Course: XML Essentials

Pierre Genevès CNRS

University Grenoble Alpes

Programming Languages



Even more "DSLs", e.g. SPARQL (2008), Julia (2009), XQuery (2010) Frequent renewal, rapid evolution

Data

• Often, data are more important than programs

Data

- Often, data are more important than programs
- One reason for this is that data often have a much longer life cycle than programs
- Examples: banking, aeronautical technical documentation, etc.

Example: Aeronautical Technical Documentation

In aeronautics, it is common to find products with life cycles that last for several decades, *e.g.* the B-52:



$\rightarrow\,$ How to ensure long-term access to data?

 \rightarrow How to design systems such that manipulated data can still be read 15 or 50 years later?

→ How to ensure long-term access to data?

 \rightarrow How to design systems such that manipulated data can still be read 15 or 50 years later?

 $\rightarrow~$ An old concern...



La pierre de Rosette (The Rosetta Stone).

\rightarrow How to ensure long-term access to data?

- $\rightarrow\,$ How to design systems such that manipulated data can still be read 15 or 50 years later?
- $\rightarrow~$ An old concern...
 - Can we really do better with computers?



La pierre de Rosette (The Rosetta Stone).

\rightarrow How to ensure long-term access to data?

- ightarrow How to design systems such that manipulated data can still be read 15 or 50 years later?
- $\rightarrow~$ An old concern...
 - Can we really do better with computers?
- \rightarrow A computer museum? \bigcirc



La pierre de Rosette (The Rosetta Stone).

What has not changed for 50 years in information representation?

What has not changed for 50 years in information representation?



- Often, data must be sent to a third-party program/person
- Data must be made explicit (*e.g.* storage in files)
- Naïve¹ approach for defining a *file format*:
 - Define (binary?) representation for data + instructions, e.g. records
 - Write file format spec (v1.0?) + implement parser

¹But used to be very widespread...

- Often, data must be sent to a third-party program/person
- Data must be made explicit (*e.g.* storage in files)
- Naïve¹ approach for defining a *file format*:
 - Define (binary?) representation for data + instructions, e.g. records
 - Write file format spec (v1.0?) + implement parser



¹But used to be very widespread...

- Often, data must be sent to a third-party program/person
- Data must be made explicit (*e.g.* storage in files)
- Naïve¹ approach for defining a *file format*:
 - Define (binary?) representation for data + instructions, e.g. records
 - Write file format spec (v1.0?) + implement parser



MachinCompany File format X Parser P_X





¹But used to be very widespread...

- Often, data must be sent to a third-party program/person
- Data must be made explicit (e.g. storage in files)
- Naïve¹ approach for defining a *file format*:
 - Define (binary?) representation for data + instructions, e.g. records
 - Write file format spec (v1.0?) + implement parser



MachinCompany File format X Parser P_X







TrucMuche SA File format *Z* Parser *P*_Z

¹But used to be very widespread...

- Often, data must be sent to a third-party program/person
- Data must be made explicit (*e.g.* storage in files)
- Naïve¹ approach for defining a *file format*:
 - Define (binary?) representation for data + instructions, e.g. records
 - Write file format spec (v1.0?) + implement parser



- Problems: exchanging data \rightarrow exchanging programs!
 - this approach cannot scale (and costs \$^{\$\$\$})
- $\rightarrow~$ Need for normalization of data exchange

¹But used to be very widespread...

Motivation for XML:

To have one language to describe and exchange data

$\mathsf{XML} = \mathsf{Data}$

Pierre Genevès CNRS pierre.geneves@inria.fr

•••

Nabil Layaïda INRIA nabil.layaida@inria.fr

Text file

XML = Data + Structure



XML = Data + Structure



Tags describe structure, independently from processors (tags are **not** implicit parameters for a given processor, *e.g.* tags are **not** intended for describing presentation)

XML = Data + Structure



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

XML Documents

- Ordinary text files (UTF8, UTF16, ...)
- Originates from typesetting/DocProcessing community
- Idea of labeled brackets ("mark up") for structure is not new (already used by Chomsky in the 1960's)
- Properly nested brackets/tags describe a tree structure
- Allows applications from different vendors to exchange data
- Standardized, extremely widely accepted
- A Lingua franca for structured data exchange...

Standards for Data Exchange



XML History

Ancestors

- 1974 SGML (Charles Goldfarb at IBM Research)
- 1989 HTML (Tim Berners-Lee at CERN, Geneva)
- 1994 Berners-Lee founds World Wide Web Consortium (W3C)
- 1996 XML (W3C draft, v1.0 in 1998) http://www.w3.org/TR/REC-xml/

Initial W3C Goals for XML²

"The design goals for XML are:

- 1. XML shall be straightforwardly usable over the Internet.
- 2. XML shall support a wide variety of applications.
- 3. XML shall be compatible with SGML.
- 4. It shall be easy to write programs which process XML documents.
- 5. The number of optional features in XML is to be kept to the absolute minimum, ideally zero.
- 6. XML documents should be human-legible and reasonably clear.
- 7. The XML design should be prepared quickly.
- 8. The design of XML shall be formal and concise.
- 9. XML documents shall be easy to create.
- 10. Terseness is of minimal importance."

²http://www.w3.org/TR/WD-xml-961114

XML is a Data Exchange Format

- Contra.. extremely verbose, lots of repetitive markup, large files
- **Pro..** answers ambitious goals:
 - $\rightarrow\,$ long-standing (mark-up does not depend on the system where it was created nor on processings)
 - \rightarrow One of the pillars of the web
 - \rightarrow We have **A STANDARD**!
 - $\rightarrow\,$ If you use XML properly, you will never need to write a parser again

XML is a Meta-Language

- XML makes it possible to create Markup-Languages
- Instead of writing a parser, you simply fix your own "XML Dialect"

by describing all "admissible structures" :

- allowed element names
- how they can be assembled together
- maybe even the specific data types that may appear inside

You do this using an XML Type definition language such as DTD³, XML Schema (W3C), or Relax NG (Oasis).

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

³Part of the XML Recommendation: http://www.w3.org/TR/REC-xml/16

XML Document Type Definition

• The XML Recommendation⁴ includes an XML type definition language for specifying **document types**: DTD

people colleague	$\rightarrow \rightarrow$	(colleague friend)* name, affil ⁺ , email	people	
friend	\rightarrow	name, affil*, phone*, email?	name affil email	
Document T	уре			

- Each element is associated with its content model: a reg. expr. (, | ? * +)
- A document type (a set of such associations + a particular root element) describes a set of valid documents used by an organisation

⁴http://www.w3.org/TR/REC-xml/

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- \rightarrow Entity References
- \rightarrow Namespaces
- \rightarrow Text (a specific node kind)

<article class="story" id="news7" datetime="2017-08-08">...</article>

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- \rightarrow Entity References
- \rightarrow Namespaces
- \rightarrow Text (a specific node kind)

<article class="story" id="news7" datetime="2017-08-08">..</article>

<!-- a comment here -->

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- \rightarrow Entity References
- \rightarrow Namespaces
- \rightarrow Text (a specific node kind)

<article class="story" id="news7" datetime="2017-08-08">..</article>

<!-- a comment here -->

- \rightarrow Attributes
- \rightarrow Comments
- → Processing Instructions <?php sql("SELECT * FROM .") ... ?)>
- \rightarrow Entity References
- \rightarrow Namespaces
- \rightarrow Text (a specific node kind)

<article class="story" id="news7" datetime="2017-08-08">..</article>

```
<!-- a comment here -->
```

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- → Entity References
- \rightarrow Namespaces
- \rightarrow Text (a specific node kind)

<?php sql("SELECT * FROM .") ... ?)>
DTD: <!ENTITY notice "All rights...">
instance: Copyright: ¬ice;

<article class="story" id="news7" datetime="2017-08-08">..</article>

```
<!-- a comment here -->
```

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- → Entity References
- \rightarrow Namespaces

```
<?php sql("SELECT * FROM .") ... ?)>
DTD: <!ENTITY notice "All rights...">
instance:  Copyright: &notice;
```

 \rightarrow Text (a specific node kind)

<article class="story" id="news7" datetime="2017-08-08">..</article>

```
<!-- a comment here -->
```

- \rightarrow Attributes
- \rightarrow Comments
- \rightarrow Processing Instructions
- → Entity References
- → Namespaces

<?php sql("SELECT * FROM .") ... ?)>
DTD: <!ENTITY notice "All rights...">
instance: Copyright: ¬ice;

 \rightarrow Text (a specific node kind)

XML Today

"There is essentially no computer in the world, desktop, handheld, or backroom, that doesn't process XML sometimes..."
Some Widespread XML Dialects...

- XHTML (W3C) the XML version of HTML
- SVG (W3C) Animated Vector Graphics
- SMIL (W3C) Synchronized Multimedia Documents, and MMS
- MathML (W3C) Mathematical formulas
- XForms (W3C) Web forms
- Fix, FPML Financial structured products, transactions ...
- CML Chemical molecules

• ...

- X3D (Web3D) 3D Graphics
- XUL (Mozilla), MXML (Macromedia), XAML (Microsoft) Interface Definition Languages
- SOAP (RPC using HTTP), WSDL (W3C), WADL (Sun) Web Services
- RDF (W3C), OWL (W3C) Metadata/Knowledge in the Semantic Web

Outline of the Sequel

- Two notions of correctness:
 - Well-formedness
 - Validity
- Defining your own classes of documents
 - DTDs, XML Schemas
 - Modeling trees and graphs
- Parsing (with or without validation)

XML Defines 2 Levels of Correctness

1. Well-formed XML (minimal requirement)

- The **flat text format** seen on the **physical** side, *i.e.* a set of (UTF8/16) character sequences being well-formed XML
- Ensures data correspond to logical tree-like structures (applications that want to analyse and transform XML data in any meaningful manner will find processing flat character sequences hard and inefficient)
- 2. Valid XML (optional, stricter requirement)
 - More often than not, applications require the XML input **trees** to conform to a **specific XML dialect**, defined by *e.g.* a DTD

Well-Formed XML





| Characters | < | > | " | , | & |
|------------|---|---|---|---|---|
| Entities | < | > | " | ' | & |

- Proper nesting of opening/closing tags
- Shortcut: <e/> for <e></e>
- Every attribute must have a (unique) value
- A document has one and only one root
- No ambiguity between structure and data
- $\rightarrow\,$ Any XML processor considers well-formed XML as a logical tree structure which is:
 - ordered (except attributes!)
 - finite (leaves are empty elements or character data)
- $\rightarrow\,$ It must stop for not well-formed XML.

Valid XML

- The header of a document may include a reference to a DTD: <!DOCTYPE root PUBLIC "public-identifier" "uri.dtd">
- A document with such a declaration must be valid wrt the declared type
- $\rightarrow~$ The parser will validate it

```
Example
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.1//EN"
"http://www.w3.org/TR/xhtml11/DTD/xhtml11.dtd">
<html>
...
</html>
```

Why Validate?

- A document type is a contract between data producers and consumers
- Validation allows:
 - ightarrow a producer to check that he produces what he promised
 - $\rightarrow\,$ a consumer to check what the producer sends
 - \rightarrow a consumer to protect his application
 - ightarrow leaving error detection up to the parser
 - $\rightarrow\,$ simplifying applications (we know where to find relevant information in valid documents)
 - $\rightarrow\,$ delivering high-speed XML throughput (once the input is validated, a lot of runtime checks can be avoided)

Zoom on Document Type Definition (DTD)

• Any element e to be used in the XML dialect needs to be introduced via

<!ELEMENT e cm>

| Content model cm | Valid content | |
|---|--|--|
| ANY | arbitrary well-formed XML content | |
| EMPTY | no child elements allowed (attributes OK) | |
| Reg. exp. over tag names,
#PCDATA, and constructors
 , + * ? | order and occurence of child elements and text content must match the regular expression | |

• Example (XHTML 1.0 Strict DTD): <!ELEMENT img EMPTY>

Reg. Exp. in DTD Content Models

| Reg. Exp. | Semantics |
|----------------|--------------------------------------|
| tagname | element named tagname |
| #PCDATA | text content (parsed character data) |
| c_1, c_2 | c_1 directly followed by c_2 |
| $c_1 \mid c_2$ | c_1 or, alternatively, c_2 |
| c ⁺ | c, one or more times |
| <i>c</i> * | c, zero or more times |
| <i>c</i> ? | optional <i>c</i> |

Example: recipes.xml (fragment)

| ELEMENT</th <th>recipe</th> <th><pre>(title,ingredient*,preparation,comment?,nutrition)></pre></th> <th></th> | recipe | <pre>(title,ingredient*,preparation,comment?,nutrition)></pre> | |
|--|--------|---|--|
| | | | |

<!ELEMENT title (#PCDATA)>

1

2

3

4

- <!ELEMENT ingredient (ingredient*,preparation)?>
- <!ELEMENT preparation (step*)>

Declaring Attributes

- Using the DTD ATTLIST declaration, validation of XML documents is extended to attributes
- The ATTLIST declaration associates a list of attribute names *a_i* with their owning element e:

<!ATTLIST e

>

$$\begin{array}{cccc} a_1 & \tau_1 & d_1 \\ & \dots \\ a_n & \tau_n & d_n \end{array}$$

- \rightarrow The attribute types τ_i define which values are valid for attributes a_i .
- \rightarrow The defaults d_i indicate if a_i is required or optional (and, if absent, if a default value should be assumed for a_i).
- $\rightarrow\,$ In XML, attributes of an element are unordered. The ATTLIST declaration prescribes no order of attribute usage.

• Via attribute types, control over the valid attribute values can be exercised:

Attribute Type τ_i	Semantics
CDATA	character data (no < but <,)
$(v_1 v_2 v_n)$	enumerated literal values
ID	value is document-wide unique identifier for owner element
IDREF	references an element via its ID attribute

____ Example: academic.xml _____

```
<!ELEMENT academic (Firstname, Middlename*, Lastname)>
<!ATTLIST academic
   title (Prof|Dr) #REQUIRED
    team CDATA #IMPLIED
>
   <academic title="Dr" team="Tyrex"> ... </academic></academic></academic>
```

• Attribute defaulting in DTDs:

1

2

3

4

Attribute Default d _i	Semantics
#REQUIRED	element must have attribute a_i
#IMPLIED	attribute <i>a_i</i> is optional
v (a value)	attribute a_i is optinal, if absent, default value v for a_i is assumed
#FIXED v	attribute a_i is optional, if present, must have value v

```
Example: contacts.xml

<!ELEMENT contact (name, email+, phone*)>

<!ATTLIST contact

emailMode (text|xhtml) "text" <!--send safely-->

>
```

Crossreferencing via ID and IDREF

- Well-formed XML documents essentially describe tree-structured data
- Attributes of type ID and IDREF may be used to encode graph structures in XML. A validating XML parser can check such a graph encoding for consistent connectivity.
- To establish a directed edge between two XML nodes a and b:



- 1. attach a unique identifier to node b (using an ID attribute), then
- 2. refer to b from a via this identifier (using an IDREF attribute).
- 3. For an outdegree > 1 (see below), use an IDREFS attribute.



Graphs in XML – An Example

```
Graph.xml -
      <?xml version="1.0"?>
 1
 \mathbf{2}
      <!DOCTYPE graph [
 3
        <!ELEMENT graph (node+) >
        <!ELEMENT node ANY > <!-- attach arbitrary data to a node -->
 4
        <! ATTLIST node
 5
          id ID #REQUIRED
 6
          edges IDREFS #IMPLIED > <!-- we may have nodes with outdegree 0 -->
 7
      1>
 8
9
10
      <graph>
11
        <node id="A">a</node>
        <node id="B" edges="A C">b</node>
12
        <node id="C" edges="D">c</node>
13
        <node id="D">d</node>
14
        <node id="E" edges="D D">e</node>
15
      </graph>
16
```



Example (Character references in "ComicsML")

	ComicsML.dtd	(fragment)
1	strip [</th <th>-</th>	-
2		
3	ELEMENT character (#PCDATA)	
4	ATTLIST character</th <th></th>	
5	id ID	#REQUIRED >
6	ELEMENT bubble (#PCDATA)	
7	ATTLIST bubble</th <th></th>	
8	speaker IDREF	#REQUIRED
9	to IDREFS	#IMPLIED
10	tone (angry question)	#IMPLIED >
1]>	

Validation results (message generated by Apache's Xerces):

- Setting attribute to some random non-existent character identifier: ID attribute 'yoda' was referenced but never declared
- Using a non-enumerated value for attribute tone: Attribute 'tone' does not match its defined enumeration list

A Real-world DTD – GraphML

- GraphML⁵ has been designed to provide a convenient file format to represent arbitrary graphs
- Graphs (element graph) are specified as lists of nodes and edges
- Edges point from source to target
- Nodes and edges may be annotated using arbitrary description and data
- Edges may be directed (+ attribute edgedefault of graph)





¹⁰

⁵http://graphml.graphdrawing.org/

```
GraphML.dtd (main part) ____
        <!-- GRAPHML DTD (flat version) ---->
 1
 2
        <!ELEMENT graphml ((desc)?,(key)*,((data)|(graph))*)>
 3
 4
        <!ELEMENT locator EMPTY>
 \mathbf{5}
        <!ATTLIST locator
 6
                                        #FIXED
                                                  "http://www.w3.org/TR/2000/PR-xlink-20001220/"
                  xmlns:xlink
                               CDATA
 7
                  xlink:href
                               CDATA
                                        #REQUIRED
 8
                  xlink:type
                               (simple) #FIXED
                                                  "simple">
 9
10
        <!ELEMENT desc (#PCDATA)>
11
12
                           ((desc)?,((((data)|(node)|(edge)|(hyperedge))*)|(locator)))>
        <!ELEMENT graph
13
        <!ATTLIST graph
14
                             ID
                                                   #IMPLIED
                  id
15
                  edgedefault (directed|undirected) #REQUIRED>
16
17
        <!ELEMENT node
                        (desc?,(((data|port)*,graph?)|locator))>
18
        <!ATTLIST node
19
                  id
                            ID
                                   #REQUIRED>
20
        <!ELEMENT port ((desc)?,((data)|(port))*)>
21
22
        <! ATTLIST port
23
                  name
                         NMTOKEN #REQUIRED>
^{24}
25
        <!ELEMENT edge ((desc)?,(data)*,(graph)?)>
26
        <!ATTLIST edge
27
                  id
                            ID
                                         #IMPLIED
^{28}
                            IDREF
                                         #REQUIRED
                  source
29
                  sourceport NMTOKEN
                                         #IMPLIED
30
                  target
                            IDREF
                                         #REQUIRED
                  targetport NMTOKEN
31
                                         #IMPLIED
32
                  directed
                            (true|false) #IMPLIED>
33
        <!ELEMENT key (#PCDATA)>
34
35
        <!ATTLIST key
36
                  id TD
                                                                   #REQUIRED
37
                  for (graph|node|edge|hyperedge|port|endpoint|all) "all">
38
39
        <!ELEMENT data (#PCDATA)>
40
        <!ATTLIST data
41
                           IDREF
                  key
                                       #REQUIRED
42
                  id
                           TD
                                       #TMPLTED>
```

Concluding Remarks

• DTD syntax:

- \checkmark **Pro:** compact, easy to understand
- × Con: ?

Concluding Remarks

• DTD syntax:

- \checkmark **Pro:** compact, easy to understand
- × Con: not in XML!
- DTD functionality:
 - × no fine-grained types (everything is character data; what about, *e.g.* integers?)
 - × no further occurence constraints (e.g. cardinality of sequences)

 $\rightarrow\,$ DTD is a very simple but quite limited type definition language

XML Schema

- With XML Schema⁶, W3C provides an XML type definition language that goes beyond the capabilities of the "native" DTD concept:
 - XML Schema descriptions are valid XML documents themselves
 - XML Schema provides a rich set of built-in data types
 - Users can extend this type system via user-defined types
 - XML element (and attribute) types may even be derived by inheritance

XML Schema vs. DTDs

 $\rightarrow~$ Why would you consider its XML syntax as an advantage?

⁶http://www.w3.org/TR/×mlschema-0/

Some XML Schema Constructs

_____ Declaring an element ___

<re><rsd:element name="author"/>

No further typing specified: the author element may contain string values only.

Absence of minOccurs/maxOccurs implies exactly once.

_____ Declaring a typed element <xsd:element name="year" type="xsd:date"/>

Content of year takes the format YYYY-MM-DD. Other simple types: string, boolean, number, float, duration, time, AnyURI, ...

Simple types are considered atomic with respect to XML Schema (e.g., the YYYY
part of an xsd:date value has to be extracted by the XML application itself).

• Non-atomic complex types are built from simple types using type constructors.

	Declaring sequenced content
1	<rsd:complextype name="Characters"></rsd:complextype>
2	<rpre><rsd:sequence></rsd:sequence></rpre>
3	<rpre><xsd:element <="" minoccurs="1" name="character" pre=""></xsd:element></rpre>
4	maxOccurs="unbounded"/>
5	
6	
7	<rsd:complextype name="Prolog"></rsd:complextype>
8	<xsd:sequence></xsd:sequence>
9	<re><rsd:element name="series"></rsd:element></re>
D	<re><rsd:element name="author"></rsd:element></re>
1	<re><rsd:element name="characters" type="Characters"></rsd:element></re>
2	
3	
4	<rpre><rsd:element name="prolog" type="Prolog"></rsd:element></rpre>

An xsd:complexType may be used anonymously (no name attribute).

 With attribute mixed="true", an xsd:complexType admits mixed content. • New complex types may be **derived** from an existing (base) type.

1	<rsd:element name="newprolog"></rsd:element>
2	<rsd:complextype></rsd:complextype>
3	<rsd:complexcontent></rsd:complexcontent>
4	<re><rsd:extension base="Prolog"></rsd:extension></re>
5	<rest:element name="colored" type="xsd:boolean"></rest:element>
6	sd:extension>
7	
8	
9	

• Attributes are declared within their owner element.

</r>

Other xsd:attribute modifiers: use (required, optional, prohibited), fixed, default.

• The validation of an XML document against an XML Schema goes as far as peeking into the **lexical representation** of simple typed values.

Restricting the value space of a simple type (enumeration)
<rsd:simpletype name="Car"></rsd:simpletype>
<restriction base="xsd:string"></restriction>
<re><rsd:enumeration value="Audi"></rsd:enumeration></re>
<rsd:enumeration value="BMW"></rsd:enumeration>
<re><xsd:enumeration value="VW"></xsd:enumeration></re>

• Other facets: length, maxInclusive (upper bound for numeric values)...

40 / 63

Other XML Schema Concepts

- Fixed and default element content,
- support for null values,
- uniqueness constraints, arbitrary keys (specified via XPath)

Intermediate Outline

- Motivations and key principles behind XML
- Languages for defining sets of documents (tree languages)
- Processing XML Documents
- Parsing
 - Two radically different approaches: DOM and SAX
 - Advantages and drawbacks

XML Processing Model

Validation is good

- Validation is better than writing code
- Remember the promise:

"you will never have to write a parser again"

- $\rightarrow\,$ instead, you will need to encode a grammar...
- Virtually all XML applications operate on the logical tree view which is provided to them by an XML parser
- An XML parser can be validating or non-validating
- XML parsers are widely available (*e.g.* Apache's Xerces).
- How is the XML parser supposed to communicate the XML tree structure to the application?

XML Parsers

- Two different approaches:
 - 1. Parser stores document into a fixed (standard) data structure (*e.g.* DOM)

```
parser.parse("foo.xml");
doc = parser.getDocument();
```

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

Next slides on DOM & SAX by Marc H. Scholl (Uni KN)...

DOM—Document Object Model

- With **DOM**, W3C has defined a **language-** and **platform-neutral** view of XML documents.
- DOM APIs exist for a wide variety of—predominantly object-oriented—programming languages (Java, C++, C, Perl, Python, ...).
- The DOM design rests on two major concepts:
 - An XML Processor offering a DOM interface parses the XML input document, and constructs the complete XML document tree (in-memory).
 - The XML application then issues DOM library calls to explore and manipulate the XML document, or generate new XML documents.



• (The complete DOM interface is too large to list here.) Some methods of the principal DOM types *Node* and *Document*:

DOM Type	Method		Comment
Node	nodeName	:: DOMString	redefined in subclasses, <i>e.g.</i> , tag name for <i>Element</i> , "#text" for <i>Text</i> nodes
	parentNode	:: Node	
	firstChild	:: Node	leftmost child node
	nextSibling	:: Node	returns NULL for root element or last child or attributes
	childNodes	:: NodeList	see below
	attributes ownerDocument	:: NameNodeMap :: Document	see below
	replaceChild	:: Node	replace new for old node, returns old
Document	createElement createComment getElementsByTagName	:: Element :: Comment :: NodeList	creates element with given tag name creates comment with given content list of all <i>Elem</i> nodes in document or- der

Example: C++/DOM Code

content.cc (2) V/ Xerces C++ DOM API support (1) 231 #include <dom/DOM.hpp> void content (DOM_Document d) 2 243 #include <parsers/DOMParser.hpp> 254 26 collect (d.getChildNodes ()); 5 void collect (DOM NodeList ns) 276 287 DOM Node n: 29int main (void) 8 30 9 for (unsigned long i = 0: 31 XMLPlatformUtils::Initialize (): 10 i < ns.getLength (); 32 i++){ DOMParser parser: 12 n = ns.item(i):34 DOM_Document doc; 13 35 parser.parse ("foo.xml"); 14 switch (n.getNodeType ()) { 36 doc = parser.getDocument (); 15 case DOM Node::TEXT NODE: 37 cout << n.getNodeValue ().transcode (); 16 38 17 break: 39 content (doc): 18 case DOM_Node::ELEMENT_NODE: 40 collect (n.getChildNodes ()); 19 41return 0; 20 3 42 21 3 22

Now: Find all occurrences of Dogbert speaking (attribute speaker of element bubble) ...

	dogbert.cc (1)
1	// Xerces C++ DOM API support
2	<pre>#include <dom dom.hpp=""></dom></pre>
3	<pre>#include <parsers domparser.hpp=""></parsers></pre>
4	
5	void dogbert (DOM_Document d)
6	-{
7	DOM_NodeList bubbles;
8	DOM_Node bubble, speaker;
9	DOM_NamedNodeMap attrs;
10	
11	<pre>bubbles = d.getElementsByTagName ("bubble");</pre>
12	
13	<pre>for (unsigned long i = 0; i < bubbles.getLength (); i++) {</pre>
14	<pre>bubble = bubbles.item (i);</pre>
15	
16	attrs = bubble.getAttributes ();
17	if (attrs != 0)
18	<pre>if ((speaker = attrs.getNamedItem ("speaker")) != 0)</pre>
19	if (speaker.getNodeValue ().
20	compareString (DOMString ("Dogbert")) == 0)
21	<pre>cout << "Found Dogbert speaking." << endl;</pre>
22	}
23	}

	dogbert.cc (2)
int main (void)	
{	
XMLPlatformUtils::Initialize	();
DOMParser parser;	
DOM_Document doc;	
<pre>parser.parse ("foo.xml");</pre>	
<pre>doc = parser.getDocument ();</pre>	
dogbert (doc);	
_	
return 0;	
}	

DOM—A Memory Bottleneck

• The two-step processing approach (① parse and construct XML tree, ② respond to DOM property function calls) enables the DOM to be "**random access**":

The XML application may explore and update any portion of the XML tree at any time.

• The inherent memory hunger of the DOM may lead to

```
heavy swapping activity
(partly due to unpredictable memory access patterns, madvise() less
helpful)
```

or

even "out-of-memory" failures.

(The application has to be extremely careful with its own memory management, the very least.)

Numbers

DOM and random node access

Even if the application touches a single element node only, the DOM API has to maintain a data structure that represents the **whole XML input** document (all sizes in kB):⁶

XML size	DOM process size DSIZ	$\frac{\text{DSIZ}}{\text{XML size}}$	Comment
7480	47476	6.3	(Shakespeare's works) many elements con- taining small text fragments
113904	552104	4.8	(Synthetic eBay data) elements containing relatively large text fragments

 $^{^{\}rm 6}{\rm The}$ random access nature of the DOM makes it hard to provide a truly "lazy" API implementation.

To remedy the memory hunger of DOM-based processing ...

- Try to **preprocess** (*i.e.*, filter) the input XML document to reduce its overall size.
 - Use an XPath/XSLT processor to preselect *interesting* document regions,
 - ▶ ♀ no updates to the input XML document are possible then,
 - Y make sure the XPath/XSLT processor is *not* implemented on top of the DOM.

Or

• Use a **completely different** approach to XML processing (\rightarrow **SAX**).

SAX—Simple API for XML

- **SAX**⁷ (**Simple API for XML**) is, unlike DOM, *not* a W3C standard, but has been developed jointly by members of the XML-DEV mailing list (*ca.* 1998).
- SAX processors use **constant space**, regardless of the XML input document size.
 - Communication between the SAX processor and the backend XML application does *not* involve an intermediate tree data structure.
 - Instead, the SAX parser sends events to the application whenever a certain piece of XML text has been recognized (*i.e.*, parsed).
 - The backend acts on/ignores events by populating a callback function table.

⁷http://www.saxproject.org/
Sketch of SAX's mode of operations



- A SAX processor reads its input document sequentially and once only.
- No memory of what the parser has seen so far is retained while parsing. As soon as a significant bit of XML text has been recognized, an event is sent.
- The application is able to act on events in parallel with the parsing progress.

SAX Events

• To meet the constant memory space requirement, SAX reports **fine-grained parsing events** for a document:

Event	reported when seen	Parameters sent
startDocument	xml? ⁸	
startElement	$\langle t a_1 = v_1 \dots a_n = v_n \rangle$	$t, (a_1, v_1), \ldots, (a_n, v_n)$
endElement		t Unicodo buffor ptr. longth
comment	c	C
processingInstruction	t pi?	t, pi
	÷	

 $^{^{8}\}textbf{N.B.}$: Event startDocument is sent even if the optional XML text declaration should be missing.

	dilbert.xml
1	xml encoding="utf-8"? *1
2	<bubbles> *2</bubbles>
3	Dilbert looks stunned *3
4	<bubble speaker="phb" to="dilbert"> *4</bubble>
5	Tell the truth, but do it in your usual engineering way
6	so that no one understands you. \star_5
7	*6
8	*7 *8

Event ^{9 10})	Parameters sent
*1	startDocument	
*2	startElement	t = "bubbles"
*3	comment	$c = "_Dilbert looks stunned_"$
*4	startElement	<pre>t = "bubble", ("speaker","phb"), ("to","dilbert")</pre>
*5	characters	<i>buf</i> = "Tell theunderstands you.", <i>len</i> = 99
*6	endElement	t = "bubble"
*7	endElement	t = "bubbles"
*8	endDocument	

 $^9\text{Events}$ are reported in document reading order $\star_1,\,\star_2,\,\ldots,\,\star_8.$

¹⁰**N.B.**: Some events suppressed (white space).

SAX Callbacks

• To provide an efficient and tight **coupling** between the SAX frontend and the application **backend**, the SAX API employs function callbacks^{.11}



Before parsing starts, the application registers function references in a table in which each event has its own slot:

Event	Callback		Event	Callback
startElement endElement	? ?	SAX register(startElement, startElement ()) SAX register(endElement, endElement ())	startElement endElement :	<pre>startElement () endElement ()</pre>

- 2 The application alone decides on the implementation of the functions it registers with the SAX parser.
- Reporting an event *_i then amounts to call the function (with parameters) registered in the appropriate table slot.

¹¹Much like in event-based GUI libraries.



In Java, populating the callback table is done via implementation of the SAX ContentHandler interface: a ContentHandler object represents the callback table, its methods (*e.g.*, public void endDocument ()) represent the table slots.

Example: Reimplement *content.cc* shown earlier for DOM (find all XML text nodes and print their content) using SAX (pseudo code):

content (File f)

// register the callback, // we ignore all other events SAX register (characters, printText); SAX parse (f); return: printText ((Unicode) buf, Int len)

Int i;

```
foreach i \in 1 \dots len do

\  \  print (buf[i]);

return;
```

SAX and the XML Tree Structure

 Looking closer, the order of SAX events reported for a document is determined by a preorder traversal of its document tree¹²:



N.B.: An *Elem* [*Doc*] node is associated with two SAX events, namely *startElement* and *endElement* [*startDocument*, *endDocument*].

¹²Sequences of sibling *Char* nodes have been collapsed into a single *Text* node.

Challenge

• This **left-first depth-first** order of SAX events is well-defined, but appears to make it hard to answer certain queries about an XML document tree.

Sollect all direct children nodes of an *Elem* node.

In the example on the previous slide, suppose your application has just received the *startElement*(t = "a") event \star_2 (*i.e.*, the parser has just parsed the opening element tag <a>).

With the remaining events $\star_3 \ldots \star_{16}$ still to arrive, can your code detect all the immediate children of *Elem* node a (*i.e.*, *Elem* nodes b and c as well as the *Comment* node)?

The previous question can be answered more generally:

SAX events are sufficient to rebuild the complete XML document tree inside the application. (Even if we most likely don't want to.)

SAX-based tree rebuilding strategy (sketch):

- [startDocument] Initialize a stack S of node IDs (e.g. ∈ Z). Push first ID for this node.
- [startElement] Assign a new ID for this node. Push the ID onto S.¹³
- [characters, comment, ...] Simply assign a new node ID.
- [endElement, endDocument]
 Pop S (no new node created).

Invariant: The **top of** *S* holds the identifier of the current **parent node**.

 13 In callbacks (2) and (3) we might wish to store further node details in a table or similar summary data structure.

Final Remarks on SAX

• For an XML document fragment shown on the left, SAX might actually report the events indicated on the right:

	XML fragm	nent			XML + SAX events
1	<affiliation></affiliation>			1	<affiliation>*1</affiliation>
2	AT&T Labs			2	$AT \star_2 \& \star_3 T$ Labs
3				3	$*_4$ $*_5$
	-	*1 *2 *3 *4	<pre>startElement(affiliation) characters("\nAT",5) characters("&", 1) characters("T_Labs\n", 7) andElement(offiliation)</pre>		(affiliation) \n_AT",5) &",1) I_Labs\n",7) affiliation)



White space is reported.

Multiple *characters* **events** may be sent for text content (although adjacent).

(Often SAX parsers break text on entities, but may even report each character on its own.)

Concluding Remarks

We have seen:

- Motivation for XML (where XML originates from and what it is aimed for)
- What is fundamental with XML:
 - $1. \ \text{standard for structured information} \\$
 - 2. independence from processors
 - 3. well-formed documents can be processed as trees
 - 4. users may agree on a dialect and save coding effort, they can exchange valid documents and the schemas...
- How the XML meta-language works, how to define your own XML dialect (using a schema language e.g. DTD, XML Schema...)
- How/when to use the 2 different kinds of XML parsers (DOM, SAX)
- $\rightarrow\,$ Welcome to the world of tree-structured information!