

# Part IV

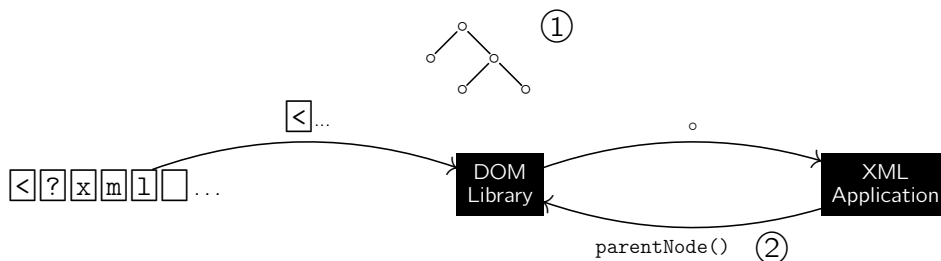
## DOM—Document Object Model

### Outline of this part

- 1 DOM Level 1 (Core)
- 2 DOM Example Code
- 3 DOM—A Memory Bottleneck

## DOM—Document Object Model

- With **DOM**, W3C has defined a **language-** and **platform-neutral** view of XML documents, much like the XML Information Set.
- 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:
  - ① An **XML Processor** offering a DOM interface parses the XML input document, and constructs the **complete XML document tree** (in-memory).
  - ② The **XML application** then issues DOM library calls to **explore** and **manipulate** the XML document, or **generate** new XML documents.



- The DOM approach has some obvious advantages:
  - Once DOM has build the XML tree structure, (tricky) issues of XML grammar and syntactical specifics are void.
  - **Constructing** an XML document using the DOM instead of serializing an XML document manually (using some variation of `print`), ensures **correctness** and **well-formedness**.
    - No missing/non-matching tags, attributes never owned by attributes, ...
  - The DOM can simplify document **manipulation** considerably.
    - Consider transforming

**Weather forecast (English)**

```

1 <?xml version="1.0"?>
2 <forecast date="Thu, Nov 8">
3   <condition>sunny</condition>
4   <temperature unit="Celsius">-1</temperature>
5 </forecast>
```

into

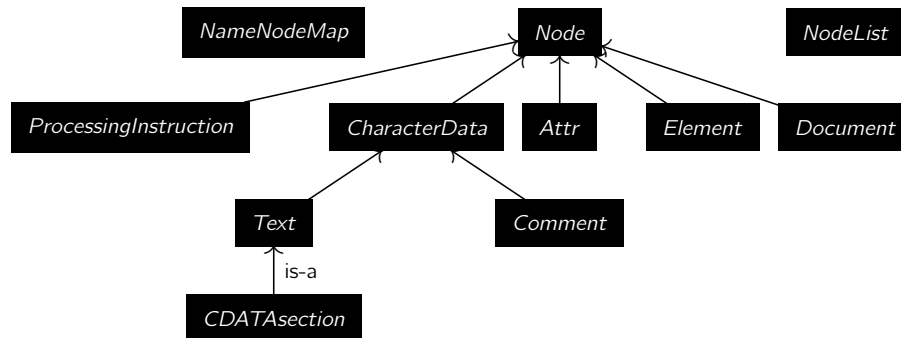
**Weather forecast (German)**

```

1 <?xml version="1.0"?>
2 <vorhersage datum="Do, 8. Nov">
3   <wetterlage>neblig</wetterlage>
4   <temperatur einheit="Celsius">-1</temperatur>
5 </vorhersage>
```

## DOM Level 1 (Core)

- To operate on XML document trees, DOM Level 1<sup>2</sup> defines an inheritance hierarchy of node objects—and methods to operate on these—as follows (excerpt):



- Character strings (DOM type *DOMString*) are defined to be encoded using UTF-16 (e.g., Java DOM represents type *DOMString* using its *String* type).

<sup>2</sup><http://www.w3.org/TR/REC-DOM-Level-1/>

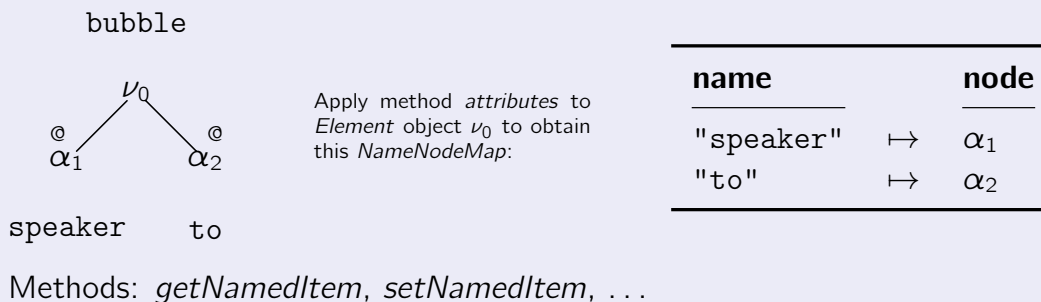
- (The complete DOM interface is too large to list here.) Some methods of the principal DOM types *Node* and *Document*:

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

## Some DOM Details

- Creating an element (or 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* (node sequence) makes up for the lack of collection datatypes in most programming languages.  
Methods: *length*, *item* (node at specific index position).
- DOM type *NamedNodeMap* represents an *association table* (nodes may be accessed by name).

### Example:



## DOM Example Code

- The following slide shows C++ code written against the Xerces C++ DOM API<sup>3</sup>.
- The code implements a variant of the *content :: Doc*  $\rightarrow$  (*Char*):
  - Function *collect()* decodes the UTF-16 text content returned by the DOM and prints it to standard output directly (*transcode()*, *cout*).

### N.B.

- A W3C DOM node type named  $\tau$  is referred to as *DOM $\tau$*  in the Xerces C++ DOM API.
- A W3C DOM property named *foo* is—in line with common object-oriented programming practice—called *getFoo()* here.

<sup>3</sup><http://xerces.apache.org/xerces-c/>

## Example: C++/DOM Code

```

11      content.cc (1)
12 void collect(DOMNodeList *ns)
13 {
14     DOMNode *n;
15
16     for (unsigned long i = 0;
17         i < ns->getLength();
18         i++){
19         n = ns->item(i);
20
21         switch (n->getNodeType())
22         {
23             case DOMNode::TEXT_NODE:
24                 cout << XMLString::transcode(
25                     n->getNodeValue());
26                 break;
27             case DOMNode::ELEMENT_NODE:
28                 collect(n->getChildNodes());
29         }
30     }
31 }

32      content.cc (2)
33 void content(DOMDocument *d)
34 {
35     collect(d->getChildNodes());
36 }
37
38 int main(void)
39 {
40     XMLPlatformUtils::Initialize();
41
42     XercesDOMParser *parser =
43         new XercesDOMParser();
44     DOMDocument *doc;
45
46     parser->parse("dilbert.xml");
47     doc = parser->getDocument();
48
49     content(doc);
50
51     return 0;
52 }

```

**Now:** Find all occurrences of Pointy-Haired Boss (phb) speaking (attribute speaker of element bubble) ...

```

12      dogbert.cc (1)
13 void dogbert(DOMDocument *d)
14 {
15     DOMNodeList *bubbles;
16     DOMNode *bubble, *speaker;
17     DOMNamedNodeMap *attrs;
18
19     bubbles = d->getElementsByTagName(XMLString::transcode("bubble"));
20
21     for (unsigned long i = 0; i < bubbles->getLength(); i++) {
22         bubble = bubbles->item(i);
23
24         attrs = bubble->getAttributes ();
25         if (attrs != NULL)
26             if ((speaker = attrs->getNamedItem(XMLString::transcode("speaker")))
27                 != NULL)
28                 if (XMLString::compareString(speaker->getNodeValue(),
29                     XMLString::transcode("phb")) == 0)
30                     cout << "Found Pointy-Haired Boss speaking." << endl;
31     }
32 }

```

```
dogbert.cc (2)
32 int main (void)
33 {
34     XMLPlatformUtils::Initialize ();
35
36     XercesDOMParser *parser = new XercesDOMParser();
37     DOMDocument *doc;
38
39     parser->parse("dilbert.xml");
40     doc = parser->getDocument();
41
42     dogbert(doc);
43
44     return 0;
45 }
```

## 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
  - ① heavy **swapping** activity  
(partly due to unpredictable memory access patterns, `madvise()` less helpful)  
or
  - ② 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):<sup>4</sup>

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

<sup>4</sup>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,
  -  make sure the XPath/XSLT processor is *not* implemented on top of the DOM.

Or

- Use a **completely different** approach to XML processing (→ **SAX**).