

INTERNET & WEB APPLICATION DEVELOPMENT SWE 444

Fall Semester 2008-2009 (081)

Module 4 (V): Document Type Definitions (DTDs)

Dr. El-Sayed El-Alfy

Computer Science Department King Fahd University of Petroleum and Minerals alfy@kfupm.edu.sa

Objectives/Outline

- Objectives
 - Understand the role of DTDs
 - Learn how to write and use DTDs
- Outline
 - What are DTDs?
 - Why DTDs?
 - DTD Syntactic Elements
 - ELEMENT
 - ATTRIBUTE
 - ENTITY
 - Types
 - Examples
 - Validation

KFUPM-081© Dr. El-Alfy

SWE 444 Internet & Web Application Development

What is DTD?

- Document Type Definition (DTD) is a grammar that describes the structure of a class of XML documents
 - structure of the documents is described via
 - element and attribute-list declarations.
- Element declarations
 - oname the allowable set of elements within the document, and
 - specify whether and how declared elements and runs of character data may be contained within each element.
- Attribute-list declarations
 - name the allowable set of attributes for each declared element, including the type of each attribute value, if not an explicit set of valid value(s).
- DTDs are written in EBNF-like notation
 - EBNF (Extended Backus-Naur Form) used to express contextfree grammars, i.e. a formal way to describe computer programming languages and other formal languages

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.3

Why DTDs?

- XML documents are designed to be processed by computer programs
 - If you can put just any tags in an XML document, it's very hard to write a program that knows how to process the tags
 - A DTD specifies what tags may occur, when they may occur, and what attributes they may (or must) have
- Applications can use a standard DTD to verify that the data you receive from the outside world is valid
- With a DTD, independent groups of people can agree to use a standard DTD for interchanging data
 - produce consistent XML documents
- You can also use a DTD to verify your own data

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Parsers

- An XML parser is an API that reads the content of an XML document
 - Currently popular APIs are DOM (Document Object Model) and SAX (Simple API for XML)
- A validating parser is an XML parser that compares the XML document to a DTD and reports any errors

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.5

Building Blocks of XML Documents

- From a DTD point of view, all XML documents (and HTML documents) are made up by the following building blocks:
 - Elements -- the main building blocks of both XML and HTML documents
 - Attributes -- provide extra information about elements, placed inside the opening tag of an element, come in namevalue pairs
 - Entities -- are a way to include fixed text
 - PCDATA Parsed Character DATA (the text that WILL be parsed by a parser; tags inside the text will be treated as markup and entities will be expanded)
 - CDATA -- Character DATA (i.e. the text found between the start tag and the end tag of an XML element and will NOT be parsed by a parser; tags inside the text will NOT be treated as markup and entities will not be expanded)

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

An XML example

- > An XML document contains (and the DTD describes):
 - Elements, such as novel and paragraph, consisting of tags and content
 - Attributes, such as number="1", consisting of a name and a value
 - Entities (not used in this example)

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.7

A DTD example

- A novel consists of a foreword and one or more chapters, in that order
 - Each chapter must have a number attribute
- A foreword consists of one or more paragraphs
- A chapter also consists of one or more paragraphs
- A paragraph consists of parsed character data (text that cannot contain any other elements)

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

DTD ELEMENT

- ➤ In a DTD, elements are declared with an ELEMENT declaration.
- > General Syntax

```
<! ELEMENT el ement-name category>
or
```

<! ELEMENT element-name (element-content)>

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

40

Elements without children

- ➤ The syntax is <! ELEMENT name category>
 - name is the element name used in start and end tags
 - category should be EMPTY
 - Example
 - In the DTD: <! ELEMENT br EMPTY>
 - In the XML, an empty element may not have any content between the start tag and the end tag
 - An empty element may (and usually does) have attributes

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Elements with unstructured children

- The syntax is <! ELEMENT name category>
 - category may be ANY
 - This indicates that any content -- character data, elements, even undeclared elements -- may be used
 - Since the whole point of using a DTD is to define the structure of a document, ANY should be avoided wherever possible
 - category may be (#PCDATA), indicating that only parsed character data may be used
 - Example
 - In DTD: <! ELEMENT paragraph (#PCDATA)>
 - In XML: <paragraph>A shot rang out! </paragraph>
 - Notes:
 - · The parentheses are required!
 - In (#PCDATA), whitespace is kept exactly as entered
 - · Elements may not be used within parsed character data
 - · Entities are character data, and may be used

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.11

Elements with children

- An element may have one or more children
- > The syntax is

<!ELEMENT element-name (child1)>

Or

<!ELEMENT element-name (child1, child2, ...)>

- > Example
 - <! ELEMENT note (to, from, heading, body)>
- Parentheses are required, even if there is only one child
- A space must precede the opening parenthesis
- \circ Commas (,) between elements mean that all children must appear, and must be in the order specified
- All child elements must themselves be declared
- · Children may have children
- " | " separators means any one child may be used
- Suffixes can be used, e.g.
 - <! ELEMENT novel (foreword, chapter+)>
- Parentheses can be used for grouping, e.g. <! ELEMENT novel (foreword, (chapter+|section+))>

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

```
ELEMENT descriptions
Suffixes:
              optional (zero or one)
                                     foreword?
                                     chapter+
              one or more
              zero or more
                                     appendi x*
Separators
                                     foreword?, chapter+
              both, in order
                                     section|chapter
▶ Grouping
                                     (section|chapter)+
   ()
              grouping
         KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development
                                                           4.13
```

Elements with mixed content

- #PCDATA describes elements with only character data
- > #PCDATA can be used in an "or" grouping, e.g.
 <! ELEMENT note (#PCDATA|message) *>
 - This is called mixed content
 - · Certain (rather severe) restrictions apply:
 - #PCDATA must be first
 - The separators must be " | "
 - The group must be starred (meaning zero or more)

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Names and namespaces

- All names of elements, attributes, and entities, in both the DTD and the XML, are formed as follows:
 - The name must begin with a letter or underscore
 - The name may contain only letters, digits, dots, hyphens, underscores, and colons
- The DTD doesn't know about namespaces -- as far as it knows, a colon is just part of a name
 - The following are different (and both legal):
 - <! ELEMENT chapter (paragraph+)>
 - <! ELEMENT myBook: chapter (myBook: paragraph+)>
 - · Avoid colons in names, except to indicate namespaces

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.15

An expanded DTD example

```
<! DOCTYPE novel [
    <! ELEMENT novel
        (foreword, chapter+, biography?, criticalEssay*)>
    <! ELEMENT foreword (paragraph+)>
        <! ELEMENT chapter (section+|paragraph+)>
        <! ELEMENT section (paragraph+)>
        <! ELEMENT biography(paragraph+)>
        <! ELEMENT criticalEssay (section+)>
        <! ELEMENT paragraph (#PCDATA)>
]>
```

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Attributes and entities

- In addition to elements, a DTD may declare attributes and entities
- An attribute describes information that can be put within the start tag of an element
 - In XML:

```
<car name= "Toyota" model = "2001"></car>
```

• In DTD:

```
<! ATTLIST car
name CDATA #REQUIRED
model CDATA #IMPLIED >
```

- An entity describes text to be substituted
 - In XML: ©ri ght;
 - In DTD: <! ENTITY copyright "Copyright KFUPM">

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.17

Attributes

> The format of an attribute is:

```
<! ATTLIST element-name
    name type requirement
    name type requirement>
```

where the *name-type-requirement* may be repeated as many times as desired

- · Only spaces separate the parts, so careful counting is essential
- The element-name tells which element may have these attributes
- The name is the name of the attribute
- Each attribute has a type, such as CDATA (character data)
- Each attribute may be required, optional, or "fixed"
- In XML, attributes may occur in any order

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

1 10

Important attribute types

- > There are ten attribute types
- > These are the most important ones:
 - CDATA The value is character data
 - (man|woman|chi|ld) The value is one from this list
 - ID The value is a unique identifier
 - ID values must be legal XML names and must be unique within the document
 - NMTOKEN The value is a legal XML name
 - · This is sometimes used to disallow whitespace in the name
 - · It also disallows numbers, since an XML name cannot begin with a digit
- > The other seven, less frequently used, are:
 - IDREF, IDREFS, NMTOKENS, ENTITY, ENTITIES, NOTATION, xml:

 $KFUPM-081 @ Dr. El-Alfy \qquad SWE \ 444 \ Internet \ \& \ Web \ Application \ Development$

4.19

Requirements

- Recall that an attribute has the form
 - <! ATTLIST element-name name type requirement>
- > The requirement is one of:
 - A default value, enclosed in quotes
 - Example: <! ATTLI ST degree CDATA "PhD">
 - #REQUI RED
 - · The attribute must be present
 - #I MPLI ED
 - · The attribute is optional
 - #FIXED "value"
 - · The attribute always has the given value
 - · If specified in the XML, the same value must be used

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Entities

- Entities are a way to include fixed text (sometimes called "boilerplate")
- > Example of use in the XML:

```
This document is &copyright; 2002.
```

There are exactly five predefined entities:

```
<, &gt;, &amp;, &quot; and &apos;
```

- Additional entities can be defined in the DTD:
 - <! ENTI TY copyright "Copyright KFUPM">
- Entities can be defined in another document:
 <! ENTITY copyright SYSTEM "MyURI">
- Entities should not be confused with character references, which are numerical values between & and #
 - Example: &233#; or &xE9#; to indicate the character é

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.21

Another example: XML

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

The DTD for this example

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.23

Inline DTDs

If a DTD is used only by a single XML document, it can be put directly in that document, as follows:

```
<?xml version="1.0">
<!DOCTYPE myRootElement [
    <!-- DTD content goes here -->
]>
<myRootElement>
    <!-- XML content goes here -->
</myRootElement>
```

An inline DTD can be used only by the document in which it occurs

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

External DTDs

- External DTDs are almost always preferable to inline DTDs, since they can be used by more than one document
- ➤ The file extension for an external DTD is .dtd
 - External DTDs can only be referenced with a URL
- > An external DTD (a DTD that is a separate document) is declared with a SYSTEM or a PUBLIC command:
 - <!DOCTYPE myRootElement SYSTEM
 "http://www.mysite.com/mydoc.dtd">
 - The name that appears after DOCTYPE (in this example, myRootElement) must match the name of the XML document's root element
 - Use SYSTEM for external DTDs that you define yourself, and use PUBLIC for official, published DTDs

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

4.25

Validators

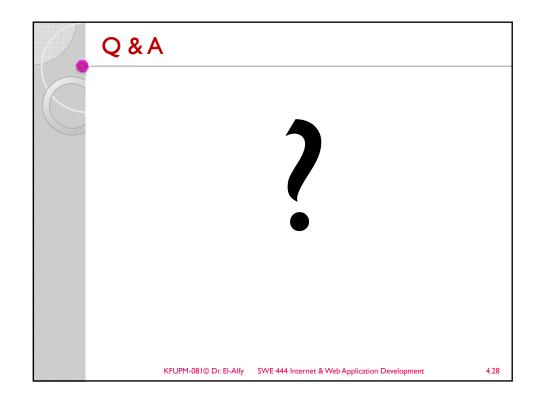
- Opera 5 and Internet Explorer 5+ can validate your XML against an internal DTD
 - IE provides (slightly) better error messages
 - Opera apparently just ignores external DTDs
 - · IE considers an external DTD to be an error
- jEdit with the XML plugin will check for wellstructuredness and (if the DTD is inline) will validate your XML each time you do a Save
 - http://www.jedit.org/
- Validate [Using Inline DTD] http://www.stg.brown.edu/service/xmlvalid/

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development

Limitations of DTDs

- > DTDs are a very weak specification language
 - You can't put any restrictions on element contents
 - It's difficult to specify:
 - · All the children must occur, but may be in any order
 - · This element must occur a certain number of times
 - There are only ten data types for attribute values
- > But most of all: DTDs aren't written in XML!
 - If you want to do any validation, you need one parser for the XML and another for the DTD
 - This makes XML parsing harder than it needs to be
 - There is a newer and more powerful technology: XML Schemas
 - · However, DTDs are still very much in use

 $KFUPM-081 @ Dr. El-Alfy \qquad SWE \ 444 \ Internet \ \& \ Web \ Application \ Development$



References

- > Some useful links with examples and other resources:
 - W3School DTD Tutorial
 - http://www.w3schools.com/dtd/default.asp
 - MSXML 4.0 SDK
 - http://www.topxml.com
 - http://www.xml.org
 - http://www.xml.com
 - Several online presentations

KFUPM-081© Dr. El-Alfy SWE 444 Internet & Web Application Development