

Objectives

After completing this chapter, you should be able to:

- draw the XDM (XPath/XQuery Data Model) Tree representation for a given XML document.
- explain the most important XDM node types and their essential properties.
- define “document order”.
- mention some details, in which XML data files with the same XDM tree might differ.

Unessential Details (5)

Problem with pattern-Facet:

- In XDM, only the internal representation of data values is stored (i.e. the value itself, not the lexical representation).

Thus, the differences between distinct lexical representations of the same value are considered unessential (e.g. 3, +3, 03).

- However, the `pattern`-facet refers to the lexical representation.
- It is not guaranteed that when a value is printed, a lexical representation is constructed that satisfies the `pattern`.

Type System (1)

- XDM is based on the type system of XML Schema.
- However, a few new types were introduced in the type hierarchy of XML Schema.
- Some of these types are needed for documents (or parts of documents) that where not validated.
 - Partial validation occurs in case of wildcards with skip or lax mode. In XDM, every value has a type.
- The new types are in the XML Schema namespace, they will probably be added in the next version.

Type System (3)

- New types, continued:
 - **dayTimeDuration**: Derived from **duration** with a pattern that permits only days, hours, minutes, seconds (including fractional seconds).

In this way, it can be represented as the total number of seconds.

- **yearMonthDuration**: Derived from **duration** by permitting only year and month components.

In this way, it can be internally represented as the total number of months. It might be inconsistent with XML Schema that it does not have the prefix "g".

Basic Definitions (3)

Atomic Value:

- An atomic value is a value in the value space of an atomic type and is labelled with the name of that type.
- The standard uses this notation in an example:

```
xs:anyURI("http://www.example.com/catalog.xml")
```

This basically looks like a call to a constructor function for the type. It is understood that `xs` is bound to `http://www.w3.org/2001/XMLSchema`.

- Working with pairs of type ID and binary value is common in dynamically typed (“untyped”) programming languages (similar to “variant record”).

Whitespace (3)

- If a schema is used for validation (i.e. the XDM is constructed from a PSVI), text nodes that consist entirely of whitespace are removed if they are children of an element node whose “content-type” property is not `"mixed"`.
- In short:
 - If a validation is done (with respect to a DTD or a schema), there will be no text nodes for whitespace between element content.
 - If the XDM instance is built without validation, such text nodes are constructed.

Whitespace (4)

Some related Remarks:

- XSLT permits to define a set of elements for which pure whitespace child nodes should be removed:

```
<xsl:strip-space elements="..." />
```

- With `xml:space="preserve"` one can specify in the XML data file that whitespace should be preserved.

The other value is "default" which means that the application can do what it wants with the whitespace. This attribute is already introduced in the XML standard. For validating parsers, it must be declared.

Example (5)

- **E1**: Element node (**STUDENT**).
 - `node-kind = "element"`
 - `parent = D`
 - `children = (T1, E2, T2, E3, T3)`
 - `node-name = "STUDENT"`
 - `attributes = (A)`

- **E2**: Element node (**FIRST**).
 - `node-kind = "element"`
 - `parent = E1`
 - `children = (T4)`
 - `node-name = "FIRST"`

Example (6)

- **E3**: Element node (**LAST**).
 - `node-kind = "element"`
 - `parent = E1`
 - `children = (T5)`
 - `node-name = "LAST"`
- **T1**: Text node (whitespace after `<STUDENT ...>`).

This node appears only if the XDM instance is constructed without validation (or if `STUDENT` has a mixed content model). The same applies to T2 and T3.

 - `node-kind = "text"`
 - `parent = E1`
 - `string-value = "\n` " (note: `\n` is not XML syntax)

Example (7)

- **T2:** Text node (whitespace after `</FIRST>`).
 - `node-kind = "text"`
 - `parent = E1`
 - `string-value = "\n "`
- **T3:** Text node (whitespace after `</LAST>`).
 - `node-kind = "text"`
 - `parent = E1`
 - `string-value = "\n"`
- **T4:** Text node (contents of `<FIRST>`).
 - `node-kind = "text"`
 - `parent = E2`
 - `string-value = "Ann"`

Example (8)

- **T5**: Text node (contents of `<LAST>`).
 - `node-kind = "text"`
 - `parent = E3`
 - `string-value = "Smith"`
- **A**: Attribute node (for `SID="101"`)
 - `node-kind = "attribute"`
 - `parent = E1`
 - `node-name = "SID"`
 - `string-value = "101"`
- In addition, there are three namespace nodes (one attached to each element node), see below.

Namespace Nodes (1)

- Although the example contains no explicit namespaces, the prefix `xml` is always bound to

`http://www.w3.org/XML/1998/namespace`

- For each namespace declaration, a namespace node is attached to each element node that is in scope of that namespace declaration.

I.e. not only to the element that explicitly contains the namespace declaration, but also to all descendants, as long as the same prefix is not bound to another URI (or to the empty URI which “undefines” the prefix).

Namespace Nodes (2)

- Namespace nodes cannot be shared between elements.

They have a link to a specific element node in the `parent` property.

- Thus, the example contains already three namespace nodes.

- E.g. **N1**: Namespace node for **STUDENT**-Element:

- `node-kind = "namespace"`

- `parent = E1`

- `node-name = "xml"`

- `string-value =`

- `"http://www.w3.org/XML/1998/namespace"`

Namespace Nodes (3)

- Namespace nodes are accessible in XPath 1.0 by the namespace axis.
 - In XPath 2.0, use of the namespace axis is deprecated. In XQuery 1.0, it does not exist.
 - Instead, one should use XPath functions.
 - These functions do not permit access to the node identity or parent node of a namespace node.
 - Then, namespace nodes can be shared between element nodes.

Namespace Nodes (4)

- Because of the problem with namespace nodes, XDM has two alternative accessor functions (both correspond to the property “`namespaces`”):

- `namespace-nodes`: This returns a sequence of namespace nodes.

If namespace nodes are not needed, this accessor does not have to be implemented.

- `namespace-bindings`: This returns the namespace declarations valid at an element node as a set of prefix/URI pairs.

The standard says that the representation is implementation-dependent, but declares the return type as sequence of `xs:string`.

Document Order (1)

- The document order is a total order on nodes.
- Within a tree, the root node is the first node.

This actually follows from the other rules.

- Every node occurs before all its children and descendants.
- The relative order of siblings is the order in which they occur in the `children` property of their parent.
- Children and descendants of a node occur before following siblings.

Document Order (2)

- Namespace nodes immediately follow the element node with which they are associated.

The relative order of namespace nodes is implementation defined, but stable (i.e. if two namespace nodes of the same element node are compared several times, the result is always the same).

- After the namespace nodes, (or the element node, if there are no namespace nodes), the attribute nodes immediately follow.

The relative order of attribute nodes is implementation defined, but stable.

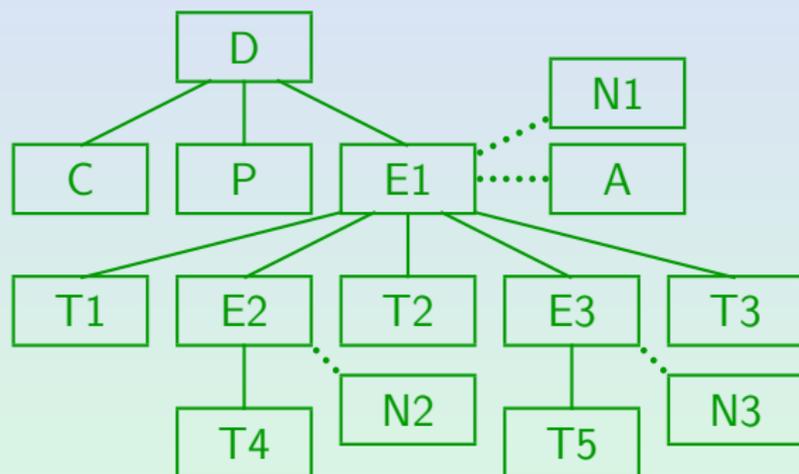
Document Order (3)

- Alternative definitions:
 - The document order is simply the sequence of a pre-order traversal of the tree, with the namespace and attribute nodes listed immediately after their element node.
 - The document order is simply the order of the begin of the start of a node value in the XML document (assuming that namespaces are defined before other attributes).

Document Order (4)

- The relative order of nodes of different trees is implementation-defined (but stable) with the following restriction:
 - If one node of tree T_1 appears between one node of tree T_2 , all nodes of tree T_1 must appear before all nodes of tree T_2 .
- I.e. the document order on the nodes of several trees can be derived from some order on the trees and the order of the nodes within each tree.

Document Order (5)



- Document order: D, C, P, E1, N1, A, T1, E2, N2, T4, T2, E3, N3, T5, T3 (unique in this example).

Exercise

Please draw the XDM tree:

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<!DOCTYPE html PUBLIC
  "-//W3C//DTD XHTML 1.0 Strict//EN"
  "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
<html xmlns="http://www.w3.org/1999/xhtml">
  <head>
    <title>My first XHTML document</title>
  </head>
  <body>
    <h1>Greeting</h1>
    <p>Hi, <a href="http://www.w3.org">W3C</a>!</p>
  </body>
</html>
```

String/Typed Value (1)

- Three important accessor functions are:

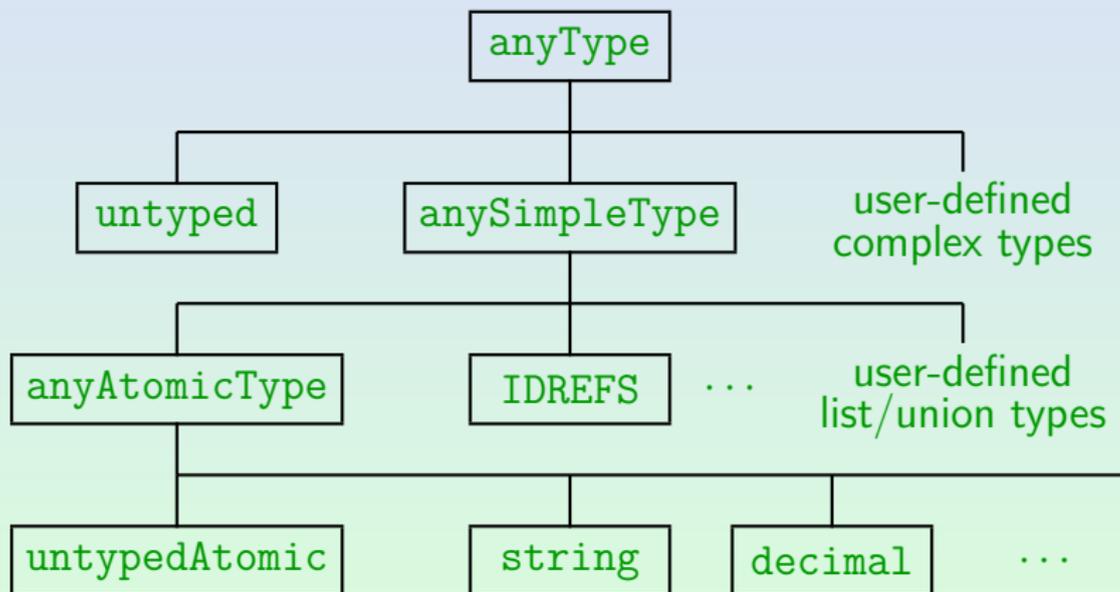
- `string-value`, with return type `xs:string`.
- `typed-value`, with static (declared) return type `xs:anyAtomicType*` (a sequence of atomic values).

The dynamic type of the actually returned values may be more specific. The sequence that can possibly be returned is used for list types.

- `type-name`, with return type `QName?` (i.e. a qualified name or the empty sequence).

This gives the real (dynamic) type of the value that `typed-value` returns. Remember that atomic values have a type attached.

String/Typed Value (2)



String/Typed Value (3)

- Document, element, and attribute nodes have properties “string-value”, “typed-value”, “type-name”.
- However, the corresponding accessor functions are also important for the other node kinds:
 - For text, comment, and processing-instruction nodes, they return the value of the “content” property,
 - for namespace nodes, they return the value of the “uri” property.

For these four kinds of nodes, the **typed-value** is the same as the string-value, the dynamic type of **typed-value** is **xs:string**.

String/Typed Value (4)

- For document nodes,
 - the string-value is simply the content of all text nodes that are descendants of the document node, concatenated in document order.

I don't know why this is a property and not simply computed by the accessor function.

- The typed-value is the same, but with the type `xs:untypedAtomic`.
- In the above example, `string-value` of the document node is `"AnnSmith"` (if validation is done) or `"\n Ann\n Smith\n"` (without validation).

String/Typed Value (5)

- For an attribute node,
 - **string-value** is the normalized attribute value.

The normalization is defined in the XML standard (Section 3.3.3). Basically, all whitespace characters are translated into space characters (CR-LF line ends are translated to a single space). If the data type of the attribute is known and is different from **CDATA**, also leading and trailing spaces are removed, and sequences of space characters are translated into a single space. If a schema is used for validation, the schema normalized value is used (see facet **whiteSpace**). In this case, any lexical representation of the typed value can be returned (see below).

- If no schema is available, the **typed-value** is the **string-value** as an **xs:untypedAtomic**.

String/Typed Value (7)

- String value of element nodes:
 - If the element is declared with a simple type as content, the **string-value** is the schema normalized value of the concatenation of all text node children.

If a schema is used, elements with a simple type as content are treated equivalently to attributes.
Actually, any lexical representation of the typed value can be returned (see below).
 - Otherwise, it is the concatenation of all descendant text nodes (in document order).

String/Typed Value (10)

- An implementation may store the string-value, the typed-value, or both.
- If it stores only the typed-value, it may return any lexical representation of this value as string value.

For example, suppose that `SID` is declared as `integer`. Consider the input:
`<STUDENT SID=" 00101">...</STUDENT>`. The typed-value of the attribute `SID` is the integer 101. An implementation may return `"101"` or `"00101"` as string-value (or other equivalent representations).

- If it stores only the string-value, it must convert the string to the correct member union type.

Although type-name is probably only the union type. Also determining the namespace URI for QName and NOTATION might not be trivial.

Base URI (2)

- If a processing instruction has a base URI different from its parent, it is difficult/impossible to keep this in the external representation.

For all other nodes, the `xml:base` attribute can be used. In this case, one would have to write to that URI, which might be impossible or at least unwanted. The problem is that one cannot use `xml:base` in processing instructions.

This shows that the XML standards do not fit completely together. (newer standards must live with design decisions done in older standards, already in the SGML standard). Things would probably become simpler and more consistent if a complete redesign were done).

- The document node has also a document-uri.

This is an absolute URI that should be used to reload the document if necessary.

