

Moving on...

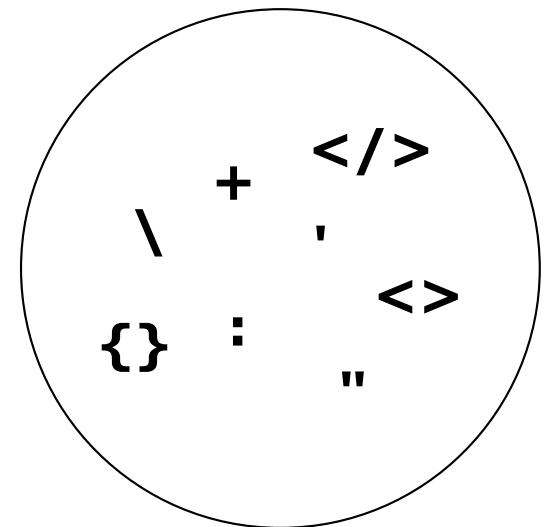
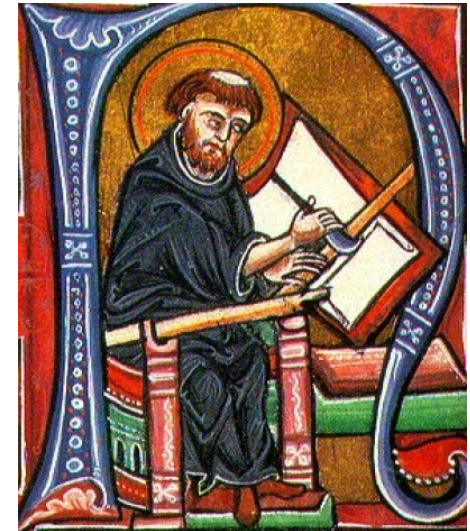
- UML Classes => XML documents
- UML Use Cases => Xquery/Xpath queries



<LectureTopic>
Fundamentals of e**X**tensible**M**arkup**L**anguage
</LectureTopic>

What is Markup

- Sequence of characters within a text or word processing file to define
 - Print properties
 - Display properties
 - Document's logical structure
- Markup indicators are often called "tags"
 - Examples
 - RTF
 - EDIFACT
 - XML



Mark Up: RTF

```
\li0\ri0\sb240\sa60\keepn\widctlpar\aspalpha\aspnum\faauto\outlinelevel2 \a
djustright \rin0\lin0\itap0
\b\f1\fs26\lang2057 \langfe1033 \cgrid\langnp2057 \langfenp1033
{\lang6153 \langfe1033 \langnp6153 Entity Relationship Diagram
\par }\pard\plain \s1\ql
\li0\ri0\sb240\sa60\keepn\widctlpar\aspalpha\aspnum\faauto\outlinelevel0 \a
djustright \rin0\lin0\itap0 \cbpat17
\b\f1\fs24\lang2057 \langfe1033 \kerning32 \cgrid\langnp2057 \langfenp1033
{\lang6153 \langfe1033 \langnp6153 Entity Type
\par }\pard\plain \ql
\li0\ri0\widctlpar\aspalpha\aspnum\faauto\adjustright \rin0\lin0\itap0
\fs24\lang2057 \langfe1033 \cgrid\langnp2057 \langfenp1033
{\b\fs20\ul\lang6153 \langfe1033 \langnp6153
Def.:}{ \b\fs20\lang6153 \langfe1033 \langnp6153 }{
\fs20\lang6153 \langfe1033 \langnp6153 An object or co ncept that is
identified by the enterprise as having an independent existence.
\par }\pard\plain \s1\ql
\li0\ri0\sb240\sa60\keepn\widctlpar\aspalpha\aspnum\faauto\outlinelevel0 \a
djustright \rin0\lin0\itap0 \cbpat17
\b\f1\fs24\lang2057 \langfe1033 \kerning32 \cgrid\langnp2057 \langfenp1033
{\lang6153 \langfe1033 \langnp6153 Entity
\par }\pard\plain \ql
```



Mark Up: EDIFACT

```
' 'ED2' 'OPENET:1111111:OVT':003705655815:OVT'ABC1234567'0'TYP:ORDERS'N
RQ:1' '
UNA:+. ?
'UNB+UNOC:2+003705655815:30+1111111:30+980729:2233+4++ORDERS911+++KKK
KATE+1'UNH+
1+ORDERS:001:911:UN:FI0030'BGM+640+1234567'DTM+4:19981201:102'DTM+2:199
90101:102'DTM+2:9901:616'RFF+BC:123'RFF+VN:123456'NAD+BY+003705655815:1
00'
NAD+SE+11111111:92'NAD+PL+53432::92++KAUPPA:KAUPUNKI+KATU
9+KAUPUNKI++00007'NAD+CN+-::ZZ++TERMINAALI+OVI 42+TOINEN
KAUPUNKI++00069'UNS+D'LIN+1++23442423234
:EN'PIA+5+3244:MF'PIA+5+2341234324:ZBU'PIA+5+234243:ZCG'IMD+F+8+-
::91:KUKKAPUR
KKI:SAVI'QTY+21:8:KPL'FTX+AAA+++T.HARMAA:V[RI'FTX+AAA+++10:KOKO'PRI+NTP
:7.23:+
RP:7.32:PE'TAX+7+VAT+++::22.00'LIN+2++543434554345:EN'PIA+5+535:MF'PIA
+5+45:
PCE`UNT+38+2'UNZ+2+4'
' 'EOF' ' '9'
```



Mark Up: XML

```
<fragment>
  <section>
    <title>Introduction</title>
    <para>Since the emergence of <acronym refid="xml">XML</acronym> in
    early 1998 and it's subsequent adoption across diverse application
    domains, one of the key benefits it enabled was the separation of
    content and presentation <bibref refloc="Bos97"/>. <acronym
    refid="xml">XML</acronym> borrowed this model (along with other
    important concepts) from the <acronym.grp><acronym
    refid="sgml">SGML</acronym><expansion id="sgml">Standard
    Generalised Markup Language</expansion></acronym.grp>. An
    <acronym refid="sgml">SGML</acronym> document consists of
    logically structured content and uses a separate file (style
    sheet) to specify how the content should be formatted for
    [...]
    <figure id="img1">
      <title>ePublishing Components</title>
      <graphic href="02-04-03-fig01.jpg" width="321" height="214"/>
    </figure>
  </section>
</fragment>
```



What is SGML?

- Standard Generalised Mark-Up Language
- ISO standard since 1986
- Meta-language for defining document mark-up vocabularies
- Uses logical mark-up (structure, content) instead physical (how document looks on printed page)
- Platform-, system-, vendor- and version-independent documents
- Very powerful, but contains a number of complex features

```
<!DOCTYPE anthology [  
  <!ELEMENT anthology - - (poem+)      >  
  <!ELEMENT poem - - (title?, stanza+) >  
  <!ELEMENT title - O (#PCDATA)      >  
  <!ELEMENT stanza - O (line+)      >  
  <!ELEMENT line O O (#PCDATA)      >  
]]>  
  
<anthology>  
  <poem>  
    <title> The SICK ROSE  
    <stanza>  
      <line>O Rose thou art sick.</line>  
      <line>The invisible worm,</line>  
      [...]   
    </stanza>  
    <stanza>  
      <line>Has found out thy bed</line>  
      <line>Of crimson joy:</line>  
      [...]   
    </stanza>  
  </poem>  
</anthology>
```



What is HTML?

- HTML, the *de facto* standard for publishing Web content, is an SGML vocabulary
- Supporting full SGML on the Web was too difficult so HTML made some simplifications
 - not extensible
 - limited structure
 - not content oriented
 - cannot be validated
- HTML is a simple language to understand and use
- Most of the content available on the Web has been created with HTML

```
<html>
  <head>
    <title>The SICK ROSE</title>
  </head>

  <body>
    <h1>The SICK ROSE</h1>

    <p>
      O Rose thou art sick.<br />
      The invisible worm,<br />
      [...]
    </p>

    <p>
      Has found out thy bed<br />
      Of crimson joy:<br />
      [...]
    </p>
  </body>
</html>
```



What is XML?

- eXtensible Markup Language
- XML is a simplified subset of SGML
- Can also be used to define document markup **vocabularies** (e.g. XHTML)
 - These can have a strictly defined structure (DTD)
- Retains the powerful features of SGML (extensibility, structure, validation)
- Ignores the complex features of SGML and is therefore easier to use and implement
- XML documents look similar to HTML documents
- Separates structure and presentation (like SGML)



XML Example

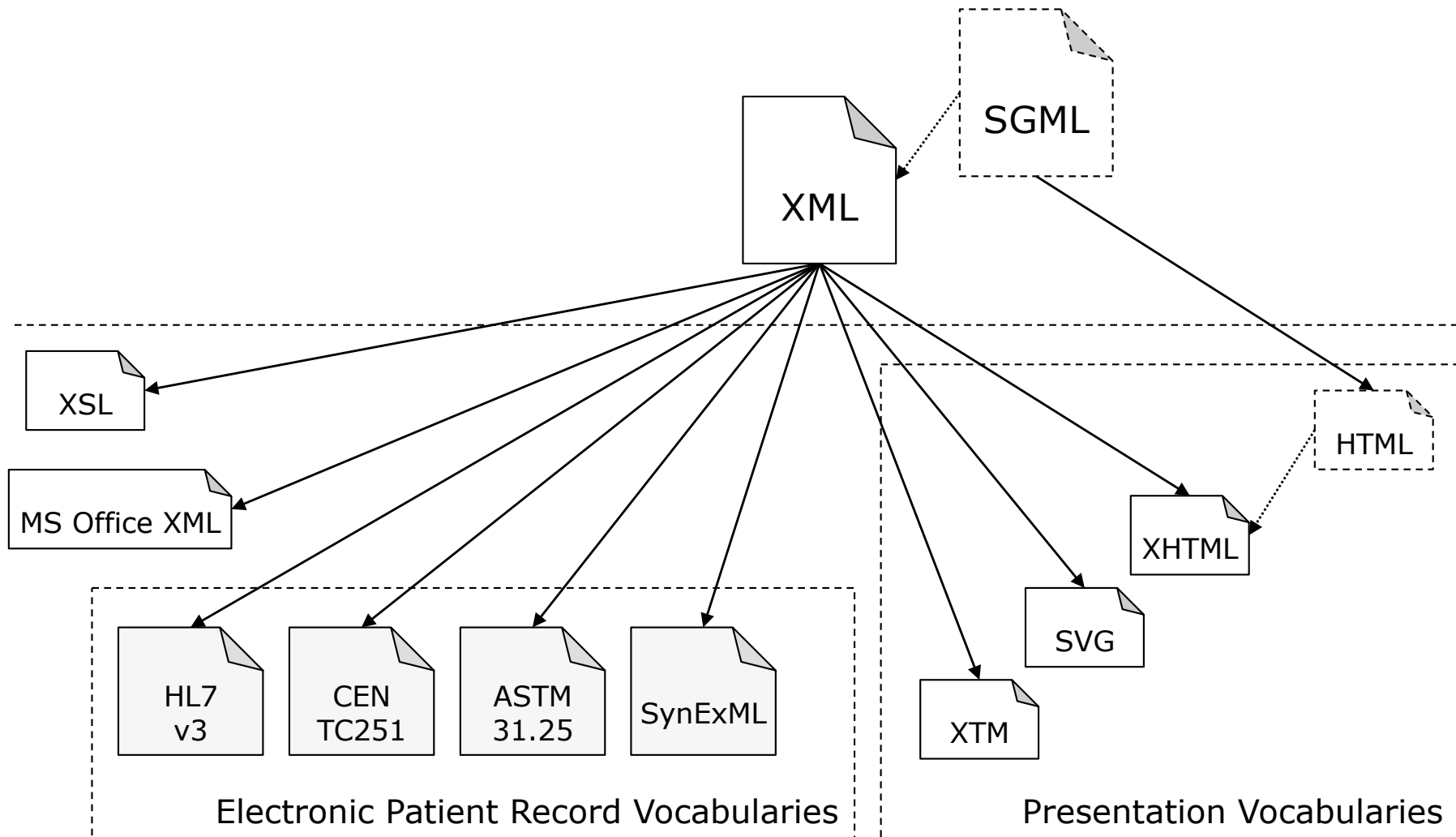
```
<?xml version='1.0' encoding='ISO-8859-1' standalone='yes' ?>
<doc type="book" isbn="1-56592-796-9" xml:lang="en">
  <title>A Guide to XML</title>
  <author>Norman Walsh</author>
  <chapter>
    <title>What Do XML Documents Look Like?</title>
    <paragraph>If you are [...]</paragraph>
    <ol>
      <item>
        <paragraph>The document begins [...]</paragraph>
      </item>
      <item>
        <paragraph>Empty elements have [...]</paragraph>
        <paragraph>In a very [...]</paragraph>
      </item>
    </ol>
    <section>[...]</section>
    [...]
  </chapter>
  <chapter>[...]</chapter>
</doc>
```



XML vocabularies you have come across? – student pov



Meta Language vs. Vocabulary

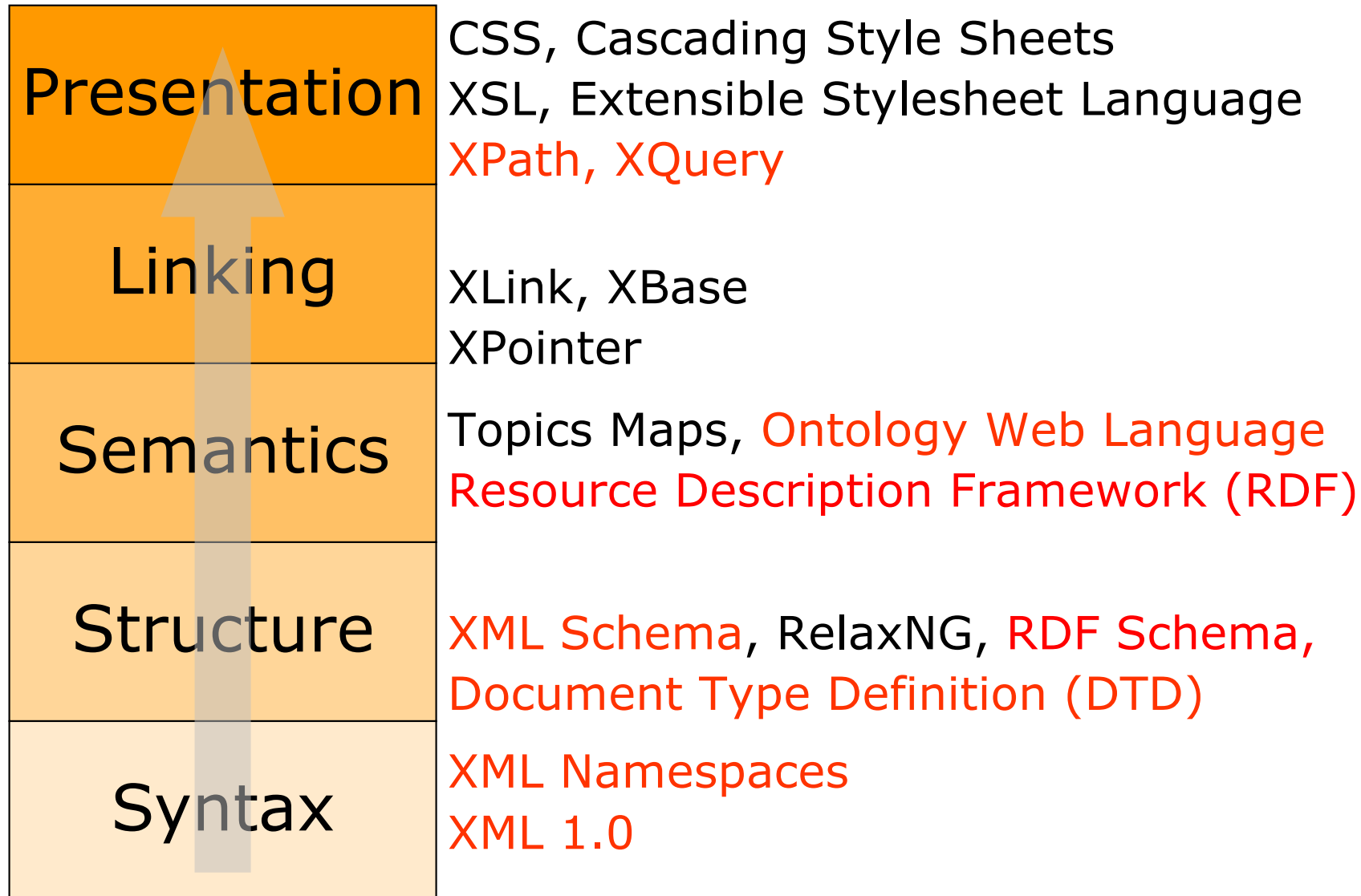


Why is the emergence of XML an important development?

- XML is a tool for defining vocabularies
 - XML vocabularies are easy to read
 - XML is self describing
 - Parse tree embedded in document
 - Grammar for language referenced via DTD/Schema
- XML vocabularies are easy for computers to process, exchange and display
 - XML tools are ubiquitous, free and conform to established standards
 - Natural affinity with Object serialization
 - Data source neutral



XML technologies



XML 1.0

- The XML 1.0 specification describes the syntax for XML documents (elements and attributes) and DTDs
- An XML document is a hierarchical data structure using self-definable tags
 - e.g. `<doc><author>[...]</author></doc>`
- There are many other technologies related to XML
- XML is

A simple common layer for tree structures in a character stream.



Design goals of XML 1.0 specification

1. XML shall be straightforwardly usable over the Internet.
2. XML shall support a wide variety of applications.
3. XML shall be compatible with SGML.
4. It shall be easy to write programs which process XML documents.
5. The number of optional features in XML is to be kept to the absolute minimum, ideally zero.
6. XML documents should be human-legible and reasonably clear.
7. The XML design should be prepared quickly.
8. The design of XML shall be formal and concise.
9. XML documents shall be easy to create.
10. Terseness in XML markup is of minimal importance.



Physical Parts of XML documents

Physical parts of XML documents

- XML Declaration
- Elements
- Attributes
- Document Type Declaration
- Entities
- Processing Instructions
- Comments
- Character Data Sections

- XML Namespaces



XML Declaration

- Placed at the start of an XML document
- Informs XML software of
 - the version of XML the document conforms to
 - the character encoding scheme used in the document
 - whether or not a set of external declarations affect the interpretation of this document

```
<?xml version="1.0" ?>

<?xml
  version="1.0" encoding="UTF-
8" ?>

<?xml
  version="1.0"
  encoding="UTF-8"
  standalone="yes" ?>

<!DOCTYPE person [
  <!ELEMENT person (name, adult,
nationality)>
]>
<person> and so on </person>

<?xml version="1.0">
<!DOCTYPE person SYSTEM
  'person.dtd'>
<person> and so on </person>
```



Elements

- Define logical structure and sections of XML documents
- Four different content types:
 - Data content
 - Element content
 - Mixed content
 - Empty.
- Each element must be completely enclosed by another element, except for the root
- **Important Note**
 - Any XML name must start with a letter, underscore but after that can include also digits, fullstops, hyphens.
Don't start with colon due to namespaces
Don't include spaces

```
<?xml version="1.0" ?>

<doc>
  <title>Java Gently</title>
  <author>Judy Bishop</author>
  <publisher name='HH' />
  <chapter>
    <thetext> this is <b>bold</b>
    bold </bold> text </thetext>
  </chapter>
</doc>
```



Attributes

- Provides additional information about an element
- Attributes are contained within the **start-tag**
- Consists of a name and associated value separated by an equals sign
- **The attribute value must always be enclosed by quotes**
- The order of attributes is insignificant

```
<?xml version="1.0" ?>

<doc type="book"
      isbn="0-201-71050-1">

  <title>Java Gently</title>
  <author>Judy Bishop</author>
  <chapter>
    <paragraph type="abstract">
      In this book ...
    </paragraph>
  </chapter>

</doc>
```



ELEMENT vs. ATTRIBUTE

- Lexically little difference,
- application specific,
- no hard/fast rules available.

ELEMENT

- Constituent data,
- Used for content,
- White space can be ignored or preserved
- Nesting allowed (child elements),
- Convenient for large values, or binary entities.

ATTRIBUTE

- Inherent data,
- Used for meta-data,
- No further nesting possible (atomic data),
- Default values,
- Minimal datatypes



Entities

- Storage units for repeated text
 - Defined in a DTD
- Character entities are used to insert characters that cannot be typed directly
- XML contains a number of 'built-in' entities
 - "
 - '
 - <
 - >
 - &

```
<math>
  5 &lt; 6 and 6 &gt; 5
</math>

<copyright>
  &copyright-notice;
</copyright>

<bullet>
  XML contains a number
  of &apos;built-in&apos;
  entities
  <list>
    <item>&quot;</item>
    <item>&apos;</item>
    <item>&lt;</item>
    <item>&gt;</item>
    <item>&amp;</item>
  </list>
</bullet>
```



Character Data Sections

- Data which is to be parsed is called PCDATA
- An XML parser will not treat the contents of a CDATA section as markup
 - Used to simplify mark-up by escaping a selection of text
- Entity references are not resolved
- Useful for including source code in XML

```
<![CDATA[
```

```
You don't need to escape  
special characters in CDATA  
sections, such as <, >, &, ,  
' and ".
```

```
]]>
```

```
<![CDATA[<<< STOP now >>>]]>
```

```
<![CDATA[<?xml version='1.0'?>
```

```
<person>
```

```
<name>Mike</name>
```

```
<age>24</age>
```

```
</person>]]>
```



Processing Instructions

- Pass additional information to application (e.g. parser)
- Application-specific instructions
- Consists of a PI Target and PI Value
- Processed by applications that recognise the PI Target

```
<?xml version="1.0" ?>

<?xml-stylesheet
  type='text/css'
  href='style.css'?>

<?xml-stylesheet
  type='text/xsl'
  href='style.xsl'?>

<?myapp filename='test.txt'?>
```



Comments

- Used to comment XML documents
- Not considered to be part of an XML document
- An XML parser is not required to pass comments to higher-level applications

```
<!-- one-line comment -->
```

```
<!--  
This  
is a  
multi-line comment  
-->
```



Well formed XML

- XML Declaration required
- At least one element
 - Exactly one root element
- Empty elements are written in one of two ways:
 - Closing tag (e.g. "
</br>")
 - Special start tag (e.g. "
")
- For non-empty elements, closing tags are required
- Attribute values must always be quoted
- Start tag must match closing tag (name & case)
- Correct nesting of elements
 - Example **incorrect nesting** and **incorrect case**

```
<full_name>  
<first_name>  
John </Full_name>  
</first_name>
```



Exercise

- Spot the 6 deliberate errors in the XML opposite

```
<?xml version="1.0"!>
<!DOCTYPE catalog SYSTEM "books.dtd">
<catalog>
  <book id='bk101' type='softback'>
    <author>Gambardella, Matthew</author>
    <title>XML Developer's Guide</title>
    <genre>Computer</genre>
    <price>44.95</price>
    <publish_date>2000-10-01</publish-date>
    <description>An in-depth look at creating applications with XML.
  </book></description>
  <book id='bk102' >
    <author nationality=irish>Jenkins, Fred</Author>
    <title>XML Technology Guide</title>
    <price>50.00</price>
    <publish date>2000-10-01</publish date>
    <description>An in-depth look at using XML
technologies.</description>
    <stocked_by>Easons</stocked_by>
    <stocked_by>Amazon</stocked_by>
  </book type='hardback'>
</catalog>
```



```

<?xml version="1.0"!>
<!DOCTYPE catalog SYSTEM "books.dtd">
<catalog>
  <book id='bk101' type='softback'>
    <author>Gambardella, Matthew</author>
    <title>XML Developer's Guide</title>
    <genre>Computer</genre>
    <price>44.95</price>
    <publish_date>2000-10-01</publish-date>
    <description>An in-depth look at creating applications with XML.
  </book></description>
  <book id='bk102' >
    <author nationality=irish>Jenkins, Fred</Author>
    <title>XML Technology Guide</title>
    <price>50.00</price>
    <publish date>2000-10-01</publish date>
    <description>An in-depth look at using XML
technologies.</description>
    <stocked_by>Easons</stocked_by>
    <stocked_by>Amazon</stocked_by>
  </book type='hardback'>
</catalog>

```



Exercise and Hand up

- Write down a “well formed” XML snippet, using elements and/or attributes, describing:
 - Your name (distinguishing first, middle, surname)
 - Student ID
 - Favourite music groups
 - County
 - Expected date of graduation

XML Declaration required

Exactly one root element

Empty elements are written in one of two ways:

Closing tag or Special start tag

For non-empty elements,

closing tags are required

Attribute values must always be quoted

Start tag must match closing tag (name & case)

Correct nesting of elements



Observations or problems typically exposed

- No XML declaration
- No Root in the document!

```
<?xml version="1.0" ?>
```

```
<student_info>  
  <student_name student-id="1234 ">  
    <firstname> Declan </firstname>  
    <surname>O' Sullivan</surname>  
  </student_name>  
  <fav_music>  
    <band> U2 </band>  
    <band> beatles </band>  
  </fav_music>  
</student_info>
```

- Quotes used in element data
 <band> "U2" </band>
- Avoidance of attributes



XML as a tree structure

```
<ASSESSMENTS>
```

```
  <STUDENT name = "Smith">
```

```
    <MARK theCourse = "CS2011">
```

```
      75 </MARK>
```

```
    <MARK theCourse = "CS2012">
```

```
      99 </MARK>
```

```
  </STUDENT> ...
```

```
  <COURSE name = "CS2011",  
  takenBy = "Smith, Jones, ... ">
```

```
    <MARK> 60 </MARK>
```

```
  </COURSE> ...
```

```
</ASSESSMENTS>
```

Describes
mark for
individual
student

Describes
average
mark for
course

Element
ASSESSMENTS

element
STUDENT

element
COURSE

ETC

attribute
name
4BA1

attribute
name
Smith

element
MARK

attribute
theCourse
CS2011

text
75



BaseX Introduction

- A light-weight, high-performance and scalable **XML Database** engine and **XPath/XQuery Processor**.
- Interactive and user-friendly **GUI frontend**
- Different programming APIs to connect to BaseX XML database
 - REST-Style Web API
 - Variety of Client APIs for different programming languages
 - See <http://docs.basex.org/wiki/Developing>
- **YOUR ACTION:** Download Core Package Java BaseX to your laptop or your U: drive or to D: drive on PC (<http://basex.org/products/download/all-downloads>)



BaseX: Database functionality

Creating Database

Open Database

Buttons for Different Visualisations of XML
that is in the Database

The screenshot shows the BaseX 7.7.1 - XQuerySample interface. The top toolbar contains buttons for creating a new database (star icon), opening a database (folder icon), and other functions. The 'Open Database' label points to the folder icon. The 'Buttons for Different Visualisations of XML that is in the Database' label points to a group of icons (tree, table, grid, etc.) which are circled in black. Below the toolbar is a command input field with a 'Command' dropdown and a 'Run' button (green play icon). The main editor area shows the XQuery `//entry` in the query input field, with a label 'Used for Filtering views' pointing to the filter icon (funnel) in the toolbar. Below the editor is a visualization area showing a tree structure of XML data. The tree has a root node 'reviews' which branches into three 'entry' nodes. Each 'entry' node has three children: 'title', 'price', and 'review'. The 'title' node has a child 'text()' with the value 'Data on the Web'. The 'price' node has a child 'text()' with the value '34.95'. The 'review' node has a child 'text()' with the value 'A very good discussion of semi-structured database systems and XML.'. Below the visualization is a 'Result' pane showing the XML output:

```
<entry>
  <title>Data on the Web</title>
  <price>34.95</price>
  <review>A very good discussion of semi-structured database
    systems and XML.</review>
</entry>
<entry>
```

Note that when XML is loaded in the database, you no longer need doc("filename") in the XQuery



Step 1: Create Database with some XML file

Creating Database



BaseX 7.7.1 - XQuerySample

3 Results

Command

Editor

```
//entry
```

OK 1 : 8

Diagram illustrating the XML structure:

```
graph TD
    reviews[reviews] --> entry1[entry]
    reviews --> entry2[entry]
    reviews --> entry3[entry]
    entry1 --> title1[title]
    entry1 --> price1[price]
    entry1 --> review1[review]
    entry2 --> title2[title]
    entry2 --> price2[price]
    entry2 --> review2[review]
    entry3 --> title3[title]
    entry3 --> price3[price]
    entry3 --> review3[review]
    title1 --> text1[text()]
    price1 --> text2[text()]
    review1 --> text3[text()]
    title2 --> text4[text()]
    price2 --> text5[text()]
    review2 --> text6[text()]
    title3 --> text7[text()]
    price3 --> text8[text()]
    review3 --> text9[text()]
```

Result

```
<entry>
  <title>Data on the Web</title>
  <price>34.95</price>
  <review>A very good discussion of semi-structured database
    systems and XML.</review>
</entry>
<entry>
```



Step 3: Open Database

Open Database

The screenshot shows the BaseX 7.7.1 - XQuerySample application window. An arrow points to the 'Open Database' button in the toolbar. The interface is divided into three main sections: Editor, Result, and a tree diagram.

Editor: The top section contains a toolbar with icons for file operations and a text area with the XQuery `//entry`. Below the text area is a status bar showing 'OK' and '1 : 8'.

Result: The bottom section displays the XML result of the query. It shows a tree structure with three `entry` elements under a `reviews` root. Each `entry` element has three children: `title`, `price`, and `review`. The `review` element is a text node containing the text 'A very good discussion of semi-structured database systems and XML.'

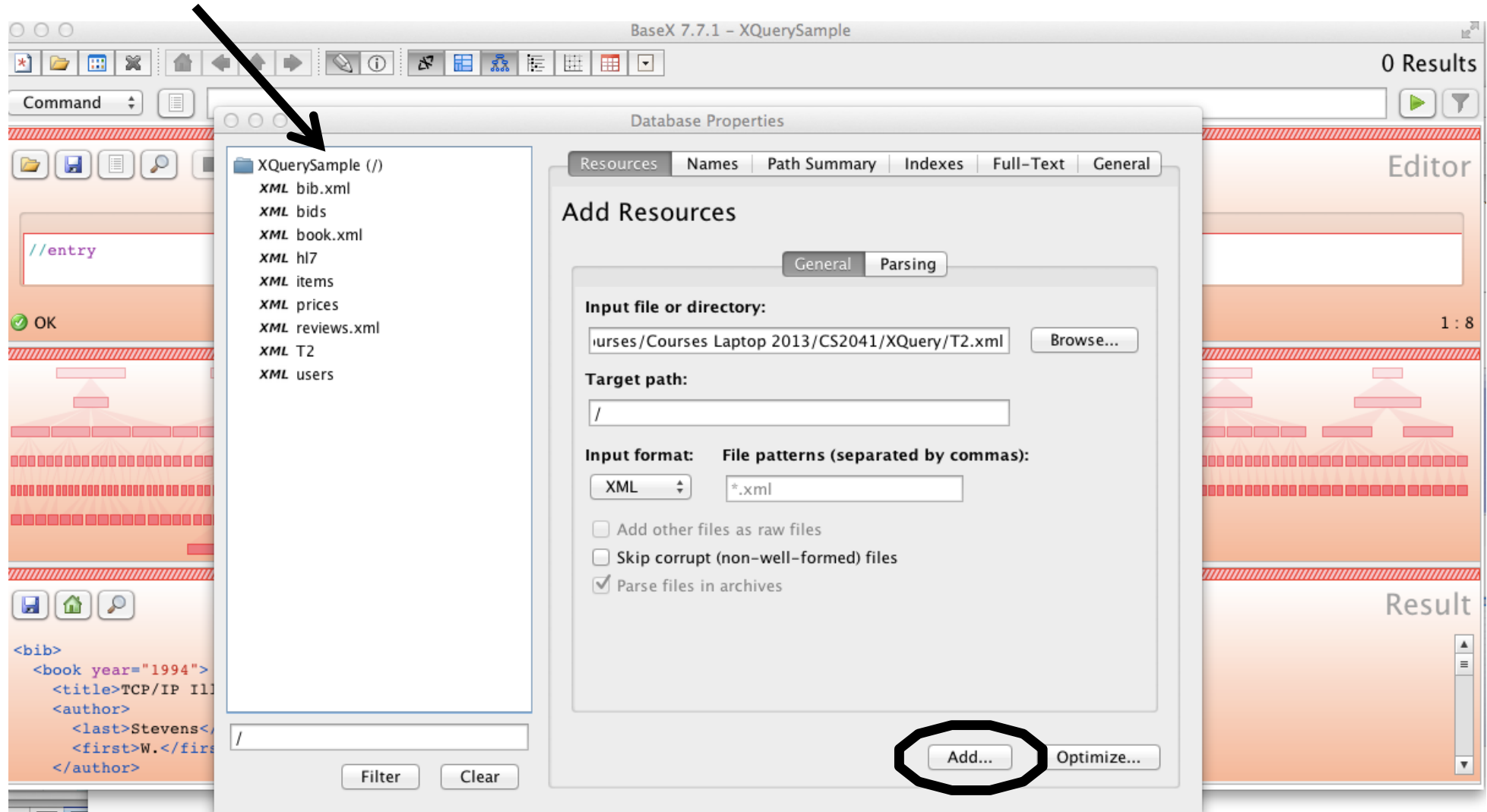
Tree Diagram: The tree diagram visualizes the XML structure. The root node is `reviews`, which has three children: `entry`, `entry`, and `entry`. Each `entry` node has three children: `title`, `price`, and `review`. The `review` node is a text node containing the text 'A very good discussion of semi-structured database systems and XML.'

Note that when XML is loaded in the database, you no longer need doc("filename") in the XQuery

Step 4: Add more XML files

Menu "Database" then item "Properties"

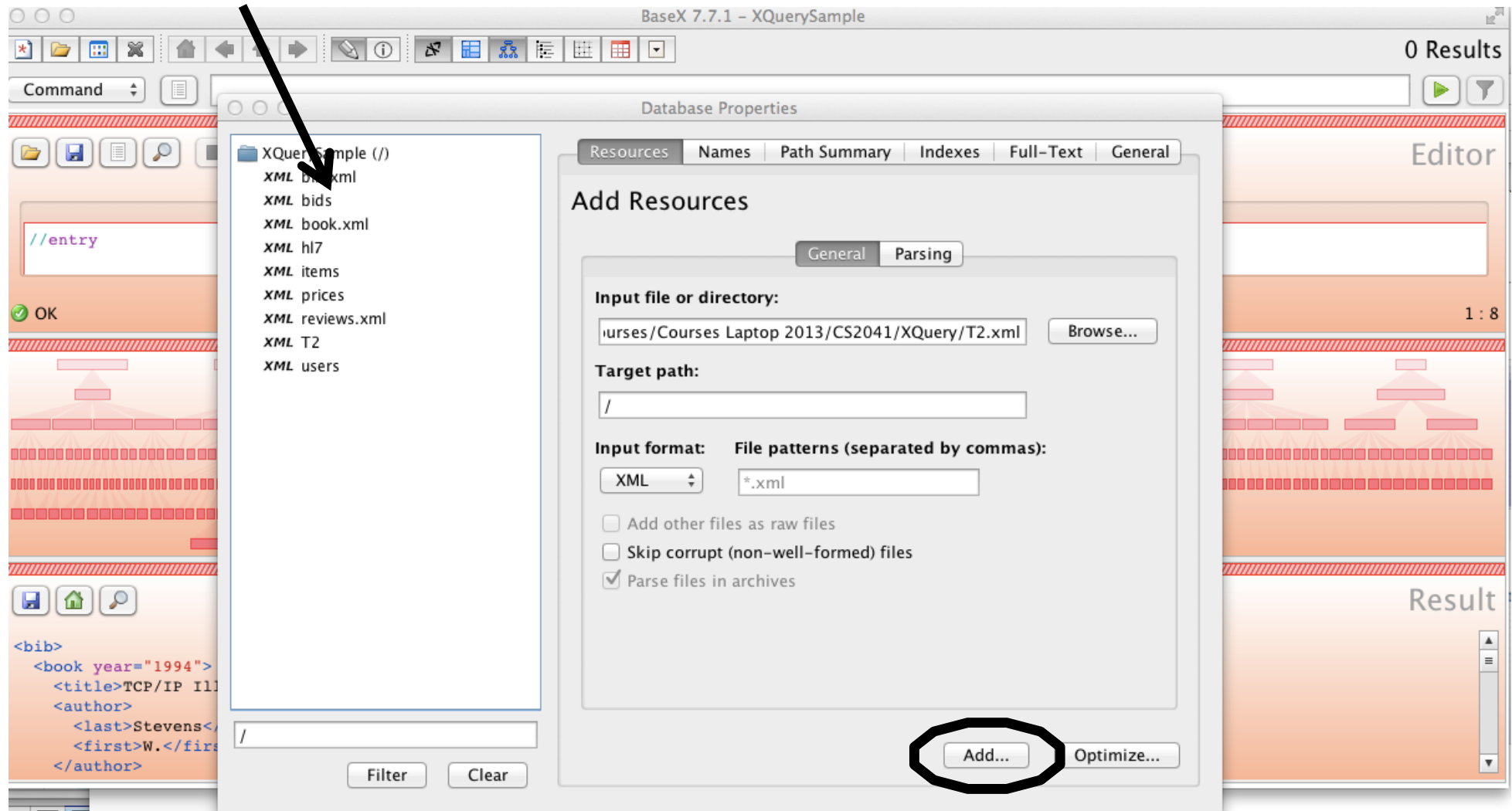
XML files that were added to the Database using Add button below



Note that when XML is loaded in the database, any changes to XML file has to be RELOADED

To Delete XML file from database so you can load new version

Select XML file, right click, select Delete



Note that when XML is loaded in the database, any changes to XML file has to be RELOADED

Step 5: Visualise Data

Buttons for Different Visualisations of XML
that is in THIS Database


The screenshot shows the XQuerySample application window. The title bar reads "XQuerySample 7.7.1 - XQuerySample". The toolbar contains various icons for file operations and data visualization. A black oval highlights a group of icons: a tree structure, a table, a grid, and a bar chart. An arrow points from the text "Buttons for Different Visualisations of XML that is in THIS Database" to this group of icons. Below the toolbar is a "Command" input field. The main area is divided into three sections: "Editor", "Result", and "Diagram". The "Editor" section shows the XML query `//entry`. The "Result" section displays the XML output:

```
<entry>
  <title>Data on the Web</title>
  <price>34.95</price>
  <review>A very good discussion of semi-structured database
    systems and XML.</review>
</entry>
<entry>
```

. The "Diagram" section shows a tree structure of the XML data, with nodes for "entry", "title", "price", "review", and "text()".




BaseX.org

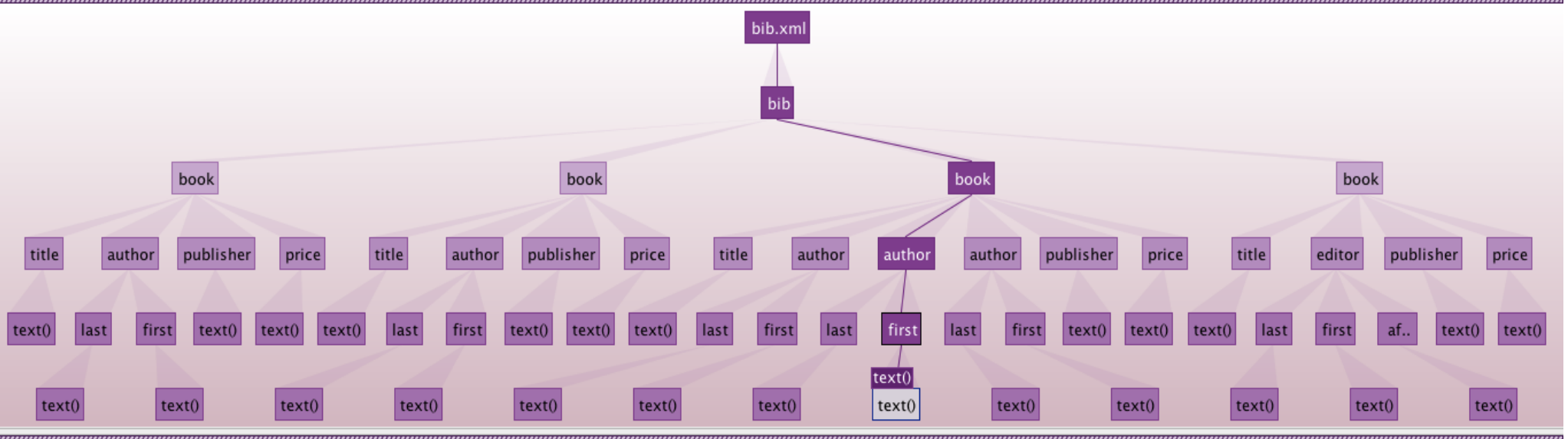


Result

```
<title>TCP/IP Illustrated</title>
<author>
  <last>Stevens</last>
  <first>W.</first>
</author>
<publisher>Addison-Wesley</publisher>
<price>65.95</price>
</book>
<book year="1992">
  <title>Advanced Programming in the Unix environment</title>
  <author>
    <last>Stevens</last>
    <first>W.</first>
  </author>
  <publisher>Addison-Wesley</publisher>
  <price>65.95</price>
</book>
<book year="2000">
  <title>Data on the Web</title>
```



Editor 1:1



```
graph TD
    bib[bib.xml] --> bib[bib]
    bib --> book1[book]
    bib --> book2[book]
    bib --> book3[book]
    bib --> book4[book]
    book1 --> title1[title]
    book1 --> author1[author]
    book1 --> publisher1[publisher]
    book1 --> price1[price]
    title1 --> text1_1[text()]
    author1 --> last1[last]
    author1 --> first1[first]
    last1 --> text1_2[text()]
    first1 --> text1_3[text()]
    publisher1 --> text1_4[text()]
    price1 --> text1_5[text()]
    book2 --> title2[title]
    book2 --> author2[author]
    book2 --> publisher2[publisher]
    book2 --> price2[price]
    title2 --> text2_1[text()]
    author2 --> last2[last]
    author2 --> first2[first]
    last2 --> text2_2[text()]
    first2 --> text2_3[text()]
    publisher2 --> text2_4[text()]
    price2 --> text2_5[text()]
    book3 --> title3[title]
    book3 --> author3[author]
    book3 --> author3_2[author]
    book3 --> publisher3[publisher]
    book3 --> price3[price]
    title3 --> text3_1[text()]
    author3 --> last3[last]
    author3 --> first3[first]
    last3 --> text3_2[text()]
    first3 --> text3_3[text()]
    author3_2 --> first4[first]
    first4 --> text3_4[text()]
    publisher3 --> text3_5[text()]
    price3 --> text3_6[text()]
    book4 --> title4[title]
    book4 --> editor4[editor]
    book4 --> publisher4[publisher]
    book4 --> price4[price]
    title4 --> text4_1[text()]
    editor4 --> last4[last]
    editor4 --> first4_2[first]
    last4 --> text4_2[text()]
    first4_2 --> text4_3[text()]
    publisher4 --> af4[af..]
    af4 --> text4_4[text()]
    price4 --> text4_5[text()]
    text3_4 --> text3_4_2[text()]
    text4_3 --> text4_3_2[text()]
    text4_5 --> text4_5_2[text()]
```