### **COMP9321 Web Application Engineering**

### Extensible Markup Language (XML)

Dr. Basem Suleiman Service Oriented Computing Group, CSE, UNSW Australia

Semester 1, 2016, Week 4

http://webapps.cse.unsw.edu.au/webcms2/course/index.php?cid=2442

• COMP9321, 16s1, Week 4

## Acknowledgement/Contributions

### Service Oriented Computing Group, CSE, UNSW Australia

- Dr. Helen Paik
- Prof. Boualem Bentallah
- Dr. Srikumar Venugopal
- Dr. Moshe Chai Barukh
- Dr. Amin Beheshti
- Dr. Basem Suleiman
- Many others from service oriented computing group

### **Extensible Markup Language (XML)**

- XML is a markup language much like HTML
  - HTML:
    - HTML was designed to display data.
    - HTML tags are not predefined.
  - XML:
    - XML was designed to describe data.
    - XML tags are not predefined.
- XML is designed to be self-descriptive
- XML is a W3C Recommendation
- XML is derived from SGML (ISO 8879).
  - (Standard Generalized Markup Language) is a standard for how to specify a document markup language or tag set.

## **Extensible Markup Language (XML)**

- XML originally designed to meet the challenges of large-scale electronic publishing.
- XML separates presentation issues from the actual data.
- XML plays an increasingly important role in the exchange of a wide variety of data on the Web and elsewhere.
  - Needs a communication protocol?
  - e.g. SOAP stands for Simple Object Access Protocol
    - SOAP is based on XML
    - SOAP is a W3C recommendation
    - SOAP uses XML Information Set for its message format.

### **Extensible Markup Language (XML)**

<!DOCTYPE html> <html> <body> <h2>COMP9321 Students</h2> Student#: 50001 Name: Adam B. Program: 8543 Stage: 1 Student#: 50002 Name: Alex C. Program: 3978 Stage: 3 COMP9321 Students . . . Student#: 50001 Name: Adam B. </body> Program: 8543 </html> Stage: 1 Student#: 50002 Name: Alex C. Program: 3978 Stage: 3

<?xml version="1.0" encoding="UTF-8"?>
<comp9321\_students>
<student>
<id>>50001</id>
<name>Adam B.</name>
<program>8543</program>
<stage>1</stage>
</student>
<id>>50002</id>
<name>Alex C.</name>
<program>3978</program>
<stage>3</stage>
</student>
</comp9321\_students>

# XML: Tags, tags

Consider the following snippet of information from a staff list:

LName	Title	FName	School	Campus	Room
Edgar	Miss	Pam	Optometry	KG	B501
Edmond	Dr	David	Information Systems	GP	S842
Edmonds	Dr	lan	Physical Sciences	GP	M206

### In XML ...

<phonebook> <entry> <lastname>Edgar</lastname> <title>Miss</title> <fistname>Pam <school>Optometry</school> <campus>KG</campus> <room>B501</room></fistname></entry></phonebook>	<pre><entry>   <lastname>Edmond</lastname> &gt; <title>Dr</title>   <fistname>David   <school>Information Systems</school> &gt; <campus>GP</campus>   <room>S842</room> </fistname></entry> </pre>
<room>B501</room>	
	Entry continues

# Why XML?

### Early Web

- Used to publish documents to be read by humans
- HTML was designed for the purpose

### Today's Web

- Many business activities are performed on the Web
  - Dynamic interactions:
    - Web app  $\Leftrightarrow$  people / Web app  $\Leftrightarrow$  Web app
  - Web becomes a platform for data exchange
  - XML provides a simple, cross-platform data format
- Web contains vast amount of data published in HTML format
  - Many programs process or analyse such data
  - $\bullet$  HTML changes ... (when data inside does not)  $\to$  the program that reads the HTML page must change too
  - XML provides a long-term, reliable data format for publishing

# Why XML?

Benefits of using XML in document (data) exchange

- Self-describing, modular and portable data
- A common, widely accepted data representation language
- Standard supports available for creating/parsing XML docs
- Standard supports for checking validity of data
- Efficient search of business information
  - standard support for querying XML docs
  - quick and simple search (XPath)
  - more comprehensive keyword + structure based search possible as well (XQuery)
- Extensible document descriptions
  - XML is flexible (cf. relational tables)!
  - reuse, adaptation of existing documents

## Separating the Content from Presentation



HTML was designed to display data.

- CSS: Cascading Style Sheets
- CSS defines how HTML elements are to be displayed
- All formatting could be removed from the HTML document, and stored in a separate CSS file.



XML was designed to describe data.

## Separating the Content from Presentation

```
<!DOCTYPE html>
<html>
```

```
<head>
<style>
body {background-color: green;}
h2 {color: red;}
</style>
</head>
```

<body>

<h2>COMP9321 Students</h2>

```
Student#: 50001
Student#: 50001
Name: Adam B.
Program: 8543
Stage: 1
```

```
Student#: 50002
Name: Alex C.
Program: 3978
Stage: 3

</body>
</html>
```



# **XML** Applications

Like any other good inventions, XML is now used for things that are far beyond its creators original imagination.

- A set of 'tags' that are developed for specific types of documents.
- e.g., Chemical Markup Language (CML)

```
<atom id="caffeine_karne_a_1">
<float builtin="x3" units="A">-2.8709</float>
<float builtin="y3" units="A">-1.0499</float>
<float builtin="z3" units="A">0.1718</float>
<string builtin="elementType">C</string>
</atom>
```



# **XML** Applications

### RSS : Really Simple Syndication

With RSS it is possible to distribute up-to-date web content from one web site to thousands of other web sites around the world.

- RSS is written in XML
- RSS allows you to syndicate your site content
- RSS defines an easy way to share and view headlines and content
- RSS files can be automatically updated
- RSS allows personalized views for different sites
- RSS is useful for web sites that are updated frequently, like:
  - e.g. News sites, Companies, and Calendars.

Without RSS, users will have to check your site daily for new updates.

# **XML** Applications

### RSS : Really Simple Syndication

Really Simple Syndication (RSS)	
<rss version="0.91"> <channel> <title>CNN.com</title></channel></rss>	RSS FEEDS
<item> <title>July ends with 76 killed</title> <link/>http://www.cnn.com//story.html <description>Three U.S. soldiers were</description> </item>	<ul> <li>Opinion</li> <li>The Nation</li> <li>Business</li> <li>Sport</li> <li>Media</li> <li>Travel</li> </ul>

Without RSS, users will have to check your site daily for new updates.

## XML is ...

- Is a Language
  - there is a grammar, and it can be parsed by machines
- Is a Markup Language
  - XML looks a bit like HTML (tags).
  - But it describes what things are, not what they are supposed to do
- Is eXtensible
  - you can define more words and add to the language
- XML is for structuring data.
- XML is for describing data.
- XML is text, but isn't meant to be read.
- XML is verbose by design.
- XML is a family of technologies.
- XML is license-free, platform-independent and well-supported.
- XML is NOT a programming language.
  - it is not something you can 'compile'

### Quick XML syntax

- All XML documents must have 'a' root element
- All XML elements must have a closing tag
  - $\bullet~$  Empty element tags end with />
- XML tags are case sensitive (NAME vs. Name)
- All XML elements must be nested (<q>??)
- Element Naming
  - letters, numbers, and other characters
  - must not start with a number, '. (period)' or '- (hyphen)'
  - must not start with 'xml' (or XML or Xml ..)
  - cannot contain spaces
- Attribute values must always be quoted (single or double)
- Comments in XML: <!-- This is a comment -->

## The XML Family

XML: a markup language used to describe information.

# The XML Family

XML: a markup language used to describe information.

#### **Document Object Model (DOM )**

- XML DOM defines a standard for accessing and manipulating XML documents.
- The DOM presents an XML document as a tree-structure.
- The DOM is a W3C standard.
- The DOM is separated into 3 different parts / levels:
  - Core DOM standard model for any structured document
  - XML DOM standard model for XML documents
    - A standard object model for XML
    - A standard programming interface for XML
    - Platform- and language-independent
  - HTML DOM standard model for HTML documents

# The XML Family

XML: a markup language used to describe information.

#### **Document Object Model (DOM):**

- XML DOM defines a standard for accessing and manipulating XML documents.
- The DOM presents an XML document as a tree-structure.
- The DOM is a W3C standard.
- The DOM is separated into 3 different parts / levels:
  - Core DOM standard model for any structured document
  - XML DOM standard model for XML documents
    - A standard object model for XML
    - A standard programming interface for XML
    - Platform- and language-independent
  - HTML DOM standard model for HTML documents

## XML Document: Bookstore Example

1	xml-version="1.0"-encoding="UTF-8"?					
2						
3 -	<pre>&gt; <bookstore></bookstore></pre>					
4	_					
5 -	<book category="Cooking"></book>					
6	<title-lang="en">Everyday Italian</title-lang="en">					
7	<author>Giada De Laurentiis</author>					
8	<year>2005</year>					
9	<price>30.00</price>					
10						
11						
12 -	<book category="Web Technologies"></book>					
13	<title-lang="en">Extensible Markup Language (XML) </title-lang="en">					
14	<author>W3C</author>					
15	<year>1998</year>					
16	<price>0.00</price>					
17						
18						
19 -	<book<category="web development"=""></book<category="web>					
20	<title-lang="en">JavaServer Pages (JSP)</title-lang="en">					
21	<author>Jason Hunter</author>					
22	<year>2000</year>					
23	<price>35.00</price>					
24						
25						
26						
27						

Using xml-formatter, Web Toolkit Online http://www.webtoolkitonline.com/xml-formatter.html

## XML Tree Structure: Bookstore Example



Source: http://www.w3schools.com/xml/xml\_tree.asp

## The XML Family – Basics

- XML documents are formed as **element trees**
- XML prolog: XML version and encoding character <?xml version="1.0" encoding="UTF-8"?>
- An XML tree starts at a root element <bookstore> and branches from the root to child elements (e.g., <book>)
- All elements can have sub elements (child elements, e.g., <author>)
- The terms parent, child, and sibling describes the relationships between elements
- Parent have children. Children have parents. Siblings are children on the same level (brothers and sisters)
- All elements can have text content (e.g., Everyday Italian) and attributes (e.g., <a href="https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript">https://www.selementscondingsuperscript@sleetscondingsuperscript@sle
- Elements has starting and closing tags (**<price> </price>**)

## The XML Family – Basics

- XML tags are case sensitive. <Student> is different from <student>
- All elements **must** be properly nested within each other
- Attribute values must always be quoted <student id="3230813">
- Comments <!-- XML comments -->
- White spaces are preserved in XML (in contrast to HTML)
- XML Stores **New Line** as LF (carriage return and line feed)

#### Entity References

- Some characters are reserved and have special meaning in XML, e.g., <, >, &
- Must be escaped using entity references

<Age> should be &lt; 25 </Age>

Character	Entity reference
< (less than)	<
> (greater than)	>
& (ampersand)	&
<b>'</b> (apostrophe)	'
" (quotation)	"

## XML – Document Type Definition

- An XML document with correct syntax is called "Well Formed"
- Errors (incorrect syntax) application processing will trigger errors
- Well Formed XML document  $\rightarrow$  it has valid XML syntax rules
- Well formed XML document  $\rightarrow$  "valid" XML document ?
- A valid XML document must be:
  - Well formed AND
  - Conform to Document Type Definition (DTD)
- Document Type Definition (DTD)
  - Defines the structure and the legal elements and attributes of an XML document
  - DTD or XML Schema (XML alternative to DTD)
  - Internal DTD declaration or external DTD declaration (.dtd)

# **DTD** Parsing

#### In DTD

- Parsed Character DATA (PCDATA)
  - Text that WILL be parsed by a parser. tags inside the text will be treated as markup and entities will be expanded
- Character Data (CDATA)
  - Text that will NOT be parsed by a parser tags inside the text will NOT be treated and entities will not be expanded.

### Phonebook.xml with Internal DTD

```
<?xml version="1.0"?>
<!DOCTYPE Phonebook [
  <!ELEMENT Phonebook (Entry)+ >
  <!ELEMENT Entry (LastName, FirstName, School, Campus, Room, Extension)>
  <!ELEMENT LastName (#PCDATA)>
  <!ELEMENT FirstName (#PCDATA)>
  <!ELEMENT School (#PCDATA)>
  <!ELEMENT Campus (#PCDATA)>
  <!ELEMENT Room (#PCDATA)>
  <!ELEMENT Extension (#PCDATA)>
  <!ATTLIST LastName Title (Miss | Ms | Mrs | Mr | Dr | Prof) #REQUIRED>
1>
<Phonebook>
   <Entry>
      <LastName Title="Miss">Edgar</LastName>
      <FirstName>Pam</FirstName>
      <School>Optometry</School>
      <Campus>GP</Campus>
      <Room>B501</Room>
      <Extension>5695</Extension>
   </Entry> <!-- more entries not shown ... -->
</Phonebook>
```

### Phonebook.xml with External DTD. Phonebook.dtd

#### Phonebook.xml

<?xml version="1.0"?>

<!DOCTYPE Phonebook SYSTEM "Phonebook.dtd">

<Phonebook>

<Entry>

```
<LastName Title="Miss">Edgar</LastName>
```

```
<FirstName>Pam</FirstName> <!-- rest of the entries -->
```

</Phonebook>

#### Phonebook.dtd

ELEMENT</th <th>Phonebook</th> <th>(Entry+</th> <th>-) &gt;</th> <th></th> <th></th> <th></th> <th></th> <th></th> <th></th>	Phonebook	(Entry+	-) >						
ELEMENT</td <td>Entry (Las</td> <td>tName,</td> <td>First</td> <td>Vame,</td> <td>Schoo</td> <td>l,Camp</td> <td>us,</td> <td>Room,</td> <td>Extension)&gt;</td>	Entry (Las	tName,	First	Vame,	Schoo	l,Camp	us,	Room,	Extension)>
ELEMENT</td <td>LastName</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	LastName	(#PCDAT	(A)>						
ELEMENT</td <td>FirstName</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	FirstName	(#PCDAT	(A)>						
ELEMENT</td <td>School</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	School	(#PCDAT	(A)>						
ELEMENT</td <td>Campus</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	Campus	(#PCDAT	(A)>						
ELEMENT</td <td>Room</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	Room	(#PCDAT	(A)>						
ELEMENT</td <td>Extension</td> <td>(#PCDAT</td> <td>(A)&gt;</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>	Extension	(#PCDAT	(A)>						
ATTLIST</td <td>LastName T</td> <td>itle (M</td> <td>liss  </td> <td>Ms  </td> <td>Mrs  </td> <td>Mr   1</td> <td>Dr  </td> <td>Prof)</td> <td>#REQUIRED&gt;</td>	LastName T	itle (M	liss	Ms	Mrs	Mr   1	Dr	Prof)	#REQUIRED>

### **CDATA Section**

Sometimes the character data of an element might contain too many characters that need to be escaped (e.g., chunk of other XML parts or HTML code).

CDATA section lets you enclose the character data as literal.

#### Example

Everything between <!CDATA[ and ]]> is treated as raw characters, not markups.

<sup>•</sup> COMP9321, 16s1, Week 4

## **Defining XML Content: Elements**

#### A Book <book> <author> <name>J.K. Rowling</name> </author> <detail> <series>Seventh</series> <title>Harry Potter and the Deathly Hallows</title> </detail> </book>

#### **Creating Elements:**

<!ELEMENT book (author, detail)> <!ELEMENT author (name)> <!ELEMENT name (#PCDATA)> <!ELEMENT detail (series, title)> <!ELEMENT series (#PCDATA)> <!ELEMENT title (#PCDATA)>

## **Defining XML Content: Modifiers**

### A Book

```
<book><br/>
<author> <!- more than one authors? -><br/>
<name>E. Harold</name> <name>S. Means</name><br/>
</author><br/>
<detail> <!- not every book is in a series -><br/>
<title>XML in a Nutshell</title><br/>
</detail> </book>
```

- ? : optional element (only once)
- 2 + : mandatory element (1 or more)

```
<!ELEMENT book (author, detail*)>
<!ELEMENT author (name+)>
<!ELEMENT name (#PCDATA)>
```

```
<!ELEMENT detail (series?, title)>
<!ELEMENT series (#PCDATA)>
<!ELEMENT title (#PCDATA)>
```

# Defining XML Content: Choices, Empty

#### **Element Choices**

<!ELEMENT newbooks (book+)> <!ELEMENT book (author+, detail\*)> <!ELEMENT author (name | penname)> <!ELEMENT name (#PCDATA)> <!ELEMENT penname (#PCDATA)> <!ELEMENT detail ((series?, title) | (publisher, release\*))> <!ELEMENT series (#PCDATA)> <!ELEMENT title (#PCDATA)> <!ELEMENT title (#PCDATA)> <!ELEMENT publisher (#PCDATA)>

#### Empty Element Content

<!ELEMENT BR EMPTY>

<BR/> is called an empty element

# Defining XML Content: Mixed content, Any

#### Mixed content: mixture of elements and text

<!ELEMENT message (#PCDATA | bold | italic)\*>

<message>You <italic> really <bold>must</bold> try this delicious <bold>new</bold> recipe for <italic>pudding </message>

#### ANY : Any predefined element could be included

<!ELEMENT book (author+, description, detail\*)> <!ELEMENT author (name+)> <!ELEMENT name (#PCDATA)> <!ELEMENT description ANY> <!ELEMENT detail (series?, title)> <!ELEMENT series (#PCDATA)> <!ELEMENT title (#PCDATA)>

## **Defining XML Content: Creating Attributes**

#### < book >

```
<author period="classical" category="children">
```

```
<name type="normal">J.K. Rowling</name>
```

```
</author>
```

```
<title>Harry Potter and the Half-Blood Prince</title>
```

</book>

### **Creating Attributes:**

<!ELEMENT book (author, title)> <!ELEMENT author (name+)> <!ELEMENT name (#PCDATA)> <!ELEMENT title (#PCDATA)> <!ATTLIST name type (normal | penname) "normal" #REQUIRED> <!ATTLIST author period CDATA #REQUIRED category CDATA #IMPLIED>

### **Defining XML Content: Creating Attributes**

#### Default values for attributes:

The default postcode in an address is to be 4001, state must be QLD.

 $<\!! {\sf ATTLIST}$  Address Postcode CDATA "4001" State CDATA  $\# {\sf FIXED} "\, {\sf QLD}" \!>$ 

The above definition has the following effects on the source doc.

- <Address />  $\rightarrow$  <Address Postcode="4001" State="QLD">
- <Address Postcode="4010" State="QLD"/>  $\rightarrow$  (no error)
- <Address Postcode="4001" State="NSW"/> → (error)

### **XML Custom Entities**

Sometimes it might be desirable to construct a document from several (other) XML documents:

<!DOCTYPE sql [ <!ELEMENT sql (select, from)> <!ELEMENT select (col+)> <!ATTLIST select order CDATA #REQUIRED> <!ELEMENT col (#PCDATA)> <!ELEMENT from (table+)> <!ELEMENT table (#PCDATA)> <!ENTITY select SYSTEM "select.xml"> <!ENTITY from SYSTEM "from.xml">]> <sql>&select;&from;</sql>

Where "select.xml" contains:	And "from.xml" contains:		
<select order="cost"></select>	<from></from>		
<col/> CarNr	Cars		
<col/> Make			
<col/> Cost			

### **Parameter Entities**

It might be a good idea to fragment the DTD in the same way that the document content is partitioned:

```
<!DOCTYPE sql [
<!ELEMENT sql (select, from)>
<!ENTITY % seldef SYSTEM "select.dtd">
%seldef;
<!ENTITY % fromdef SYSTEM "from.dtd">
%fromdef;
<!ENTITY select SYSTEM "from.dtd">
<!ENTITY select SYSTEM "select.xml">
<!ENTITY from SYSTEM "from.xml">]>
<sql>&select;&from;</sql>
```

Where select.dtd is defined as:	and "from.dtd" is:			
ELEMENT select (col+)	ELEMENT from (table+)			
ATTLIST select order CDATA #REQUIRED	ELEMENT table (#PCDATA)			
ELEMENT col (#PCDATA)				

### Well-formedness and Validity of XML

Well-formedness Rules:

- Open and close all tags
- Empty-element tags end with />
- There is a unique root element
- Elements may not overlap
- Attribute values are quoted
- $\circ$  < and & are only used to start tags and entity references, respectively
- Only the five predefined entity references are used

Validity Rules:

- Well-formed
- Must have a Document Type Definition (DTD)
- Must comply with the constraints specified in the DTD
## Limitations of DTD

- Limited support for constraining attribute values
- No limits on character data
- Limited support for namespaces
- No self-documentation, does not have XML syntax

### The XML Family – XML Schema

#### XML Schema (or SML Schema Definition XSD)

- is an XML-based alternative to DTD.
- describes the structure of an XML document.
- defines elements and attributes that can appear in a document
- defines data types for elements and attributes
- defines default and fixed values for elements and attributes
- defines the child elements, their orders, etc.
- XML Schemas are much more powerful than DTDs.
- The XML Schema language is also referred to as XML Schema Definition (XSD).

### The XML Family – XML Schema

#### XML Schema (orXSD)

- W3C's recommendation for replacing DTD with features such as:
- Simple and complex data types
- Type derivation and inheritance
- Namespace-aware element and attributes
- Limits on number of appearances by an element
- Combining with regular expressions for finer control over document structure
- Most importantly, XML Schemas are well-formed XML documents themselves

## XML Namespaces

• XML elements can have any names.

• What if a name could mean two different things (ie., name clash)?

The following two XML documents that describe student information.

From University X:	From University Y:
<student></student>	<student></student>
<id>12345</id>	<id>534-22-5252</id>
<name>Jeff Smith</name>	<name>Bob Citizen</name>
<language>C#</language>	<language>Spanish</language>
<rating>9.5</rating>	<rating>3.2</rating>

How could a program distinguish the different elements?

# XML Namespaces (example)



# XML Namespaces

A namespace is a set of *names* in which all names are unique.



- The name 'title' can now be identified as: Book.title, Project.title, Employee.title ...
- These names are called "qualified names".
- XML namespaces give elements and attributes a unique name across the Internet.
- XML namespaces enable programmers to process the tags and attributes they care about and ignore those that don't matter to them.

# Previous examples can now be ...

The previous examples can now have qualified names:





# XML Namespace Syntax

- wmlns:<prefix>='namespace identifier'
- eg., <books xmlns:xdc="http://www.xml.com/books">
- not a normal XML attribute (treated differently)
- the URI must be unique, but may not represent a 'useful' resource
- the prefix is by convention or author's choice

Consider the following XML document: painting.xml

### XML Schema Definition (XSD)

- a recommendation of the World Wide Web Consortium (W3C)
- specifies how to formally describe the elements in an Extensible Markup Language (XML) document.

Types and Declarations

- Simple Types Basic data types such as strings, integers, boolean, etc.
- Complex Types Composed of simple types. Consists of an arrangement of elements and attributes.
- Element declarations Associates an element name to an instance of a simple or complex type
- Attribute declaration Associates an attribute name with an instance of a simple type.

# Simple Types

Simple Type Declaration

• Cannot contain elements or attributes

• Can be pre-defined or user-defined. Examples:

<element name='sname' type="string" />
<element name='age' type="integer" />
<element name='course' type='string'/>

Which defines:

<sname>John Doe</sname> <age>24</age> <course>COMP9321</course>

### Attributes

Attributes are simple types as well (however, simple types themselves CANNOT contain attributes)

```
<attribute name="currency" type="string"/>
```

Attributes can have default or fixed values and can be optional or required

<attribute name="currency" type="string" default="EUR"/>
<attribute name="currency" type="string" fixed="AUD"/>
<attribute name="currency" type="string" use="required"/>

- <attribute name="attribute-name" type="attribute-type" use="uu"/>
- Built-in data types
  - String, decimal, integer, boolean, date, time
- A default value is automatically assigned to the attribute when no other value is specified
- A fixed value is also automatically assigned to the attribute, and you cannot specify another value
- Attributes are optional by default, otherwise use="required"

# **Type Restrictions**

Restrictions are used to specify a range of acceptable values for XML elements or attributes. Examples:

```
<simpleType name="nameString">
<restriction base="string">
<pattern value="([a-zA-Z])+"/>
</restriction>
</simpleType>
<element name='sname' type="tns:nameString" />
```

```
<simpleType name="ageNum">
<restriction base="integer">
<minInclusive value="1"/>
<maxInclusive value="80"/>
</restriction>
</simpleType>
<element name='age' type="tns:ageNum" />
```

```
<simpleType name="courseString">
<restriction base="string">
<enumeration value="COMP9321"/>
<enumeration value="COMP9322"/>
<enumeration value="COMP9323"/>
</restriction>
</simpleType>
<element name='course' type="tns:courseString"/>
```

More examples : http://www.w3schools.com/xml/schema\_facets.asp

# **Type Restrictions**

han
)
lue)
d
t

More examples : http://www.w3schools.com/xml/schema\_facets.asp

# **Complex Types**

Complex Types can be empty or composed of only elements, or only text, or a mix of both elements and text.

The number of elements are controlled by the *indicators* as below: Order indicators:

- All all elements specified in the type can occur in any order but must occur only once
- Choice either one or the other element must be present
- Sequence all elements must occur in the order specified. You can have more than one element.

Occurrence indicators:

- maxOccurs
- minOccurs

# **Complex Types**

```
<complexType name="student">
<sequence>
<element name='sname' type="tns:nameString" />
<element name='age' type="tns:ageNum" />
<element name='course' type="tns:courseString"/>
</sequence>
</complexType>
```

which defines:

<student> <sname>John Doe</sname> <age>24</age> <course>COMP9321</course> </student>

# XML Example - Creating XSD

#### Shiporder.xml

#### <?xml version="1.0" encoding="UTF-8"?> 1 k?xml version="1.0" encoding="UTF-8" ?> 1 2 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"> 2 -3 -<shiporder orderid="889923"</pre> 3 -<xs:element name="shiporder"> xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" 4 4 -<xs:complexType> 5 xsi:noNamespaceSchemaLocation="shiporder.xsd"> 5 -<xs:sequence> 6 <orderperson>John Smith</orderperson> <xs:element name="orderperson" type="xs:string"/> 6 7 -<shipto> 7 -<xs:element name="shipto"> 8 <name>01a Nordmann</name> 8 -<xs:complexType> 9 <address>Langgt 23</address> 9 -<xs:sequence> 10 <city>4000 Stavanger</city> 10 <xs:element name="name" type="xs:string"/> <country>Norway</country> 11 11 <xs:element name="address" type="xs:string"/> 12 </shipto> <xs:element name="city" type="xs:string"/> 12 13 -<item> 13 <xs:element name="country" type="xs:string"/> 14 <title>Empire Burlesque</title> 14 </xs:sequence> <note>Special Edition</note> 15 15 </xs:complexType> <quantity>1</quantity> 16 16 </xs:element> <price>10.90</price> 17 17 -<xs:element name="item" maxOccurs="unbounded"> 18 </item> 18 -<xs:complexType> 19 -<item> 19 -<xs:sequence> <title>Hide your heart</title> 20 <xs:element name="title" type="xs:string"/> 20 <quantity>1</quantity> 21 21 <xs:element name="note" type="xs:string" minOccurs="0"/> <price>9.90</price> 22 <xs:element name="quantity" type="xs:positiveInteger"/> 22 </item> 23 23 <xs:element name="price" type="xs:decimal"/> 24 </shiporder> 24 </xs:sequence> 25 25 </xs:complexType> 26 </xs:element> 27 </xs:sequence> <xs:attribute name="orderid" type="xs:string" use="required"/> 28 29 </xs:complexType> </r></r></xs:element></r> 30

Source: http://www.w3schools.com/xml/schema\_example.asp

Shiporder.xsd

### The XML Family – XSL and XSLT

XML: a markup language used to describe information.

**DOM**: a programming interface for accessing and updating documents.

**DTD and XML Schema:** describes the structure and content of XML documents. **XSLT:** 

- XSL stands for eXtensible Stylesheet Language, and is a style sheet language for XML documents.
- CSS = Style Sheets for HTML
- XSL = Style Sheets for XML
  - XSL describes how the XML document should be displayed!
- XSLT (XSL Transformations) a language for transforming XML documents.

*e* C:\Users\Amin\student.xml

P → C @ C:\Users\Amin\Dropbox\\_U...×

File Edit View Favorites Tools Help

¢

id	name	program	stage			
50001	Adam B.	8543	1			
50002	Alex C.	3978	3			
Eile	student	<b>xml - Note</b> t View He	pad	-		×
xm<br xm<br href <com <st <st <st <st <st <st <st <st< th=""><td><pre>il version il-styles il-styles</pre></td><th>n="1.0" er heet type= t.xsl"?&gt; udents&gt;  m B.8543/stage&gt;  x C.3978/stage&gt; tudents&gt;</th><th>e&gt; gram&gt; gram&gt;</th><td>y="UTF /xsl"</td><td>-8 ?</td><th>.&gt; .</th></st<></st </st </st </st </st </st </st </com 	<pre>il version il-styles il-styles</pre>	n="1.0" er heet type= t.xsl"?> udents>  m B.8543/stage>  x C.3978/stage> tudents>	e> gram> gram>	y="UTF /xsl"	-8 ?	.> .

🧾 student.xsl - Notepad 🗕 🗖	×
<u>F</u> ile <u>E</u> dit F <u>o</u> rmat <u>V</u> iew <u>H</u> elp	
<pre>k?xml version="1.0" encoding="UTF-8"?&gt; <xsl:stylesheet version="1.0" xmlns:xsl="http://www.w3.org/1999/XSL/Transform"></xsl:stylesheet></pre>	^
<pre><xs1:template match="/"></xs1:template></pre>	
<h2>comp9321_students</h2>	
id	
name	
program	
stage	
<xsl:for-each select="comp9321_students/student"></xsl:for-each>	
<xsl:value-of select="id"></xsl:value-of>	
<xsl:value-of select="name"></xsl:value-of>	
<xsl:value-of select="program"></xsl:value-of>	
	~

# The XML Family - XPath

XML: a markup language used to describe information.

**DOM**: a programming interface for accessing and updating documents.

DTD and XML Schema: describes the structure and content of XML documents.XSLT: a language for transforming XML documentsXPath:

- XPath (XML Path language) is a language for finding information in an XML document.
- XPath contains a library of standard functions
- XPath is a major element in XSLT
- XPath is also used in XQuery, XPointer and XLink
- XPath is a W3C recommendation



# The XML Family – XPath Examples

)r

g

tε

Х

V,

 $a^{\dagger}$ 

. . .

- <?xml version="1.0" ....>
- <comp9321\_students>
- <student>

D

D

- <id>50001</id>
- Х <name>Adam B.</name>
- X <program>8543</program> <stage>1</stage>
  - </student>
  - <student>
  - <id>50002</id>
  - <name>Alex C.</name>
  - <program>3978</program>
  - <stage>3</stage>
  - </student>
  - </comp9321\_students>

- d some **XPath** expressions:
- /comp9321\_student/student[1] es
  - Selects the first "student" element that is the child
  - of the "comp9321 student" element
- a /comp9321\_student/student[last()]
  - Selects the last "student" element that is the child of the "comp9321 student" element
  - /comp9321\_student/student[position()<3]</pre>
  - Selects the first two "student" element that is the
  - child of the "comp9321\_student" element

#### /comp9321\_student/student[stage>2]

• Selects all the "student" elements of the "comp9321 student" element that have a "stage" element with a value greater than 2.

# The XML Family

XML: a markup language used to describe information.
DOM: a programming interface for accessing and updating documents.
DTD and XML Schema: describes the structure and content of XML documents.
XSLT: a language for transforming XML 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.

#### An XML document is a tree ...

An XML document is a tree ...



### Attributes in XML tags

LName	Title	FName	School	Campus	Room
Edgar	Miss	Pam	Optometry	KG	B501

#### Phonebook with Attributes

```
<Phonebook>
<Entry entryNumber="001">
<Name Title="Miss">
<Last>Edgar</Last>
<First>Pam</First>
</Name>
<School Campus="KG">Optometry</School>
<Room Building="B" Level="5">01</Room>
</Entry>
</Phonebook>
```

#### • Attribute order is not significant

### Attributes in XML tags

LName	Title	FName	School	Campus	Room
Edgar	Miss	Pam	Optometry	KG	B501

#### Phonebook with many attributes ...

```
<Phonebook>
<Entry entryNumber="001">
<Name Title="Miss" LName="Edgar" FName="Pam"/>
<Location Campus="KG" School="Optometry" Building="B" Room="501"/>
</Entry>
</Phonebook>
```

• Avoid using too many (loses structure, more parsing effort ...)

### Parsing XML documents with Java

Inside an XML file, there are ...

- Markup:
  - Tags, Entity References, Comments Processing Instructions, DTD declarations, XML declaration, and CDATA Section Delimiters
- and Character Data which includes everything else
  - Parsed Character Data (PCDATA): character data left after entity references are replaced with their text
  - e.g. Given <PUBLISHER>A &amp; M Records</PUBLISHER>, the parsed character data is A & M Records

```
<?xml version="1.0" encoding="UTF-8">
<!DOCTYPE SONG SYSTEM "song.dtd">
<SONG xmlns="http://www.cafeconleche.org/namespace/song">
<TITLE>New Slang</TITLE>
<ARTWORK ALT="Garden State" WIDTH="100" HEIGHT="200"/>
<ARTIST>The Shins</ARTIST>
<!-- The publisher is actually Polygram but I needed an example of a general entity reference. -->
<PUBLISHER>A &amp; M Records</PUBLISHER>
<LENGTH>3:29</LENGTH>
<YEAR>2004</YEAR>
</SONG>
```

# Parsing XML documents with Java

What do you mean by 'parsing (or processing) XML docs'?

- Parsing makes an interface available to your application that needs to make use of the document
- Through the interface, you can modify, retrieve the document contents.



What if the interface provided by the parser is parser-specific? → your application will have to be 'parser-specific'. Obviously ... we want "STANDARD"!



#### SAX and DOM as the Standard Interfaces

- SAX the Simple API for XML
- DOM the Document Object Model

Why two standards?  $\rightarrow$  trade-off between control and performance

- DOM gives you a tree structure
  - you have a complete control over the structure
  - ie., traverse the tree, modify structure, etc.
  - the tree gets stored in memory all at once
- SAX lays out the document in time, as a sequence of 'events'
  - events are associated with each tag (open/close), each tag body, etc.
  - you will write event handlers (ie., you can ignore certain events)
  - It requires much less memory
  - It becomes difficult to use if processing an element depends on earlier/later elements

## Document Object Model (DOM)

- DOM is an API for HTML and XML documents, its specification is developed by W3C (http://www.w3.org/DOM/DOMTR)
- It defines the logical structure of documents and the way a document is accessed and manipulated.



http://java.sun.com/j2se/1.5.0/docs/api/org/w3c/dom/package-summary.html • COMP9321, 16s1, Week 4

# Dealing with Nodes in DOM

In DOM, XML Documents are treated as a tree of nodes

**Types of Nodes**: Twelve different kinds of node are defined by the W3C DOM standard.

- elements
- attributes
- text
- CDATA
- entity reference
- entity
- processing instruction
- comment
- Document
- Document type
- Document fragment
- Notation

## An example XML here ...

```
<?xml version="1.0"?>
 <bibliography>
 <book>
 <author>North, Ken</author>
 <title>Database magic with Ken North</title>
 <address>Upper Saddle River, NJ</address>
 <publisher>Prentice Hall</publisher>
<year>1999</year>
<isbn>0136471994</isbn>
</book>
<book>
<author> Elmasri, Ramez, and Shamkant B. Navathe</author>
<title>Fundamentals of database systems</title>
<address>Reading, Mass.</address>
<publisher>Addison-Wesley Pub Co</publisher>
<year>2000 </year>
<edition>3rd</edition>
 <isbn>0201542633</isbn>
</book>
<book>
<author>Feiler, Jesse</author>
<title>Database-driven Web sites</title>
<address>San Francisco</address>
<publisher>Morgan Kaufmann</publisher>
<year>1999</year>
<isbn>0122513363</isbn>
</book>
 <!-- more book -->
</bibliography>
\% Text of an element node is stored in a text node.
```

# DOM for XML

- The root node has no parent but every other node has exactly one parent node
- A node can have any number of children
- Two nodes that have the same parent are called *siblings*
- A node with no children is called a *leaf node*
- Among siblings, the node that appears first in sequence is called the first child and the node that appears last is the last child
- Important: Text of an element node is stored in a separate text node.

# Using a DOM Parser (eg., Apache Xerces)

```
import org.w3c.dom.*;
import org.apache.xerces.parsers.DOMParser;
public class DOMCountNames {
   public static void main(String[] args) {
         try {
    DOMParser parser = new DOMParser();
    parser.parse(args[0]);
    Document doc = parser.getDocument();
    // do something ..
         catch(Exception e){
            e.printStackTrace(System.err);
 }
```

#### The DOMParser Class:

- DOMParser class is derived from the XMLParser class
- parse() method parses the input source given by a system identifier
- Document getDocument() method returns the document itself

#### **Document Interface Methods**

Once you have the **Document** object, you can:

- Attr createAttribute(String name): Creates an attribute
- *Element* createElement(String tagName): Creates an element
- Text createTextNode(String data): Creates a Text Node
- Element getDocumentElement(): Gets the root element of the document
- *Element* getElementById(String elementId): Get the element by ID
- NodeList getElementsByTagName(String tagname): Returns a NodeList of all the elements with a given tag name

#### NodeList Interface Methods:

- int getLength(): Gets the number of nodes in this list
- Node item(int index): Gets the item at the specified index value in the collection

• COMP9321, 16s1, Week 4

# Examples of Node Properties (XML), p.9.25



document.getElementsByTagName("person")[1] document.getElementsByTagName("person")[1].parentNode document.getElementsByTagName("person")[1].childNodes document.getElementsByTagName("person")[1].firstChild document.getElementsByTagName("person")[1].lastChild document.getElementsByTagName("person")[1].previousSibling document.getElementsByTagName("person")[1].lastChild.firstChild

last

person

persons

first

first

last

person

Li

#### Count/Print the number of 'book' elements

```
import org.w3c.dom.*;
import org.apache.xerces.parsers.DOMParser;
public class DOMCountNames
   public static void main(String[] args) {
      try{
 DOMParser parser = new DOMParser();
 parser.parse(args[0]);
 Document doc = parser.getDocument();
 NodeList nodelist = doc.getElementsByTagName("book");
 System.out.println(args[0] + " has " + nodelist.getLength()
                    + " <book> elements.");
      catch(Exception e){
 e.printStackTrace(System.err);
```

#### Compiling and Running

}

}

% javac -classpath ":xerces.jar" DOMCountNames.java % java -classpath ":xerces.jar" DOMCountNames books.xml

### Dealing with Nodes in DOM

There is a large range of methods that can be applied to the nodes:

#### Node Interface Methods

getNodeName()
getNodeType()
getFirstChild()
getNextSibling()
insertBefore(...)
appendChild(...)
normalize()
hasAttributes()

getNodeValue()
getParentNode()
getLastChild()
getAttributes()
replaceChild(...)
hasChildNodes()
isSupported(...)
getLocalName()

setNodeValue(...)
getChildNodes()
getPreviousSibling()
getOwnerDocument()
removeChild(...)
cloneNode(...)
getNamespaceURI()

Node interface documentationhttps://www.w3.org/2003/01/dom2-javadoc/org/w3c/dom/Node.html
## Dealing with Nodes

Consider the following program:

```
Document doc = parser.getDocument();
Element docRoot = doc.getDocumentElement();
String docRootName = docRoot.getTagName();
System.out.println("Doc root: "+docRootName);
int i =
          0:
for (Node node = docRoot.getFirstChild();
     node != null:
     node = node.getNextSibling()) {
  if (node.getNodeType()==Node.ELEMENT_NODE) {
    System.out.println(i+": " + node.getNodeType()
                              + node.getNodeName());
  }
  else {
    System.out.println(i+": " + node.getNodeType());
  i++;
```

The method getNodeType() returns a number in the range 1 to 12. Thus we can tell which kind of node we are dealing with. What will the output of this program be?

## The Element interface

This interface outlines operations that are specific to elements:

- getTagName(): This method returns the name of the tag associated with the element.
- getAttribute(name): This method returns a string containing the value of an attribute:
  - name is the name of the attribute.

Element Interface documentation – full list of methods and operations

https://www.w3.org/2003/01/dom2-javadoc/org/w3c/dom/Element.html

## More with DOM ...

Heaps of hands on tutorials on the web ...

- DOM and Javascript: e.g., http://www.sitepoint.com/print/xml-javascript-mozilla https://developer.mozilla.org/en/The\_DOM\_and\_JavaScript
- DOM, XML, Javascript and Ajax: e.g., http://www.w3schools.com/Ajax/ajax\_intro.asp

## References

- <u>http://www.w3.org/XML/</u>
- XML in a nutshell, Chapters 9 and 10
- <u>http://www.ibm.com/developerworks/library/xml-schema/</u>
- <u>http://www.w3schools.com/xml/</u>
- Some examples in these notes are originated from Dr. David Edmond from QUT, Brisbane

