

Part III

Well-Formed XML

Outline of this part

1 Formalization of XML

- Elements
- Attributes
- Entities

2 Well-Formedness

- Context-free Properties
- Context-dependent Properties

3 XML Text Declarations

- XML Documents and Character Encoding
- Unicode
- XML and Unicode

4 The XML Processing Model

- The XML Information Set
- More XML Node Types

Formalization of XML

- We will now try to approach XML in a slightly more formal way. The nuts and bolts of XML are pleasingly easy to grasp.
- This discussion will be based on *the* central XML technical specification:
 - **Extensible Markup Language (XML) 1.0 (Second Edition)**
W3C Recommendation 6 October 2000
(<http://www.w3.org/TR/REC-xml>)



Visit the W3C site

This lecture does *not* try to be a “guided tour” through the XML-related W3C technical documents (*boring!*).

Instead we will cover the basic principles and most interesting ideas. Visit the W3C site and use the original W3C documents to get a full grasp of their contents.

Elements

- The **element** is the main markup construct provided by XML.
 - Marked up document region (**element content**) enclosed in matching **start** end **closing (end) tags**:
 - **start tag**: `<t>` (*t* is the **tag name**),
 - matching **closing tag**: `</t>`

Well-formed XML (fragments)

```
1 <foo> okay </foo>
2 <This-is-a-well-formed-XML-tag.> okay
3 </This-is-a-well-formed-XML-tag.>
4 <foo>okay</foo>
```

Non-well-formed XML

```
1 <foo> oops </bar>
2 <foo> oops </Foo>
3 <foo> oops ... (EOT)
```

- Element content may contain document characters as well as **properly nested elements**, so-called **mixed content**):

Well-formed XML

```
1 <foo><bar>
2     <baz> okay </baz>
3     </bar>
4     <ok> okay </ok> still okay
5 </foo>
```

Non-well-formed XML

```
1 <foo><bar> oops </foo></bar>
2 <foo><bar> oops </bar><bar> oops </foo></bar>
```



Check for proper nesting

Which data structure would you use to straightforwardly implement the check for proper nesting in an XML parser?

- Element content may be **empty**:
 - The fragments `<t> </t>` and `<t/>` are well-formed XML and considered equivalent.
- Element nesting establishes a **parent–child relationship between elements**:
 - In the XML fragment `<p><c>... </c>...<c'>... </c'> </p>`,
 - element p is *the* **parent** of elements c, c' ,
 - elements c, c' are **children** of element p ,
 - elements c, c' are **siblings**.
- There is exactly one element that encloses the **whole** XML content: the **root element**.

Non-well-formed XML

```
1 <one>
2   one eins un
3 </one>
4 <two> two zwei deux </two>
```

Attributes

- Elements may further be classified using **attributes**:
(It is common practice to denote an attribute named a by $@a$ in written text (**attribute** a).)

`<t a="..." b='...' ...> ...</t>`

- An attribute **value** is restricted to character data
(attributes may *not* be nested),
- attributes are *not* considered to be children of the containing element
(instead they are **owned** by the containing element).

Well-formed XML (fragment)

```
1 <price currency="US$" multiplier='1'>
2   23.45
3 </price>
4 <price>
5   <currency>US$</currency>
6   <multiplier>1</multiplier>
7   23.45
8 </price>
```

Entities

- In XML, document **content** and **markup** are specified using a *single set of characters*.
- The characters { <, >, &, ", ' } form pieces of XML markup and may instead be denoted by **predefined entities** if they actually represent content:

Character	Entity
<	<
>	>
&	&
"	"
'	'

Well-formed XML

```

1 <operators>Valid comparison operators
2   are &lt;; =, &amp; &gt;.</operators>
```

- The XML entity facility is actually a versatile recursive **macro** expansion machinery (more on that later).

Well-Formedness

- The W3C XML recommendation is actually more formal and rigid in defining the syntactical structure of XML:

*“A textual object is **well-formed** XML if,*

- ① *Taken as a whole, it **matches the production labeled document.***
- ② *It meets all the **well-formedness constraints** given in this [the W3C XML Recommendation] specification. . . .”*

Well-formedness #1: Context-free Properties

- 1 All **context-free** properties of well-formed XML documents are concisely captured by a **grammar** (using an EBNF-style notation).
 - **Grammar**: system of **production (rule)s** of the form

$$lhs ::= rhs$$

Excerpt of the XML grammar

[1]	<i>document</i>	::=	<i>prolog element Misc*</i>
[2]	<i>Char</i>	::=	\langle a Unicode character \rangle
[3]	<i>S</i>	::=	$(\text{'\u0027'} \text{'\u005C'} \text{'\u005C'} \text{'\u005C'} \text{'\u005C'})^+$
[4]	<i>NameChar</i>	::=	<i>Letter Digit '.' '-' '_' ':'</i>
[5]	<i>Name</i>	::=	$(\text{Letter} \text{'\u0027'} \text{'\u005C'})^+ (\text{NameChar})^*$
[10]	<i>AttValue</i>	::=	$\text{'\"'} ([\text{'\u0026'}] \text{Reference})^* \text{'\"'}$ $ \text{'\u0027'} ([\text{'\u0026'}] \text{Reference})^* \text{'\u0027'}$
[14]	<i>CharData</i>	::=	$[\text{'\u0026'}]^*$
[22]	<i>prolog</i>	::=	<i>XMLDecl? Misc*</i>
[23]	<i>XMLDecl</i>	::=	$\text{'<?xml' VersionInfo EncodingDecl? S? '?>'}$
[24]	<i>VersionInfo</i>	::=	$S \text{'version' Eq ('\"' VersionNum '\"' '\u0027' VersionNum '\u0027')}$
[25]	<i>Eq</i>	::=	$S? \text{'='} S?$
[26]	<i>VersionNum</i>	::=	$([\text{a-zA-Z0-9}_. :] \text{'-'})^+$
[27]	<i>Misc</i>	::=	<i>S</i>
[39]	<i>element</i>	::=	<i>EmptyElemTag</i> $ $ <i>S</i> <i>Tag content ETag</i>
[40]	<i>S</i> <i>Tag</i>	::=	$\text{'<' Name (S Attribute)* S? '>'}$
[41]	<i>Attribute</i>	::=	<i>Name Eq AttValue</i>
[42]	<i>ETag</i>	::=	'</' Name S? '>'
[43]	<i>content</i>	::=	$(\text{element} \text{CharData} \text{Reference})^*$
[44]	<i>EmptyElemTag</i>	::=	$\text{'<' Name (S Attribute)* S? '>'}$
[67]	<i>Reference</i>	::=	<i>EntityRef</i>
[68]	<i>EntityRef</i>	::=	'\u0026' Name ';'
[84]	<i>Letter</i>	::=	$[\text{a-zA-Z}]$
[88]	<i>Digit</i>	::=	$[\text{0-9}]$

N.B.

- The numbers in $[\cdot]$ refer to the corresponding productions in the W3C XML Recommendation.

Expression...	... denotes	
r^*	$\epsilon, r, rr, rrr, \dots$	zero or more repetitions of r
r^+	rr^*	one or more repetitions of r
$r?$	$r \mid \epsilon$	optional r
$[abc]$	$a \mid b \mid c$	character class
$[\^abc]$		inverted character class

Remarks

Rule...	...implements this characteristic of XML:
[1]	an XML document contains exactly one root element
[10]	attribute values are enclosed in " or '
[22]	XML documents may include an optional declaration prolog
[14]	characters < and & may <i>not</i> appear literally in element content
[43]	element content may contain character data and entity references as well as nested elements
[68]	entity references may contain arbitrary entity names (other than lt, amp, ...)
⋮	⋮

- As usual, the XML grammar may systematically be transformed into a program, an **XML parser**, to be used to check the syntax of XML input.

Parsing XML

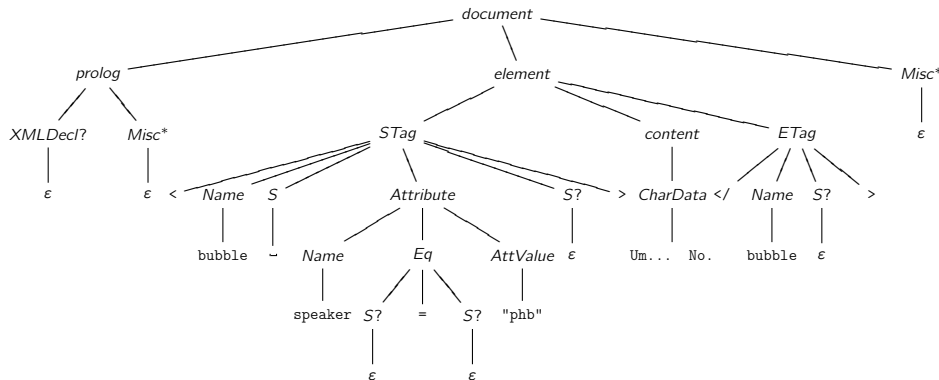
- 1 Starting with the symbol *document*, the parser uses the $lhs ::= rhs$ rules to expand symbols, constructing a **parse tree**.
- 2 The leaves of the parse tree are characters which have no further expansion.
- 3 The XML input is **parsed** successfully if it perfectly matches the parse tree's **front** (concatenate the parse tree leaves from left to right¹).

¹**N.B.:** $x\epsilon y = xy$.

Example 1

Parse tree for XML input

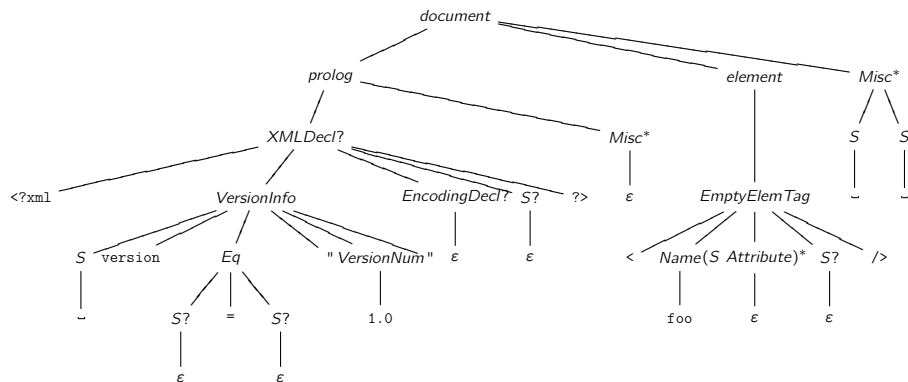
`<bubble speaker="phb">Um... No.</bubble>` :



Example 2

Parse tree for the “minimal” XML document

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



Well-formedness #2: Context-dependent Properties

- The XML grammar cannot enforce *all* XML **well-formedness constraints (WFCs)**.
- Some XML WFCs depend on
 - ① what the XML parser has **seen before** in its input, or
 - ② on a **global state**, *e.g.*, the definitions of user-declared entities.
- These WFCs cannot be checked by simply comparing the parse tree front against the XML input (**context-dependent WFCs**).

Sample WFCs

WFC	Comment
(2) Element Type Match	The <i>Name</i> in an element's end tag must match the element name in the start tag.
(3) Unique Att Spec	No attribute name may appear more than once in the same start tag or empty element tag.
(5) No < in Attribute Values	The replacement text of any entity referred to directly or indirectly in an attribute value (other than <) must not contain a <.
(9) No Recursion	A parsed entity must not contain a recursive reference to itself, either directly or indirectly.

All 10 XML WFCs are given in <http://www.w3.org/TR/REC-xml>.

How to implement the XML WFC checks?

Devise methods—besides parse tree construction—that an XML parser could use to check the XML WFCs listed above.

Specify *when* during the parsing process you would apply each method.

The XML Text Declaration `<?xml...?>`

- Remember that a well-formed XML document may start off with an optional header, the **text declaration** (grammar rule [23]).
 - N.B.** Rule [23] says, *if* the declaration is actually there, no character (whitespace, etc.) may precede the leading `<?xml`.

The leading `<?xml`

Can you imagine why the XML standard is so rigid with respect to the placement of the `<?xml` leader of the text declaration?

- An XML document whose text declaration carries a *VersionInfo* of `version="1.0"` is required to conform to W3C's XML Recommendation posted on October 6, 2000 (see <http://www.w3.org/TR/REC-xml>).

XML Documents and Character Encoding

- For a computer, a character like X is nothing but an 8 (16/32) bit **number** whose value is *interpreted* as the character X when needed (e.g., to drive a display).
- Trouble is, a large number of such *number* → *character* mapping tables, the so-called **encodings**, are in parallel use today.
- Due to the huge amount of characters needed by the global computing community today (Latin, Hebrew, Arabic, Greek, Japanese, Chinese . . . languages), **conflicting intersections** between encodings are common.

Example:

0xa4 0xcb 0xe4 0xd3	<u>iso-8859-7</u> →	☐, Λ δ Σ
0xa4 0xcb 0xe4 0xd3	<u>iso-8859-15</u> →	€ Ë ä Ó

Unicode

- The **Unicode** (<http://www.unicode.org/>) Initiative aims to define a new encoding that tries to embrace all character needs.
- The Unicode encoding contains characters of “all” languages of the world, plus scientific, mathematical, technical, box drawing, ... symbols (see <http://www.unicode.org/charts/>).
- Range of the Unicode encoding: 0x0000–0x10FFFF (16 × 65536 characters).
 - Codes that fit into the first 16 bits (denoted U+0000–U+FFFF) have been assigned to encode the most widely used languages and their characters (**Basic Multilingual Plane, BMP**).
 - Codes U+0000–U+007F have been assigned to match the 7-bit ASCII encoding which is pervasive today.

UTF-32

- Current CPUs operate most efficiently on **32-bit words (16-bit words, 8-bit bytes)**.
- Unicode thus developed **Unicode Transformation Formats (UTF)** which define how a Unicode character code between $U+0000$ – $U+10FFFF$ is to be mapped into a 32-bit word (16-bit words, 8-bit bytes).

UTF-32 (map a Unicode character into a 32-bit word)

- 1 Map any Unicode character in the range $U+0000$ – $U+10FFFF$ to the corresponding 32-bit value $0x00000000$ – $0x0010FFFF$.
- 2 **N.B.** For each Unicode character encoded in UTF-32 we waste at least 11 zero bits.

UTF-16

... map a Unicode character into one or two 16-bit words

- 1 Apply the following mapping scheme:

Unicode range	Word sequence
U+000000–U+00FFFF	□□□□□□□□□□□□□□□□
U+010000–U+10FFFF	110110□□□□□□□□□□ 110111□□□□□□□□□□

- 2 For the range U+000000–U+00FFFF, simply fill the □ positions with the 16 bit of the character code.
(Code ranges U+D800–U+DBFF and U+DC00–U+DFFF are unassigned!)
- 3 For the U+010000–U+10FFFF range, subtract 0x010000 from the character code and fill the □ positions using the resulting 20-bit value.

Example

Unicode character U+012345 ($0x012345 - 0x010000 = 0x02345$):
 UTF-16: 11011000000001000 1101111101000101

UTF-8

N.B. UTF-16 is designed to facilitate efficient and robust decoding:

- If we see a leading 11011 bit pattern in a 16-bit word, we know it is the first **or** second word in a UTF-16 multi-word sequence.
- The sixth bit of the word then tells us if we actually look at the first or second word.

UTF-8 (map a Unicode character into a sequence of 8-bit bytes)

- UTF-8 is of special importance because
 - (a) a stream of 8 bit bytes (*octets*) is what flows over an IP network connection,
 - (b) text-processing software today is built to deal with 8 bit character encodings (*iso-8859-x*, ASCII, *etc.*).

UTF-8 encoding

- 1 Apply the following mapping scheme:

Unicode range	Byte sequence
U+000000–U+00007F	0□□□□□□□
U+000080–U+0007FF	110□□□□□ 10□□□□□□
U+000800–U+00FFFF	1110□□□□ 10□□□□□□ 10□□□□□□
U+010000–U+10FFFF	11110□□□ 10□□□□□□ 10□□□□□□ 10□□□□□□

- 2 The spare bits (□) are filled with the bits of the character code to be represented (rightmost □ is least significant bit, pad to the left with 0-bits).

Examples:

- Unicode character U+00A9 (© sign):
UTF-8: 11000010 10101001 (0xC2 0xA9)
- Unicode character U+2260 (math relation symbol ≠):
UTF-8: 11100010 10001001 10100000 (0xE2 0x89 0xA0)

Advantages of UTF-8 encoding

N.B. UTF-8 enjoys a number of highly desirable properties:

- For a UTF-8 multi-byte sequence, the **length of the sequence** is equal to the number of leading 1-bits (in the first byte), *e.g.*:

11100010 10001001 10100000

(Only single-byte UTF-8 encodings have a leading 0-bit.)

- **Character boundaries** are simple to detect (even when placed at some arbitrary position in a UTF-8 byte stream).
- UTF-8 encoding does not affect (binary) sort order.
- Text processing software which was originally developed to work with the pervasive 7-bit ASCII encoding remains functional. This is especially true for the C programming language and its string (`char[]`) representation.

C and UTF-8

Can you explain the last points made?

XML and Unicode

- A conforming XML parser is required to correctly process UTF-8 and UTF-16 encoded documents (The W3C XML Recommendation predates the UTF-32 definition).
- Documents that use a different encoding *must* announce so using the XML text declaration, *e.g.*

```
<?xml version="1.0" encoding="iso-8859-15"?>
or <?xml version="1.0" encoding="utf-32"?>
```

- Otherwise, an XML parser is encouraged to **guess** the encoding while reading the very first bytes of the input XML document:

Head of doc (bytes)	Encoding guess
0x00 0x3C 0x00 0x3F	UTF-16 (<i>big-endian</i>)
0x3C 0x00 0x3F 0x00	UTF-16 (<i>little-endian</i>)
0x3C 0x3F 0x78 0x6D	UTF-8 (or ASCII, iso-8859-*: <i>erroneous</i>)

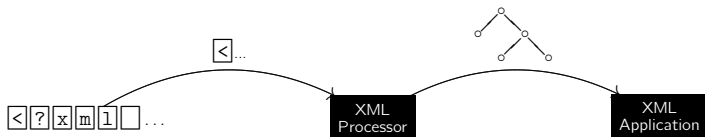
(Notice: < = U+003C, ? = U+003F, x = U+0078, m = U+006D)

The XML Processing Model

- On the **physical** side, XML defines nothing but a **flat text format**, *i.e.*, defines a set of (UTF-8/16) character sequences being well-formed XML.
- Applications that want to analyse and transform XML data in any meaningful manner will find processing flat character sequences hard and inefficient.
- The **nesting** of XML elements and attributes, however, defines a **logical tree-like structure**.

XML Processors

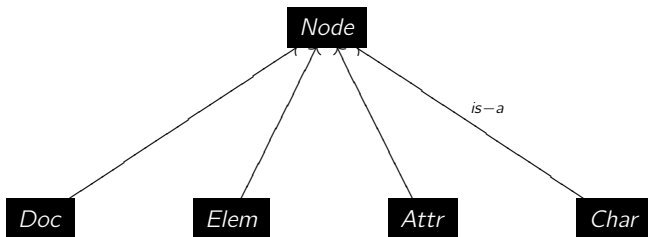
- Virtually all **XML applications operate on the logical tree view** which is provided to them through an **XML Processor** (*i.e.*, the XML parser):



- XML processors are widely available (*e.g.*, Apache's Xerces see <http://xml.apache.org/>).
- How is the XML processor supposed to **communicate the XML tree structure** to the application ... ?

The XML Information Set

- Once an XML processor has checked its XML input document to be well-formed, it offers its application a set of document **properties** (functions).
- The application calls property functions and thus explores the input XML tree as needed.
- An XML document tree is built of different types of **node objects**:



- The set of properties of all document nodes is the document's **Information Set** (see <http://www.w3.org/TR/xml-infoset/>).

Node properties

Node Object Type	Property		Comment
<i>Doc</i>	<i>children</i>	:: <i>Doc</i> → <i>Elem</i>	root element
	<i>base-uri</i>	:: <i>Doc</i> → <i>String</i>	
	<i>version</i>	:: <i>Doc</i> → <i>String</i>	<?xml version="1.0"?>
<i>Elem</i>	<i>localname</i>	:: <i>Elem</i> → <i>String</i>	
	<i>children</i>	:: <i>Elem</i> → (<i>Node</i>)	* ¹
	<i>attributes</i>	:: <i>Elem</i> → (<i>Attr</i>)	
	<i>parent</i>	:: <i>Elem</i> → <i>Node</i>	* ²
<i>Attr</i>	<i>localname</i>	:: <i>Attr</i> → <i>String</i>	
	<i>value</i>	:: <i>Attr</i> → <i>String</i>	
	<i>owner</i>	:: <i>Attr</i> → <i>Elem</i>	
<i>Char</i>	<i>code</i>	:: <i>Char</i> → <i>Unicode</i>	a single character
	<i>parent</i>	:: <i>Char</i> → <i>Elem</i>	

- Read symbol :: as “has type”.
- For any node type τ , (τ) denotes an ordered **sequence** of type τ .

 Make sense of the types of the *Elem* properties *children* (*¹) and *parent* (*²)!

Information set of a sample document

Document δ_0 (weather forecast)

```

1 <?xml version="1.0"?>
2 <forecast date="Thu, Oct 30">
3   <condition>foggy</condition>
4   <temperature unit="Celsius">-1</temperature>
5 </forecast>
  
```

$children(\delta_0) = \varepsilon_1$	$code(\gamma_4) = U+0066 \text{ 'f'}$
$base-uri(\delta_0) = \text{"file:/..."}$	$parent(\gamma_4) = \varepsilon_3$
	⋮
$version(\delta_0) = \text{"1.0"}$	
$localname(\varepsilon_1) = \text{"forecast"}$	$code(\gamma_8) = U+0079 \text{ 'y'}$
$children(\varepsilon_1) = (\varepsilon_3, \varepsilon_9)$	$parent(\gamma_8) = \varepsilon_3$
$attributes(\varepsilon_1) = (\alpha_2)$	$localname(\varepsilon_9) = \text{"temperature"}$
$parent(\varepsilon_1) = \delta_0$	$children(\varepsilon_9) = (\gamma_{11}, \gamma_{12})$
$localname(\alpha_2) = \text{"date"}$	$attributes(\varepsilon_9) = (\alpha_{10})$
$value(\alpha_2) = \text{"Thu, Oct 30"}$	$parent(\varepsilon_9) = \varepsilon_1$
	⋮
$localname(\varepsilon_3) = \text{"condition"}$	
$children(\varepsilon_3) = (\gamma_4, \gamma_5, \gamma_6, \gamma_7, \gamma_8)$	
$attributes(\varepsilon_3) = ()$	
$parent(\varepsilon_3) = \varepsilon_1$	

N.B. Node objects of type *Doc*, *Elem*, *Attr*, *Char* are denoted by δ_i , ε_i , α_i , γ_i , respectively (subscript i makes object identifiers unique).

Working with the Information Set

- The W3C has introduced the XML Information Set to aid the specification of further XML standards.
- We can nevertheless use it to write simple “programs” that explore the XML tree structure. The resulting code looks fairly similar to code we would program using the DOM (Document Object Model, see next chapter).

Example: Compute the list of **sibling** *Elem* nodes of given *Elem* ε (including ε):

siblings (ε) :: *Elem* \rightarrow (*Elem*)

Node ν ;

$\nu \leftarrow \text{parent}(\varepsilon)$;

if $\nu = \delta_{\square}$ then

 | // ν is the Doc node, i.e., ε is the root element
 | return (ε);

else

 | return *children* (ν);

Another Example

Return the text **content** of a given *Doc* δ (the sequence of all Unicode characters δ contains):

```
content ( $\delta$ ) :: Doc  $\rightarrow$  (Unicode)
return collect ((children ( $\delta$ )));
```

```
collect ( $\nu s$ ) :: (Node)  $\rightarrow$  (Unicode)
Node  $\nu$ ;
(Unicode)  $\gamma s$ ;
 $\gamma s \leftarrow ()$ ;
foreach  $\nu \in \nu s$  do
  if  $\nu = \gamma_{\square}$  then
    // we have found a Char node ...
     $\gamma s \leftarrow \gamma s + (\text{code } (\nu))$ ;
  else
    // otherwise  $\nu$  must be an Elem
    node
     $\gamma s \leftarrow \gamma s + \text{collect } (\text{children } (\nu))$ ;
return  $\gamma s$ ;
```

- Example run: $\text{content } (\delta_0) = ('f', 'o', 'g', 'g', 'y', '-', '1')$.

“Querying” using the Information Set

- Having the XML Information Set in hand, we can analyse a given XML document in arbitrary ways, *e.g.*
 - 1 In a given document (comic strip), find **all** *Elem* nodes with local name *bubble* owning an *Attr* node with local name *speaker* and value "Dilbert".
 - 2 List **all** scene *Elem* nodes **containing** a bubble spoken by "Dogbert" (*Attr* *speaker*).
 - 3 Starting in panel number 2 (no *Attr*), find **all** bubbles **following** those spoken by "Alice" (*Attr* *speaker*).
- Queries like these are quite common in XML applications. An XML standard exists (**XPath**) which allows to specify such **document path traversals** in a declarative manner:
 - 1 `//bubble[./@speaker = "Dilbert"]`
 - 2 `//bubble[@speaker = "Dogbert"]/../..`
 - 3 `//panel[@no = "2"]//bubble[@speaker = "Alice"]/following::bubble`

More XML node types ...

- The XML standard defines a number of additional node types that may occur in well-formed documents (and thus in their XML Information Set).
- **CDATA** nodes (embed *unparsed* non-binary character data)

CDATA

```
1 <source>
2   <![CDATA[ May use <, >, and & and
3     anything else freely here ]]>
4 </source>
```

- **Comment** nodes (place comments in XML documents)

Comment

```
1 <proof>
2   <!-- Beware! This has not been properly
3     checked yet... -->
4   ...
5 </proof>
```

... more XML node types

- **PI** nodes (embed *processing instructions* in XML documents)

PI

```
1 <em>
2   <b>Result:</b>
3   <?php sql ("SELECT * FROM ...") ...?>
4 </em>
```

- For a complete list of node types see the W3C XML Recommendation (<http://www.w3.org/TR/REC-xml>).