

Chapter 2: Well-Formed XML

1

Chapter 2 Objectives

- Understand components of an XML document
- How to create XML elements using start-tags and end-tags
- How to further describe elements with attributes
- How to declare your document as being XML
- How to send instructions to applications that are processing the XML document
- Which characters aren't allowed in XML – and how to use them in your documents anyway!
- What is well-formedness anyway?

2

Components

- **Elements**
- **Attributes**
- **Declaration**
- **Instructions**
- **Comments**
- **CDATA**

3

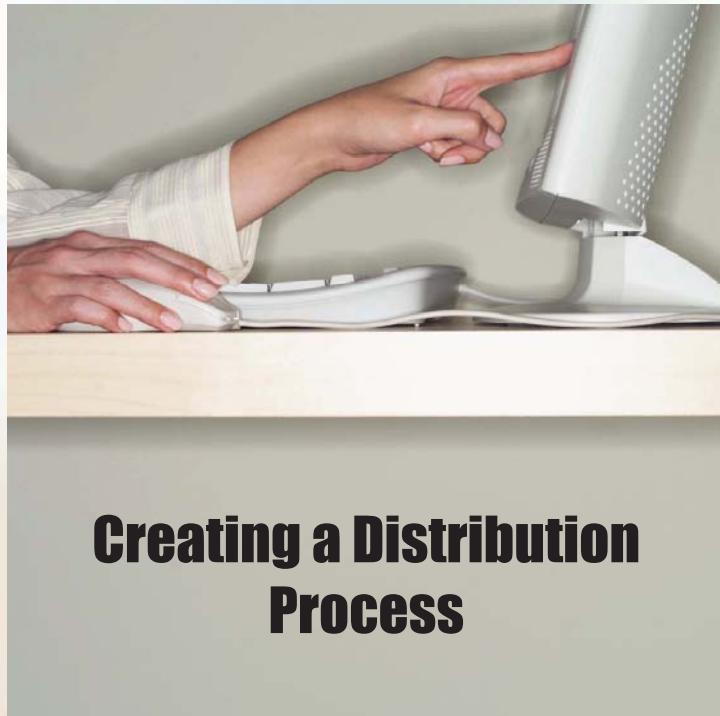
Tags and Elements

```
<name>
<first>John</first>
<middle>Fitzgerald Johanson</middle>
<last>Doe</last>
</name>
```

- **<first>** is a start-tag
- **</first>** is an end-tag
- **<first>John</first>** is an element
- John is the content of the element
- Analogy: Sandwich

4

Try It Out



5

Rules for Elements

Every Start Tag Must Have An End Tag

Bad Example

<P>Here is some text in an HTML paragraph.

Here is some more text in the same paragraph.

<P>And here is some text in another HTML
paragraph.</P>

Good Example

<P>Here is some text in an HTML paragraph.</P>

<P>Here is some more text in the same paragraph.</P>

<P>And here is some text in another HTML
paragraph.</P>

6

Rules for Elements

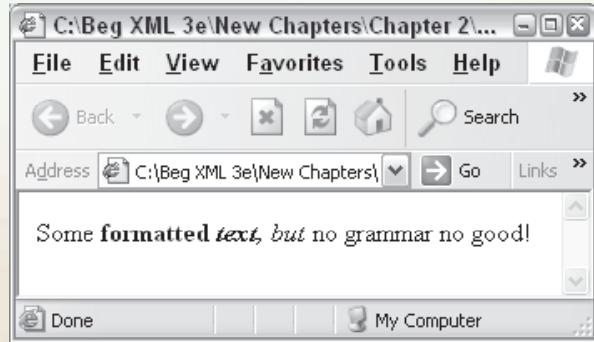
Elements Must Be Properly Nested

Bad Example

```
<P>Some <STRONG> formatted <EM>text</STRONG>, but</EM>  
no grammar no good!</P>
```

Good Example

```
<P>Some  
  <STRONG>  
    Formatted  
    <EM>  
      Text  
    </EM>  
  </STRONG>  
  <EM>  
    , but  
  </EM>  
No grammar no good!  
</P>
```



7

Rules for Elements

Can Have Only One Root Element

Bad Example

```
<name>John</name>  
<name>Jane</name>
```

Good Example

```
<names>  
<name>John</name>  
<name>Jane</name>  
</names>
```

8

Rules for Elements

Elements Must Obey XML Naming Conventions

- Names can start with letters, no numbers
- After first character, numbers, hyphens and periods are allowed
- Names can't contain spaces
- There are *reserved* characters like ":"
- Names can't start with the letters "xml", "XML", or "Xml", or any other combination
- No spaces after the "<", but before the ">" if desired

9

Rules for Elements

Case Sensitivity

These are two different elements

```
<name>John</name>
<NAME>John</NAME>
```

10

Rules for Elements

Whitespace in PCDATA

Whitespace stripping takes place in HTML...

```
<P>This is a paragraph.  
  &nbsp;&nbsp;&nbsp;&nbsp;It has a whole  
  bunch<BR>&nbsp;&nbsp;of space.</P>
```

...but not in XML

Example

```
<tag>This is a paragraph.  It has a whole  
Bunch of space.</tag>
```

This is a paragraph.   It has a whole
Bunch of space.

11

Rules for Elements

End-of-Line Whitespace

- Windows uses both the line feed and the carriage return
- UNIX uses only line feed
- XML parsers will convert all Windows “line feed and carriage returns” to just line feed characters to standardize end-of-line logic

12

Rules for Elements

Readability Whitespace

```
<Tag>  
  <AnotherTag>This is some XML</AnotherTag>  
</Tag>
```

This is known as extraneous whitespace in the markup.

13

Attributes

- Simple name/value pairs associated with an element
- name="value" or name='value'
- Can have more than one attributes for an element
- Order is not important
- Elements' order is important!

```
<name nickname="Shiny John">  
  <first>John</first>  
  <middle>Fitzgerald Johansen</middle>  
  <last>Doe</last>  
</name>
```

(Adding attributes to orders ...)

14

Elements or Attributes?

- What is the difference between

```
<name nickname="Shiny John">  
</name>
```

and

```
<name>  
  <nickname>Shiny John</nickname>  
</name>
```

- Element can contain elements
 - Not the case for attribute
- Order
- Part of data or property of data
- Personal preference / style

15

Comments

```
<name nickname='Shiny John'>
```

```
  <first>John</first>
```

```
  <!--John lost his middle name in a fire-->
```

```
  <middle></middle>
```

```
  <last>Doe</last>
```

```
</name>
```

(Adding comments to orders ...)

16

Empty Elements

```
<name nickname='Shiny John'>
  <first>John</first>
  <!--John lost his middle name in a fire-->
  <middle></middle>
  <last>Doe</last>
</name>
```

can also use self-closing tag <middle/>

ok to have attributes

```
<name nickname="Shiny John"/>
```

17

XML Declaration

Define the Encoding and Standalone Attributes

```
<?xml version='1.0' encoding='UTF-16' standalone='yes'?>
<name nickname='Shiny John'>
  <first>John</first>
  <!--John lost his middle name in a fire-->
  <middle/>
  <last>Doe</last>
</name>
```

(Adding declaration to orders ...)

18

Processing Instructions

```
<?xml version='1.0'?>
<name nickname='Shiny John'>
    <first>John</first>
    <!--John lost his middle name in a
        fire-->
    <middle/>
    <?nameprocessor SELECT * FROM blah?>
    <last>Doe</last>
</name>
```

19

Processing Instructions

```
<?xml version='1.0'?>
<name nickname='Shiny John'>
    <first>John</first>
    <!--John lost his middle name in a fire-->
    <middle/>
    <?nameprocessor SELECT * FROM blah?>
    <last>Doe</last>
</name>
```

PITarget *instruction*

Is the XML declaration a PI? **NO!**

(Adding PI to orders ...)

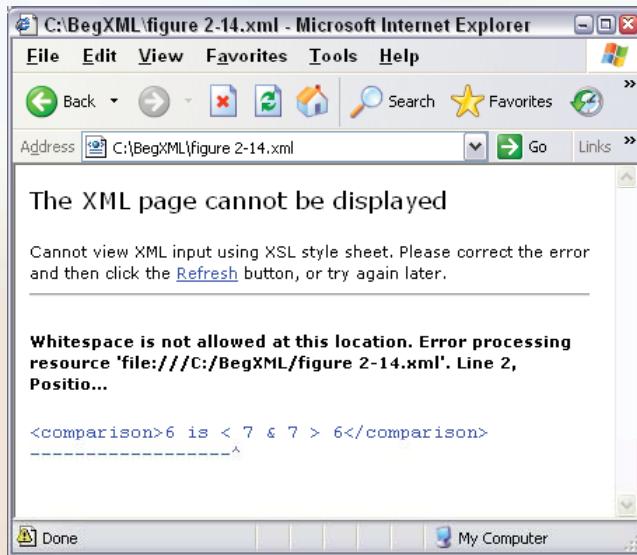
20

Illegal PCDATA Characters

- PCDATA: Parsed character data

```
<!--This is not well-formed XML!-->  
<comparison>6 is < 7 & 7 > 6</comparison>
```

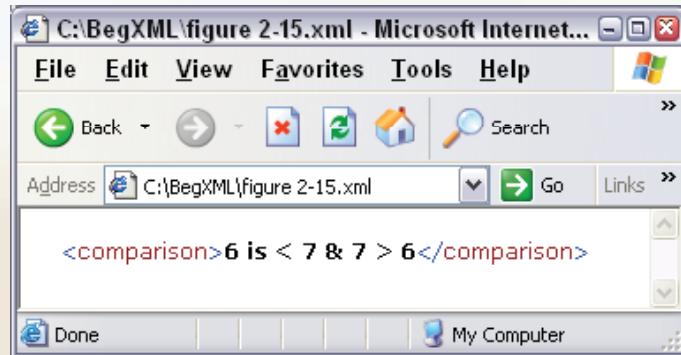
Ouch!



21

Escaping Characters

```
<comparison>  
 6 is &lt; 7 &amp;gt; 7 &gt; 6  
</comparison>
```



22

Escaping Characters

- Entity reference

- & &
- < <
- > >
- &apos ‘
- " “
- &#n or &#xn
 - © or &xA9

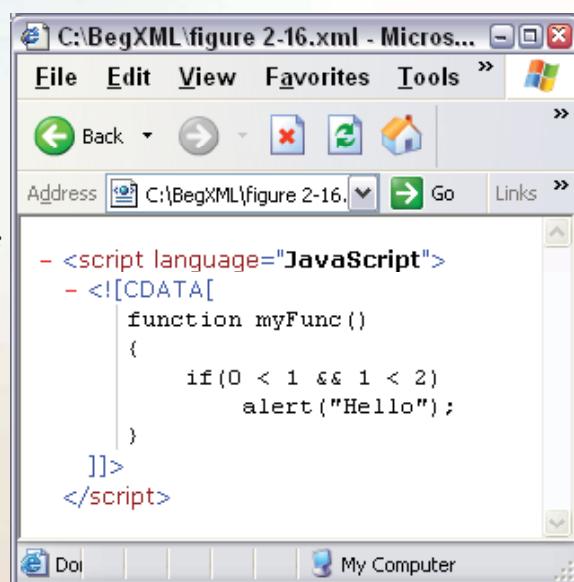
©

23

CDATA Sections

- CDATA: Character data

```
<script language='JavaScript'>
<! [CDATA[
    function myFunc()
{
    if(0 < 1 && 1 < 2)
        alert("Hello");
}
]]>
</script>
```



- Example: HTML in XML

24

Errors in XML

There Are Errors, and Then There Are Fatal Errors

- **Errors**
 - Violations rules in recommendation
 - May recover
 - Continue processing
- **Fatal errors**
 - Draconian error handling
 - Not allowed to continue
 - XML parser should not try to “recover” errors

25

Well-Formedness

- **XML documents must have exactly one root element**
- **A start tag must have a matching end tag**
- **XML tags are case sensitive**
- **XML elements must be properly nested**
- **XML attribute values must be quoted**
- **Reserved characters must be handled properly**

26