





Semistructured Data (cont.)
 Its features facilitate: the integration of information from different sources the exchange of information between applications
 Example: company A receives data from company B A only cares about certain fields in certain types of records B's data includes: other types of records other fields within the records that company A cares about with semistructured data, A can easily recognize and ignore unexpected elements the exchange is more complicated with structured data



XML Elements
 An XML <i>element</i> is: a begin tag an end tag (in some cases, this is merged into the begin tag) all info. between them. example: name>CS 460
 An element can include other nested child elements. <course> <name>CS 460</name> <begin>1:25</begin> </course>
 Related XML elements are grouped together into <i>documents</i>. may or may not be stored as an actual text document



Att	ridutes vs. Child	Elements
	attribute	child element
number of occurrences	at most once in a given element	an arbitrary number of times
value	always a string	can have its own children

• The string values used for attributes can serve special purposes (more on this later)





Specifying a Separate Schema

- XML doesn't require a separate schema.
- However, we still need one if we want programs to:
 - easily process XML documents
 - · validate the contents of a given document
- The resulting schema can still be semistructured.
 - · for example, can include optional components
 - more flexible than ER models and relational schema

Special Types of Attributes ID an identifier that must be unique within the document (among *all* ID attributes – not just this attribute) IDREF a single value that is the value of an ID attribute elsewhere in the document IDREFs a *list* of ID values from elsewhere in the document













XPath Expressions (cont.) • Attribute names are preceded by an @ symbol: • example: //person/@pid • selects all pid attributes of all person elements • We can specify a particular document as follows: document("doc-name") path-expression • example: document("university.xml")//course/start













Predicates in XPath Expressions (cont.)
<room> <building>CAS</building><room_num>212</room_num> </room>
<office> <building>CAS</building><room_num>100</room_num> </office>
<room> <building>KCB</building><room_num>101</room_num> </room>
<office> <building>PSY</building><room_num>228D</room_num> </office>
 If there are other elements that also have nested room_num and building elements (like office elements above)
 //room_num[/building="CAS"] will get room_num children from all such elements with a building child = "CAS"
 //room[building="CAS"]/room_num will only get room_num children from room elements with a building child = "CAS"







FLWOR Expressions
<pre>for \$r in //room[contains(name, "CAS")], \$c in //course let \$e := //person[contains(@enrolled, \$c/@id)] where \$c/@room = \$r/@id and count(\$e) > 20 order by \$r/name return (\$r/name, \$c/name)</pre>
 The for clause is like the FROM clause in SQL. the query iterates over all combinations of values from its XPath expressions (like Cartesian product!) query above looks at combos of CAS rooms and courses
 The let clause is applied to each combo. from the for clause. each variable gets the <i>full set</i> produced by its XPath expr. unlike a for clause, which assigns the results of the XPath expression one value at a time

FLWOR Expressions (cont.)

```
for $r in //room[contains(name, "CAS")],
    $c in //course
let $e := //person[contains(@enrolled, $c/@id)]
where $c/@room = $r/@id and count($e) > 20
order by $r/name
return ($r/name, $c/name)
```

- The where clause is applied to the results of for and let.
- If the where clause is true, the return clause is applied.
- The order by clause can be used to sort the results.

```
Note: The Location of Predicates
 for $r in //room[contains(name, "CAS")],
     $c in //course
let $e := //person[contains(@enrolled, $c/@id)]
where $c/@room = $r/@id and count($e) > 20
order by $r/name
 return ($r/name, $c/name)

    It's sometimes possible to move components of the

  where clause up into the for clause as predicates.
• In the above query, we could move the first condition up:
 for $r in //room[contains(name, "CAS")],
     $c in //course[@room = $r/@id]
 let $e := //person[contains(@enrolled, $c/@id)]
 where count($e) > 20
 order by $r/name
 return ($r/name, $c/name)
```





Reshaping the Output We can reshape the output by constructing new elements: for \$c in //course where c/start > "11:00"return <after11-course> {\$c/name/text(), " - ", \$c/start/text()} </after11-course> the text() function gives just the value of a simple element without its start and end tags · when constructing a new element, need curly braces around expressions that should be evaluated • otherwise, they'll be treated as literal text that is the value of the new element here again, use commas to separate items because we're using text(), there are no newlines after the name and start time · we use a string literal to put something between them



for vs. let

· Here's an example that illustrates how they differ:

- for each value of \$d, the let clause assigns to \$e the *full set* of emp elements from that department
- the where clause limits us to depts with >= 10 employees
- we create a new element for each such dept.
- · we use functions on the set \$e and on values derived from it









