

XPath

Querying and Storing XML

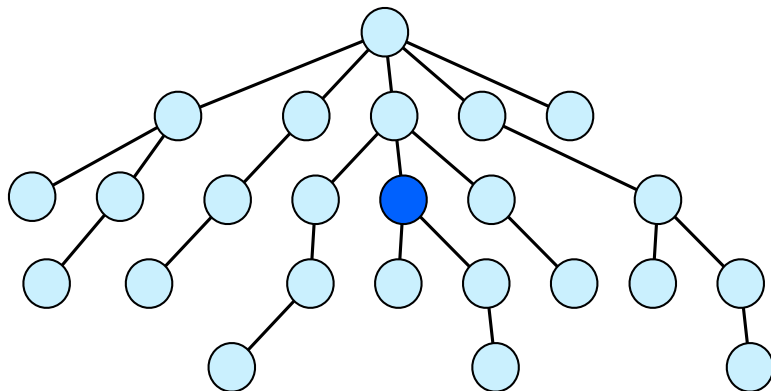
Week 2
XPath & XQuery
January 22-25, 2013

- A language of “path expressions”
- Loosely related to “file paths”
 - root (/)
 - sequential composition (p/q)
 - wildcards (*)
 - axis steps (child, parent, descendant, etc.)
- also: filters, text nodes, label tests
 - plus positional & string functions
- Used for navigation
 - component of XSLT, XQuery, etc.

QSX

January 22-25, 2013

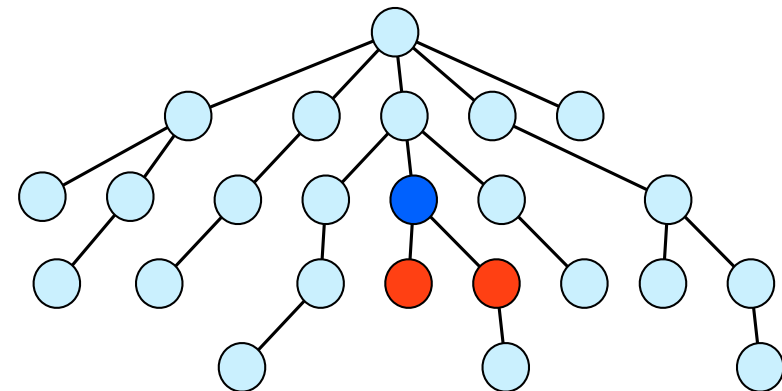
Context node (starting point)



QSX

January 22-25, 2013

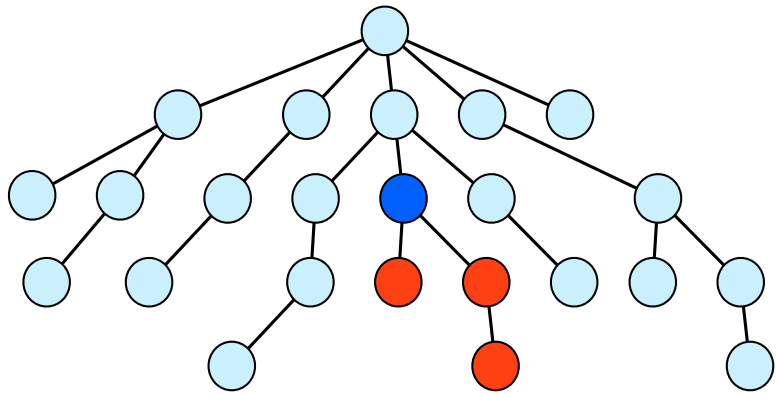
Child



QSX

January 22-25, 2013

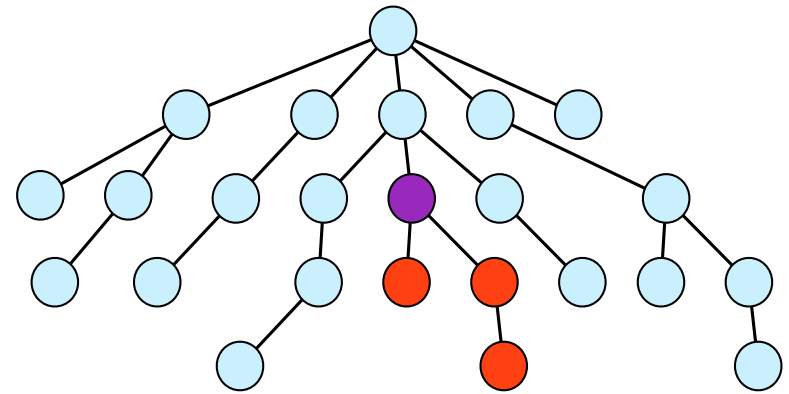
Descendant



Qsx

January 22-25, 2013

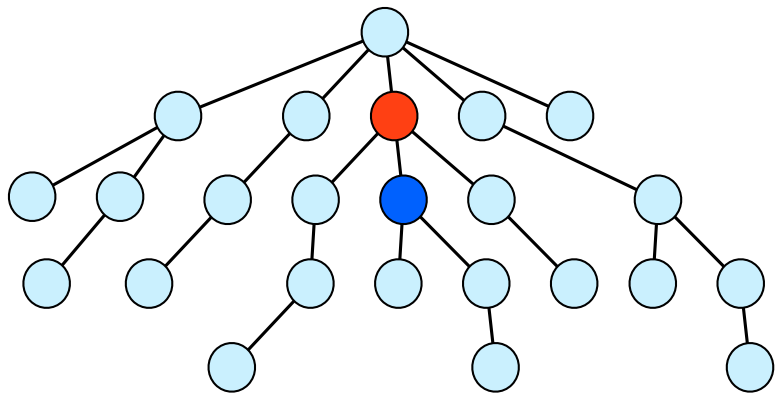
Descendant-or-self



Qsx

January 22-25, 2013

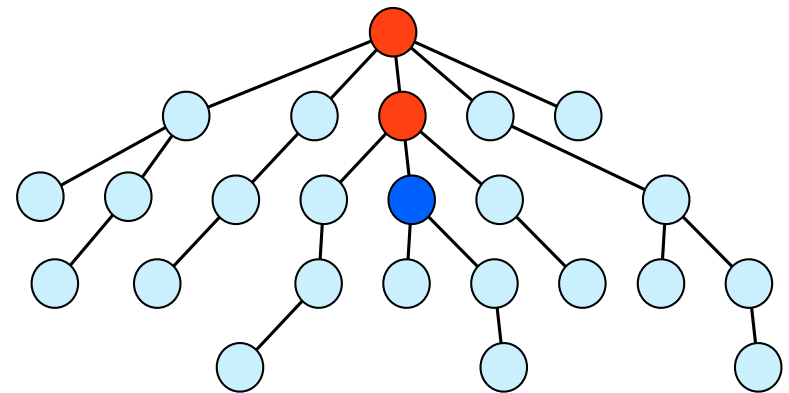
Parent



Qsx

January 22-25, 2013

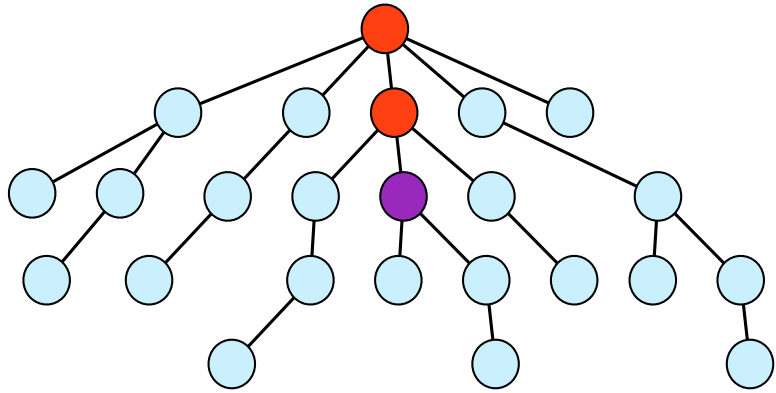
Ancestor



Qsx

January 22-25, 2013

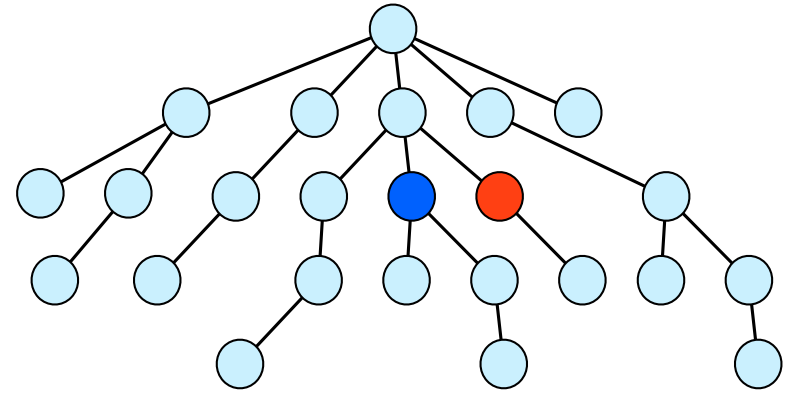
Ancestor-or-self



Qsx

January 22-25, 2013

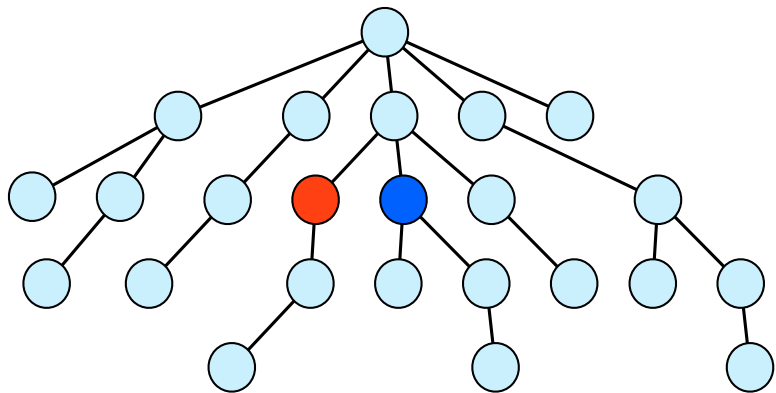
Following-sibling



Qsx

January 22-25, 2013

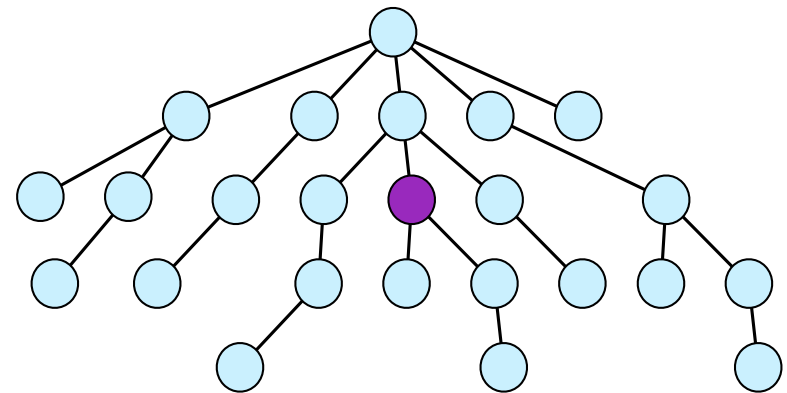
Preceding-sibling



Qsx

January 22-25, 2013

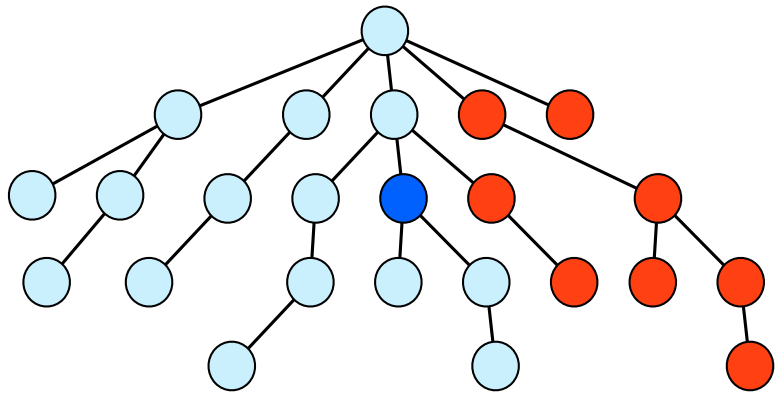
Self



Qsx

January 22-25, 2013

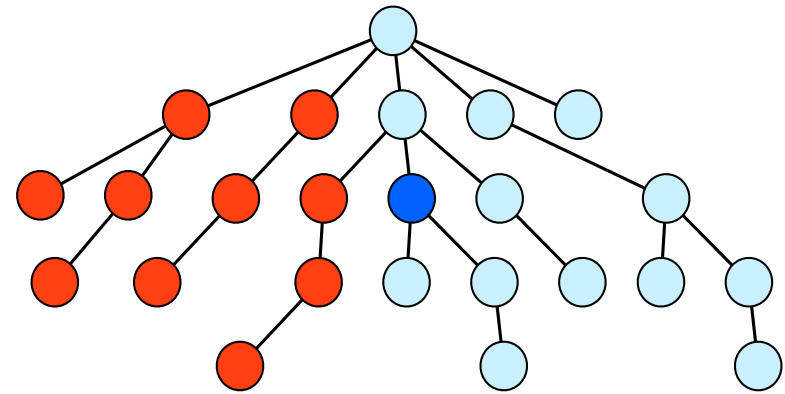
Following



Qsx

January 22-25, 2013

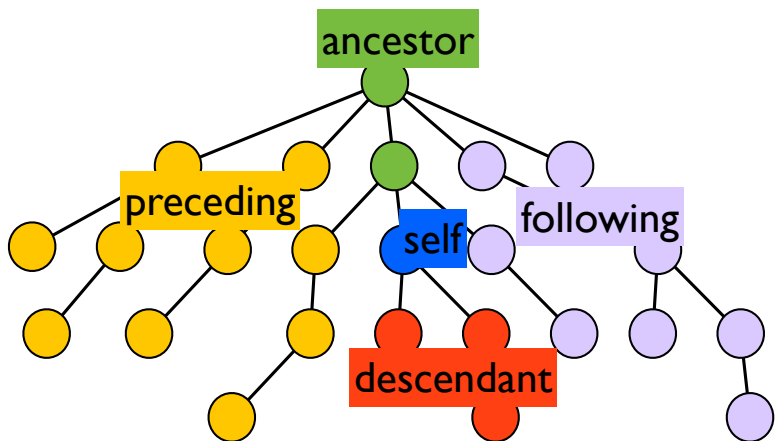
Preceding



Qsx

January 22-25, 2013

Partition



Qsx

January 22-25, 2013

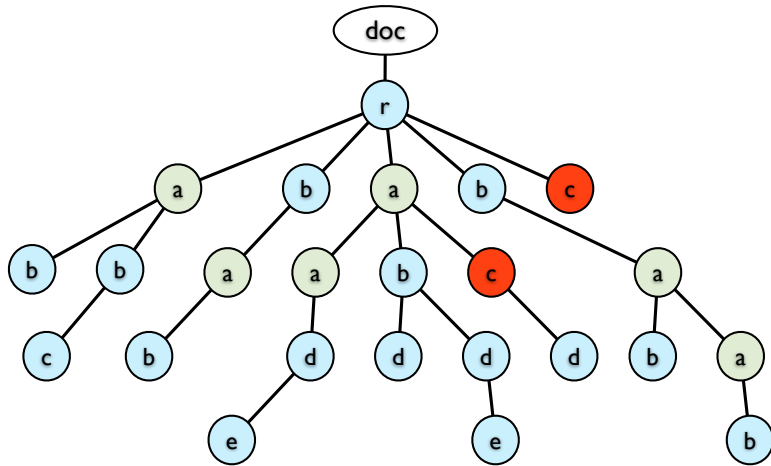
Syntax

- axis ::= child | descendant | descendant-or-self | parent | ancestor | ancestor-or-self | preceding-sibling | following-sibling | self | preceding | following
- test ::= * | text() | node() | a | b | @a | ...
- p ::= ax::tst[q] | p/p'
- q ::= p | q and q' | q or q' | not(q) | ...
- ap ::= /p

Qsx

January 22-25, 2013

Sequential composition

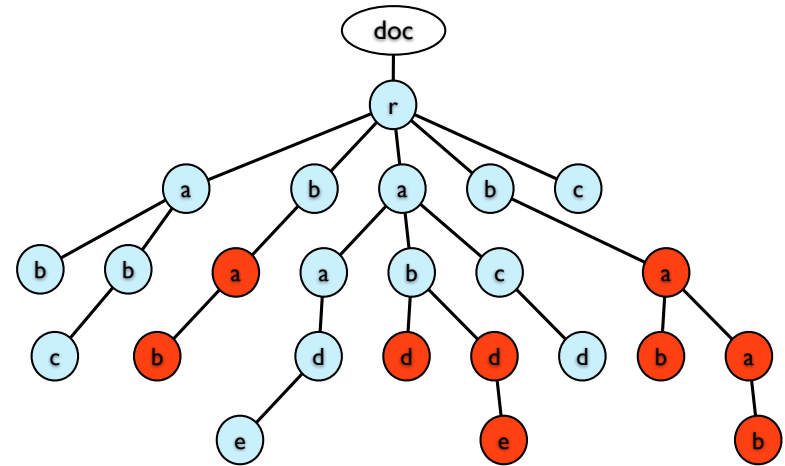


//a/following-sibling::c

Q SX

January 22-25, 2013

Sequential composition

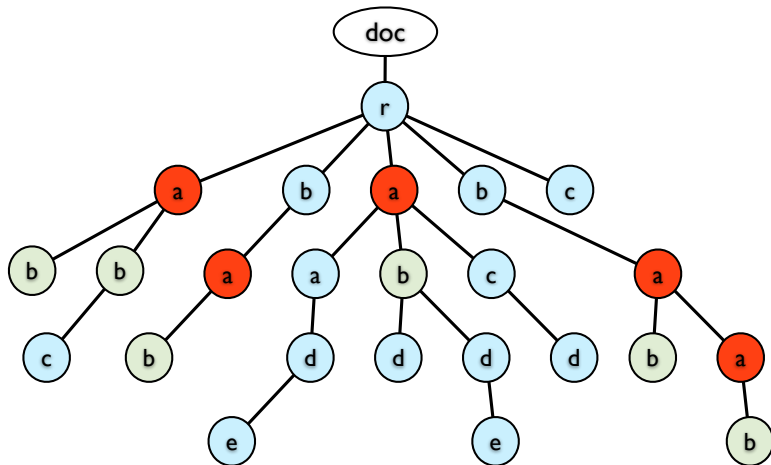


//a/following-sibling::b//*

Q SX

January 22-25, 2013

Filters

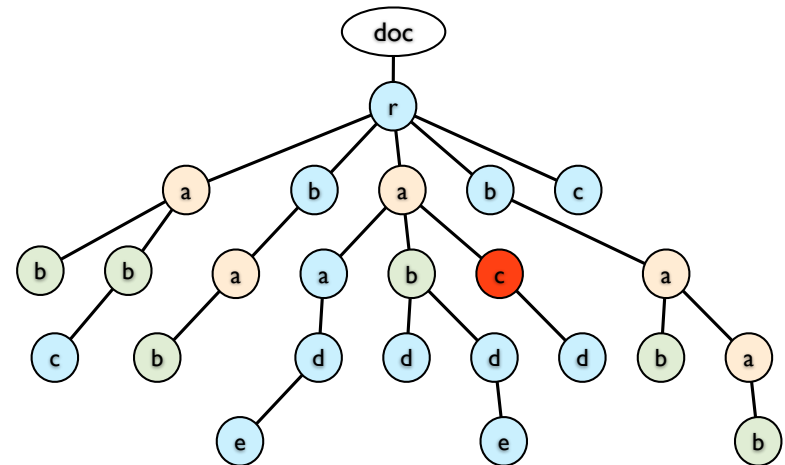


//a[b]

Q SX

January 22-25, 2013

Filters

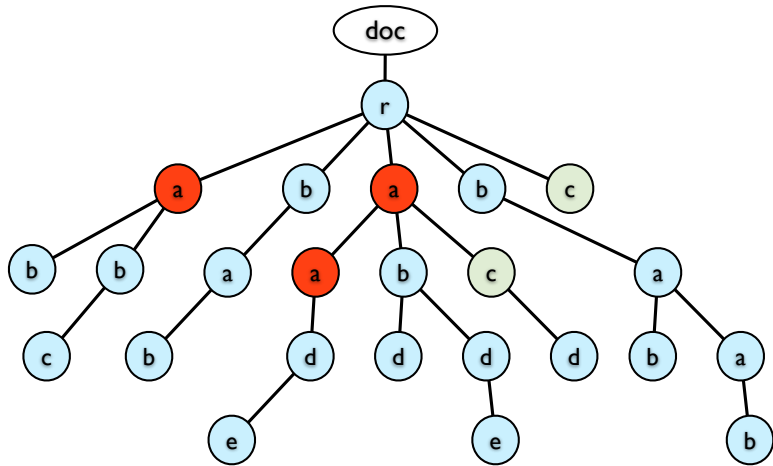


//a[b]/c

Q SX

January 22-25, 2013

Filters

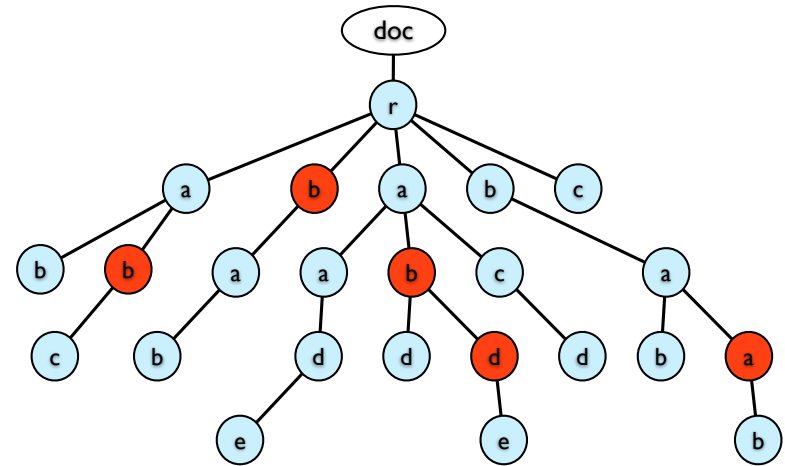


`//a[following-sibling::c]`

Q SX

January 22-25, 2013

Positional tests

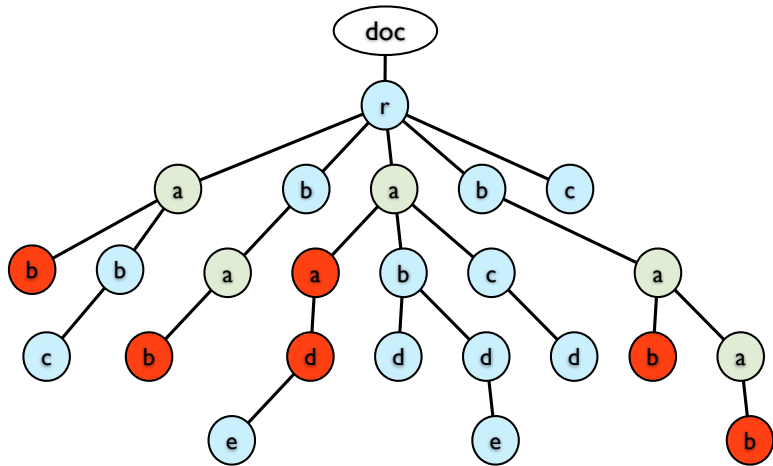


`//*[position()=2] (or just //a[2])`

Q SX

January 22-25, 2013

Positional tests

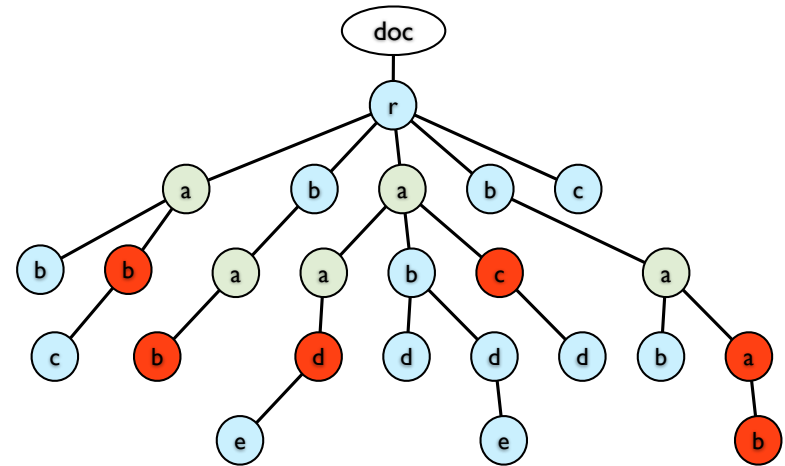


`//a/*[first()]`

Q SX

January 22-25, 2013

Positional tests



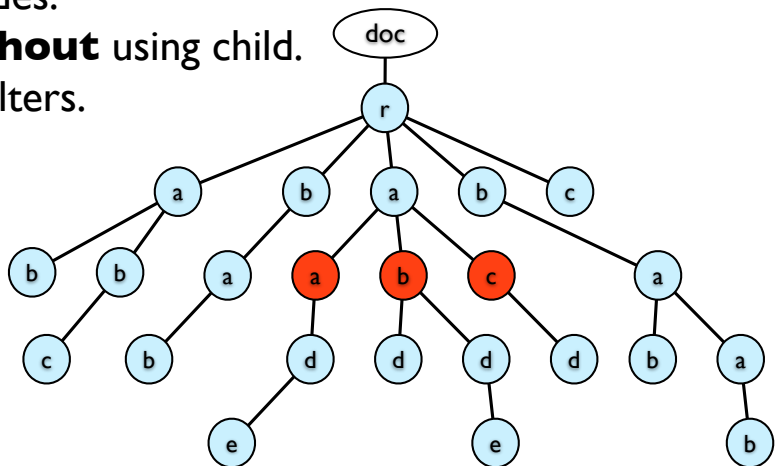
`//a/*[last()]`

Q SX

January 22-25, 2013

Quiz

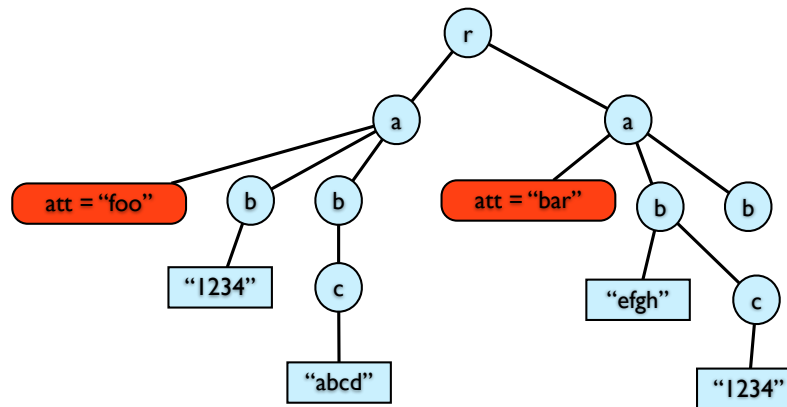
1. Write XPath to select red nodes.
2. **Without** using child.
3. Or filters.



Qsx

January 22-25, 2013

Attributes & Text

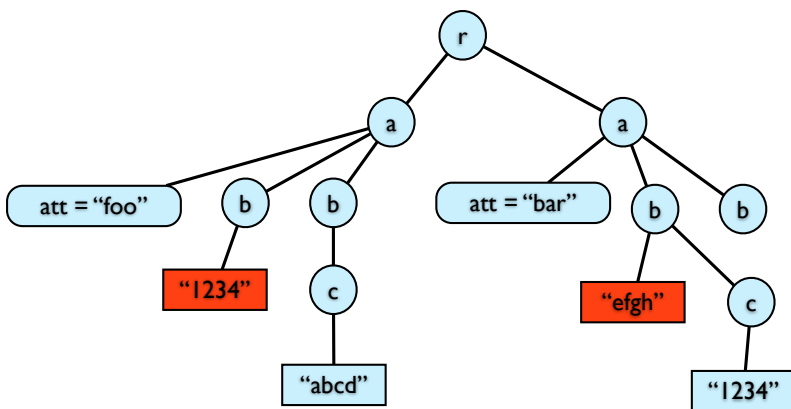


`//a/@att`

Qsx

January 22-25, 2013

Attributes & Text

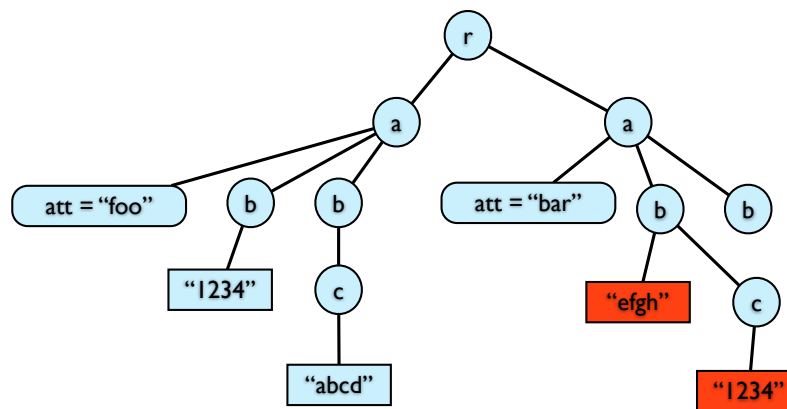


`//a/b/text()`

Qsx

January 22-25, 2013

Equality

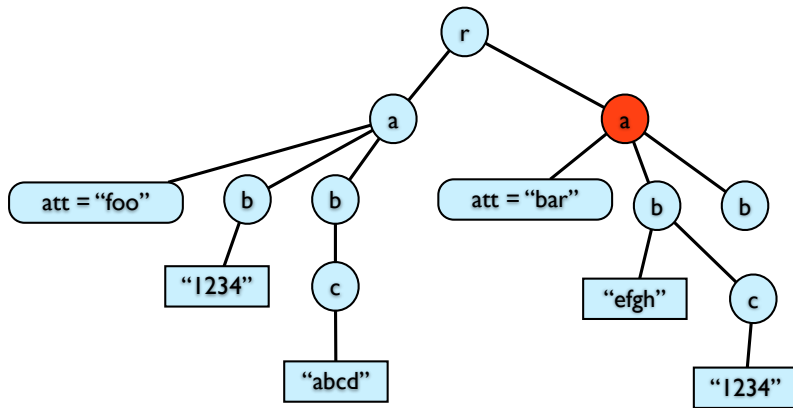


`//a[@att="bar"]//*[text()]`

Qsx

January 22-25, 2013

Equality

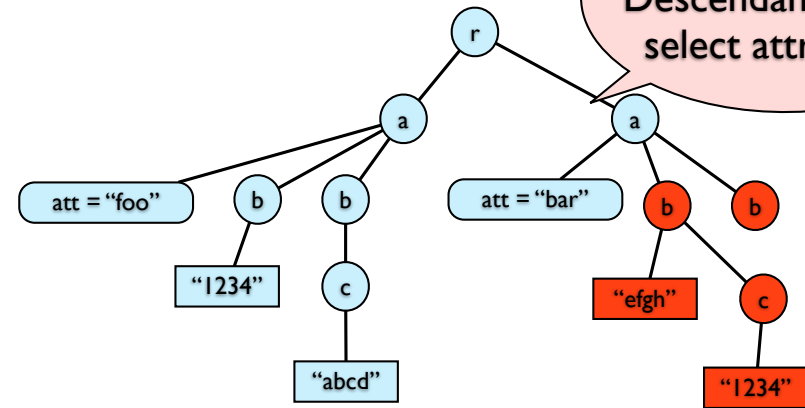


```
//a[@att="bar"]
```

Q SX

January 22-25, 2013

Equality

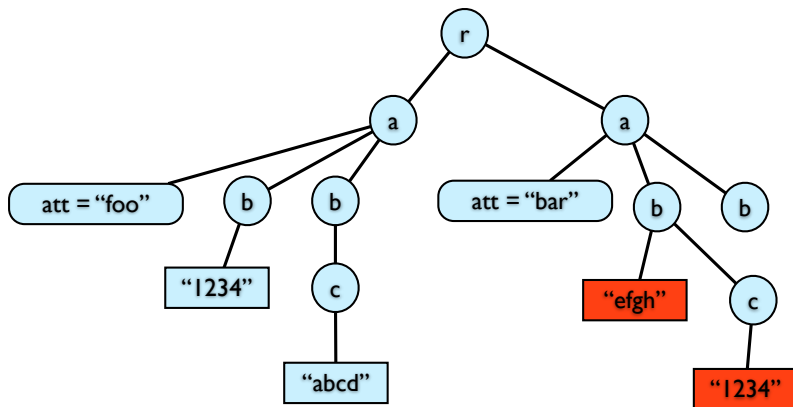


```
//a[@att="bar"]//*
```

Q SX

January 22-25, 2013

Equality

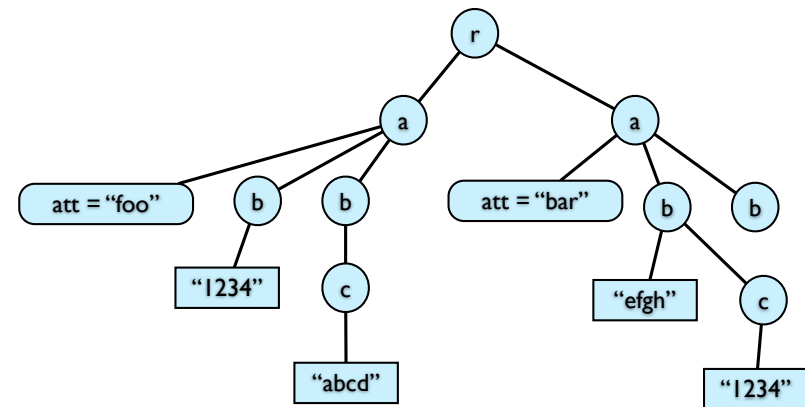


```
//a[@att="bar"]//text()
```

Q SX

January 22-25, 2013

Equality quiz

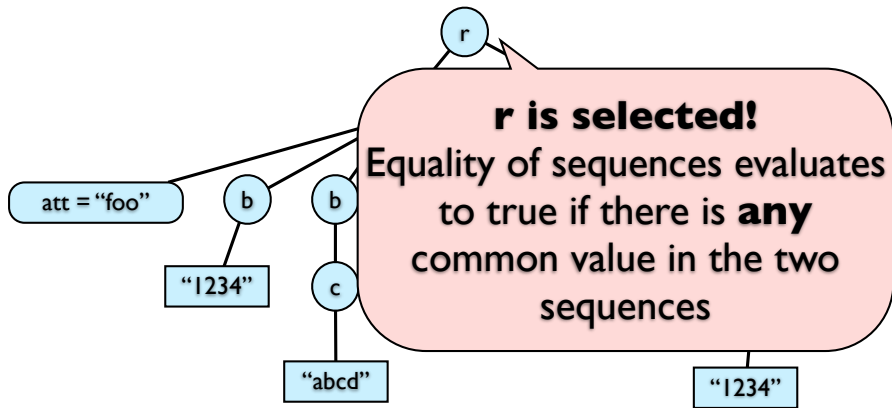


```
/r[a/b/text() = a/b/c/text()]
```

Q SX

January 22-25, 2013

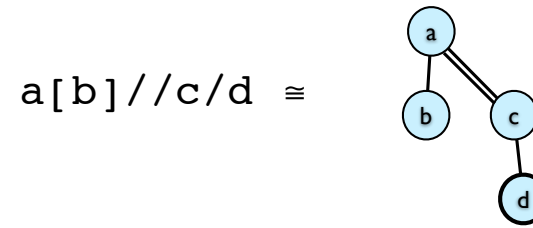
Equality quiz



`/r[a/b/text() = a/b/c/text()]`

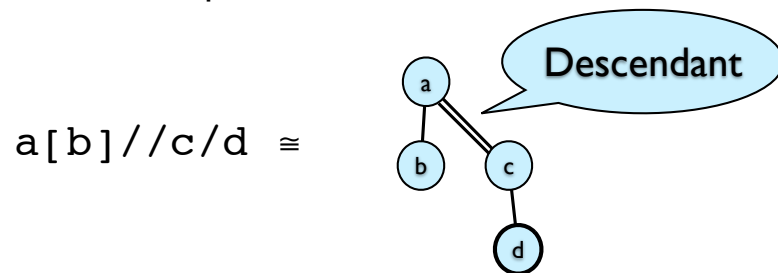
Tree patterns

- A graphical notation for (downward) XPath queries/filters



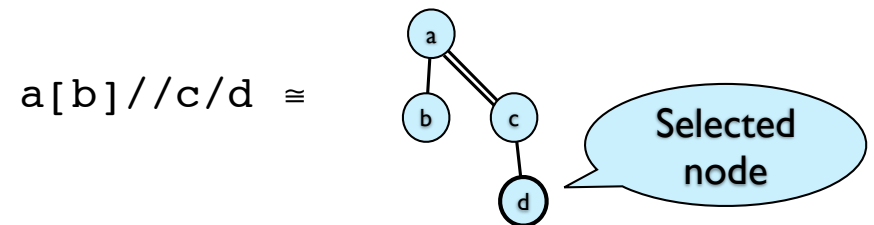
Tree patterns

- A graphical notation for (downward) XPath queries/filters



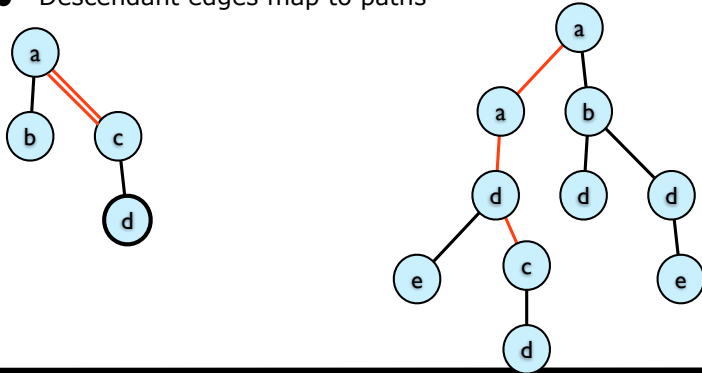
Tree patterns

- A graphical notation for (downward) XPath queries/filters



Tree pattern matching

- A function $h: P \rightarrow T$ such that:
 - Child edges map to edges
 - Descendant edges map to paths

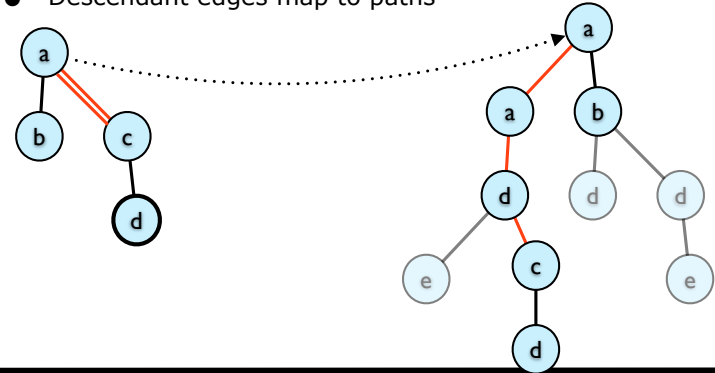


QSX

January 22-25, 2013

Tree pattern matching

- A function $h: P \rightarrow T$ such that:
 - Child edges map to edges
 - Descendant edges map to paths

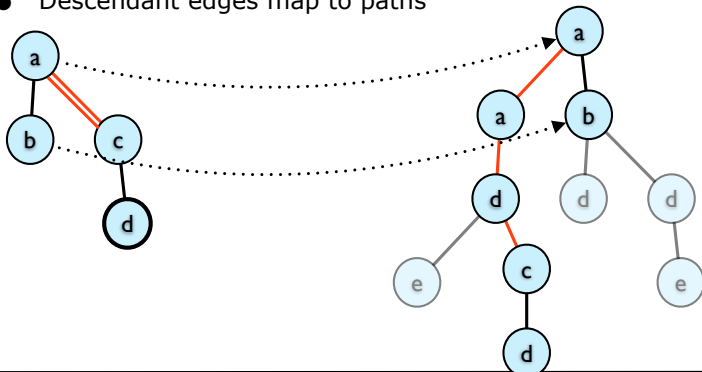


QSX

January 22-25, 2013

Tree pattern matching

- A function $h: P \rightarrow T$ such that:
 - Child edges map to edges
 - Descendant edges map to paths

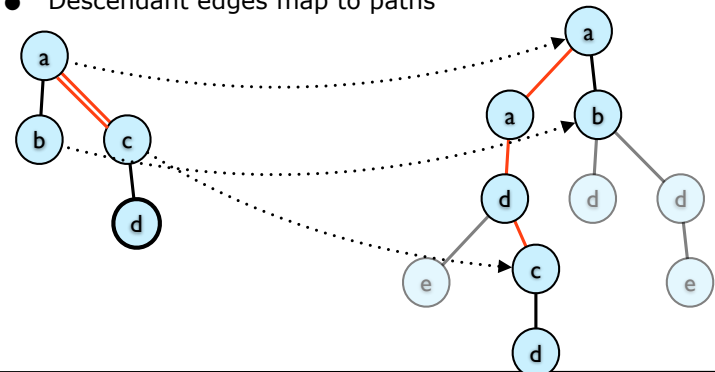


QSX

January 22-25, 2013

Tree pattern matching

- A function $h: P \rightarrow T$ such that:
 - Child edges map to edges
 - Descendant edges map to paths

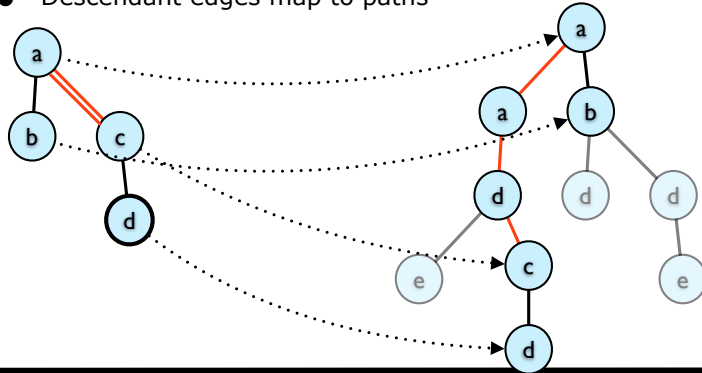


QSX

January 22-25, 2013

Tree pattern matching

- A function $h: P \rightarrow T$ such that:
 - Child edges map to edges
 - Descendant edges map to paths



Qsx

January 22-25, 2013

Semantics of XPath: steps

- $Ax[\text{self}](T) = \{(x,x) \mid x \in V\}$
- $Ax[\text{child}](T) = E$
- $Ax[\text{descendant}](T) = E^+$
- $Ax[\text{descendant-or-self}](T) = E^*$
- $Ax[\text{parent}](T) = \{(y,x) \mid (x,y) \in E\}$
- ...
- $Ax[\text{following-sibling}](T) = \{(x,y) \mid x < y\}$

Qsx

January 22-25, 2013

Semantics of XPath

- Represent tree as $T = (V, E, \lambda, <)$
 - Σ is set of possible node labels
 - $E \subseteq V \times V$ is parent/child edge relation
 - $\lambda : V \rightarrow \Sigma$ gives node labels
 - $< \subseteq V \times V$ linearly orders children of each node
- For simplicity, will ignore text nodes, attributes
 - but in general these need to be modeled too!

Qsx

January 22-25, 2013

Semantics of XPath: tests, paths & filters

- $\text{Test}[*](T) = \{x \mid x \in V\}$
- $\text{Test}[a](T) = \{x \mid x \in V, \lambda(x) = a\}$
- $\text{Path}[ax::\text{test}](T) = \{(x,y) \in Ax[ax](T) \mid y \in \text{Test}[\text{test}](T)\}$
- $\text{Path}[p/p'](T) = \{(x,z) \mid (x,y) \in \text{Path}[p](T), (y,z) \in \text{Path}[p'](T)\}$
- $\text{Path}[p[q]](T) = \{(x,y) \mid (x,y) \in \text{Path}[p](T), y \in \text{Filt}[q](T)\}$
- $\text{Filt}[p](T) = \{x \mid \exists y. (x,y) \in \text{Path}[p](T)\}$
- $\text{Filt}[q \text{ and } q'](T) = \text{Filt}[q](T) \cap \text{Filt}[q'](T)$
- $\text{Filt}[q \text{ or } q'](T) = \text{Filt}[q](T) \cup \text{Filt}[q'](T)$
- $\text{Filt}[\text{not}(q)](T) = \{x \in V \mid x \notin \text{Filt}[q](T)\}$

Qsx

January 22-25, 2013

Next time

- XQuery
 - Putting XPath to work
 - Iteration, binding, sequences, and XML construction expressions
 - Recursive functions

Q5X

January 22-25, 2013

What can XPath **not** do (well)?

- Construct new XML documents
- Combine information from different parts of document
 - Joins
- Abstraction over parts of query
 - Function definitions/recursion

Q5X

January 22-25, 2013

XQuery

Q5X

January 22-25, 2013

XQuery - a query language for XML

- Goals:
 - "SQL-like" query language for XML
 - Support query optimization
 - Support data types/XML Schema (will cover next week)
- Design:
 - Purely functional (more or less)
 - Every expression evaluates to a value (= sequence of XML trees or primitive values)
 - Extends XPath 2.0 with comprehensions, functions

Q5X

January 22-25, 2013

A first example

A first example

```
for $x in document("books.xml")/books/book
where $x/author="Abiteboul"
return <result>
    <title>{$x/title/text()}</title>
    <year>{$x/year/text()}</year></result>
```

Q SX

January 22-25, 2013

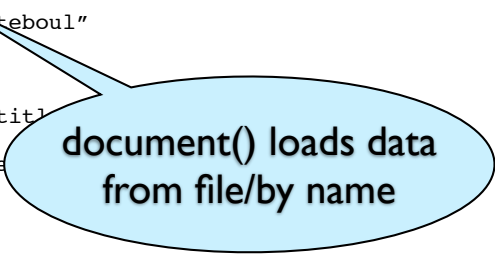
Q SX

January 22-25, 2013

A first example

A first example

```
for $x in document("books.xml")/books/book
where $x/author="Abiteboul"
return <result>
    <title>{$x/title/text()}</title>
    <year>{$x/year/text()}</year></result>
```



document() loads data from file/by name

```
for $x in document("books.xml")/books/book
where $x/author="Abiteboul"
return <result>
    <title>{$x/title/text()}</title>
    <year>{$x/year/text()}</year></result>
...
<result><title>Data on the Web</title>
    <year>2000</year></result>
<result><title>Web Data Management</title>
    <year>2011</year></result>
```

Q SX

January 22-25, 2013

Q SX

January 22-25, 2013

Atomic values

- Integers 1,2,3
- Strings 'abcd', "abcd"
- Dates / times
- Other basic types from XML Schema (will cover these later)

Values

- Atomic constants (last slide)
- XML trees
 - `<elt att1=v1 ... attn=vn>...value seq...</elt>`
- Value sequences are sequences of atomic/tree values
 - `()`, `(v1,v2, ..., vn)`
 - cannot be nested, i.e., `((v1,v2), v3) = (v1,v2,v3)`
 - however, `v1` could be an element with another sequence as content
- Formally:
`v ::= c | <elt att=v ... att=v>{vs}</elt>`
`vs ::= () | (v1,...,vn)`

Variables

- In XQuery, variables always start with \$
 - \$x, \$y, \$z, \$i
- This is common in other W3C standards with human-readable syntax
- A variable denotes a value sequence (more or less)

XML constructors

- XML values can be embedded in XQuery directly
`<element att1="v1" ...>...</element>`
- Can "antiquote" to embed XQuery expressions in elements
`<element>{$x/a/b}</element>`
- Can explicitly construct elements (with arbitrary names, attributes)
`element $foo { attribute {$bar} {$baz},
text {$some_text}}`

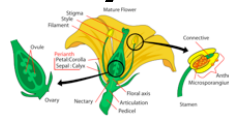
Building sequences

- Empty sequence: $()$
 - like empty list in other languages
- Sequence concatenation: (e_1, e_2)
 - evaluates e_1, e_2 to value sequences vs_1, vs_2
 - concatenates vs_1 and vs_2
- Examples: (expression equivalence)
 - $(1, 2, ()) \equiv (1, 2) \equiv ((), 1, 2) \equiv (1, (), 2)$
 - $(1, (2, 3)) \equiv (1, 2, 3) \equiv ((1, 2), 3)$
 - $((1, 2), (3, 4)) \equiv (1, 2, 3, 4)$

Qsx

January 22-25, 2013

Anatomy of a query: FLWOR



<code>for \$x in ...xpath...</code>	iterates over items in sequence
<code>let \$y := ...expression...</code>	binds variable to expression
<code>where ...condition...</code>	filters results based on boolean test
<code>order by ...ordering...</code>	orders results by key value
<code>return ...expression...</code>	constructs return values

Qsx

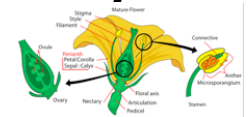
January 22-25, 2013

Reminder:
Next review
assignment due:
Monday (Jan 28) 4pm
Electronic handin only!

Qsx

January 22-25, 2013

Anatomy of a query: FLWOR



<code>for \$x in ...xpath...</code>	iterates over items in sequence
<code>let \$y := ...expression...</code>	binds variable to expression
<code>where ...condition...</code>	filters results based on boolean test
<code>order by ...ordering...</code>	orders results by key value
<code>return ...expression...</code>	constructs return values

Essentially list comprehensions (see also Haskell, Python, ...)

Qsx

January 22-25, 2013

For / comprehension

```
for $x in ...xpath...
```

- Evaluates xpath to a sequence
 - actually can be any expression
- Generates one binding of \$x for each element
- Evaluate rest of query once for each \$x-binding
- Concatenate results in order

Qsx

January 22-25, 2013

Let binding

```
let $y := ...expression...
```

- Evaluates expression to value
- Binds \$x to the value
- Evaluates rest of query with new binding

Qsx

January 22-25, 2013

Where clause

```
where ...condition...
```

- Evaluates condition expression to (Boolean) value
- If true, continue evaluating query
- If false, rest of query evaluates to ()
 - i.e., filters out results that don't satisfy condition

Qsx

January 22-25, 2013

Order by

```
order by ...ordering...
```

- Orders results of rest of query by key
- Key specification is defined in terms of values available so far
- can specify increasing or decreasing
 - many other options

Qsx

January 22-25, 2013

Return

```
return ...expression...
```

- Ends current iteration of query and generates result for it
 - unless filtered out by where-clause earlier
- Evaluates expression under current bindings

Let vs. for

- Both bind variables

```
let $x := (1,2,3)
let $y := ("a","b")
return ($x,$y)
...
(1,2,3,"a","b")
```

Let vs. for

- Both bind variables

```
for $x in (1,2,3)
let $y := ("a","b")
return ($x,$y)
...
(1,"a","b",2,"a","b",3,"a","b")
```

Let vs. for

- Both bind variables

```
let $x := (1,2,3)
for $y in ("a","b")
return ($x,$y)
...
(1,2,3,"a",1,2,3,"b")
```

Let vs. for

- Both bind variables

```
for $x in (1,2,3)
```

```
for $y in ("a","b")
```

```
return ($x,$y)
```

...

```
(1, "a", 1, "b", 2, "a", 2, "b", 3, "a", 3, "b")
```

Putting it all together

- A join: pairs of books having author in common, ordered by year of first one

```
let $books := document("books.xml")/books
```

```
for $x in $books/book, $y in $books/book
```

```
let $year := $x/year/text()
```

```
where $x/author/text() = $y/author/text()
```

```
order by $year
```

```
return <result>{$x},{ $y}</result>
```

Evaluating a join naively

- Iterates over all pairs of \$x,\$y
- Evaluates test
- Generates result for each pair satisfying test
- Problem: Quadratic.
 - Can do better using hash or sort join algorithms
 - Especially for large data
- XML databases can do this
- **Unordered** mode helps

Conditionals

```
if ...test... then ... else ...
```

```
if ...test... then ...
```

- Evaluate test
 - if true, evaluate then-branch
 - if false, evaluate else-branch
 - or () if no else-branch specified

Built-in functions

- Includes all XPath primitive functions
 - `first()`, `last()`, `position()`, `not()`, etc.
- equality: has same (strange) semantics as in XPath
 - i.e., $(1,2) = (2,3)$ evaluates to true
- Also `document(<xmlfile>)`
 - loads in an XML file and binds it to a value

Qsx

January 22-25, 2013

Aggregation and emptiness tests

`sum()`, `average()`, `min()`, `max()`,
`count()`

- calculate corresponding functions on numerical sequences (like in SQL)
 - (can also use in XPath)
- `empty()`, `exists()`
- test whether a sequence is empty or nonempty

Qsx

January 22-25, 2013

Set operations

- These are also allowed in XPath 2.0
- Union `e1 union e2`:
 - $\text{Path}[p \text{ union } p'](T) = \text{Path}[p](T) \cup \text{Path}[p'](T)$
- Intersection `e1 intersect e2`:
 - $\text{Path}[p \text{ intersect } p'](T) = \text{Path}[p](T) \cap \text{Path}[p'](T)$
- Difference `e1 except e2`:
 - $\text{Path}[p \text{ except } p'](T) = \text{Path}[p](T) \setminus \text{Path}[p'](T)$

Qsx

January 22-25, 2013

Quantifiers

some `$x` in ...`exp1`... satisfies ...`exp2`...

- true iff `exp2` evaluates to true for **some** bindings of `$x` to element of `exp1`
 - `exists(for $x in p where q return <z/>)`

every `$x` in ...`exp1`... satisfies ...`exp2`...

- true iff `exp2` evaluates to true for **all** bindings of `$x` to element of `exp1`
 - `empty(for $x in p where not(q) return <z/>)`

Qsx

January 22-25, 2013

Ordering & duplicates

- XQuery values are ordered sequences
- Can turn ordering off: `unordered {...}`
 - which enables more optimizations
- Or require it: `ordered {...}`
- Can also eliminate duplicates
 - `fn:remove-duplicates()`
- This happens automatically with some operations
 - such as `union`

QSX

January 22-25, 2013

User-definable functions

- Can define functions to abbreviate parts of queries

```
define function f($x,$y) {
  for $z in $x/a, $w in $y/b
  where $z/text() = $y/text()
  return <result>{$z}{$w}</result>
}
```

QSX

January 22-25, 2013

Quiz

- Starting with XML that lists cities, states and **optional** nicknames:

```
<cities><city><name>New York City</name>
  <state>NY</state>
  <nickname>The Big Apple</nickname>
</city> ...
</cities>
```
- 1. Write query that **ignores** state and lists city **by nickname if any**; otherwise **uses the name**
- 2. Write query that produces a list of states, each containing a list of city names in that state.
- 3. ... And gives a count of the number of cities in each state.

QSX

January 22-25, 2013

Functions can be recursive!

- example: recursive parts query

```
define function totalcost($x) {
  for $y in $x/part
  return $x/price + totalcost($y)
}
```

QSX

January 22-25, 2013

Turing-completeness

- Due to recursive functions, XQuery is a fully Turing-complete language
 - even without arithmetic
 - can simulate tape, arithmetic using trees
 - Big contrast to SQL, which lacks general recursion
- Can write whole Web applications using XQuery + web server interface library
- In practice, XQuery engines focus optimization effort on FLWOR queries

Qsx

January 22-25, 2013

Types

- XQuery has a native regular expression-based type system
 - Basic idea: if $\$x : \text{element } \{(a, (b, a)^*, c)\}$
 - then for $\$y \text{ in } \$x \text{ return } \$y : (a|b|c)^*$
- We will cover types and regular expression typing in more detail later
 - including XML Document Type Definitions, XML Schemas
 - and more precise systems for path/query typing

Qsx

January 22-25, 2013

Semantics

- XQuery Formal Semantics
 - uses operational rules to explain meaning of XQuery expressions
 - Also formalizes typing rules
- Will also look at this in more detail later
 - needed for proving correctness of type systems, optimizations

Qsx

January 22-25, 2013

Next time

- XSLT
- Type systems, XML DTDs
- Review assignment (due **Monday 4pm**):
 - XSLT overview
 - Read about XML Schemas
 - Read "Keys for XML"

Qsx

January 22-25, 2013