# eXtensible Markup Language (XML)
## Basic Concepts

**Giuseppe Della Penna**

Università degli Studi di L'Aquila

*dellapenna@univaq.it*

*http://www.di.univaq.it/gdellape*

University of L'Aquila
Computer Science Department

# Notes to the English Version

*These slides contain an English translation of the didactic material used in the Web Engineering course at University of L'Aquila, Italy.*

*The slides were initially written in Italian, and the current translation is the first result of a long and complex adaptation work.*
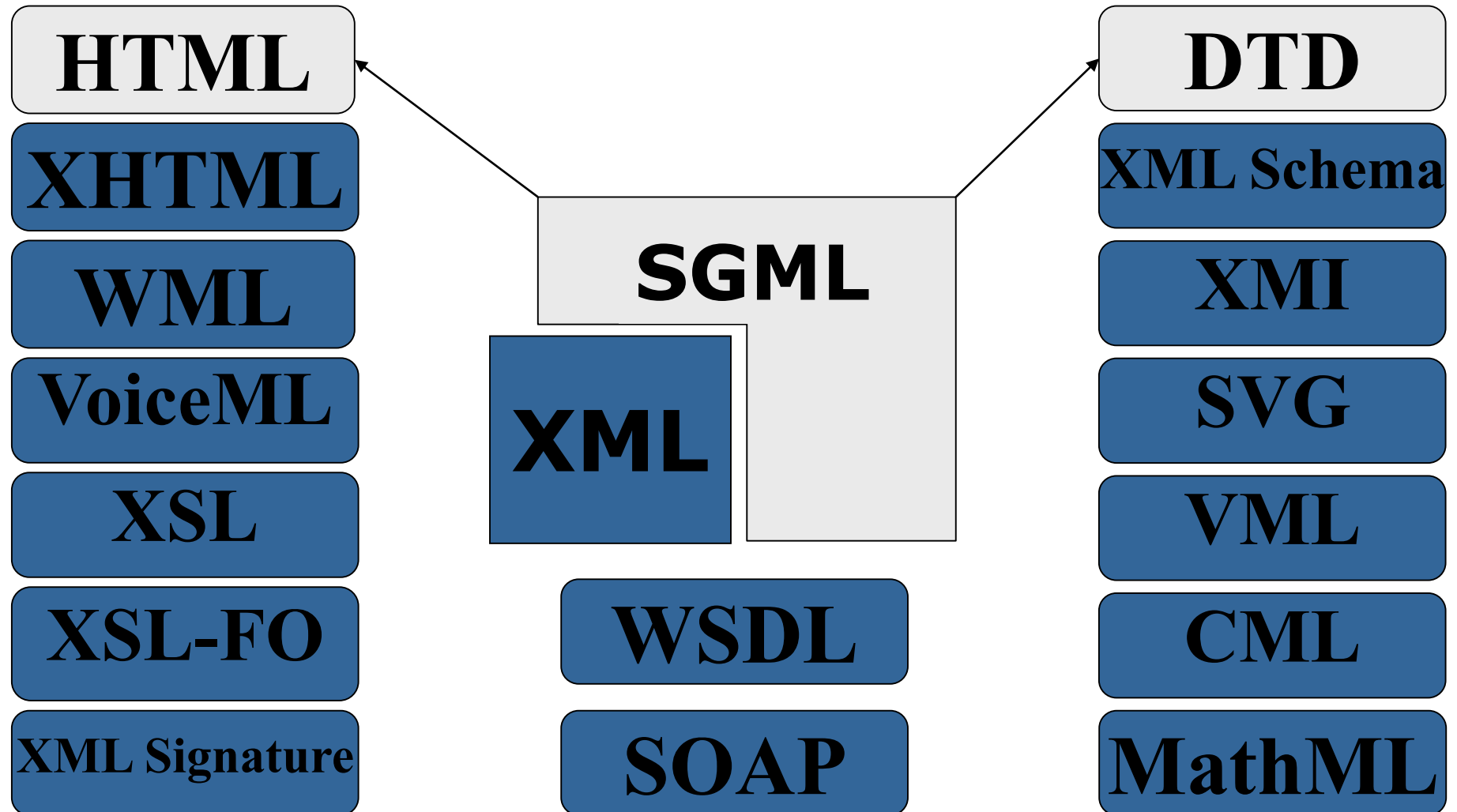
*Therefore, the slides may still contain some errors, typos and poorly readable statements.*

*I'll do my best to refine the language, but it takes time.*

*Suggestions are always appreciated!*

# Origins of XML

- XML is a *metalanguage*, i.e., a language that is used to create other languages.

- In particular, XML defines the basic rules for creating **markup languages,** i.e., languages whose content (text) is structured by special delimiters called *tags*.

- XML derives from **SGML,** another well metalanguage used mainly in professional settings (e.g., publishing).

- Compared to SGML, XML has been greatly **simplified** and small **extensions** have been added to make it more user friendly.

# The Family of XML

# XML Pros

- XML allows developers to easily create **ad-hoc languages** to contain **structured information.**

- XML is completely **text-based,** so it is humanreadable and can be easily hand-edited. It Supports UNICODE, so it is suitable for all types of languages.

- The structures defined with XML are useful for creating platform-independent and self-descriptive **data structures**.

- The **automatic processing** of an XML language is particularly simple and efficient. The strict syntactic rules of XML-based languages make them very suitable for automatic processing.

- Since XML is actually written as plain text, XML data can be easily and safely **transported using the HTTP protocol** through firewalls (SOAP, web services).

# XML Cons

- XML documents, because of their textual structure and tags, tend to be much more large than the corresponding binary format.

- XML manipulation libraries are not as fast as the ad-hoc parsers written for specific formats, especially the binary ones.

- In general, therefore, the use of XML is more expensive in terms of necessary resources (network, memory and CPU time required for decoding it, etc..)

# XML Applications

- Despite the (few) disadvantages seen, the use of XML is widespread and growing:
  - **Web Services**
    - SOAP, WSDL, ...
  - **Science**
    - MathML, CML,...
  - **Web and Publishing**
    - XHTML, WML, VoiceML, XSL, XSL-FO, ...
  - **Multimedia**
    - SMIL, SVG,...
  - **Definition of formal structures**
    - XMLSchema, XMI,...
  - **Security**
    - XML Encryption, XML Signature

University of L'Aquila
Computer Science Department

# An XML Document

```xml
<?xml version="1.0" encoding="iso-8859-1"?>
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN"
  "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd">
<html xmlns="http://www.w3.org/1999/xhtml">
    <head>
        <meta http-equiv="Content-Language" content="it"/>
        <title>Sito Ufficiale dei Corsi di Laurea  in Informatica – Universit&agrave; dell'Aquila ::
        <?php echo $pageTitle; ?>
        </title>
        <link rel="stylesheet" media="print" href="css/stile_stampa.css" type="text/css"/>
        <link rel="stylesheet" media="screen" href="css/stile_grafico.css" type="text/css" title="deep blue"/>
        <!--[if lte IE 6]>
        <link rel="stylesheet" media="all" type="text/css" href="css/ie6_hacks.css" />
        <![endif]-->
        <link rel="SHORTCUT ICON" href="favicon.ico" type="image/x-icon"/>
        <script type="text/javascript">
        //<![CDATA[
        ...
        //]]></script>
        <link rel="alternate" type="application/rss+xml" title="RSS Feed" href="..."/>
    </head>
    <body>
    ...
    </body>
</html>
```

# The Structure of an XML document

- An XML document consists of a *prologue* and a *body*
- The body of the document may contain:
  - **text,**
  - **tags** (element delimiters),
  - **annotations** (comments)
  - **processing instructions** (instructions for external automatic processors)
  - **entities** (similar to macros)
  - In addition, tags may contain **attributes** and **namespaces.**

# Prologue: XML Declaration

<?xml version="1.0" encoding="ISO-8859-1"?>

- The first line of the prologue is the **XML declaration,** which is *mandatory* and must appear at *the very beginning of the document*.
- The expression "<?xml" is called **the opening tag** of the XML declaration. The statement is closed by the symbol "?>".
- Within the statement, there are two expressions of the form **name = "value".** This type of notation is used to define an **attribute** contained in the tag. An attribute refines or extends the meaning of a tag, and it is widely used in XML.
- The attributes of the XML declaration are:
  - **version:** (required) indicates the version of XML used.
  - **encoding:** (optional) is the name of the character encoding used in the document (default: UTF-8 or 16, that is, 8 or 16-bit Unicode, ISO-8859-1 is the most suitable for western European characters)
  - **Standalone:** (optional) if true *yes* indicates that the file does not refer to other external files. (Default: *no*)

# Prologue: DOCTYPE Declaration

- XML documents can (and should) be associated to **a formal specification that defines the language used** in the document and its syntax rules.
- The default XML way to create this specification is the *document type definition* **(DTD)**
- If a document has an associated DTD, you must include a DOCTYPE declaration in the prologue that declares the association. This statement inherits the syntax of the corresponding SGML one.
- However, there are other formalisms for the definition of XML languages, such as *schemas,* which use different association methods.
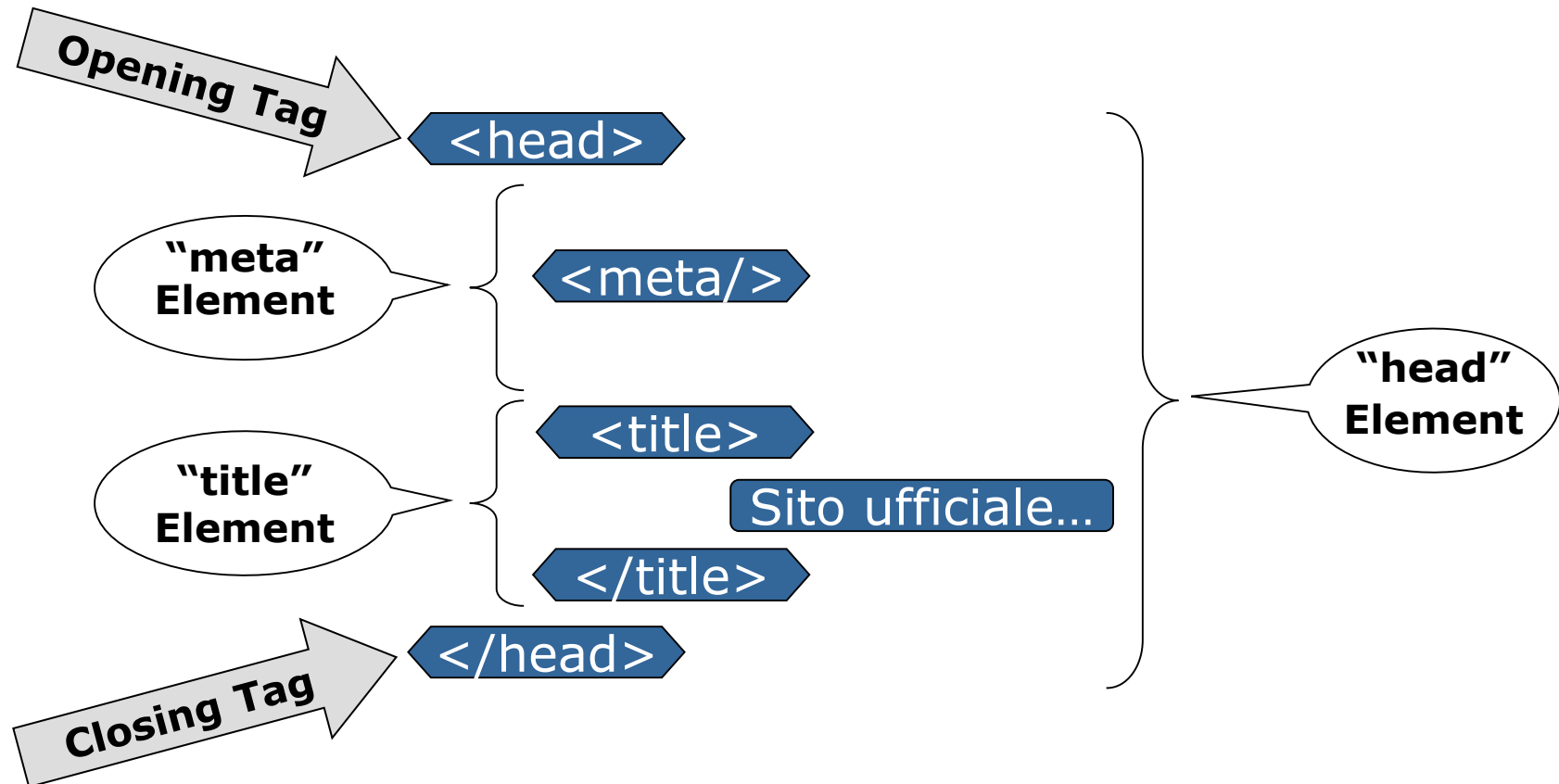
# Prologue: DOCTYPE Declaration

<!DOCTYPE *RootElement ExternalDTDReference* [*InternalDTDSubset* ]>

- The declaration begins with the tag "<!DOCTYPE" and is closed by the symbol ">". Inside there are the following:
- **RootElement** *(mandatory)* is the document root element name, i.e., the name of the tag that will contain the entire document.
- **ExternalDTDReference** *(optional)* points to a file that contains the DTD itself, and may be:
  - **SYSTEM** *"uri"*, an *uri* which identifies an external file.
  - **PUBLIC** *"pubid" "uri,"* where *pubid* is a unique identifier for the DTD and *uri* points to a file that contains it.
- **InternalDTDSubset** *(optional)* is a DTD, or a DTD fragment, which can be specified directly within the document.

# Elements

- **Elements** are the base of the structure of XML documents.
- An element is a **piece of data,** *limited and identified* (by name) by a *tag*.
- The content of an element is anything that appears between its opening tag and its closing tag.
- Elements can be nested, i.e., elements may be part of the contents of an outer element.

# Elements



Opening Tag

`<head>`

"meta" Element

`<meta/>`

"title" Element

`<title>`

Sito ufficiale…

`</title>`

`</head>`

Closing Tag

"head" Element

# Elements: Basic Rules

- Element names are **case-sensitive.**
- **Each element must be closed,** that is, its closing tag must appear before the end of the document.
- In the case of nested elements, **end tags must appear in reverse order of opening,** i.e., the element contents cannot "overlap".
- Every XML document must have **a unique "root" element,** where all the others are nested.

# Elements: Syntax

**①** *<name>*

  …

**②** *</name>*

**③** *<name/>*
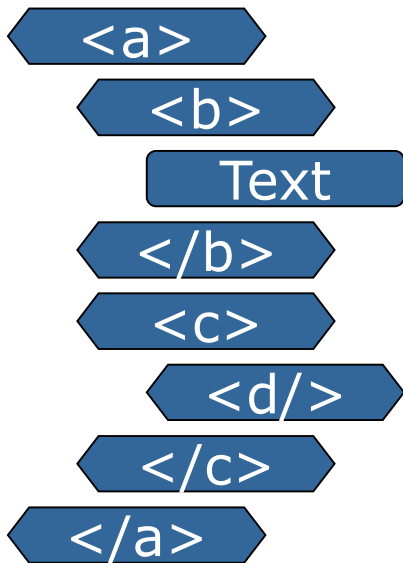
```
<a>
  <b/>
  <c>
    Text
    <d>
      <e/>
    </d>
  </c>
</a>
```

- The opening tag of an element has the form shown in (1), where *name* is the name of the element. The corresponding closing tag is shown in (2)
- Finally, some elements may be empty: in this case you can omit the closing tag writing the opening in the abbreviated form shown in (3).

University of L'Aquila
Computer Science Department

# Hierarchy of Elements

<a>
<b>
Text
</b>
</a>

- Nested elements create the **tree structure** of XML documents.
- Within this structure it is possible to define some useful "relationships":
    - a is the root node
    - b and c are the children of a, the text is the child of b, d is the child of c
    - c is the father of d, b is the father of the text, a is the father of b and c
    - b and c are brothers
    - b, c, d and the text are descendants of a, d is a descendant of c, the text is a descendant of b
    - a is a ancestor of b, c, d and text, b is an ancestor of the text, c is an ancestor of d.

# Attributes

- Attributes allow you to specify **properties of the elements,** modifying or better defining their meaning.

- Attributes are **inserted within the opening tag** of the elements.

- The **order** in which attributes appear in the opening tag is not significant.

- The value of an attribute should be simple: otherwise it is better to use a nested element to contain it.

# Attributes: Basic Rules

- Attribute names are **case-sensitive.**
- The same element cannot contain two attributes with the same name.
- Attributes with no value (only name) are not allowed.
- The attribute value must be specified **between single or double quotes.**
- The value may contain **entity references,** but no other structure (XML elements, processing instructions, etc.).

# Attributes: Syntax

**①** *<name attribute="value">*

**②** *<name attr1="val1" attr2="val2">*

**③** *<name attribute=' "value" '>*

```
<a x="txt" y="2">
 <c> Text
  <d>
   <e z="abc123"/>
  </d>
 </c>
</a>
```

- The basic syntax for an attribute inserted in the opening tag of an element is shown in (1)

- To specify multiple attributes it is enough to separate them with one or more spaces as shown in (2)

- To include quotation marks in an attribute value, you must use quate different from the one used to surround the value itself (3)

# Namespaces

- Namespaces are used to **declare membership of elements and attributes to a particular XML language,** providing a semantics.

- They are particularly useful if **multiple languages are mixed in the same document,** with possible name collisions.

- The **namespace declarations** are inserted in the opening tags, similar to an attribute, and are valid within the element and its contents.

# Namespaces: Syntax

**①**  *<name* xmlns:*prf="uri">*

**②**  *<name* xmlns=*"uri">*

**③**  *<name xmlns="uri" xmlns:prf="uri">*

- The explicit namespace declaration (1), inserted into an opening tag, indicates that all the elements whose name is prefixed by *"prf"* (namespace prefix) will be considered as belonging to the namespace identified by *uri*.

- The special standard namespace declaration (2) indicates the namespace of all the elements with no explicit namespace prefix.

- In each element it is possible to declare multiple explicit namespace prefixes, but only ane default namespace (3)

- *URI used in such declarations are only conventional identifiers associated with different namespaces, and do not point to any particular internet resource.*

# Namespaces: Examples

```
<a xmlns="ns1" xmlns:html="ns2">
  <b/>
  <html:p><html:b>testo</html:b>
    <c xmlns="ns3"><d/></c>
    <d/>
    <e xmlns:xsl="ns4" xsl:attr="val">
      <xsl:f>testo</xsl:f>
    </e>
  </html:p>
</a>
```

- Understanding namespaces is important in order to manage complex XML documents and their semantics.

- In this example:

- The namespace "ns1" contains the elements a, b, d, e.

- The namespace "ns2" contains elements html:p, html:b.

- The namespace "ns3" contains the elements c and d.

- The namespace "ns4" contains the xsl:attr attribute and the xsl:f element

- Note that there are two elements *d* in the document, belonging to different namepsaces!

# Entities

- **In XML parlance documents are composed by a set of *entities*.**
    - Each character is a *character entity*, each tag is an entity and the document itself is an entity.
- **Each entity, except for the document and the external DTD, has a name.**
- **The entities are divided into *parsed* and *unparsed*:**
    - Each *parsed* entity has a corresponding textual value. The XML parser replaces the entity with its value when it parses the document.
    - An *unparsed* entity, however, is not replaced by the parser, and can have even a binary value, accessible via the *notations*.

# Entities (parsed): Syntax

**①** &name;

**②** &#number;

**③** &#x*number;*

&gt;      →      >
&lt;      →      <
&quot;   →      "
&amp;   →      &
&#32;   →      [space]
&#x20;  →      [space]

- **General entities,** which can represent any string, are defined in the DTD and the XML document refers to them using the syntax (1), where *name* is the name of the entity.

- **Character entities,** which represent single UNICODE characters, are referred with the syntax (2), where *number* is the decimal code for the Unicode character, or with the syntax (3), where *number* is the hexadecimal code for the Unicode character.

# Entities: Use

- Parsed entities are a handy way to insert strings in the document referring to an external definition, instead of writing them explicitly.

- They are useful if there are **characters that cannot directly** typed, or to **expand strings used frequently,** or to write **characters that are not explicitly allowed in a context,** such as quatation marks or the '<' and '>' symbols.

# Text

- The text that can be inserted in XML documents **includes all the characters defined in UNICODE.**

- You can insert special or reserved characters using **character entities.**

- You can insert predefined strings using **general entities.**

- *<u>You can not explicitly use</u>* the characters '>', '<' and '&', for which you should always use the corresponding character entities.

# CDATA sections

```
<![CDATA[
<< &goofy;
Text only!<
>>
]]>
```

- CDATA sections explicitly define **areas where there is only text.**

- within CDATA sections the parser does not look for elements, attributes, entities, and other XML structures

- The opening tag of a CDATA section is the string "<![CDATA[", while the closing tag of "]]>", which obviously can not appear in the content.

# Processing Instructions

<?*target  data* ?>

- The Processing Instructions (PIs) are used to **pass extra information to programs that manipulate the XML file** and can appear anywhere in the document.
- The general form of a PI has an opening tag like £<? t*arget*" where *target* identifies which application will process the instruction, and a closing tag "?>". Note that the XML declaration is nothing more than a processing instruction!
- Inside the tag you can write any type of textual data. The only rule is that the data cannot contain the sequence "?>". The two examples below are respectively (1) the PI that associates an XSL style sheet to a document and (2) a PHP script.

**<?xml-stylesheet type="text/xsl" href="sms_pdf.xslt"?>**
**<?php echo "hello" ?>**

# Comments

<!-- This is a XML (and SGML) comment -->

- Comments are useful to humans, and are ignored by XML manipulation programs.

- Comments may appear anywhere except within the value of an attribute.

- The comments follow the syntax of SGML, and are identical to those used, for example, in HTML.

- The opening tag of a comment is the sequence "<!--", and the closing tag is the sequence "-->"

- The content of the comment is generic text, which should not contain the closing sequence.

# Validation of XML Documents

- **An XML document is *well formed*** if it respects the general syntax rules seen in the previous slides.
- **An XML document is** well formed and *valid* if it meets the syntactic and semantic rules contained in the associated DTD. **A document with no DTD is never valid.**
- There are **validating** and **not validating parsers.** The latter may ignore any DTD, except for the definition of general entities.

# References

- **XML specification from the W3C**
  http://www.w3c.org/TR/XML//