



JACOBS  
UNIVERSITY

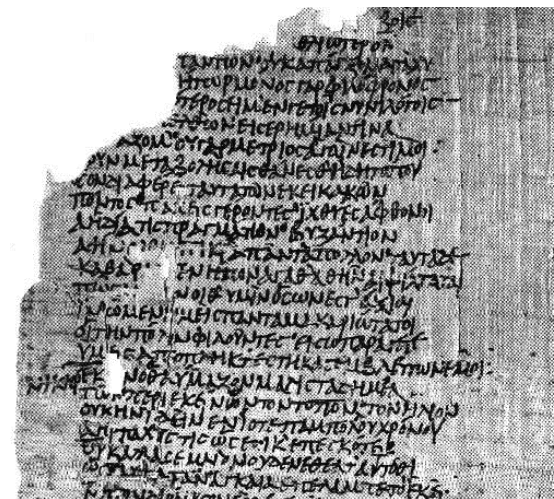
# XML

Ramakrishnan & Gehrke, Chapter 20

[w3schools.com](http://w3schools.com)

# So What's Wrong with HTML?

- "Web services" ultimately means: **programs communicate** via Web!
- Assume end user is not a human, but a program (ex: automated orders)
  - what can it recognize in HTML?
- Actually, HTML does not help:
  - Freedom to hide semantics in layout conventions (bold-face for name...)
  - No **semantic** document structure –  
how to locate address in letter? Check validity?
  - Only one fixed HTML definition – cannot define document types
  - No support for reuse – cannot identify address field across documents
  - Navigation cumbersome – e.g., link into document requires modification of this doc!
- XML to catch semantics of different document types – "semantic Web"



## ■ XML = eXtensible Markup Language

- is not a protocol (uses HTTP!), is not a database
- but is a flexible mechanism for defining domain-specific data exchange formats
- Designed to allow easy implementation

```
<molecule>
  <weight>234.5</weight>
  <spectra>...</spectra>
  <figures>...</figures>
</molecule>
```

## ■ "Extensible": meta language for defining new markup languages

- Each language defines aka "document type", still leaving large degree of variability to single document instances
  - *DTD = Document Type Definition*
- "application of XML" = use XML to define such a new markup language
  - *Example: XHTML = redefinition of HTML in XML*
- Automatic validity checking against definition
- All ASCII, only references to binary data

# Comparison: Degrees of Freedom

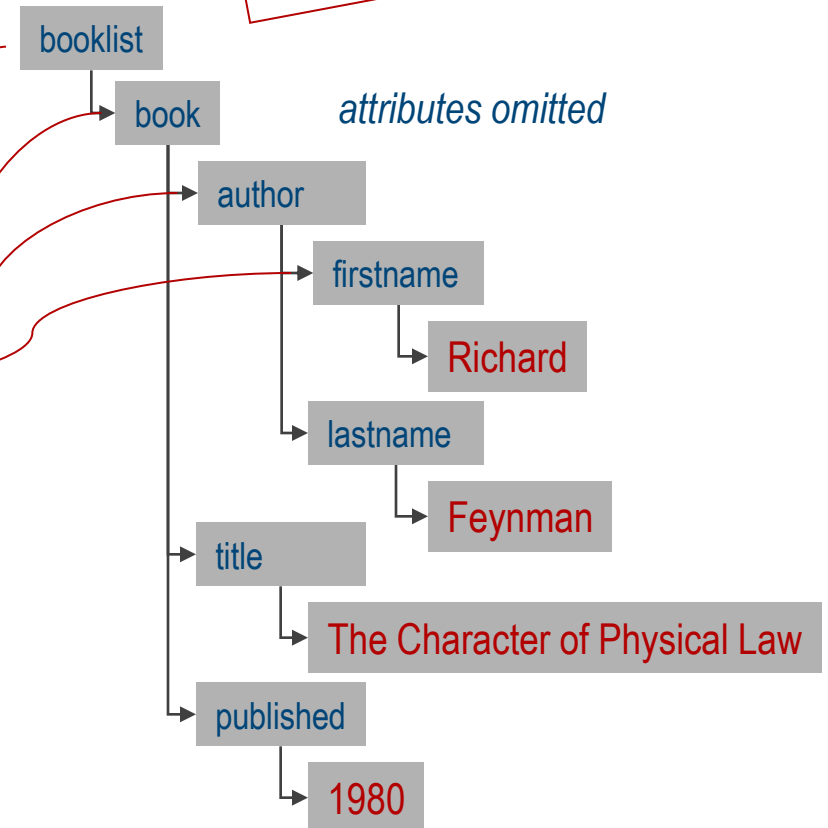
	SGML	XML	HTML
<b>SGML declaration</b>	variable	fixed	fixed
<b>DTD</b>	variable	variable	fixed
<b>Document</b>	variable	variable	variable

# XML Documents Describe Trees

- Tree nodes  
= [ element | attribute | text ] info set items

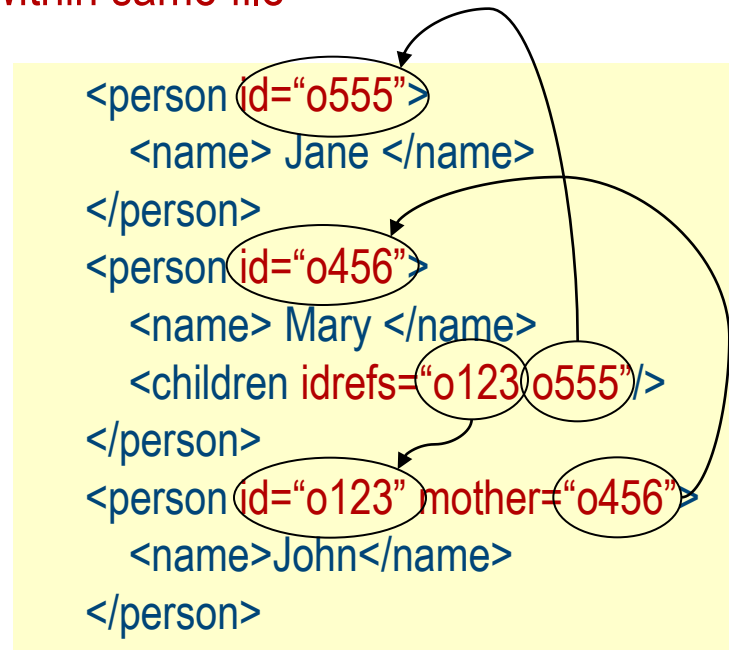
[www.w3schools.com](http://www.w3schools.com)

```
<?xml version="1.0" encoding="UTF-8">
<!DOCTYPE booklist SYSTEM "library.std">
<booklist>
  <book genre="Science" format="Hardcover">
    <author>
      <firstname>Richard</firstname>
      <lastname>Feynman</lastname>
    </author>
    <title>The Character of Physical Law</title>
    <published>1980</published>
  </book>
</booklist>
```



# OIDs and References

- OLD = object identifier
  - Symbolic name for elements **within same file**



- XML just syntax, **referential integrity in no way guaranteed!**

- Problem:
  - multiple markup and vocabularies
  - semantics valid only within XML document
    - $name1 = name2 \iff object(name1) = object(name2)$
- Namespace = collection of names, identified by a URI reference
  - XML vocabularies from different DTDs
  - Ex:

```
<?xml version="1.0"?>
<book:book xmlns:book='urn:loc.gov:books' xmlns:isbn='urn:ISBN:0-395-36341-6'>
  <book:title>Cheaper by the Dozen</book:title>
  <isbn:number>1568491379</isbn:number>
</book:book>
```

# XML DTD: Document Type Definition

- Two ways to define XML named structures (i.e., types):
  - DTD – simpler, less powerful
  - XML Schema – more powerful & complex
- Embedded in document, external reference, or both
- DTD = BNF grammar describing constraints on structure & content
  - what elements and attributes are required / optional
  - like a schema but not really
  - Defines formal structure of the language,



# DTD Example

```
<?xml version='1.0'?>
<!ELEMENT Basket (Cherry+, (Apple | Orange)*) >
  <!ELEMENT Cherry EMPTY>
    <!ATTLIST Cherry flavor CDATA #REQUIRED>
  <!ELEMENT Apple EMPTY>
    <!ATTLIST Apple color CDATA #REQUIRED>
  <!ELEMENT Orange EMPTY>
    <!ATTLIST Orange location 'Florida'>
<!ELEMENT>
```

```
<Basket>
  <Cherry flavor='good'/>
  <Apple color='red'/>
  <Apple color='green'/>
</Basket>
```

```
<Basket>
  <Apple/>
  <Cherry flavor='good'/>
  <Orange/>
</Basket>
```

- **Well-formedness:** minimal set of requirements for a "good" XML document
  - Either DTD declaration or standalone="yes"
  - Exactly one document element
    - *Only comments and PIs outside root element*
    - *Only comments and PIs outside document element (and declaration of course)*
  - Correct cascading of elements (nesting)
    - *For each start tag there is an end tag*
    - *All attribute values in quotes or apostrophies*
    - *No element has two attributes with the same name*
    - *No metatags inside element/attribute values*

- W3C Recommendation (ie: std), 2012
  - Schema for XML document instances, expressed in XML
  - extensible, built-in data type support, modular through namespaces

■ Ex:

```
<?xml version="1.0"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema"
targetNamespace="http://www.w3schools.com"
xmlns="http://www.w3schools.com"
elementFormDefault="qualified">

  <xs:element name="note">
    <xs:complexType>
      <xs:sequence>
        <xs:element name="to" type="xs:string"/>
        <xs:element name="from" type="xs:string"/>
        <xs:element name="heading" type="xs:string"/>
        <xs:element name="body" type="xs:string"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>

</xs:schema>
```

„complex“ = contains other elements

„simple“ = no sub-elements

```
<?xml version="1.0" encoding="UTF-8"?>
<note xmlns="http://w3schools.com"
  xmlns:xsi="http://w3.org/2001/XMLSchema-instance"
  xsi:schemalocation="http://www.w3schools.com note.xsd">
  <to>sample recipient</to>
  <from>sample sender</from>
  <heading>as per phone call</heading>
  <body>Dear X, confirming our phone agreement. Yours, Y</body>
</note>
```

# XML Schema: Some Details

- Simple element of name **xxx** and type **yyy**:

```
<xs:element name="xxx" type="yyy"/>
```

- xs:string, xs:integer, xs:boolean, xs:date, xs:time, ....

```
<lastname>Refsnes</lastname>  
<age>36</age>  
<dateborn>1970-03-27</dateborn>
```

- Attribute **xxx** of type **yyy**:

```
<xs:attribute name="xxx" type="yyy"/>
```

```
<lastname lang="EN">Smith</lastname>
```

- Complex element **employee** with **firstname**, **lastname**:

```
<xs:element name="employee">  
  <xs:complexType>  
    <xs:sequence>  
      <xs:element name="firstname" type="xs:string"/>  
      <xs:element name="lastname" type="xs:string"/>  
    </xs:sequence>  
  </xs:complexType>  
</xs:element>
```

```
<employee>  
  <firstname>John</firstname>  
  <lastname>Smith</lastname>  
</employee>
```

# More on Simple & Complex

Yes, confusing



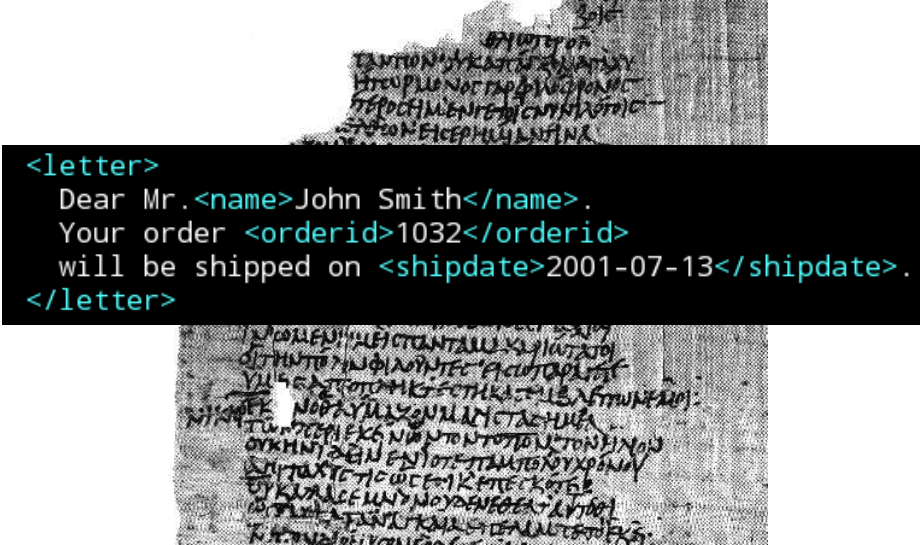
JACOBS  
UNIVERSITY

- **xs:complexContent** → elements & attributes & text outside elements
- **xs:simpleContent** → no sub-elements, only text & attributes
- **xs:complexType** → non-atomic structure
- Indicators define use of elements
  - xs:sequence, xs:choice, ...
  - xs:minOccurs, xs:maxOccurs

```
<xs:sequence>
  <xs:element name="full_name" type="xs:string"/>
  <xs:element name="child_name" type="xs:string"
    minOccurs="0" maxOccurs="5"/>
</xs:sequence>
```

```
<full_name>Tove Refsnes</full_name>
<child_name>Hege</child_name>
<child_name>Stale</child_name>
<child_name>Jim</child_name>
<child_name>Borge</child_name>
```

- Document semantics **captured**
  - Named „markups“ indicate meaning
  - Powerful document structuring, incl. nested elements, attributes, ...
  - Strong typing
  - Automatic validation & processing

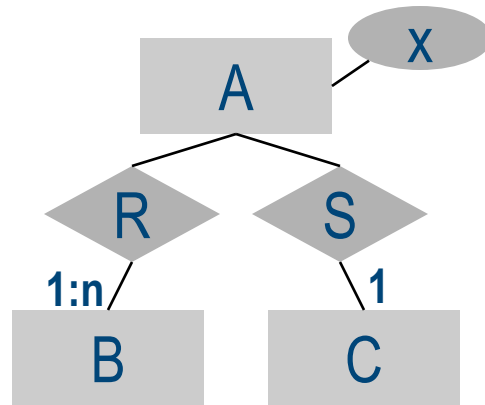


```
<letter>  
  Dear Mr.<name>John Smith</name>.  
  Your order <orderid>1032</orderid>  
  will be shipped on <shipdate>2001-07-13</shipdate>.  
</letter>
```

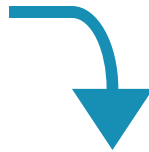
- Drawbacks:
  - inefficient (2x tag names!)
  - DTD rudimentary; XML Schema powerful, but sometimes weird

# Establishing a DTD: Some Practice

- Example 1: From ER



- Exploit constraints!
- Can sometimes leave out relationships



- Example 2: from sample XML

```
<A x="123"/>
  <B /> <B />
  <C>son montuno</C>
</A>
```

- never sure if miniworld caught completely!



# XML Domain Standards

- Many domain-specific schemas existing
- MathML
- Chemical Markup Language
- OpenGIS family of geo service standards: WMS, WFS, WCS, ...
- MusicML
- BIPS (Bank Internet Payments System)
- ...



- XML allows to define **document types** (i.e., data exchange formats)
  - In terms of **infoset items**: elements, attributes, text, (references), ...
  - Becoming *de facto* standard also for, e.g., configuration files (but is no database!)
  - DTD vs XML Schema vs NG Relax vs Schematron vs ...
- Have seen **transformation from ER**
  - To generate XML connectors suitable for the miniworld of database applications
- Many facets not covered, such as:
  - XSLT: transforming XML to, e.g., HTML
  - XML DOM

# Practising DOM Trees

- Let's consider XHTML
- DOM tree?
  - Attributes prefixed with „@“
- References?

```
<html>
<head> <title>Friendly Homepage</title> </head>
<body>
Hello World.
<a name= "picture" > <img src= "world.jpg" /> </a>
Click <a href= "#picture">here</a> for a picture.
</body>
</html>
```

- We have seen the content models
  - elements only; text only; empty; mixed content (text and elements)
  - ...do you remember how to specify them?
- create a DTD / XML Schema for a person's address
  - think of optional parts
  - why not simply use #CDATA for the DTD address?
- think about the zip code
  - what can you do with a DTD / with XML Schema?
  - what would you like to do?
- what could empty elements be good for?