

Empty Elements

- In XML, for every opening tag there must be a corresponding closing tag.

In contrast, SGML has “tag minimization” rules that permit to leave out tags that can be uniquely reconstructed by the SGML parser.

- Since `
</br>` does not look very nice, and the closing tag contains no additional information, empty element tags were introduced in XML: “`
`” is equivalent to “`
</br>`”.

This was one of the few points where XML was not a subset of SGML at the beginning. Of course, SGML was/will be extended to permit this syntax, too.

Line Ends

- In SGML, line ends (record boundaries) directly after a start tag or directly before an end tag are ignored (i.e. at the start or end of the content).

- In XML, line ends or empty space is not ignored.

The parser passes it to the application, which can of course ignore it. E.g. a validating parser, which knows that an element contains pure element content (not mixed content) will ignore whitespace between the elements.

- In XML, line ends are normalized to a line feed.

Even on a Windows system (which uses CR, LF for line ends), the XML application receives LF (ASCII 10) from the parser.

Attributes (1)

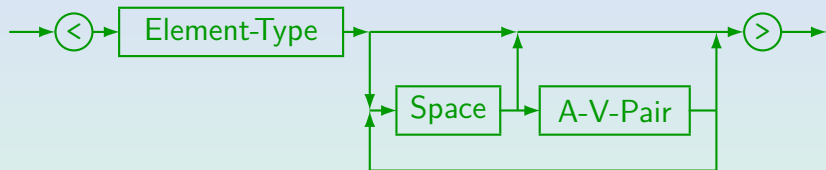
- In the start tag, attribute-value pairs can be optionally specified.
- E.g. in XHTML, links to other documents are marked with the element `a` (“anchor”):

```
XML was developed by the  
<a href="http://www.w3.org">W3C</a>.
```

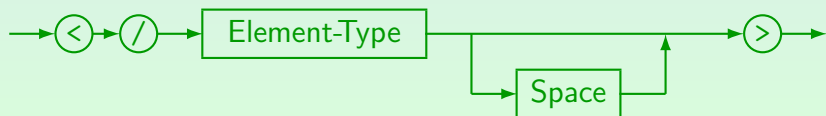
- The text of the reference is given in the element content, the URI of the referenced web page is specified in the attribute “`href`”.

Attributes (2)

Start-Tag:

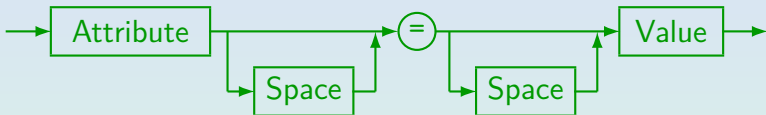


End-Tag:

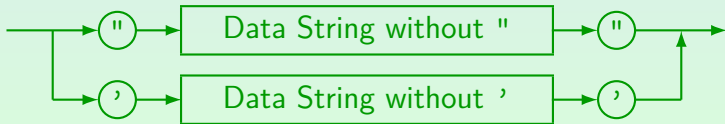


Attributes (4)

A-V-Pair:



Value:



Attributes (5)

- Attribute values can be enclosed in " or ' .
The other sign can appear inside the string.

If one needs both quotation marks, one must use an entity or character reference (see below).

- Attribute values cannot contain elements.
- The character "<" is forbidden in attribute values.

If necessary, one can include it with a character reference or an entity reference. Excluding "<" in attribute values helps to detect errors earlier (such as a missing quote). To make this clear: "<" is not forbidden in the internal value of an attribute (which an XML parser can pass to the application), it is only forbidden in the external representation. But it never creates elements in attribute values.

Attributes (6)

- The character “&” is treated special in attribute values (character/entity reference, see below).
- Attribute values can extend over multiple lines. The parser replaces tabs and line ends in the attribute value by a space.

Depending on the type of the attribute, white space may be normalized: It is then removed at the beginning and at the end of the attribute value, and several consecutive spaces are merged into one. However, this does not happen for normal “CDATA” attributes.

- The sequence in which several attribute-value-pairs are listed in a tag is not important.

Character References (1)

- One must distinguish between
 - the repertoire of characters used internally (e.g. data passed from XML parser to application)
 - the encoding of these characters in bytes for exchanging documents (external representation).
- Internally, XML uses the Unicode character set.
- For exchanging documents, one can e.g. use the ISO 8859-1 (ISO Latin 1) character codes, which contains only a subset of all Unicode characters.

Other encodings contain e.g. cyrillic or japanese characters.

Character References (2)

- The XML declaration at the beginning of the XML file defines the encoding (see below).

The encoding can also be specified in the HTTP protocol.

- The characters in the ISO Latin 1 character set are also contained in the Unicode character set and have the same numeric codes in both character sets.

I.e. Unicode is upward compatible to ISO Latin 1. However, the encoding as sequence of bytes is different. At the beginning, Unicode character numbers had 16 Bit, now there are 17 planes of 16 Bit each. With the UTF-8 encoding of Unicode, at least the 7-bit ASCII characters have the same encoding in ASCII, ISO Latin 1, and Unicode. However, for German national characters (ä, ö, ü, etc.) this is no longer true: UTF-8 uses two bytes for them. See Slide 36.

Character References (3)

- Characters that cannot be directly entered, can be written as a “character reference” using their numeric code:

`ä`

is an “ä”. Hexadecimal notation can also be used:

`ä`

- The numbers refer to the repertoire (i.e. Unicode), not to the encoding for exchange.

ASCII and ISO Latin 1 codes can be used since Unicode is upward compatible.

- Character references can also be used to “escape” characters that otherwise would have special meaning in SGML/XML.

The result of a character reference is always treated as data.

E.g. if a double quote (ASCII 34) needs to be included in an attribute value that is enclosed in double quotes, one can write it as “"”.

Short Digression: UTF-8 (1)

- The character “ä” has the number (“code point”) `U+00E4` (228) in Unicode.
- It has the same number in ISO Latin 1, and this encoding uses one byte per character, therefore it is represented by the byte `0xE4`.
- UTF-8 uses one byte per character only for characters with codes up to `U+007F` (127).

I.e. when the first bit of a byte is 0, this byte encodes a character by itself.
This ensures that UTF-8 is upwards compatible to ASCII.

Short Digression: UTF-8 (2)

- For characters with code point above **U+007F** (including **ä**), multi-byte sequences are needed in UTF-8:
 - For n -byte sequences, the first byte starts with n **1**-bits followed by a **0**-bit.

E.g. the prefix 110 in the first byte means that this byte sequence consists of two bytes.
 - All other bytes start with **10**.

If one jumps somewhere into an UTF-8 byte stream, one can detect the start of character encodings (bytes starting not with 10).
- Unicode reserves 17 planes of 16 bit each (up to **U+10FFFF**), which requires max. 4 bytes in UTF-8.

Short Digression: UTF-8 (3)

- A two byte sequence consists of bytes **110xxxxx** and **10yyyyyy**, representing the code point **xxxxxyyyyyy** and is used for the range **U+0080** to **U+7FF**.
- The character “ä” with code **U+00E4** (**11100100**) is encoded as **11000011** (**0xC3**) and **10100100** (**0xA4**).
- If an editor assumes that the text is encoded in ISO Latin 1, it will display the two characters **Ã ¤**.

The second is a “general currency symbol”: U+00A4 (= 164).
- Each code point must be encoded in the shortest byte sequence (e.g. ASCII characters as one byte).

Comments (1)

- Comments can be used to enter notes or explanations for a reader of the SGML/XML source file into the document.
- Comments are ignored by programs that process an SGML/XML file. E.g. they might not appear in the formatted output.

The XML standard permits that an XML parser passes comments to the application program, but it does not require this.

- A comment in SGML/XML has the form

```
<!-- This is a comment -->
```

Comments (2)

- Comments can extend over several lines.

I.e. they do not have to be closed on the same line.

- Within a comment, it is forbidden to write two consecutive hyphens “--”.

In SGML, the comment actually extends from “--” to “--”. However, it can only be used in a markup declaration, which starts with “<!” and ends with “>”.

- Tags within a comment are permitted, but confuse many browsers.

Browsers try to correct syntax errors. When they see a tag, they might assume that the author forgot the “end of comment” mark.

Comments (3)

- Comments can be used anywhere in the document outside other markup.
- They cannot be used within tags.
- In SGML (but not in XML), comments “`-- ... --`” can appear in markup declarations at places permitted by the grammar.

In modern programming languages, whitespace including comments is allowed between tokens. SGML/XML are different: maybe because they are languages for writing documents, not programs, maybe they are a bit outdated in this aspect.

- XML supports only “`<!-- ... -->`”.

Exercise

Please find syntax errors:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<GradesDB3>
  <student sid='101' first='Ann' last="Smith"
    <result cat='H' eno = 1 points=10>
    <result cat='H' eno ='1' points='8'>
    <result cat='M" eno ="1" points='12'
      >
  </ student >
  <!------- Exercises ----->
  <ex cat='H' eno='1' note='<em>difficult</em>''
    Rel&#97 tional Algebra</ex>
</Grades-DB>
```

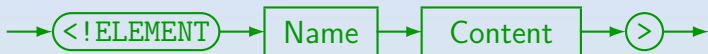

Example

Simple DTD for a HTML-Subset:

```
<!ELEMENT html (head, body)>
<!ELEMENT head (title)>
<!ELEMENT title (#PCDATA)>
<!ELEMENT body ((#PCDATA|p|em|ul)*)>
<!ELEMENT p ((#PCDATA|em|br)*)>
<!ELEMENT em (#PCDATA)>
<!ELEMENT br EMPTY>
<!ELEMENT ul (li+)>
<!ELEMENT li ((#PCDATA|p|em|ul|br)*)>
```


Element-Type Declarations (2)

Element-Type Declaration:



White space is required between “`<!ELEMENT`” and the name, and between the name and the content specification. It is permitted but not required between content specification and the “`>`”.

Names in XML must start with a letter, an underscore “`_`” or a colon “`:`”, and can otherwise contain letters, digits, periods “`.`”, hyphens “`-`”, underscores “`_`”, colons “`:`”, or certain special Unicode characters. Names starting with “`xml`” in any capitalization are reserved, the colon is treated specially by the XML namespace standard.

The element type declaration is SGML is more complicated: There, also specifications for markup minimization are required (if the markup minimization parameter `OMITTAG` is set), “exclusions” and “inclusions” are possible, several element types can be declared together, etc.

Content Specifications (1)

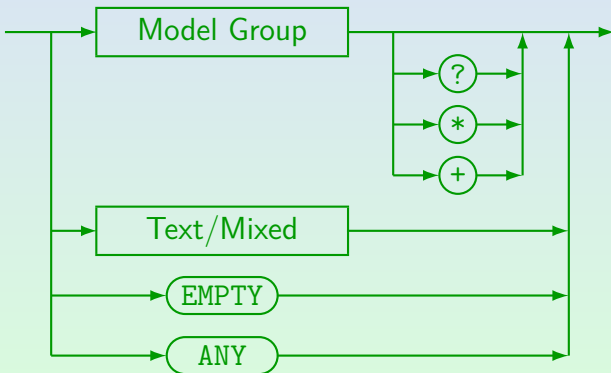
- The building blocks of content specifications are
 - Names X of element types: This pattern matches exactly one element of type X , i.e. basically `<X>...</X>`.
 - The keyword **#PCDATA**: Pure textual data without tags (but possibly character/entity references).
 - #PCDATA stands for “Parsed Character Data” (it is checked that there are no nested tags). SGML (not XML) has also CDATA (treats `<` as text).
- One can specify the optionality/multiplicity of elements and groups by attaching occurrence indicators:
 - **A?**: Optional, non repeatable (0 or 1 time).
 - **A***: Optional, repeatable (0 or more times).
 - **A+**: Required, repeatable (1 or more times).

Content Specifications (2)

- Content specifications can be connected with
 - $(A_1 | \dots | A_n)$: “Alternative”/“Choice” (one of the A_i).
E.g. $(A | B)$: “A or B”. The content must match A or match B.
 - (A_1, \dots, A_n) : “Sequence” (all A_i in the given sequence).
E.g. (A, B) : “First A, then B” (“A followed by B”). A prefix of the content must match A, the rest B.
SGML (not XML) has also $(A \& B)$: “A and B”. A and B must both appear, but in arbitrary sequence. This is equivalent to $((A, B) | (B, A))$.
- The A_i are
 - An element type (possibly with $?$ / $*$ / $+$).
 - **#PCDATA** (in XML with restrictions, see below).
 - A nested model group (possibly with $?$, $*$, $+$).

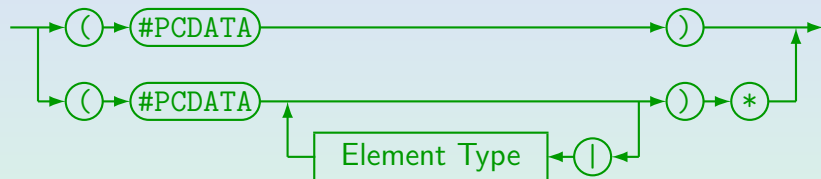
Content Specifications (4)

Content:



Content Specifications (5)

Text/Mixed:



- In XML, the only content models that can contain `#PCDATA` are (SGML has no such restriction):
 - `(#PCDATA)`
 - `(#PCDATA | Element-Type | ... | Element-Type)*`

Attribute Declarations (1)

- Example (symbol used for marking list items):

```
<!ATTLIST UL type (disc|square|circle) #IMPLIED>
```

In HTML 4.01 Strict this attribute was removed.

- Several attributes (of one element type) can be declared in a single `ATTLIST` command.
- E.g. some attributes of images in HTML:

```
<!ATTLIST IMG src CDATA #REQUIRED  
alt CDATA #REQUIRED  
width CDATA #IMPLIED  
height CDATA #IMPLIED>
```


Exercise (1)

Please find syntax errors:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<!DOCTYPE GradesDB4 [
    <!-- contains syntax errors -->
    <!ELEMENT GradesDB4 (STUDENT, RESULT)*>
    <!ELEMENT STUDENT RESULT+>
    <!ATTLIST STUDENT FIRST CDATA #REQUIRED
                LAST CDATA #REQUIRED>
    <!ELEMENT RESULT #EMPTY>
    <!ATTLIST RESULT EX_ID IDREF #REQUIRED
                POINTS NMTOKEN #OPTIONAL>
    <!ELEMENT EXERCISE #PCDATA>
    <!ATTLIST EXERCISE ID ID #REQUIRED>
]> <!-- continued on next slide -->
```

Exercise (2)

Please validate against DTD on last slide:

```
<GradesDB4>
  <student sid='101' first='Ann' last='Smith'>
    <email>smith@acm.org</email>
    <result ex_id='H1' points='A+'/>
    <result ex_id='2' points='8'/>
    <result ex_id='M1' points='12 points'/>
  </student>
  <student first='Maria' last='Brown'/>
  <exercise id='H1'>ER</exercise>
  <exercise id='2'>SQL</exercise>
</GradesDB4>
```

Exercise (3)

Please develop a DTD for this document:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<smallbusiness>
  <product id='P01' name='Apple' price='0.40'>
    Really <em>delicious</em>Apples!</product>
  <product id='P02' name='Banana' price='0.50'>
    The best bananas!</product>
  <order id='R100' customer='Ann Smith'>
    <item prodid='P01' />
    <item prodid='P02' quantity='5' /> </order>
  <order id='R100' customer='Maria Brown'>
    <item prodid='P01' quantity='3' /> </order>
</smallbusiness>
```


DOCTYPE Declaration (1)

- One usually refers at the beginning of the document to the corresponding DTD:

```
<?xml version="1.0"?>  
<!DOCTYPE EMAIL SYSTEM "mail.dtd">  
<EMAIL>  
    ...  
</EMAIL>
```

- The file “`mail.dtd`” contains the declaration of elements, attributes, and entities as described above.

```
<!ELEMENT EMAIL (TO, FROM, DATE, SUBJECT?,  
                CONTENTS)>  
    ...
```


XML Declaration (1)

- XML documents should start with an XML declaration that specifies at least the XML version:

```
<?xml version="1.0"?>
```

For SGML processors, the XML declaration is simply a processing instruction.

- Version “1.0” is still the most widely used version, but there is now also a version “1.1”.

There are new editions of the W3C recommendation for XML 1.0. but they only clarify/correct a few points. The W3C recommendation for XML 1.0 was published on February 10, 1998. The second edition was published on October 6, 2000. The third edition of XML 1.0 was published on February 4, 2004, together with the first edition of XML 1.1. The current, fourth edition of XML 1.0 was published together with the second edition of XML 1.1 on August 16, 2006, both were edited in place on September 29, 2006. [<http://www.w3.org/XML/>].

