

XML

Extensible Markup Language

Prof. Cesare Pautasso

<http://www.pautasso.info>

cesare.pautasso@unisi.ch

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

1

XML is the future (1999)

XML will be **foundation** for
future Web standards

XML will become the **language** for
international desktop and Web publishing

XML will become the
universal data exchange format
between heterogeneous environments

XML will replace all
existing Word Processing **storage** formats

Predictions by John Bosak (Sun), 1999

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

2

Contents

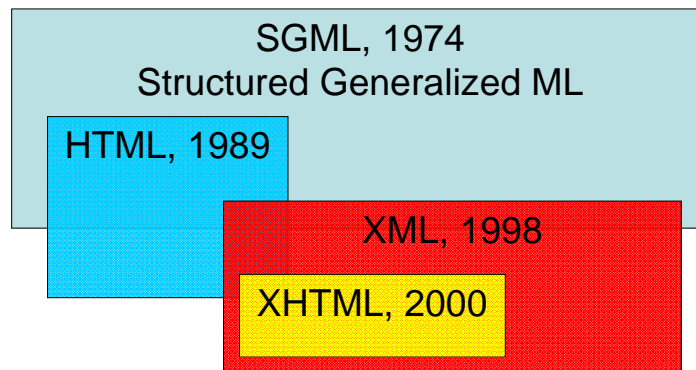
- Motivation: **Universal Syntax**
- History
- HTML vs. XML
- What is the Extensible Markup Language?
 - XML Syntax
 - XML Structure
- A critical look at XML
- XML Technology Landscape

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

3

XML History



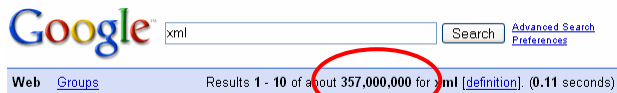
7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

4

(X)HTML Example

- Write a program that downloads some amazon.com web page and extracts the title, author and price of a book (given its ISBN number)
- Write a program that asks google.com how many Web pages are there about a certain topic

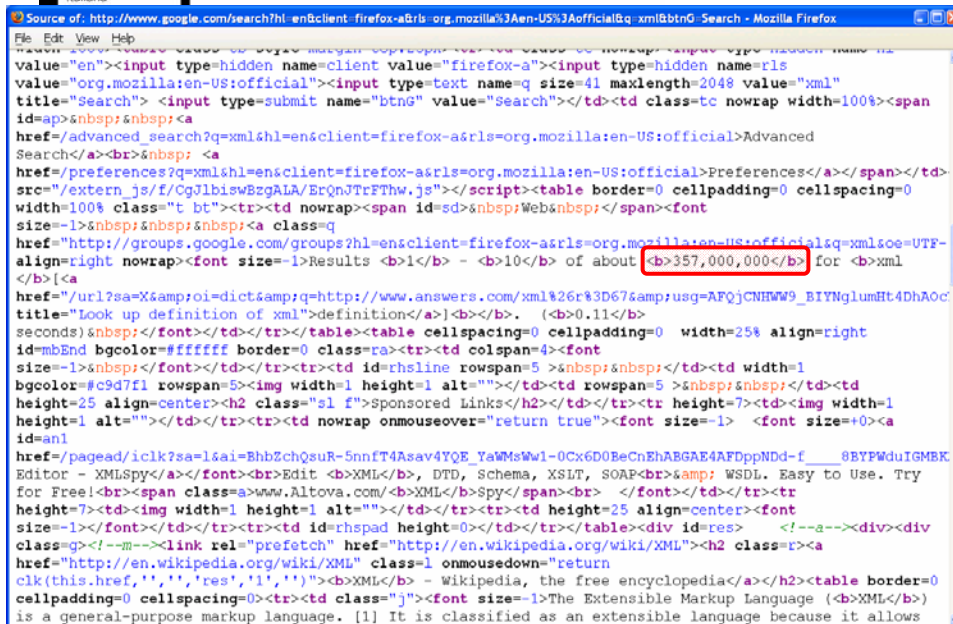


7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

5

HTML Scraping



XML Example

- Would not be better if Google would return the search results in a format that would make it easier for a program to figure it out?

```
<search keyword="xml ">
  <results total="357,000,000">
    <site url="www.xml.com">XML.com</site>
  </results>
</search>
```

- Represent and structure the data using the “most appropriate” markup tags so that we can give it some conventional semantic meaning and write programs that can “understand” it.

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

7

HTML not enough

- Forging syntax check
 - No formal validation “at compile-time” of content correctness
 - Difficult to parse from your own programs (need to handle lots of garbage)
- Fixed set of tags – limited extensibility:
 - Users cannot customize the syntax
 - Extended by browser makers with lots of interoperability problems
- Predefined semantics of tags
 - Documents written in the HyperText Markup Language are very good for representing HyperText pages, but very bad for everything else!

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

8

XHTML vs. XML

- XHTML – hypertext markup language that follows the XML syntax rules
- XHTML Tags describe **Web pages** (hypertext nodes)
- XHTML Documents displayed by browsers for people to read
- XML – generic *meta-language* that defines syntax rules for a whole family of markup languages
- Can be used to describe the structure and markup any type of content
- XML documents typically processed automatically by software, not read by people

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

9

XML Syntax

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

10

XML Syntax

- Elements (enclosed in tags)
 - Attributes (of elements)
 - Text (inside elements)
 - Namespace Declarations (special attributes)
 - Prolog `<?xml version="1.0" encoding="UTF-8" ?>`
 - Entities (&) Character Data (CDATA)
 - Comments `<!-- -->` (à la SGML)
 - Processing Instructions `<? ?>` (SGML legacy)
-
- Note: **Order** of elements in document is preserved, order of attributes within element is not

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

11

XML Elements Syntax

- Markup Elements enclosed in Tags
`<element>...</element>`
`<element/>` (empty element tag)
 - Elements tags can be nested and mixed with text
`<element>...`
`<subElement>...</subElement>...`
`</element>`
- Warning: make sure nesting is well formed!
`</element>_{<element>}`

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

12

XML Attribute Syntax

- Element attributes

```
<el ement attribute="...">...</el ement>
```

```
<el ement attribute="..." />
```

- Do we really need attributes?

```
<el ement>
  <attri bute>...</attri bute>
</el ement>
```

Attributes represent “hidden” data about the element, while element text is “visible”

- Warning: Attribute names must be unique!

```
<el ement a=" 1" a=" ABC" ></el ement>
```

7.11.2007

Fall Semester 2007

Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

13

XML Entities

- Special characters can be encoded (or escaped) using character entities
- In general, any Unicode character can be entered with the notation:

Character	Entity
<	&l t;
>	>
&	&
"	"
'	'

&#N; (N is decimal)

&#xM; (M is hex)

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

14

XML Prolog Syntax

- Header found at the beginning of all XML documents:

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

- Specify the encoding usually Unicode (UTF-8, UTF-16) or ISO- 8859-1
- Also XML version “1.1” exists, not yet widely used.

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

15

XML CDATA

- **Character Data** Sections are used to stop parsing the XML markup and treat the enclosed data as is:

```
<element>
```

```
  <![CDATA[<not An Element>]]>
```

```
</element>
```

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

16

XML Whitespace

- Should whitespace (spc, tabs, cr, lf) be ignored?
- It depends:
 - Data-oriented XML – ignore whitespace (more compact documents, faster processing, less readability)
 - Document-oriented XML – whitespace is part of the content and must be preserved:

```
<chapter xml:space="preserve">
  <section name="Introduction">
    Once upon a time...
  </section>
</chapter>
```

xml:space="default"
Application decides
how to handle it

7.11.2007

Software

17

XML Namespaces: Problem

- Problem:
 - What happens if two XML documents are pasted together?
 - How to distinguish “tag vocabularies” of different data sources?
- Solution:
 - Use **namespaces** to qualify and distinguish the elements of an XML subtree so that tag name collisions can be avoided.

```
<col or>
  <bl ue/>
  <orange/>
  <red/>
  <bl ack/>
  <whi te/>
</col or>

<frui t>
  <appl e/>
  <orange/>
  <ananas/>
  <banana/>
</frui t>
```

Example: mix XHTML tags
with your own XML data
to store “rich text” descriptions

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

18

Working with Namespaces

```
<ns:element
  xml ns:ns="URI" ...>...
</ns:element>
```

- Declare a namespace prefix **ns** uniquely identified by **URI** using the special attribute **xml ns**
- The namespace declaration scope contains all attributes and all children elements

```
<mix xmlns:color="URI 1"
      xmlns:fruit="URI 2">
  <color:orange/>
  <fruit:orange/>
</mix>
```

- Prefixes are used to avoid repeating long URI for each element/attribute
- Elements without prefix belong to the **default** namespace (declared with **xml ns="Default URI"**)

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

19

XML Document Correctness

Well formed documents

- Verify the basic XML **syntax** constraints
- Only well formed documents can be parsed

Valid documents

- Well formed documents that satisfy additional **structural** constraints defined in a “schema”
- Non valid documents can still be processed even if no guarantees can be made on their “meaning”

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

20

Summary: Well Formed XML Syntax Rules

1. XML is case sensitive
2. Single Root Element for a document
3. Elements enclosed in angle brackets `</>`
`<el ement></el ement>` or `<el ement />`
4. Correct element nesting (XML tree)
5. Attributes must have unique names
6. Attribute values must be "quoted"
7. Document begins with XML prolog `<?xml ...?>`
8. Namespaces prefixes must be declared before use

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

21

XML Structure Document Type Definition (DTD) XML Schema

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

22

XML Structure

- Easy to start writing XML documents, simply make it well formed and make up the tags as you go.
- Challenges:
 - How to ensure multiple documents use the same set of tags and have the same structure?
 - How to write an application if we cannot be sure of which tags are used in the documents?
 - How to agree in advance on a specific XML format for exchanging a document between programs?
- Solution: use a “schema” to restrict and constrain the tags and the structure of an XML document

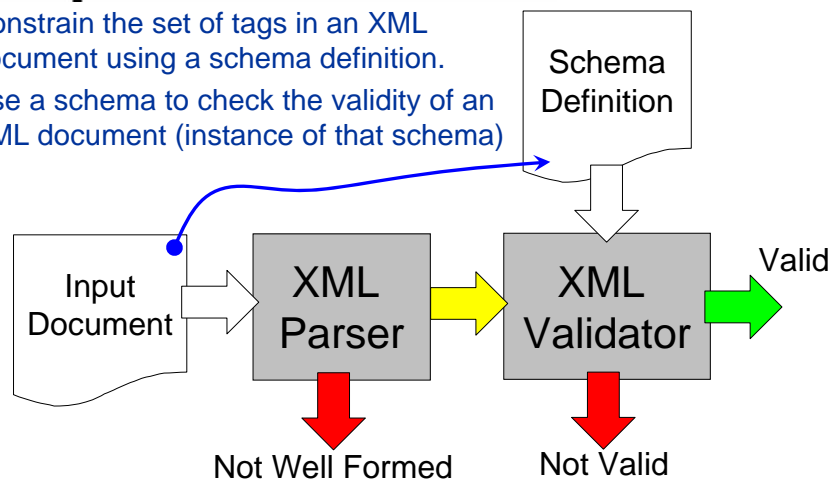
7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

23

Schemas

- Constrain the set of tags in an XML document using a schema definition.
- Use a schema to check the validity of an XML document (instance of that schema)



7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

24

Defining XML Structure

- Different schema languages are available to specify new XML-based languages.
- **Document Type Definition (DTD)** – part of the XML 1.0 specification. A relic from SGML, it does not use XML syntax, nor it supports namespaces (which were invented later)
To address some of the limitations of DTD, other languages for constraining the structure of XML documents have been introduced:
- **XML Schema** – the new XML type system from W3C. Uses XML syntax. Very complex data modeling language.
- **RelaxNG** – a simpler schema language from OASIS with both XML syntax and a simplified notation.

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

25

Document Type Definitions DTD

- The Document Type Definition (DTD) specifies the tree structure of XML documents:
 - Nesting of Elements
 - Attributes of Elements
 - Values of Attributes
- A DTD contains a set of definitions of:
 1. Element Type
 2. Attribute List

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

26

Element Type Definition

`<! ELEMENT name content >`

The element with tag name contains certain elements, text, or a combination thereof.

`<! ELEMENT name EMPTY >`

The element name is empty: `<name/>`

`<! ELEMENT name (#PCDATA) >`

The element name only contains character data:

`<name>Text without other elements</name>`

`<! ELEMENT name (mix) >`

The element name can contain a mixture of elements and text:

`<name>Text with <other/> elements</name>`

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

27

Choice Element Type Definition

```
<! ELEMENT name (first | last) >
  <! ELEMENT first (#PCDATA) >
  <! ELEMENT last (#PCDATA) >
```

The element name can contain a **choice** between two child elements: **first** xor **last**:

`<name><first>Cesare</first></name>`

`<name><last>Pautasso</last></name>`

Invalid (name contains both elements)

`<name><last>Pautasso</last><first>Cesare</first></name>`

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

28

Sequence Element Type Definition

```
<! ELEMENT book (title, author, price)>
  <! ELEMENT title (#PCDATA)>
  <! ELEMENT author (first | last)>
  <! ELEMENT price (#PCDATA)>
```

The element `book` can contain a **sequence** of `title`, `author`, `price` elements (found in this order)

```
<book> <title>XML for dummies</title>
  <author>
    <first>Ed</first>
    <last> Tittle</last>
  </author>
  <price>$0.68</price>
</book>
```

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

29

Repeating elements (cardinality)

```
<! ELEMENT book (title, author+, edition*)>
  <! ELEMENT title (#PCDATA)>
  <! ELEMENT author (first | last)>
  <! ELEMENT price (#PCDATA)>
  <! ELEMENT edition (#PCDATA | year? | price?)>
```

The element `book` can contain a sequence of **exactly one** `title`, **at least one** `author`, an **optional** `price`, and **zero or more** `edition` elements (which may include **optional** `year` and `price` elements).

```
<book>
  <title>XML for dummies</title>
  <author>Ed Tittle</author>
  <author>Ramesh Chandak</author>
  <edition>Paperback <year>1998</year></edition>
  <edition>Hardcover <price>$0.68</price></edition>
</book>
```

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

30

Element Type Definition Summary

`<! ELEMENT name content >`

content	
ANY	any contents
EMPTY	<name/>
name	an element
,	Sequence (Concatenation) of elements
	Choice (Any Order) between elements
?	Optional
*	Zero or More
+	One or More

If no cardinality is specified, then the child element **must** be present only **once**

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

31

Attribute List Definition

`<! ELEMENT element content >`

`<! ATTLIST element attribute type options >`

Define the **attribute** for an element. Attributes have types.
Options are used to control their presence and default values.

`<! ATTLIST book number ID #REQUIRED >`

`<! ATTLIST price currency CDATA CHF >`

`<! ATTLIST author prefix (Prof|Dr) #IMPLIED >`

`<! ATTLIST edition number CDATA #REQUIRED >`

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

32

Attribute Types

<! ATTLIST element attribute type options>

Attribute type	
CDATA	Character Data (String)
ID	Unique value within document to identify element
IDREF	Reference to other element ID
IDREFS	Multiple space-separated Ref.
(value ...)	Enumeration

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

33

Attribute Default Options

<! ATTLIST element attribute type options>

Attribute options	
#REQUIRED	Attribute must be present
#IMPLIED	Optional attribute w/o default
value	Optional attribute w/ default
#FIXED value	Set value no matter if attribute is present or not

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

34

Local embedded definitions:

```
<!DOCTYPE root [definitions]>
```

External definitions:

Use this one!

```
<!DOCTYPE root SYSTEM "url/filename">
```

```
<!DOCTYPE root PUBLIC "FPI">
```

```
<!DOCTYPE root PUBLIC "FPI" "url">
```

(FPI = Formal Public Identifier)

Example:

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN"
"http://www.w3.org/TR/html4/strict.dtd">
```

Yes, but DTDs are about defining **document** structure and not modeling data types!

- Limited data types for element content and attributes
 - only CDATA, ID/IDREF or Enumeration
- Does not support Namespaces
- Definitions not specified using XML syntax
- References too simple
 - Cannot constrain ID uniqueness within a document subtree
 - Cannot have primary keys across multiple attributes
 - Does not support referential integrity
- Context insensitive definitions
 - Cannot constrain structure based on parents of an element
 - Cannot have elements present if an attribute has a certain value or if also another element is present (co-occurrence)

DTD vs. XML Schema

- ~~Limited data types for element content and attributes~~
 - ~~only CDATA, ID/IDREF or Enumeration~~
 - ~~Does not support Namespaces~~
 - ~~Definitions not specified using XML syntax~~
 - ~~References too simple~~
 - Cannot constrain ID uniqueness within a document subtree
 - Cannot have primary keys across multiple attributes
 - ~~Does not support referential integrity~~
 - Context insensitive definitions
 - Cannot constrain structure based on parents of an element
 - Cannot have elements present if an attribute has a certain value or if also another element is present (co-occurrence)
- Rich built-in XML Data Types (44)
- Namespace Support
- XML Syntax
- Key+References

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

37

XML Schema (XSD)

- XML language to define the structure of XML documents. W3C standard from 2001.
- Schemas are composed of:
 - Data Type Definitions
 - Simple Types - Restrictions
 - Complex Types
 - Constructors: Sequence, All (=Choice)
 - Cardinality (minOccurs, maxOccurs)
 - Element/Attribute Declarations (of type)
 - Key Identifiers

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

38

XML Schema Example

```
<xs:element name="person">
  <xs:complexType>
    <xs:sequence>
      <xs:element ref="name"/>
      <xs:element ref="email" minOccurs="0" maxOccurs="unbounded"/>
      <xs:element ref="url" minOccurs="0" maxOccurs="unbounded"/>
      <xs:element ref="link" minOccurs="0" maxOccurs="1"/>
    </xs:sequence>
    <xs:attribute name="id" type="xs:ID" use="required"/>
    <xs:attribute name="note" type="xs:string"/>
    <xs:attribute name="contract" default="false">
      <xs:simpleType>
        <xs:restriction base="xs:string">
          <xs:enumeration value="true"/>
          <xs:enumeration value="false"/>
        </xs:restriction>
      </xs:simpleType>
    </xs:attribute>
    <xs:attribute name="salary" type="xs:integer"/>
  </xs:complexType>
</xs:element>
```



7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

XML Schema Example Complex Types/Element

```
<xs:element name="person">
  <xs:complexType>
    <xs:sequence>
      <xs:element ref="name"/>
      <xs:element ref="email"
        minOccurs="0" maxOccurs="unbounded"/>
      <xs:element ref="url"
        minOccurs="0" maxOccurs="unbounded"/>
      <xs:element ref="link"
        minOccurs="0" maxOccurs="1"/>
    </xs:sequence>
  </xs:complexType>
</xs:element>
```

Anonymous Type Definition

<p>minOccurs="0" maxOccurs="unbounded"/> minOccurs="0" maxOccurs="unbounded"/> minOccurs="0" maxOccurs="1"/></p> <p>Cardinality default (minOccurs="1", maxOccurs="1")</p>
--

sequence (all, in order)

all (all, in any order)

choice (only one)

any

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

40

XML Schema Example Simple Types/Attribute

```
<xs:element name="person">
  <xs:complexType>
    <xs:attribute name="id" type="xs:ID" use="required"/>
    <xs:attribute name="note" type="xs:string"/>
    <xs:attribute name="contr" default="false">
      <xs:simpleType>
        <xs:restriction base="xs:string">
          <xs:enumeration value="true"/>
          <xs:enumeration value="false"/>
        </xs:restriction>
      </xs:simpleType>
    </xs:attribute>
    <xs:attribute name="salary" type="xs:integer"/>
  </xs:complexType>
</xs:element>
```

Custom Defined Type

Built-in Types (ID, integer, string, ...)

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

41

Link XML Document to XSD

- Add a pointer in the root element of the XML document

```
<root xmlns="namespace"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="namespace url"
/>
```

- Configure namespaces in the schema

```
<schema xmlns="http://www.w3.org/2001/XMLSchema"
  targetNamespace="namespace"
/>
```

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

42

XML Critique

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

43

The X factor

*If I would invent another
programming language, its name
will contain the letter X*

Niklaus Wirth,
2001

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

44

XML as the Universal Syntax

The Extensible Markup Language is used in all aspects of a Web application:

- Data Representation
 - regular structure (like in database management)
 - semi-structured
 - unstructured textual documents (e.g., XHTML pages)
- Meta-Data
 - Define rules on how the data should be structured
- Code
 - Specify how data should be processed/transformed

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

45

Text vs. Structured Data

Text

The Conductor says the
Cisalpino train is arriving in
Lugano in 5 minutes

```
<message from="Conductor">
The Cisalpino train will arrive in
Lugano in 5 minutes
</message>
```

XML

```
<message from="Conductor">
  <train>Cisalpino</train>
  <event>arrival</event>
  <station>Lugano</station>
  <time unit="minutes">5</time>
</message>
```

Data-oriented
XML

Document
-oriented
XML

```
<message from="Conductor">
The <train>Cisalpino</train> will arrive
in <station>Lugano</station>
in <time>5 <unit>minutes</unit></time>
</message>
```

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

46

Where is XML used?



- **More than the Web (XHTML)**
- Blogs (Rich Site Summary – RSS Feeds) 
- Document Markup Formats
 - Office Documents (OOXML/ODF, WordML) 
 - Wireless Mobile Applications (WML)
 - Chemistry (CML)
 - Theology (ThML)
 - Music (MusicML)
 - Speech (VoiceML) 
 - Graphics (Scalable Vector Graphics – SVG) 
 - Many more!
- Meta-Data (XML Schema, WSDL, Resource Description Framework – RDF, Adobe Extensible Metadata Platform – XMP) 
- Logs (Common Base Events – CBE Format)
- Configuration Files (J2EE Deployment Descriptors, ANT build scripts) 
- Communication Protocols (Web Services)
- Databases

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso



47

XML Benefits

1. Text-based (Unicode) format
 - Human readable
 - Machine readable
2. W3C Standard
 - *Independent* of Platform, OS, Vendor, Programming Language, Communication Protocol
3. Easy to start writing well formed XML documents
4. Rich Technology Tool-chain
 - More than general data representation format:
choose XML and you can reuse a rich family of data management and processing technologies

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

48

XML Limitations

1. Text-based (Unicode) format
 - **Verbose** compared to a binary or simpler textual format (like JSON)
2. Not a graph, only a tree
 - The notion of reference not well integrated
3. “Self-Describing” Data
 - XML is about syntax (well-formed) and structure (valid, according to a schema)
 - Tags do not have any implicit semantics: it depends on application interpretation

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

49

XML Technology Landscape

- Data Representation
 - XML Syntax
 - XML Information Set (InfoSet)
 - XML Namespaces
 - XML Schema, DTD
 - XLink, XPointer
- Data Processing
 - XPath
 - XSLT – Extensible Stylesheet Transformation Language
 - XQuery
 - XUpdate
- Data Processing API
 - DOM
 - SAX
 - JAXP
- Communication Protocols
 - XML Forms
 - XML Web Services (SOAP, WSDL, UDDI)
 - XML Encryption
 - XML Digital Signature

7.11.2007

Fall Semester 2007
Software Atelier III – Web Development Lab
©2007 Cesare Pautasso

50

References

- Anders Moller and Michael Schwartzbach, **An Introduction to XML and Web Technologies**, Addison-Wesley, 2006
- Elliotte Rusty Harold and W. Scott Means, **XML in a Nutshell**, O'Reilly, 3rd Ed. 2004
- John Bosak, **XML Ubiquity and the Scholarly Community**, [Computers and the Humanities](#), Volume 33, Numbers 1-2, April 1999 , pp. 199-206(8)