

XML

- Similar to HTML (Berners-Lee, CERN → W3C)
- use your own tags.
- Amount/popularity of XML data is growing steadily
(faster than computing power)

XML and Databases

Lecture 1
Introduction to XML

Sebastian Maneth
NICTA and UNSW

CSE@UNSW - Semester 1, 2010

XML and Databases

- This course will
- introduce you to the **world of XML** and to the **challenges** of dealing with XML in a RDBMS.

Some of these challenges are

Existing (DB) technology **cannot** be applied to XML data.

XML

- Similar to HTML (Berners-Lee, CERN → W3C)
- use your own tags.
- Amount/popularity of XML data is growing steadily
(faster than computing power)

XML and Databases

Lecture 1
Introduction to XML

Sebastian Maneth
NICTA and UNSW
(today: Kim Nguyen)

CSE@UNSW - Semester 1, 2010

XML and Databases

- This course will
- introduce you to the **world of XML** and to the **challenges** of dealing with XML in a RDBMS.

Some of these challenges are

Existing (DB) technology **cannot** be applied to XML data.

XML and Databases

- This course will
- introduce you to the **world of XML** and to the **challenges** of dealing with XML in a RDBMS.

Some of these challenges are

Existing (DB) technology **cannot** be applied to XML data.



- can handle huge amounts of data stored in **relations**
- storage management
- index structures
- join/sort algorithms
- ...

Databases

XML and Databases

This course will

- introduce you to the **world of XML** and to the **challenges** of dealing with XML in a RDMS.
- **Tree structured data (XML)**
- XML parsers & efficient memory representation
- Query Languages for XML (XPath, XQuery, XSLT...)
- Efficient evaluation using finite-state automata
- Mapping XML to databases
- Advanced Topics (query optimizations, access control, update languages...)

About XML

We will talk about *algorithms and programming techniques* to efficiently manipulate XML data:

- **Regular expressions**, can be used to validate XML data
- **Finite state automata** lie at the heart of highly efficient **XPath implementations**
- **Tree traversals** may be used to preprocess XML trees in order to support **XPath evaluation**, to store **XML trees** in databases, etc.
- Hacking CGI scripts
→ HTML
→ JavaScript
...

You will NOT learn about

- Hacking CGI scripts
→ HTML
→ JavaScript
...
- You NEED to program in
- Java or
→ C++

Course Organization

Lecture	Tuesday, 15:00 – 18:00 ChemicalSc M17 (ex AppliedSc) (K-F10-M17)
Lecturer Consult All email to	Sebastian Maneth Friday, 11:00-12:00 (E508, L5) cs4317@cse.unsw.edu.au

Book	None!
Suggested reading material:	Course slides of Marc Scholl, Uni Konstanz http://www.inf.uni-konstanz.de/dbis/teaching/ws05/06/database-xml/XML_DB.pdf

Book	None!
Theory / PL oriented, book draft: http://arbie.is.s.u-tokyo.ac.jp/~hahosoya/xmlbook/	

About XML

- XML is the World Wide Web Consortium's (W3C, <http://www.w3.org/>) **Extensible Markup Language**
- We hope to convince you that XML is not yet another hyped TLA, but is useful technology.
- You will become best friends with one of the *most important data structures in Computing Science*, the **tree**.
XML is all about tree-shaped data.
- You will learn to apply a number of closely related **XML standards**:

- > Representing data: **XML itself, DTD, XML Schema, XML dialects**
- > Interfaces to connect PLs to XML: **DOM, SAX**
- > Languages to query/transform XML: **XPath, XQuery, XSLT**.

- Tutors ? , Kim Nguyen
- All email to cs4317@cse.unsw.edu.au

Course Organization

Lecture	Tuesday, 15:00 – 18:00 ChemicalSc M17 (ex AppliedSc) (K-F10-M17)
Lecturer Consult All email to	Sebastian Maneth Friday, 11:00-12:00 (E508, L5) cs4317@cse.unsw.edu.au

Programming Assignments

- 5 assignments, due every other Monday. (1st is due 15th March
2nd is due 29th March ...)
- Per assignment: 10 points (total: 60 points) (+2 bonus points)
- Final Exam:**

Final exam: 40 points (must get 16/40 to pass, that is 40%)

Course Organization

Lecture	Tuesday, 15:00 – 18:00 ChemicalSc M17 (ex AppliedSc) (K-F10-M17)
Lecturer Consult All email to	Sebastian Maneth Friday, 11:00-12:00 (E508, L5) cs4317@cse.unsw.edu.au

Programming Assignments

- Tutorial** Tuesday, 12:00-14:00 @ Quadrangle G040 (K-E15-G040) -- before lecture
Wednesday, 12:00-14:00 @ Quadrangle G022 (K-E15-G022)
Wednesday, 14:00-16:00 @ Quadrangle G022 (K-E15-G022)
Thursday, 14:00-16:00 @ Quadrangle G040 (K-E15-G040)
Thursday, 16:00-18:00 @ Quadrangle G022 (K-E15-G022)
- Tutors ? , Kim Nguyen
- All email to cs4317@cse.unsw.edu.au

Outline - Assignments

You can freely choose to program your assignments in

- C / C++, or
- Java

However, your code **must compile with gcc / g++, javac**,
as installed on CSE linux systems!

Submit code (using give) by Monday 23:59 (every other week)

Assignment 4 (harder) gets four weeks / counts double (20 Points)
due date 17th May

Outline - Assignments

1. Read XML, using DOM parser. Create document statistics **13 days**
2. SAX Parse into memory structure: Tree vs DAG **2 weeks**
3. Map XML into RDBMS **2 weeks**
(+1 week break)
4. XPath evaluation over main memory structures **4 weeks**
(+ streaming support)
5. XPath into SQL Translation **2 weeks**

Outline - Lectures

- 1. Introduction to XML, Encodings, Parsers
- 2. Memory Representations for XML: Space vs Access Speed
 - (1) religious
 - (2) practical
 - (3) theoretical / mathematical
- 3. RDBMS Representation of XML
- 4. DTDs, Schemas, Regular Expressions, Ambiguity
- 5. Node Selecting Queries: XPath
- 6. Efficient XPath Evaluation
- 7. XPath Properties: backward axes, containment test
- 8. Streaming Evaluation: how much memory do you need?
- 9. XPath Evaluation using RDBMS
- 10. Properties of XPath
- 11. XSLT
- 12. XQuery
- 13. Wrap up, Exam Preparation, Open questions, etc

Outline

- 1. Three **motivations for XML**
 - (1) religious
 - (2) practical
 - (3) theoretical / mathematical
- 2. **Well-formed XML**
- 3. **Character Encodings**
- 4. **Parsers for XML**
 - parsing into DOM (Document Object Model)

XML Introduction

- Religious** motivation for XML:
- to have **one language** to speak about data.

Lecture 1

XML Introduction



29

XML Motivation (historical)

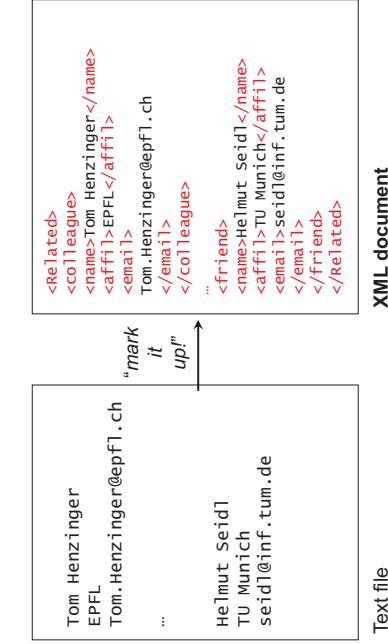
→ XML is a Data Exchange Format

- 1974 SGML (Charles Goldfarb at IBM Research)
- 1989 HTML (Tim Berners-Lee at CERN/Geneva)
- 1994 Berners-Lee founds Web Consortium (W3C)
- 1996 XML (W3C draft, v1.0 in 1998)

<http://www.w3.org/TR/REC-xml/>

27

XML = data + structure



28

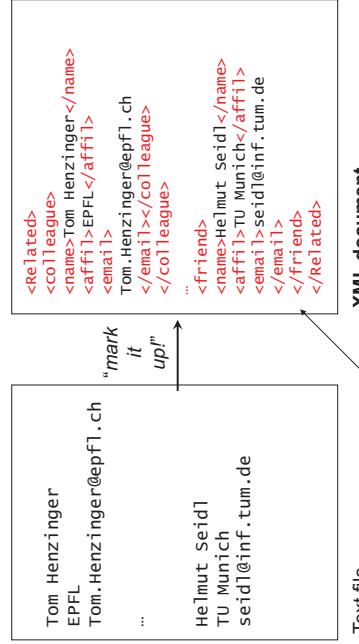
XML Motivation (historical)

→ XML is a Data Exchange Format

- 1974 SGML (Charles Goldfarb at IBM Research)
- 1989 HTML (Tim Berners-Lee at CERN/Geneva)
- 1994 Berners-Lee founds Web Consortium (W3C)
- 1996 XML (W3C draft, v1.0 in 1998)

28

XML = data + structure



Is this a good "template"? What about last/first name?
Several affil's / email's...?

XML Documents

XML Documents

- Ordinary text files (UTF-8, UTF-16, UCS-4 ...)
 - Originates from typesetting/DocProcessing community
 - Idea of labeled brackets ("mark up") for structure is not new!
(already used by Chomsky in the 1960's)
 - Brackets describe a tree structure
 - **Allows applications from different vendors to exchange data!**
 - **standardized, extremely widely accepted!**
-
- Problem** highly verbose, lots of repetitive markup, large files

... instead of writing a parser, you simply fix your own "XML dialect",
by describing all "admissible templates" (+ maybe even the specific
data types that may appear inside).

You do this, using an *XML Type definition language* such
as DTD or Relax NG (Oasis).

Of course, such type definition languages are SIMPLE, because you
want the parsers to be efficient!

They are similar to EBNF. → context-free grammar with reg. expr's in
the right-hand sides. ☺

XML Documents

XML Documents

XML Documents

- Ordinary text files (UTF-8, UTF-16, UCS-4 ...)
 - Originates from typesetting/DocProcessing community
 - Idea of labeled brackets ("mark up") for structure is not new!
(already used by Chomsky in the 1960's)
 - Brackets describe a tree structure
 - **Allows applications from different vendors to exchange data!**
 - **standardized, extremely widely accepted!**
-
- Contra..** highly verbose, lots of repetitive markup, large files

Social Implications!
All sciences (biology, geography, meteorology, astrology...) have own XML "dialects" to store their data optimally

Pro.. we have a standard! A Standard! A STANDARD!
→ ☺ You never need to write a parser again! Use XML! ☺

Example DTD (Document Type Description)

```
Related          → (coleague | friend | family)*
colleague       → (name,affil*,email*)
friend          → (name,affil*,email*)
family          → (name,affil*,email*)
name            → (#PCDATA)
...

```

Element names and their content

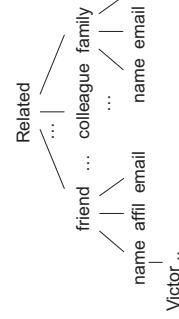
XML Documents

```
Example DTD (Document Type Description)
```

```

Related      → (colleague | friend | family)*
colleague   → (name, affil*, email*)
friend     → (name, affil*, email*)
family     → (name, affil*, email*)
name       → (#PCDATA)
...

```



XML Documents

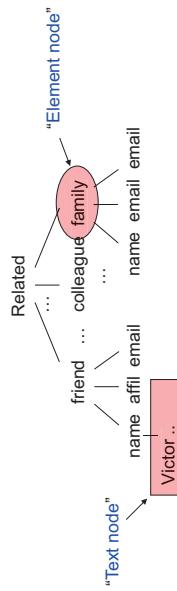
```
Example DTD
```

```

Related      → (colleague | friend | family)*
colleague   → (name, affil*, email*)
friend     → (name, affil*, email*)
family     → (name, affil*, email*)
name       → (#PCDATA)
...

```

Element names and their content



What else: (besides element and text nodes)

- attributes
- processing instructions
- comments
- namespaces
- entity references (two kinds)

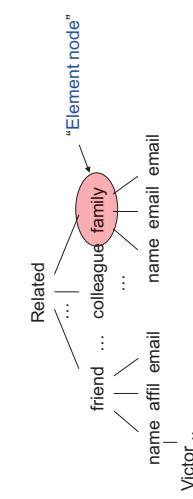
XML Documents

```
Example DTD
```

```

Related      → (colleague | friend | family)*
colleague   → (name, affil*, email*)
friend     → (name, affil*, email*)
family     → (name, affil*, email*)
name       → (#PCDATA)
...

```



What else: (besides element and text nodes)

- attributes
- processing instructions
- comments
- namespaces
- entity references (two kinds)

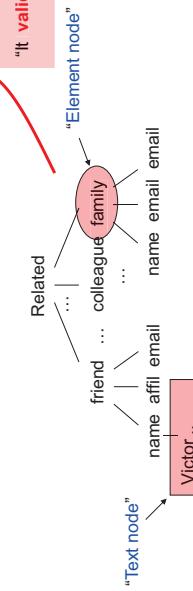
XML Documents

```
Example DTD
```

```

Related      → (colleague | friend | family)*
colleague   → (name, affil*, email*)
friend     → (name, affil*, email*)
family     → (name, affil*, email*)
name       → (#PCDATA)
...

```



What else: (besides element and text nodes)

- attributes
- processing instructions
- comments
- namespaces
- entity references (two kinds)

XML Documents

What else: (besides element and text nodes)

- attributes
- processing instructions
- comments
- namespaces
- entity references (two kinds)

<family rel="brother", age="25">

</family>

...

<name>

...

</name>

Terminology

document is valid wrt the DTD

it validates

What else: (besides element and text nodes)

- attributes
- processing instructions
- comments
- namespaces
- entity references (two kinds)

XML Documents

```
<?php sql ("SELECT * FROM ...") ...?>
See 2.6 Processing Instructions
```

What else:

- **attributes**
- processing instructions
- comments
- namespaces
- entity references (two kinds)

```
<family rel="brother",age="25">
<name>
...
</family>
```

XML Documents

```
<?php sql ("SELECT * FROM ...") ...?>
See 2.6 Processing Instructions
```

What else:

- **attributes**
- processing instructions
- comments <!-- some comment -->
- namespaces
- entity references (two kinds)

```
<family rel="brother",age="25">
<name>
...
</family>
```

<!-- the 'price' element's namespace is http://ecommerce.org/schema -->
<edi:price xmlns:edi='http://ecommerce.org/schema' units='Euro'>32.18</edi:price>

XML Documents

```
<?php sql ("SELECT * FROM ...") ...?>
See 2.6 Processing Instructions
```

What else:

- **attributes**
- processing instructions
- comments <!-- some comment -->
- namespaces
- entity references (two kinds) → character reference
Type <key><ss-than</key> (<) to save options.

```
<family rel="brother",age="25">
<name>
...
</family>
```

<!-- the 'price' element's namespace is http://ecommerce.org/schema -->
<edi:price xmlns:edi='http://ecommerce.org/schema' units='Euro'>32.18</edi:price>

XML Documents

```
<?php sql ("SELECT * FROM ...") ...?>
See 2.6 Processing Instructions
```

What else:

- **attributes**
- processing instructions
- comments <!-- some comment -->
- namespaces
- entity references (two kinds)

```
<family rel="brother",age="25">
<name>
...
</family>
```

<!-- the 'price' element's namespace is http://ecommerce.org/schema -->
<edi:price xmlns:edi='http://ecommerce.org/schema' units='Euro'>32.18</edi:price>

XML Documents

```
<?php sql ("SELECT * FROM ...") ...?>
See 2.6 Processing Instructions
```

What else:

- **attributes**
- processing instructions
- comments <!-- some comment -->
- namespaces
- entity references (two kinds) → character reference
Type <key><ss-than</key> (<) to save options.

```
<family rel="brother",age="25">
<name>
...
</family>
```

This document was prepared on &docdate; and

<!-- the 'price' element's namespace is http://ecommerce.org/schema -->
<edi:price xmlns:edi='http://ecommerce.org/schema' units='Euro'>32.18</edi:price>

Early Markup

The term markup has been coined by the **typesetting** community, not by computer scientist.

- With the advent of printing press, writers and editors used (often marginal notes to instruct printers to
 - Select certain fonts
 - Let passages of text stand out
 - Indent a line of text, etc

The markup language is designed to be easily recognizable in the actual flow of text.

The markup language is designed to be easily recognizable in the actual flow of text.

Fortran 77 source, fixed form, space characters made explicit (.)

Fortran 77

```
1 C THIS,PROGRAM,CALCULATES,THE,CIRCONFERENCE,AND,AREA,OF,HALFCIRCLE,WITH
2 C,RADIUS,R.
3 C
4 C DEFINE,YARIABLE,JAMES,
5 C,.....,RADIUS:,REAL,*,JAMES:
6 C,.....,PI:,REAL,*,PI,3.14159
7 C,.....,CIRCLE:,CIRCONFERENCE,=,2,PI,*R
8 C,.....,AREA:,REAL,*,PI,=,PI*R
9 ****
10 C
11 C,.....,RADIUS,B,CIRCLE,AREA
12 C
13 C,.....,PI,=,3.14159
14 C
15 C,SPECIFY,AVERAGE,PI,F,B:
16 C,.....,A,B,=,A,0
17 C
18 C,CALCULATIONS,=,2,*PI*R
19 C,.....,CIRCLE,=,2,*PI*R
20 C,.....,AREA,=,PI*R*R
21 C
22 C,WRITE,RESULTS:
23 C,.....,WRITE,(6,*,"1.",FOR,A,CIRCLE,OF,RADIUS,' ',B,
24 C,.....,1,THE,CIRCONFERENCE,IS,' ',CIRCLE,
25 C,.....,AND,THE,AREA,IS,' ',AREA
26 C
27 C,.....,END
```

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 20

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 22

Early Markup

Computer scientists adopted the markup idea – originally to **annotate program source code**:

- Design the markup language such that its constructs are **easily recognizable** by a machine.

→ Approaches

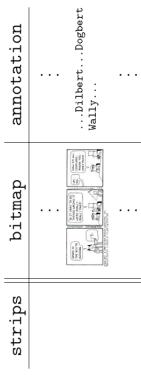
- (1) Markup is written using a **special set of characters**, disjoint from the set of characters that form the tokens of the program
- (2) Markup occurs in places in the source file where program code may not appear (**program layout**).

Example: Fortran 77 fixed form source:

- Fortran statements start in column 7 and do not exceed column 72,
- a Fortran statement longer than 66 chars may be continued on the next line if a character not in {0,1,...} is place in column 6 of the continuing line
- comment lines start with a "C" or ";" in column 1,
- Numeric labels (DO, FORMAT statements) have to be placed in columns 1-5.

An Application of Markup: A Comic Strip Finder Problem:

- **Query a database of comic strips by content.** We want to approach the system with queries like:
 - ① Find all strips featuring Dilbert but not Dogbert.
 - ② Find all strips with Wally being angry with Dilbert.
 - ③ Show me all strips featuring characters talking about XML.
- **Approach:**
- Unless we have nextⁿ generation image recognition software available, we obviously have to **annotate** the comic strips to be able to process the queries above:



Stage 1: ASCII-Level Markup

Sample Markup Application

A Comic Strip Finder

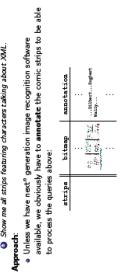
- Next 8 slides from Marc Scholl's 2005 lecture.

An Application of Markup: A Comic Strip Finder
Query: database of comic strips by content. We want to approach the system with queries like

- ① Find all strips featuring Dilbert or not Dogbert.
- ② Find all strips with Wally being angry with Dilbert.
- ③ Show me all strips featuring characters talking about XML.

Approach:

- Unless we have nextⁿ generation image recognition software available, we obviously have to **annotate** the comic strips to be able to process the queries above.



Which kind of queries may we ask now?
And what kind of software do we need to complete the comic strip finder?

Markup Languages	Sample Markup Application
Stage 2: HTML-Style Physical Markup	<pre><h1>Dilbert</h1> <h2>Panel 1</h2> Pointy-Haired Boss Speed is the key to success. <h2>Panel 2</h2> Dilbert Is it okay to do things wrong if we're really really fast? <h2>Panel 3</h2> Pointy-Haired Boss I'm... No. Wally Now I'm all confused. Thank you very much. </pre>
Stage 3: XML-Style Logical Markup	<p>We create a set of tags that is customized to reflect the content of comics, e.g.:</p> <ul style="list-style-type: none"> <character>Dilbert</character> Speed is the key to success <bubble>New types of queries may require new tags. No problem! Resulting set of tags forms a new markup language All tags need to appear in properly nested pairs (<t> ... <s> ... </t>). Tags can be freely nested to reflect the logical structure of comic content. <p>Parsing XML?</p> <p>In comparison to the stage 1 ASCII-level markup parsing, this construction of an XML</p>

Marc H. Schoell (DBIS, Uni KN) XML and Databases Winter 2005/06 24

HTML: Observations

```
----- dibert.xml -----
1 <strip>
2   <panel>
3     <speech>
4       <character>Pointy-Haired Boss</character>
5       <bubble>Speed is the key to success.</bubble>
6     </speech>
7   </panel>
8   <panel>
9     <speech>
10       <character>Dilbert</character>
11       <bubble>Is it okay to do things wrong
12         if we're really, really fast?</bubble>
13     </speech>
14   </panel>
15   <panel>
16     <speech>
17       <character>Pointy-Haired Boss</character>
18       <bubble>Uh... No.</bubble>
19     </speech>
20     <speech>
21       <character>Wally</character>
22       <bubble>I'm all confused.
23         Thank you very much.</bubble>
24     </speech>
25   </panel>
26 </strip>
```

```

<series href="http://www.dilbert.com/s/dilbert/series">
  <character id="Adam" >Author</character>
  <character id="Dilbert" >The Pointy-Haired Boss</character>
  <character id="Wally" >Dilbert's Dog, The Engineer</character>
  <character id="Alice" >Wally's Girlfriend</character>
  <character id="Maggie" >Alice's Daughter</character>
</series>

<spans>
  <span id="1" length="3">
    Spanda length="3">
  </span>
  <span id="2" length="4">
    <scene visible="phb">
      Pointy-haired Boss pointing to presentation slide.
    </scene>
    <bubbles>
      <bubble>phb>Speed is the key to success.</bubble>
      <bubbles>
        <bubble>phb>Speaker is the key to success.</bubble>
      </bubbles>
    </panels>
    <panel no="2">
      <scene visible="wally_dilbert_alice">
        Wally, Dilbert, and Alice sitting at conference table.
      </scene>
      <bubbles>
        <bubble>phb>Toons "question" toons "question" toons "question"
        <bubble>phb>Is it ok to do things wrong if we're real? , really fast?
      </bubbles>
    </panels>
    <panel no="3">
      <scene visible="wally_dilbert_alice">
        Wally turning to Dilbert, angrily
      </scene>
      <bubbles>

```

- We create a **set of tags that is customized** to represent the content of comics, e.g.:
 - <character>Dilbert</character>
 - <bubble>Speed is the key to success</bubble>
- New types of queries may require new tags: No problem for XML!
 - Resulting set of tags forms a new markup language (**XML dialect**).
- All tags need to appear in **properly nested** pairs (e.g.,
 - <t><s>...</s>...</t>...</t>).
- Tags can be freely nested to reflect the **logical structure** of the comic content.

Parsing XML?

In comparison to the stage 1 ASCII-level markup parsing, how difficult do you rate the construction of an XML parser?

Stage 4: Full-Featured XML Markup

- Although fairly simplistic, the previous stage improvement.
- XML comes with a number of additional constructs to convey even more useful information, e.g.
 - **Attributes** may be used to qualify tags (instead of
 - * <question> Is it okay ... ?</question>
 - * <angry> Now I'm ...</angry>
 - * <bubble tone="question">Is it ...?</bubble>
 - * <bubble tone="angry">Now I'm ...</bubble>
- **References** establish links internal to an XML document:
 - Establish link target:
 - * <character id="phb">The Pointy Boss</character>
 - Reference the target:
 - * <bubble speaker="phb">Speed is

Marc H. Schödl (DBIS, Uni Kln) XML and Databases Winter 2005/06 28

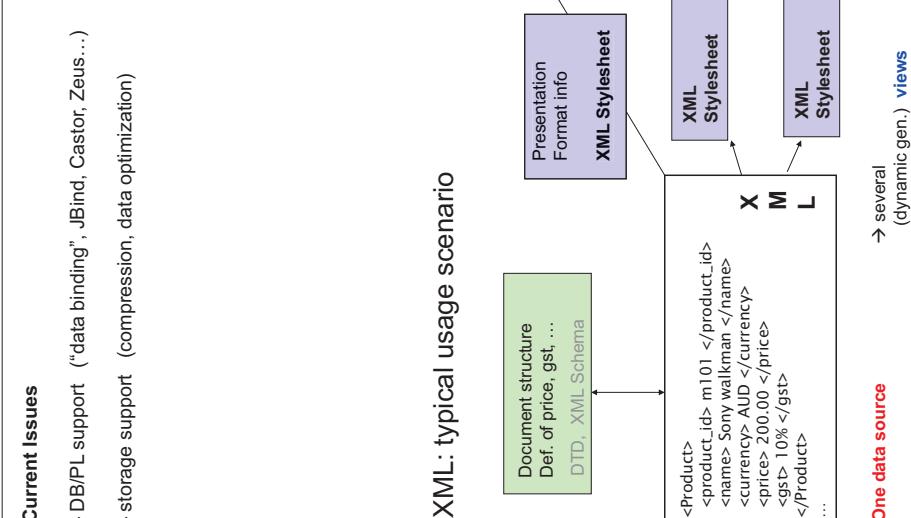
```
<?xml version="1.0" encoding="iso-8859-1"?>
<rss version="2.0" xmlns:dc="http://purl.org/dc/elements/1.1/">
  <channel>
    <title>...</title>
    <link>...</link>
    <description>...</description>
    <language>...</language>
    <copyright>...</copyright>
    <lastBuildDate>...</lastBuildDate>
    <pubDate>...</pubDate>
    <item>
      <title>...</title>
      <link>...</link>
      <description>...</description>
      <dc:creator>...</dc:creator>
      <dc:date>...</dc:date>
      <dc:source>...</dc:source>
      <dc:type>...</dc:type>
      <dc:format>...</dc:format>
      <dc:language>...</dc:language>
      <dc:subject>...</dc:subject>
      <dc:rights>...</dc:rights>
      <dc:relation>...</dc:relation>
      <dc:spatial>...</dc:spatial>
      <dc:temporal>...</dc:temporal>
    </item>
  </channel>
</rss>
```

- HTML defines a number of markup **tags**, some of which are required to match ($< t >$... $</ t >$).
 - Note that HTML tags primarily describe **physical markup** (font size, font weight, indentation, ...)
 - Physical markup is of limited use for the comic strip finder (the **tags do not reflect the structure of the comic content**).

Today, XML has many friends:



XML: typical usage scenario



XML: where is it used ?

- Document formats:**
SVG (vector images), OpenDocument Format (OpenOffice, GoogleDocs), DocBook (Text formatting), EPUB (electronic books), XHTML (web), ...
- Protocol formats:**
SOAP (WebServices), XMPP (Jabber, GoogleTalk), AJAX (custom XML messages used for interactive web sites; Facebook, Gmail, Tweeter,...), RSS feeds (used for blogs/news sites), ...
- Custom data format:**
Bio-informatics, Linguistics, Configuration files, Geographic data (OpenStreetMap, Google maps), Bibliographic Resources (Medline, ADC)

Current Issues

- DB/PL support ("data binding", JBind, Castor, Zeus...)
- storage support (compression, data optimization)

63

2. Well-Formed XML

From the W3C XML recommendation

<http://www.w3.org/TR/REC-xml/>

- "A textual object is **well-formed XML** if,
- (1) taken as a whole, it **matches the production labeled document** given in this specification .."
 - (2) it meets all the **well-formedness constraints** given in this specification .."

document = start symbol of a context-free grammar ("XML grammar")

- (1) contains the **context-free properties** of well-formed XML
- (2) contains the **context-dependent properties** of well-formed XML

- There are 10 WFCs (well-formedness constraints).
- E.g.: **Element Type Match** "The Name in an element's end tag must match the element name in the start tag."

64

2. Well-Formed XML

From the W3C XML recommendation

<http://www.w3.org/TR/REC-xml/>

- "A textual object is **well-formed XML** if,
- (1) taken as a whole, it **matches the production labeled document** given in this specification .."
 - (2) it meets all the **well-formedness constraints** given in this specification .."

document = start symbol of a context-free grammar ("XML grammar")

- (1) contains the **context-free properties** of well-formed XML
- (2) contains the **context-dependent properties** of well-formed XML

- There are 10 WFCs (well-formedness constraints).
- E.g.: **Element Type Match** "The Name in an element's end tag must match the element name in the start tag."
- **Why is this not context-free?**

65

2. Well-Formed XML

From the W3C XML recommendation

<http://www.w3.org/TR/REC-xml/>

- "A textual object is **well-formed XML** if,
- (1) taken as a whole, it **matches the production labeled document** given in this specification .."
 - (2) it meets all the **well-formedness constraints** given in this specification .."

document = start symbol of a context-free grammar ("XML grammar")

- (1) contains the **context-free properties** of well-formed XML
- (2) contains the **context-dependent properties** of well-formed XML

- There are 10 WFCs (well-formedness constraints).
- E.g.: **Element Type Match** "The Name in an element's end tag must match the element name in the start tag."
- **Why is this not context-free?**

2. Well-Formed XML

Context-free grammar in EBNF = System of production rules of the form

Lhs ::= **rhs**

Lhs a nonterminal symbol (e.g., **document**)
rhs a string over nonterminal and terminal symbols.

Additionally (EBNF), we may use regular expressions in **rhs**.
 Such as:

[1]	Char	::= prolog element Misc*
[2]	a	unique character
[3]	S	::= (' ' '\t' '\n' '\r')+
[4]	NameChar	::= (Letter digit '_' '-' '.' Namechar)*
[5]	Name	::= Letter

[22] **prolog** ::= **'<?xml version="1.0' encoding="UTF-8"?>**

[23] **XMLDecl** ::= **'<?xml version="1.0" encoding="UTF-8"?>**

[24] **versionInfo** ::= **'<versionInfo>'**

[25] **Eq** ::= **'='**

[26] **versionNum** ::= **'1.0'**

[39] **element** ::= **EmptyElementTag**

[40] **STag** ::= **StartTag**

[41] **Attribute** ::= **AttributeName EqAttributeValue**

[42] **ETag** ::= **'</' Name ' '>**

[43] **content** ::= **(Element | Reference | CharData)?***

[44] **EmptyElementTag** ::= **'<' Name EqAttributeValue '*' '/>**

[67] **Reference** ::= **EntityRef | CharRef**

[68] **EntityRef** ::= **'&' Name ';'**

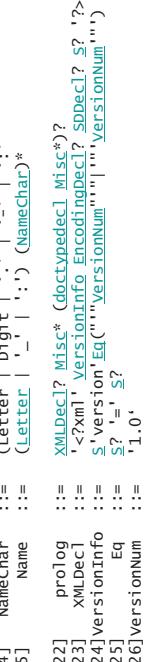
[84] **Letter** ::= **[a-zA-Z]**

[88] **digit** ::= **[0-9]**

XmL Grammar - EBNF-style

Example 1

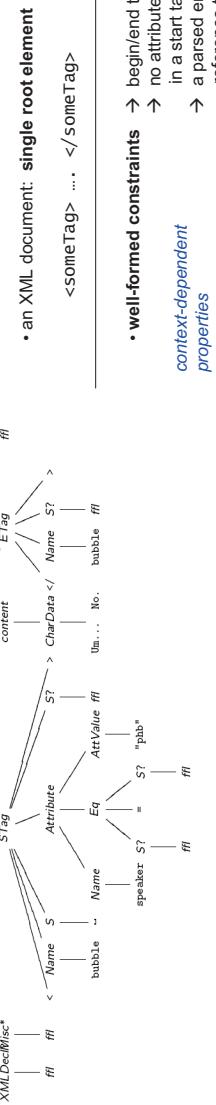
Parse tree for XML input
<bubble speaker="phb">Um... No.</bubble>



2. Well-Formed XML

Terminology

- tag names name, email, author, ...
- start tag <name>, end tag </name>
- elements <name> ... </name>, <author age="99"> ... </author>
- elements may be nested
- empty element <red></red> is abbreviated as <red/>



well-formed constraints

- an XML document: single root element
- no attribute name may appear more than once
- in a start tag or empty element tag
- a parsed entity must not contain a recursive reference to itself, either directly or indirectly

XML Grammar - EBNF-style

As usual, the XML grammar can be systematically 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. 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's leaves from left-to-right, while removing ε symbols).

[1]	r*	denoting ε, r, rr, rrr, ...
[2]	r+	denoting r ε
[3]	r?	denoting a b c
[4]	[abc]	denoting a b z
[5]	[a-z]	character class

[6] **zero or more repetitions**

[7] **one or more repetitions**

[8] **optional r**

[9] **character class**

2. Well-Formed XML

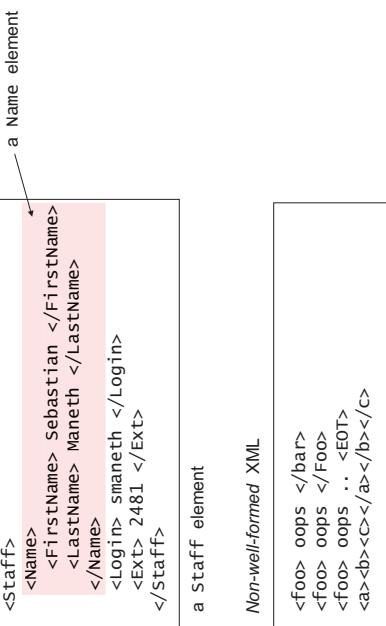
Terminology

- tag names name, email, author, ...
- start tag <name>
- elements <name> ... </name>
- elements may be nested
- empty element <red></red>

- an XML document: single root element
- no attribute name may appear more than once
- in a start tag or empty element tag
- a parsed entity must not contain a recursive reference to itself, either directly or indirectly

Character Encoding

Unicode Transformation Formats



- 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.
- Problem: many such number → character mappings, the so called **encodings** are in use today.
- Due to the huge amount of characters needed by the global computing community today (Latin, Hebrew, Arabic, Greek, Japanese, Chinese ...), **conflicting intersections** between encodings are common.

Example

```
0xcb 0xe4 0xd3 → iso-8859-7 → Æ δ Σ
0xcb 0xe4 0xd3 → iso-8859-15 → É à Ó
```

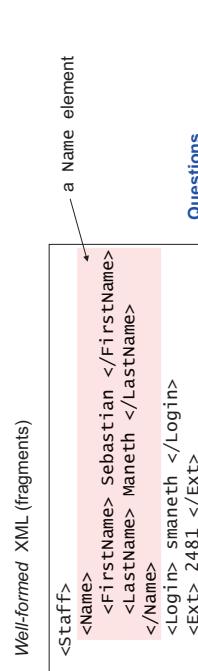
Unicode Transformation Formats

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

UTF-32

- Simply map exactly to the corresponding 32-bit value
- For each Unicode character in UTF-32: waste of at least 11 bits!

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

Questions

How can you implement the three well-formed constraints?

When, during parsing, do you apply the checks?



Unicode Transformation Formats

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

UTF-32

- Simply map exactly to the corresponding 32-bit value
- For each Unicode character in UTF-32: waste of at least 11 bits!

UTF-16

- Map a Unicode character into **one or two 16-bit words**
- U+0000 to U+FFFF map exactly to the corresponding 16-bit value
- above U+FFFF: subtract 0x10000 and then fill the □'s in
1101 10□ □□□ □□□ □□□

- E.g. Unicode character U+012345 (0x012345 – 0x10000 = 0x02345)
UTF-16: 1101 1000 0000 1000 1101 1111 0100 0101

Unicode Transformation Formats

•**UTF-16** works correctly, because the character codes between

1101 10□ □□□ □□□ and

1101 11□ □□□ □□□ (with each □ replaced by a 0)

are left unassigned in Unicode!!! (range 0xD800 – 0xDFFF is reserved)

UTF-8

Maps a unicode character into **1, 2, 3, or 4 bytes**.

Unicode range	Byte sequence
U+000000 → U+00007F	0□□□□□□
U+000080 → U+000FFF	110□□□□ 10□□□□□
U+000800 → U+0FFFFF	1110□□□□ 10□□□□□
U+010000 → U+10FFFF	11110□□□ 10□□□□□

Spare bits (□) are filled from right to left. Pad to the left with 0-bits.

E.g. U+00A9 in **UTF-8** is 11000010 10101001 10100000
U+2260 in **UTF-8** is 11100010 10001001 10100000

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

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

Encoding guess
UTF-16 (little Endian)
0x00 0x3C 0x00 0x3F
0x3C 0x00 0x3F 0x00
0x3C 0x3F 0x78 0x6D
UTF-8

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

81

83

Unicode Transformation Formats

•**UTF-16** works correctly, because the character codes between

1101 10□ □□□ □□□ and

1101 11□ □□□ □□□ (with each □ replaced by a 0)

are left unassigned in Unicode!!! (range 0xD800 – 0xDFFF is reserved)

•**UTF-16** is the **only** encoding in Java: a char is **16 bits** large, **not 8** as in C/C++

WARNING: a "single" character which uses two 16 bits fields is made up of 2 chars!

Example, the string "ö" takes two 16 bit fields:

```
String s = "\uD834\uDD1E";  
s.length(); // returns 2  
s.charAt(1); // returns \uDD1E, it's an INVALID character.  
s.codePointAt(1); // returns the integer 0xD834DD1E
```

UTF-8

Maps a unicode character into **1, 2, 3, or 4 bytes**.

Unicode range	Byte sequence
U+000000 → U+00007F	0□□□□□□
U+000080 → U+000FFF	110□□□□ 10□□□□□
U+000800 → U+0FFFFF	1110□□□□ 10□□□□□
U+010000 → U+10FFFF	11110□□□ 10□□□□□

Spare bits (□) are filled from right to left. Pad to the left with 0-bits.

E.g. U+00A9 in **UTF-8** is 11000010 10101001 10100000
U+2260 in **UTF-8** is 11100010 10001001 10100000

UTF-8

- 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)
- Character boundaries are simple to detect
- UTF-8 encoding does not affect (binary) sort order
- Text processing software designed to deal with 7-bit ASCII remains functional.

- What does "guess the encoding" mean? Under which circumstances does the parser **know** it has determined the correct encoding?
- Are there cases when it canNOT determine the correct encoding?

(especially true for the C programming language and its string (`char[]`) representation)

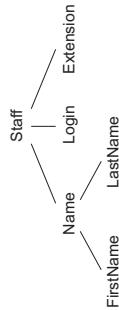
- What about efficiency of the UTF8? For different texts, compare the space requirement in UTF-8/16 and UTF-32 against each other. Which characters do you find above 0xFFFF in Unicode?
- Can you imagine a scenario where UTF-32 is *faster* than UTF-8/16?

82

84

The XML Processing Model

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



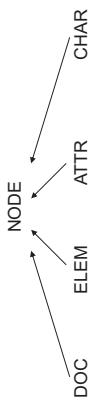
The XML Processing Model

- Virtually all **XML applications operate on the logical tree view** which is provided to them by an **XML processor** (i.e., "parse & store").
- XML processors are widely available (e.g., Apache's Xerces).

How is the XML processor supposed to communicate the **XML tree structure** to the application?

- For many PL's there are "data binding" tools.
Gives very flexible way to get PL view of the XML tree structure.

- through a fixed interface of accessor functions: **The XML Information Set**
- How is the XML processor supposed to communicate the **XML tree structure** to the application?



The accessor functions operate on different types of node objects:

The XML Processing Model

- Virtually all **XML applications operate on the logical tree view** which is provided to them by an **XML processor** (i.e., "parse & store").
- XML processors are widely available (e.g., Apache's Xerces).

How is the XML processor supposed to communicate the **XML tree structure** to the application?

- For many PL's there are "data binding" tools.
Gives very flexible way to get PL view of the XML tree structure.
- But first, let's see what the standard says...

The XML Processing Model

- Virtually all **XML applications operate on the logical tree view** which is provided to them by an **XML processor** (i.e., "parse & store").
- XML processors are widely available (e.g., Apache's Xerces).

How is the XML processor supposed to communicate the **XML tree structure** to the application?

Node type	Property	Comment
DOC	children: DOC→ELEM base-uri: DOC→STRING version : DOC→STRING	root element
ELEM	localName : ELEM→STRING children : ELEM→[NODE] attributes: ELEM→[ATTR] parent : ELEM→NODE	[...] = sequence type
ATTR	localName : ATTR→STRING value : ATTR→STRING owner : ATTR→ELEM	
CHAR	code : CHAR→UNICODE parent : CHAR→ELEM	a single character

```
<?xml version="1.0"?>
<forecast date="Thu, May 16">
  <condition sunny/>
  <temperature unit="celsius">23</temperature>
</forecast>

children(d) = [e1]                                     (= 's')
base-uri(d) = "file: /..."                            (= 's')
version(d) = "1.0"                                     (= 's')

localName(e1) = "forecast"                            (= 'y')
children(e1) = [e2,e3]                                (= 'y')
attributes(e1) = [a1]                                  (= 'y')
parent(e1) = d                                       (= 'y')

localName(a1) = "date"                               (= 'y')
value(a1) = "Thu, May 16"                            (= 'y')
owner(a1) = e1                                       (= 'y')

localName(e2) = "condition"                          (= 'y')
children(e2) = [c1,c2,c3,c4,c5]                      (= 'y')
attributes(e2) = []                                    (= 'y')
parent(e2) = e1                                      (= 'y')

code(c1) = U+0073                                     (= 's')
parent(c1)= e2                                       (= 's')
...
code(c5) = U+0079                                     (= 's')
parent(c5)= e2                                       (= 's')

localName(e3) = "temperature"                         (= 's')
children(e3) = [c6,c7]                                (= 's')
attributes(e3)=[a2]                                    (= 's')
parent(e3) = a1                                       (= 's')
. . .

```

DOM – Document Object Model

Querying the Infoset

Using the Infoset, we can analyze a given XML document in many ways.
For instance:

- Find all ELEM nodes with localname=bubble, owning an ATTR node with localname=speaker and value=Albert.
 - List all panel ELEM nodes containing a bubble spoken by "Albert"
 - Starting in panel 2 (ATTR no), find all bubbles following those spoken by "Alice"
 - Such queries appear very often and can conveniently be described using **XPath queries**:
- ```

value(a1) = "Thu, May 16"
 → //bubble[@speaker="Albert"]
localName(e2) = "condition"
 → //panel[//bubble[@speaker="Albert"]]
children(e2) = [c1,c2,c3,c4,c5]
 → //panel[@no="2"]//bubble[@speaker="Alice"]
attributes(e2) = []
 → //panel[@no="2"]//bubble[@speaker="Alice"]
parent(e2) = e1
 → //panel[//bubble[@speaker="Alice"]]/following::bubble

```

## Questions

## DOM – Document Object Model

→ Language and platform-independent view of XML  
→ DOM APIs exist for many PLs (Java, C++, C, Perl, Python, ...)

DOM relies on two main concepts

- (1) The XML processor constructs the **complete XML document tree** (in-memory)
  - (2) The XML application issues DOM library calls to **explore** and **manipulate** the XML tree, or to **generate** new XML trees.
- Advantages**
- easy to use
  - once in memory, no tricky issues with XML syntax anymore
  - all DOM trees serialize to well-formed XML (even after arbitrary updates)!

- (2) What about WHITESPACE?  
Where in an XML document does it matter, and where not?  
Where in the Infoset Example (previous slide)  
(did we do a mistake? If so, what is the correct Infoset?)

- Advantages**
- easy to use
  - once in memory, no tricky issues with XML syntax arise anymore
  - all DOM trees serialize to well-formed XML (even after arbitrary updates)!

- Disadvantage**      **Uses LOTS of memory!**

## DOM – Document Object Model

→ Language and platform-independent view of XML  
→ DOM APIs exist for many PLs (Java, C++, C, Perl, Python, ...)

DOM relies on two main concepts

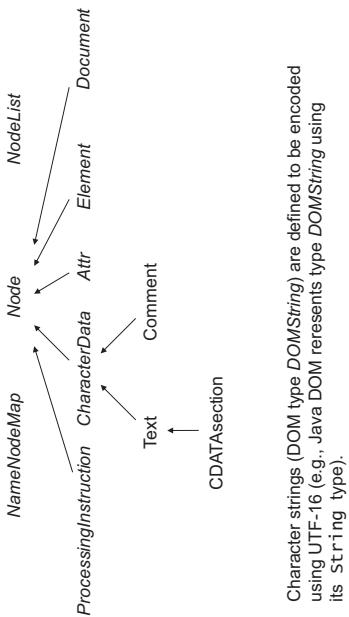
- (1) The XML processor constructs the **complete XML document tree** (in-memory)
- (2) The XML application issues DOM library calls to **explore** and **manipulate** the XML tree, or to **generate** new XML trees.

- (2) Parser triggers "events". Does not store!  
User has to write own code on how to store / process the events triggered by the parser.

- DOM = Document Object Model**

- W3C standard,  
see <http://www.w3.org/TR/REC-DOM-Level-1/>

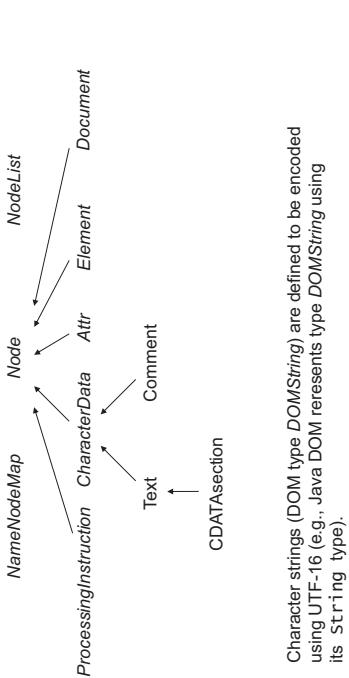
## DOM Level 1 (Core)



E.g. Find all occurrences of Dogbert speaking (attribute `speaker` of element `bubble`)

| The values of <code>nodeName</code> , <code>nodeValue</code> , and attributes vary according to the node type as follows: |                           |
|---------------------------------------------------------------------------------------------------------------------------|---------------------------|
|                                                                                                                           |                           |
| <code>nodeName</code>                                                                                                     | <code>nodeValue</code>    |
| Element                                                                                                                   | tagName                   |
| Attr                                                                                                                      | name of attribute         |
| Text                                                                                                                      | #text                     |
| CDataSection                                                                                                              | #cdata-section            |
| EntityReference                                                                                                           | name of entity referenced |
| Entity                                                                                                                    | entity name               |
| ProcessingInstruction                                                                                                     | target                    |
| Comment                                                                                                                   | #comment                  |
| Document                                                                                                                  | #document                 |
| DocumentType                                                                                                              | document type name        |
| DocumentFragment                                                                                                          | #document-fragment        |
| Notation                                                                                                                  | notation name             |

## DOM Level 1 (Core)



## DOM Level 1 (Core)

Name, value, and attributes depend on the type of the current node.

```

1 // Xerces C++ DOM API support— dogbert.cc (1)
2 #include <dom/DOM.hpp>
3 #include <parsers/DOMParser.hpp>
4
5 void dogbert (DOM::Document d)
6 {
7 DOM::NodeList bubbles;
8 DOM::Node bubble;
9 DOM::NamedNodeMap attrs;
10
11 bubbles = d.getElementsByTagName ("bubble");
12
13 for (unsigned long i = 0; i < bubbles.getLength (); i++) {
14 bubble = bubbles.item (i);
15
16 attrs = bubble.getAttributes ();
17
18 if (attrs != 0)
19 if ((speaker = attrs.getNamedItem ("speaker")) != 0)
20 if (speaker.getAttribute () ==
21 compareString (DOMString ("Dogbert")) == 0)
22 cont << "Found Dogbert speaking." << endl;
23 }
}

```

## DOM Level 1 (Core)

Some details

Creating an `element/attribute` using `createElement/createAttribute` does not wire the new node with the XML tree structure yet.

→ Call `insertBefore`, `replaceChild`, ... to wire a node at an explicit position

DOM type `NodeList` makes up to the lack of collection data types in most programming languages

DOM type `NamedNodeMap` represents an association table (nodes may be accessed by name)

Example:

Methods: `getNamedItem`, `setNamedItem`, ...

## DOM Level 1 (Core)

Some methods

DOM type Method

Node `nodeName` : `DOMString`  
`nodeValue` : `Node`  
`parentNode` : `Node`  
`firstChild` : `Node`  
`nextSibling` : `Node`  
`childNodes` : `NodeList`  
`attributes` : `NamedNodeMap`  
`ownerDocument` : `Document`  
`replaceChild` : `Node`

Creates element with given tag name

Creates comment with given tag name

Creates element with list of all Element nodes in document order

Creates comment with list of all Element nodes in document order

## DOM Level 1 (Core)

Questions

Given an XML file of, say, 50K, how large is its DOM representation in main memory?

How much larger, in the worst case, is a DOM representation with respect to the size of the XML document? (difficult!)

How could we decrease the memory need of DOM, while preserving its functionality?

# END Lecture 1

**Next**

"Event trigger" Parsers for XML:

- Build your own XML data structure  
and fill it up as the parser triggers input "events".