>
2. <name>Sanger</name>
3. <runs>87</runs>
4. <out style="caught">
5. <bowler>Jones</bowler>
6. <catcher>Kent</catcher>
7. </out>
8. </bat>

Example 6: What's XML got to do with SQL?

How might write the following as XML?

1. Select CarNr, Make, Cost
2. From Cars
3. Order by Cost

1. <sql>
2. <select order=".....">
3. <col>CarNr</col>
4. <col>Make</col>
5. <col>Cost</col>
6. </select>
7. <.....>
8. <table>Cars</table>
9. </.....>
10. </sql>

Example 6: What's XML got to do with SQL?

How might write the following as XML?

1. Select CarNr, Make, Cost
2. From Cars
3. Order by Cost

1. <sql>
2. <select order="cost">
3. <col>CarNr</col>
4. <col>Make</col>
5. <col>Cost</col>
6. </select>
7. <from>
8. <table>Cars</table>
9. </from>
10. </sql>

Example 7:

Consider the following simple program:

```
1. while (x.test()) {  
2.     y.modify();  
3.     x.meddle();  
4. }
```

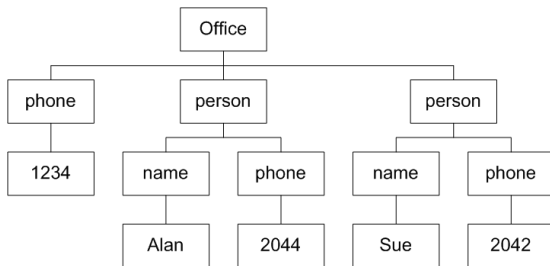
We might encode this, XML style, as follows:

```
1. <while>  
2.   <condition>  
3.     <var>x</var><method>test</method>  
4.   </condition>  
5.   <body>  
6.     <statement>  
7.       <var>y</var><method>modify</method>  
8.     </statement>  
9.     <statement>  
10.      <var>x</var><method>meddle</method>  
11.    </statement>  
12.  </body>  
13. </while>
```

NB: Most important!

An XML document is a tree

```
1. <office>
2.   <phone>1235</phone>
3.   <person>
4.     <name>Alan</name>
5.     <phone>2044</phone>
6.   </person>
7.   <person>
8.     <name>Sue</name>
9.     <phone>2043</phone>
10.  </person>
11. </office>
```



XML: at the beginning

XML design goals

- ▶ XML shall be straightforwardly usable over the Internet.
 - ▶ XML shall support a wide variety of applications.
 - ▶ XML shall be compatible with SGML.
 - ▶ It shall be easy to write programs which process XML documents.
 - ▶ The number of optional features in XML is to be kept to the absolute minimum, ideally zero.
 - ▶ XML documents should be human-legible and reasonably clear.
 - ▶ The XML design should be prepared quickly.
 - ▶ The design of XML shall be formal and concise.
 - ▶ XML documents shall be easy to create.
 - ▶ Terseness in XML markup is of minimal importance.
- ✓ <http://www.w3.org/TR/REC-xml>
- ✓ <http://www.xml.com/axml/testaxml.htm>

The XML Family

- ▶ XML: for generating languages that describe information.
- ▶ DOM: a programming interface for accessing and updating documents.
- ▶ XML schema: a language for specifying the structure and content of documents.
- ▶ XSLT: a language for transforming documents.
- ▶ XPath: a query language for navigating XML documents.
- ▶ XPointer: for identifying fragments of a document.
- ▶ XLink: generalises the concept of a hypertext link.
- ▶ XInclude: for merging documents.
- ▶ XQuery: a language for making queries across documents.
- ▶ RDF: a language for describing resources.

The XML specification: a start

The XML Specification

- ▶ The original specification was first released in 1998.
 - ▶ It specified both the physical and logical makeup of an XML document.
 - ▶ A 2nd edition was released in 2000, a 3rd edition in 2004, and a 4th in 2006.
 - ▶ None of these later editions significantly changed the original.
- ✓ `http://www.w3.org/TR/REC-xml`

Production rules

The core of the specification is a series of production rules. Each “production” has the form:

1. `symbol ::= expression`

Each production is given a number.

Selected productions

Some of the productions appearing in the XML specification are:

1. `document ::= prolog element Misc*`
2. `prolog ::= XMLDecl? Misc* (doctypeddecl Misc*)?`
3. `element ::= EmptyElemTag | STag content ETag`
4. `Misc ::= Comment | PI | S`
5. `Comment ::= '<!--' ((Char - '-') | ('-' (Char - '-')))*`
6. `'-->'`
7. `S ::= (#x20 | #x9 | #xD | #xA)+`

The production that defines an XML document suggests that there can only be one element. But we know that is not true. So ...

EBNF: Core

The rules are expressed in a version of the notation called *Extended Backus-Naur Form*. And this notation allows the following forms of expression:

<code>A?</code>	matches A or nothing; optional A.
<code>A B</code>	matches A followed by B.
<code>A B</code>	matches A or B but not both.
<code>A+</code>	matches one or more occurrences of A.
<code>A*</code>	matches zero or more occurrences of A.
<code>#xN</code>	matches the ISO/IEC 10646 character corresponding to the hexadecimal number N.
<code>(expression)</code>	expression is treated as a unit and may be combined as described in this list.

EBNF: Literals, ranges, enumerations, etc

<code>'string'</code>	matches the enclosed string exactly.
<code>"string"</code>	matches the enclosed string exactly.
<code>A - B</code>	matches any string that matches A but does not match B.
<code>[a-zA-Z]</code>	matches any Char with a value in the range(s) indicated (inclusive).
<code>[abc]</code>	matches any Char with a value among the characters enumerated. Enumerations and ranges can be mixed in one set of brackets.
<code>[^a-z]</code>	matches any Char with a value outside the range indicated.
<code>[^abc]</code>	matches any Char with a value not among the characters given. Enumerations and ranges of forbidden values can be mixed in one set of brackets.

Names and tokens

The XML specification contains the following set of production rules:

1. NameChar ::= Letter | Digit | '.' | '-' | '_' | ':'
2. | CombiningChar | Extender
3. Name ::= (Letter | '-' | ':') (NameChar)*
4. Names ::= Name (S Name)*
5. Nmtoken ::= (NameChar)+
6. Nmtokens ::= Nmtoken (S Nmtoken)*

Example 8: Names and tokens

Classify each of the following as being one of the following: Name, Names, Nmtoken or Nmtokens.

- | | |
|---------------------|-------|
| 1. xsl:value-of | |
| 2. 123ABC | |
| 3. pardon_my_french | |
| 4. _excuse_me | |
| 5. -dash-dash-dash | |
| 6. dash dash dash | |
| 7. 15 August | |
| 8. _:_1 | |

Example 8: Names and tokens

Classify each of the following as being one of the following: Name, Names, Nmtoken or Nmtokens.

- | | |
|----------------------------------|----------|
| 1. <code>xsl:value-of</code> | Name |
| 2. <code>123ABC</code> | Nmtoken |
| 3. <code>pardon_my_french</code> | Name |
| 4. <code>_excuse_me</code> | Name |
| 5. <code>-dash-dash-dash</code> | Nmtoken |
| 6. <code>dash dash dash</code> | Names |
| 7. <code>15 August</code> | Nmtokens |
| 8. <code>_: _1</code> | Name |

Conclusions

About this week's topics

- ▶ XML was designed to allow the data underlying some message or document to be expressed.
- ▶ Display or formatting information was to be removed.
- ▶ Formatting is the business of the receiving program.
- ▶ The most fundamental component of an XML document are its elements and attributes.
- ▶ Every XML document is tree-structured (hierarchical).
- ▶ The (meta-)language is expressive enough to allow a limitless range of examples.
- ▶ The rules for constructing an XML document are written in the form of a set of production rules.

✓ <http://www.w3.org/TR/REC-xml>

Next week's topics

- ▶ *XPath*: a language for querying XML documents.
- ▶ *Namespaces*: how to mix different vocabularies in the one document.