

XML and Databases

Lecture 9
Properties of XPath

Sebastian Maneth
NICTA and UNSW

CSE@UNSW -- Semester 1, 2009

Outline

1. XPath Equivalence
2. No Looking Back: How to Remove Backward Axes
3. Containment Test for XPath Expressions

A Note on Equality Test in XPath

Useful Functions (on Node Sets)

Careful with equality ("=")

XPath 2.0 has clearer comparison operators!

```
<a>  
<b>  
  <d>red</d>  
  <d>green</d>  
  <d>blue</d>  
</b>  
<c>  
  <d>yellow</d>  
  <d>orange</d>  
  <d>green</d>  
</c>  
</a>
```

XPath 1.0
Equality ("=") is based on
string value of a node!

//a[b/d = c/d] selects a-node!!!

there is a node in the node set for b/d
with same string value as a node in node set c/d

3

4

A Note on Equality Test

p1, p2 XPath (1.0) Expressions

(p1 = p2) is true if there exists a node selected by p1
that is *identical* to a node selected by p2

XPath 2.0
XQuery 1.0

```
<a>  
<b>  
  <d>red</d>  
  <d>green</d>  
  <d>blue</d>  
</b>  
<c>  
  <d>yellow</d>  
  <d>orange</d>  
  <d>green</d>  
</c>  
</a>
```

//a[b/d = c/d] selects what?

5

A Note on Equality Test

p1, p2 XPath (1.0) Expressions

(p1 = p2) is true if there exists a node selected by p1
that is *identical* to a node selected by p2

XPath 2.0
XQuery 1.0

```
<a>  
<b>  
  <d>red</d>  
  <d>green</d>  
  <d>blue</d>  
</b>  
<c>  
  <d>yellow</d>  
  <d>orange</d>  
  <d>green</d>  
</c>  
</a>
```

false (on any document)

//a[b/d = c/d] selects what?

//*[child::node()[1]
= child::node()[position=last()]]

6

A Note on Equality Test

$p1, p2$ XPath (1.0) Expressions

$(p1 == p2)$ is true if there exists a node selected by $p1$ that is *identical* to a node selected by $p2$

XPath 2.0
XQuery 1.0

XPath 1.0 simulation of (node) equality test ($==$)

Instead of $(p1 == p2)$ write:

$(count(p1 | p2) < count(p1) + count(p2))$ ☺

7

1. XPath Equivalence

$p1, p2$ XPath (1.0) Expressions

$(p1 \equiv p2)$ $p1$ "is equivalent to" $p2$ is true if, for any document D , and any context node N of D ,

$p1$ evaluated on D with context N gives the same result as $p2$ evaluated on D with context N .

Examples

$/a/**/b \equiv /a/**/b$
 $//a/b/c/././ \equiv //a[b/c/]$
 $//a[b | c] \equiv //a/*[self::b | self::c]/.$

NOT equivalent: $child::*/parent::* \neq self::*$

→ show a counter example!

8

1. XPath Equivalence

EBNF for XPaths that we want to consider now:

$path ::= path \ path \ / \ path \ path \ path \ path \ [\ qualif \] \ | \ axis \ :: \ nodetest \ | \ \perp$
 $qualif ::= qualif \ and \ qualif \ | \ qualif \ or \ qualif \ | \ (\ qualif \)$
 $path = path \ path == path \ path .$
 $axis ::= reverse_axis \ | \ forward_axis .$
 $reverse_axis ::= parent \ | \ ancestor \ | \ ancestor-or-self \ | \ preceding \ | \ preceding-sibling .$
 $forward_axis ::= self \ | \ child \ | \ descendant \ | \ descendant-or-self \ | \ following \ | \ following-sibling .$
 $nodetest ::= tagname \ | \ * \ | \ text() \ | \ node() .$

An XPath starting with "/" (root node) is called *absolute*, otherwise it is called *relative*.

9

1. XPath Equivalence

$p1, p2$ XPaths
 p arbitrary XPath
 q arbitrary qualifier

Rel→Abs If $p1 \equiv p2$, then $/p1 \equiv /p2$.

Adjunct If $p1 \equiv p2$ and p is a relative, then $p1/p \equiv p2/p$.
 If $p1 \equiv p2$ and $p1, p2$ relative, then $p/p1 \equiv p/p2$.
 If $p1 \equiv p2$, then $p1[q] \equiv p2[q]$ and $p[p1] \equiv p[p2]$.

Qualifier Flattening $p[p1/p2] \equiv p[p1][p2]$

$ancestor-or-self::n \equiv ancestor::n \ | \ self::n$
 $descendant-or-self::n \equiv descendant::n \ | \ self::n$

$p[p1 = /p2] \equiv p[p1[self::node() = /p2]]$
 $p[p1 == /p2] \equiv p[p1[self::node() == /p2]]$

10

1. XPath Equivalence

Lemma 3.2. Let m and n be node tests, i.e. m and n are tag names or one of the XPath constructs $*$, $node()$, or $text()$.

- Let a be one of the axes *parent*, *ancestor*, *preceding*, *preceding-sibling*, *self*, *following*, or *following-sibling*. Then the following holds:

$$/a::n \equiv \begin{cases} / & \text{if } a = \text{self and } n = \text{node}() \\ \perp & \text{otherwise} \end{cases}$$

- Let a be the *preceding* or *ancestor* axis. Then the following equivalences hold:

$$/child::m/a::n \equiv \begin{cases} /self::node()[child::m] & \text{if } a = \text{ancestor and } n = \text{node}() \\ \perp & \text{otherwise} \end{cases}$$

$$/child::m[a::n] \equiv \begin{cases} /child::m & \text{if } a = \text{ancestor and } n = \text{node}() \\ \perp & \text{otherwise} \end{cases}$$

11

2. No Looking Back

Dual backward forward



Thus: $dual(\text{parent}) = \text{child}$
 $dual(\text{following}) = \text{preceding}$
 etc.

Rewrite rule #1 (p, s : relative paths, ax : reverse axis)

$p[ax::m/s] \rightarrow p[/math>descendant:: $m[s]/dual(ax)::node() == self::node()]$$

12

Rewrite rule #1 (p,s: relative paths, ax: reverse axis)

$p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

any "m[s]-node" in the tree but, via dual axis, must reach context node

E.g. ax = ancestor

$p[ancestor::m] \rightarrow p[/descendant::m/descendant::node() == self::node()]$

"any m-node from which the context node can be reached via descendant, must be an ancestor of the context node."

13

Rewrite rule #1 (p,s: relative paths, ax: reverse axis)

$p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

any "m[s]-node" in the tree but, via dual axis, must reach context node

E.g. ax = preceding-sibling

$p[preceding-sibling::m] \rightarrow p[/descendant::m/following-sibling::node() == self::node()]$

"any m-node from which the context node can be reached via following-sibling, must be a preceding-sibling of the context node."

14

Rewrite rule #1 (p,s: relative paths, ax: reverse axis)

$p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

any "m[s]-node" in the tree but, via dual axis, must reach context node

E.g. ax=preceding-sibling

$p[preceding-sibling::m] \rightarrow p[/descendant::m/following-sibling::node() == self::node()]$

"any m-node from which the context node can be reached via following-sibling, must be a preceding-sibling of the context node."

Similar for **parent** and **preceding**. (ancestor-or-self not really needed. **Why?**)

15

Rewrite rule #1 (p,s: relative paths, ax: reverse axis)

$p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

→ navigation in left-hand side of equivalence
 - - - navigation in right-hand side of equivalence
 ○ root node
 ● context node
 ■ selected nodes

16

Rewrite rule #1 (p,s: relative paths, ax: reverse axis)

$p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

Removes first reverse axis inside a filter (qualifier).

Use **qualifier flattening** to replace "any" reverse axis from inside a filter.

Qualifier Flattening $p[p1/p2] \equiv p[p1[p2]]$

Similar rules for **absolute paths**:

$/p/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /p/fAx::n]$

$/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /fAx::n]$

Rewrite rules #2 and #2a

17

E.g.

`/descendant::pri ce/preceding::name`

is rewritten via Rule #2a into:

`/descendant::name[following::pri ce==/descendant::pri ce]`

Similar rules for **absolute paths**:

$/p/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /p/fAx::n]$

$/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /fAx::n]$

Rewrite rules #2 and #2a

18

E.g.

```
/descendant::pri ce/preceding::name
```

is rewritten via Rule #2a into:

```
/descendant::name[following::pri ce=/descendant::pri ce]
```

Of course, the "join" can be removed in this example:

```
/descendant::name[following::pri ce]
```

Not needed, in this example.

Similar rules for **absolute paths**:

```
/p/fAx::n/ax::m → /descendant::m[dual(ax)::n == /p/fAx::n]
```

```
/fAx::n/ax::m → /descendant::m[dual(ax)::n == /fAx::n]
```

Rewrite rules #2 and #2a

19

E.g.

```
/descendant::journal[child::title]/descendant::pri ce/preceding::name
```

becomes

```
/descendant::name[following::pri ce=
/descendant::journal[child::title]/descendant::pri ce]
```

Can you avoid the **join**, also for this example??

Similar rules for **absolute paths**:

```
/p/fAx::n/ax::m → /descendant::m[dual(ax)::n == /p/fAx::n]
```

```
/fAx::n/ax::m → /descendant::m[dual(ax)::n == /fAx::n]
```

Rewrite rules #2 and #2a

20

```
path ::= path | path / path | path / path [qualif] | axis :: nodetest | ⊥
qualif ::= qualific and qualific | qualific or qualific | ( qualific )
path = path | path == path | path .
axis ::= reverse_axis | forward_axis .
reverse_axis ::= parent | ancestor | ancestor-or-self |
preceding | preceding-sibling .
forward_axis ::= self | child | descendant | descendant-or-self |
following | following-sibling .
nodetest ::= tagname | * | text() | node() .
```

(1) $p[ax::m/s] \rightarrow p[/descendant::m[s]/dual(ax)::node() == self::node()]$

(2) $/p/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /p/fAx::n]$

(2a) $/fAx::n/ax::m \rightarrow /descendant::m[dual(ax)::n == /fAx::n]$

Rules (1),(2),(2a) suffice to remove ALL backward axes from above queries!

Why?

- Size Increase?
- How many joins?

21

2. No Looking Back

| | | | |
|--|--|---|------------|
| Dual | backward | forward | |
| parent ancestor ancestor-or-self preceding preceding-sibling | parent ancestor ancestor-or-self preceding preceding-sibling | child descendant descendant-or-self following following-sibling | not needed |

Joins (==) are expensive! (typically quadratic wrt data)

To obtain queries with fewer joins consider the **forward-axis** left of the **reverse-axis** to be removed!

New rules will be of the form

```
p/forw/back → p_new
```

```
p/forw[back] → p_new
```

22

2. No Looking Back

Interaction of **back=parent** with forward axes:

```
descendant::n/parent::m ≡ descendant-or-self::m[child::n] (3)
```

23

2. No Looking Back

Interaction of **back=parent** with forward axes:

```
descendant::n/parent::m ≡ descendant-or-self::m[child::n] (3)
```

```
child::n/parent::m ≡ self::m[child::n] (4)
```

24

2. No Looking Back

Interaction of **back=parent** with forward axes:

$descendant::n/parent::m \equiv descendant\text{-or-self}::m[child::n]$ (3)
 $child::n/parent::m \equiv self::m[child::n]$ (4)
 $p/self::n/parent::m \equiv p[self::n]/parent::m$ (5)

25

2. No Looking Back

Interaction of **back=parent** with forward axes:

$descendant::n/parent::m \equiv descendant\text{-or-self}::m[child::n]$ (3)
 $child::n/parent::m \equiv self::m[child::n]$ (4)
 $p/self::n/parent::m \equiv p[self::n]/parent::m$ (5)
 $p/following\text{-sibling}::n/parent::m \equiv p[following\text{-sibling}::n]/parent::m$ (6)

26

2. No Looking Back

Interaction of **back=parent** with forward axes:

$descendant::n/parent::m \equiv descendant\text{-or-self}::m[child::n]$ (3)
 $child::n/parent::m \equiv self::m[child::n]$ (4)
 $p/self::n/parent::m \equiv p[self::n]/parent::m$ (5)
 $p/following\text{-sibling}::n/parent::m \equiv p[following\text{-sibling}::n]/parent::m$ (6)
 $p/following::n/parent::m \equiv p/following::m[child::n]$ (7)
 $\quad | p/ancestor\text{-or-self}::*[following\text{-sibling}::n]$
 $\quad /parent::m$

27

2. No Looking Back

Interaction of **back=parent** with forward axes:

$descendant::n/parent::m \equiv descendant\text{-or-self}::m[child::n]$ (3)
 $child::n/parent::m \equiv self::m[child::n]$ (4)
 $p/self::n/parent::m \equiv p[self::n]/parent::m$ (5)
 $p/following\text{-sibling}::n/parent::m \equiv p[following\text{-sibling}::n]/parent::m$ (6)
 $p/following::n/parent::m \equiv p/following::m[child::n]$ (7)
 $\quad | p/ancestor\text{-or-self}::*[following\text{-sibling}::n]$
 $\quad /parent::m$

 $descendant::n [parent::m] \equiv descendant\text{-or-self}::m[child::n]$ (8)
 $child::n [parent::m] \equiv self::m[child::n]$ (9)
 $p/self::n [parent::m] \equiv p[self::n]/parent::m$ (10)
 $p/following\text{-sibling}::n [parent::m] \equiv p[following\text{-sibling}::n]$ (11)
 $p/following::n [parent::m] \equiv p/following::m[child::n]$ (12)
 $\quad | p/ancestor\text{-or-self}::*[parent::m]$
 $\quad /following\text{-sibling}::n$

28

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$p/descendant::n/ancestor::m \equiv p[descendant::n]/ancestor::m$ (13)
 $\quad | p/descendant\text{-or-self}::m[descendant::n]$

29

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$p/descendant::n/ancestor::m \equiv p[descendant::n]/ancestor::m$ (13)
 $\quad | p/descendant\text{-or-self}::m[descendant::n]$ (13a)
 $/descendant::n/ancestor::m \equiv /descendant\text{-or-self}::m[descendant::n]$ (13a)

30

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$$p/\text{descendant}::n/\text{ancestor}::m \equiv p[\text{descendant}::n]/\text{ancestor}::m \quad (13)$$

$$\quad | p/\text{descendant-or-self}::m[\text{descendant}::n]$$

$$/\text{descendant}::n/\text{ancestor}::m \equiv /[\text{descendant-or-self}::m[\text{descendant}::n]] \quad (13a)$$

$$p/\text{child}::n/\text{ancestor}::m \equiv p[\text{child}::n]/\text{ancestor-or-self}::m \quad (14)$$

31

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$$p/\text{descendant}::n/\text{ancestor}::m \equiv p[\text{descendant}::n]/\text{ancestor}::m \quad (13)$$

$$\quad | p/\text{descendant-or-self}::m[\text{descendant}::n]$$

$$/\text{descendant}::n/\text{ancestor}::m \equiv /[\text{descendant-or-self}::m[\text{descendant}::n]] \quad (13a)$$

$$p/\text{child}::n/\text{ancestor}::m \equiv p[\text{child}::n]/\text{ancestor-or-self}::m \quad (14)$$

$$p/\text{self}::n/\text{ancestor}::m \equiv p[\text{self}::n]/\text{ancestor}::m \quad (15)$$

32

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$$p/\text{descendant}::n/\text{ancestor}::m \equiv p[\text{descendant}::n]/\text{ancestor}::m \quad (13)$$

$$\quad | p/\text{descendant-or-self}::m[\text{descendant}::n]$$

$$/\text{descendant}::n/\text{ancestor}::m \equiv /[\text{descendant-or-self}::m[\text{descendant}::n]] \quad (13a)$$

$$p/\text{child}::n/\text{ancestor}::m \equiv p[\text{child}::n]/\text{ancestor-or-self}::m \quad (14)$$

$$p/\text{self}::n/\text{ancestor}::m \equiv p[\text{self}::n]/\text{ancestor}::m \quad (15)$$

$$p/\text{following-sibling}::n/\text{ancestor}::m \equiv p[\text{following-sibling}::n]/\text{ancestor}::m \quad (16)$$

33

2. No Looking Back

Interaction of **back=ancestor** with forward axes:

$$p/\text{descendant}::n/\text{ancestor}::m \equiv p[\text{descendant}::n]/\text{ancestor}::m \quad (13)$$

$$\quad | p/\text{descendant-or-self}::m[\text{descendant}::n]$$

$$/\text{descendant}::n/\text{ancestor}::m \equiv /[\text{descendant-or-self}::m[\text{descendant}::n]] \quad (13a)$$

$$p/\text{child}::n/\text{ancestor}::m \equiv p[\text{child}::n]/\text{ancestor-or-self}::m \quad (14)$$

$$p/\text{self}::n/\text{ancestor}::m \equiv p[\text{self}::n]/\text{ancestor}::m \quad (15)$$

$$p/\text{following-sibling}::n/\text{ancestor}::m \equiv p[\text{following-sibling}::n]/\text{ancestor}::m \quad (16)$$

$$p/\text{following}::n/\text{ancestor}::m \equiv p/\text{following}::m[\text{descendant}::n] \quad (17)$$

$$\quad | p/\text{ancestor-or-self}::*$$

$$\quad | [\text{following-sibling}::*/\text{descendant-or-self}::n]$$

$$\quad | /[\text{ancestor}::m]$$

Similar rules for **ancestor** in a filters.

34

2. No Looking Back

Interaction of **back=preceding** with forward axes:

$$p/\text{descendant}::n/\text{preceding}::m \equiv p[\text{descendant}::n]/\text{preceding}::m \quad (33)$$

$$\quad | p/\text{child}::*$$

$$\quad | [\text{following-sibling}::*/\text{descendant-or-self}::n]$$

$$\quad | /[\text{descendant-or-self}::m]$$

$$/\text{descendant}::n/\text{preceding}::m \equiv /[\text{descendant}::m[\text{following}::n]] \quad (33a)$$

$$p/\text{child}::n/\text{preceding}::m \equiv p[\text{child}::n]/\text{preceding}::m \quad (34)$$

$$\quad | p/\text{child}::*[\text{following-sibling}::n]$$

$$\quad | /[\text{descendant-or-self}::m]$$

$$p/\text{self}::n/\text{preceding}::m \equiv p[\text{self}::n]/\text{preceding}::m \quad (35)$$

$$p/\text{following-sibling}::n/\text{preceding}::m \equiv p[\text{following-sibling}::n]/\text{preceding}::m \quad (36)$$

$$\quad | p/\text{following-sibling}::*[\text{following-sibling}::n]$$

$$\quad | /[\text{descendant-or-self}::m]$$

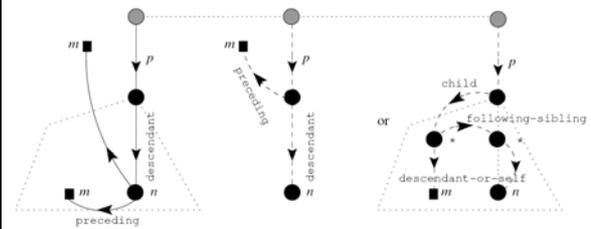
$$p/\text{following}::n/\text{preceding}::m \equiv p[\text{following}::n]/\text{preceding}::m \quad (37)$$

$$\quad | p/\text{following}::m[\text{following}::n]$$

$$\quad | p/\text{following}::n/\text{descendant-or-self}::m$$

35

Rule 33

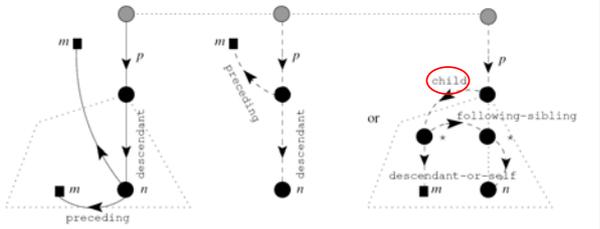


$$p/\text{descendant}::n/\text{preceding}::m \equiv p[\text{descendant}::n]/\text{preceding}::m$$

$$\quad | p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::n]/\text{descendant-or-self}::m$$

36

Rule 33



$p/\text{descendant}::n/\text{preceding}::m \equiv p[\text{descendant}::n]/\text{preceding}::m$
 $| p[\text{child}]^*[\text{following-sibling}::*/\text{descendant-or-self}::n]/\text{descendant-or-self}::m$

Wrong, I think!
Should be **descendant** instead!

37

2. No Looking Back

$/\text{descendant}::\text{price}/\text{preceding}::\text{name}$

is rewritten via Rule #2a into:

$/\text{descendant}::\text{name}[\text{following}::\text{price}]/\text{descendant}::\text{price}$

Now, let us use Rule (33a)

$/\text{descendant}::n/\text{preceding}::m \rightarrow / \text{descendant}::m[\text{following}::n]$

We obtain

$/\text{descendant}::\text{name}[\text{following}::\text{price}]$

38

$/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}/\text{preceding}::\text{name}$

becomes

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}$

Rule (33a)

$/\text{descendant}::n/\text{preceding}::m \rightarrow / \text{descendant}::m[\text{following}::n]$
 doesn't work because descendant is absolute here.

Rule (33):

$p/\text{descendant}::n/\text{preceding}::m \rightarrow p[\text{descendant}::n]/\text{preceding}::m$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::n]/\text{descendant-or-self}::m$

We obtain

$p[\text{descendant}::\text{price}]/\text{preceding}::\text{name}$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::\text{price}]/\text{descendant-or-self}::\text{name}$

39

$/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}/\text{preceding}::\text{name}$

becomes

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}$

Rule (33a)

$/\text{descendant}::n/\text{preceding}::m \rightarrow / \text{descendant}::m[\text{following}::n]$
 doesn't work because descendant is absolute here.

Rule (33):

$p/\text{descendant}::n/\text{preceding}::m \rightarrow p[\text{descendant}::n]/\text{preceding}::m$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::n]/\text{descendant-or-self}::m$

\rightarrow Rule (33a) with $n = \text{journal}[\text{child}::\text{title}][\text{descendant}::\text{price}]$

$p[\text{descendant}::\text{price}]/\text{preceding}::\text{name}$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::\text{price}]/\text{descendant-or-self}::\text{name}$

40

$/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}/\text{preceding}::\text{name}$

becomes

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}$

Rule (33a)

$/\text{descendant}::n/\text{preceding}::m \rightarrow / \text{descendant}::m[\text{following}::n]$
 doesn't work because descendant is absolute here.

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::\text{price}]/\text{descendant-or-self}::\text{name}$

\rightarrow Rule (33a) with $n = \text{journal}[\text{child}::\text{title}][\text{descendant}::\text{price}]$

$p[\text{descendant}::\text{price}]/\text{preceding}::\text{name}$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::\text{price}]/\text{descendant-or-self}::\text{name}$

41

$/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}/\text{preceding}::\text{name}$

becomes

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}$

Rule (33a)

$/\text{descendant}::n/\text{preceding}::m \rightarrow / \text{descendant}::m[\text{following}::n]$
~~doesn't work because descendant is absolute here.~~ seems it does work! ☺

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $| p/\text{child}::*[\text{following-sibling}::*/\text{descendant-or-self}::\text{price}]/\text{descendant-or-self}::\text{name}$

What about this one:

$/\text{descendant}::\text{name}[\text{following}::\text{price}]=$
 $/\text{descendant}::\text{journal}[\text{child}::\text{title}]/\text{descendant}::\text{price}$

42

Theorem

(from D. Olteanu, H. Meuss, T. Furche, F. Bry
 XPath: Looking Forward. [EDBT Workshops 2002](#): 109-127)

Given an XPath expression p that has no joins of the form $(p1 == p2)$ with both $p1, p2$ relative, an equivalent expression u without reverse axes can be computed.

Time needed: at most exponential in length of p
 Length of u : at most exponential in length of p

(moreover: no joins are introduced when computing u)

Questions

- Can you find a subclass for which Time to compute u is linear or polynomial?
- What is the problem with joins $(p1 == p2)$ for removal of reverse axes?

3. XPath Containment Test

Given two XPath expressions p, q :
 Are all nodes selected by p , also selected by q ? (on any document)
 (p "contained in" q)

Has many applications!

Boolean query

Want to select documents that "match p ".
 → If a document matches p , and p contained in q ,
 then we know the document also matches q !

→ If a document does not match q , and p contained in q ,
 then we know that document does not match p !

Applications

- Decrease online-time of publish/subscribe systems based on XPath
- Decrease query-time by making use of materialized intermediate results
- Optimization by ruling out queries with empty result set etc, etc

3. XPath Containment Test

Given two XPath expressions p, q

"0-containment" For every tree, if p selects a node then so does q .
 $p \subseteq_0 q$

"1-containment" For every tree, all nodes selected by p are also selected by q .
 $p \subseteq_1 q$

"2-containment" For every tree, and every context node N ,
 $p \subseteq_2 q$ all nodes selected by p starting from N ,
 are also selected by q starting from N .

- 1. Inclusion on Booleans
 - 2. Inclusion on Node Sets
 - 3. Inclusion on Node Relations
- } start from root

(If only child and descendant axes are allowed
 then \subseteq_1 and \subseteq_2 are the same! -- Why?)

3. XPath Containment Test

Given two XPath expressions p, q

"0-containment" For every tree, if p selects a node then so does q .
 $p \subseteq_0 q$

"1-containment" For every tree, all nodes selected by p are also selected by q .
 $p \subseteq_1 q$

Question

Given p, q and the fact $p \subseteq_1 q$,
 how can you determine from a result set of nodes for q ,
 the correct result set of nodes for p ?

3. XPath Containment Test

Given two XPath expressions p, q

Sometimes we want to test containment wrt a given DTD:

$p = /a/b//d$
 $q = /a//c$

Boolean!

Want to check if $p \subseteq_0 q$.

NO!
 a
 |
 b
 |
 d

But, what if documents are valid wrt to this DTD?

root → a^*
 a → $b^* | c^*$
 b → $d+c+$
 c → $b?c?$

2. XPath Containment Test

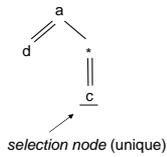
from:

T. Schwentick
 XPath query containment.
[SIGMOD Record 33\(1\)](#): 101-109 (2004)

| | |
|-------------|--|
| PSPACE | $XP(/, //, *, *)$ [21] $XP(/, //, *, *)$ (see [19]) $XP(/, //, //, //)$ [2], with fixed bounded SXICs [9] $XP(/, //) + DTDs$ [22] $XP(/, //) + DTDs$ [22] |
| coNP | $XP(/, //, //, //)$ [19] $XP(/, //, //, //, //, //, //, //)$ [22] $XP(/, //) + DTDs$ [22] $XP(/, //) + DTDs$ [22] |
| Π_2^c | $XP(/, //, //, //)$ + existential variables + path equality + ancestor-or-self axis + fixed bounded SXICs [9] $XP(/, //, //, //, //, //)$ + existential variables + all backward axes + fixed bounded SXICs [9] $XP(/, //, //, //)$ + existential variables with inequality [22] |
| PSPACE | $XP(/, //, //, //, //, //)$ if the alphabet is finite [22] $XP(/, //, //, //, //, //)$ + variables with XPath semantics [22] |
| EXPTIME | $XP(/, //, //, //)$ + existential variables + bounded SXICs [9] $XP(/, //, //, //, //, //)$ + DTDs [22] $XP(/, //, //, //)$ + DTDs [22] $XP(/, //, //, //, //, //)$ + DTDs [22] |
| Undecidable | $XP(/, //, //, //)$ + existential variables + unbounded SXICs [9] $XP(/, //, //, //)$ + existential variables + bounded SXICs + DTDs [9] $XP(/, //, //, //, //, //)$ + nodeset equality + simple DTDs [22] $XP(/, //, //, //, //, //)$ + existential variables with inequality [22] |

Pattern trees

E.g. $p = a[./d]/**/c$



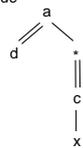
Note: child order has no meaning in pattern trees!

Test \subseteq_1 (node set inclusion) using \subseteq_0 (Boolean inclusion)

→ Simply add a new node below the selection node

New tree is Boolean (no selection node)

In a given XML tree:
pattern matches / does not match.



49

3. XPath Containment Test

4 techniques of testing XPath (Boolean) containment:

- (1) The Canonical Model Technique
- (2) The Homomorphism Technique
- (3) The Automaton Technique
- (4) The Chase Technique

50

3. XPath Containment Test

Canonical Model - XPath(/, //, [], *)

Idea: if there exists a tree that matches p but not q , then such a tree exists of size polynomial in the size of p and q .

Simple: remember, if you know that the XML document is only of height 5, then $//a/b/*c$ could be enumerated by $/a/b/*c \mid /*/a/b/*c \mid /*/*/a/b/*c \mid /*/*/*/a...$

Similarly, we try to construct a counter example tree, by replacing in p

- every $*$ by some new symbol "z"
- every $//$ by $/z, /z/z, /z/z/z, \dots, /z/z/z/z$

$N =$ length of longest $*/./$ chain in q

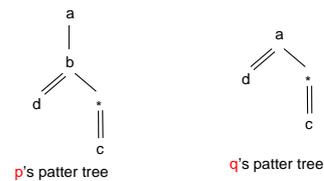
$N+1$ many z's

51

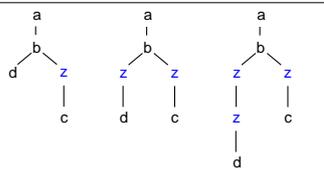
3. XPath Containment Test

Canonical Model - XPath(/, //, [], *)

Example



Test for q-match:



Formally, must test 1 and 2 more z's at right branch of each of the trees.

52

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that

- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not $*$), then $h(u)$ is also labeled by "e".

p,q expressions in XPath(/, //, [], [])

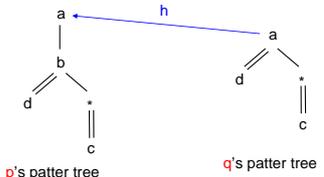
Theorem

$p \subseteq_0 q$ if and only if there is a homomorphism from Q to P .

53

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that

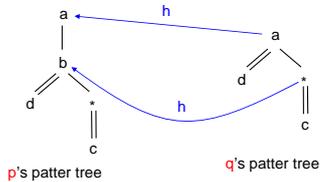


- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not $*$), then $h(u)$ is also labeled by "e".

54

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that

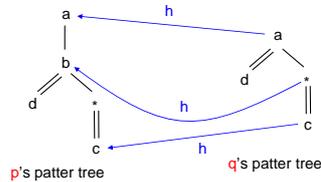


- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not *), then $h(u)$ is also labeled by "e".

55

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that

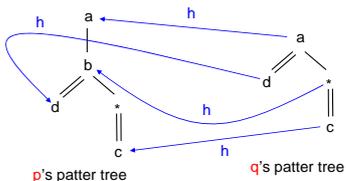


- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not *), then $h(u)$ is also labeled by "e".

56

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that



→ hom. h exists from Q to P , thus $p \subseteq_0 q$ must hold!

- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not *), then $h(u)$ is also labeled by "e".

57

3. XPath Containment Test

Homomorphism h maps each node of q 's query tree Q to a node of p 's query tree P such that

- (1) root of Q is mapped to root of P
- (2) if (u,v) is child-edge of Q then $(h(u),h(v))$ is child-edge of P
- (3) if (u,v) is descendant-edge of Q , then $h(v)$ is a "below" $h(u)$ in P
- (4) if u is labeled by "e" (not *), then $h(u)$ is also labeled by "e".

p, q expressions in $\text{XPath}(/, //, [], *)$

Theorem
 $p \subseteq_0 q$ if and only if there is a homomorphism from Q to P .

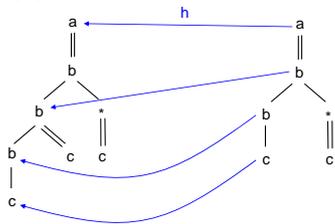
Cave If we add the star (*) then homomorphism need not exist!

→ there are $p, q \in \text{XPath}(/, //, [], *)$ such that $p \subseteq_0 q$ and there is **no** homomorphism from Q to P ☹

58

3. XPath Containment Test

$[a/b[./b[./b[d]/c]*/c]$ $[a/b[./b[d]*/c]$



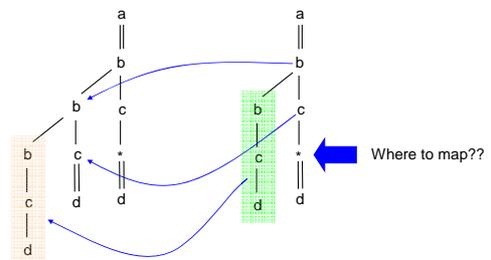
IS there a homomorphism??

Cave If we add the star (*) then homomorphism need not exist!

→ there are $p, q \in \text{XPath}(/, //, [], *)$ such that $p \subseteq_0 q$ and there is **no** homomorphism from Q to P ☹

59

$p = /a[./b[c/*//d]/b[c//d]/b[c/d]$
 $q = /a[./b[c/*//d]/b[c/d]$



Cave If we add the star (*) then homomorphism need not exist!

→ there are $p, q \in \text{XPath}(/, //, [], *)$ such that $p \subseteq_0 q$ and there is **no** homomorphism from Q to P ☹

60

$p = /a[. //b[c^*/d]/b[c/d]/b[c/d]]$
 $q = /a[. //b[c^*/d]/b[c/d]]$

Is p contained in q?
 → Test this, using the canonical model!!

Where to map??

Cave If we add the star (*) then homomorphism need not exist!
 → there are $p, q \in \text{XPath}(/, //, [], *, *)$ such that $p \subseteq_0 q$ and there is **no** homomorphism from Q to P ☹

61

Let's check the web... → **YES** p contained in q!

XPath-Containment Checker
 Implemented by Khaled Haj-Yahya (khaled.h at gmx.de)
 Supervised by B.C. Hammerschmidt (fomer)

Query p = /a[./b[c*/d]/b[c/d]/b[c/d]]
 Query q = /a[./b[c*/d]/b[c/d]]

$p \subseteq q$

XPath-Query p: /a[./b[c*/d]/b[c/d]/b[c/d]]
 XPath-Query q: /a[./b[c*/d]/b[c/d]]

Submit Query

62

3. XPath Containment Test

Automaton Technique

Recall: for any DTD there is a tree automaton which recognized the corresponding trees.

Similarly, for any $\text{XPath}(/, //, [], *, *)$ expression ex we can construct a (*non-deterministic* bottom-up) tree automaton A which accepts a tree if and only if ex matches the tree.

Theorem
 Containment test of $\text{XPath}(/, //, [], *, *)$ in the presence of DTDs can be solved in **EXPTIME**.

Exponential (deterministic) time
 Blow-up due to non-determinism of tree automaton.

BUT: no hope for improvement:
 The problem is actually **complete** for EXPTIME.

63

3. XPath Containment Test

Automaton Technique

Recall: for any DTD there is a tree automaton which recognized the corresponding trees.

Similarly, for any $\text{XPath}(/, //, [], *, *)$ expression ex we can construct a (*non-deterministic* bottom-up) tree automaton A which accepts a tree if and only if ex matches the tree.

Theorem
 Containment test of $\text{XPath}(/, //, [], *, *)$ in the presence of DTDs can be solved in **EXPTIME**.

Union of automata Intersection of automata ("product construction")

Proof Idea construct automaton for **all possible counter example trees**. Test if this automaton accepts any tree.

64

3. XPath Containment Test

Automaton Technique

Recall: for any DTD there is a tree automaton which recognized the corresponding trees.

Similarly, for any $\text{XPath}(/, //, [], *, *)$ expression ex we can construct a (*non-deterministic* bottom-up) tree automaton A which accepts a tree if and only if ex matches the tree.

Theorem
 Containment test of $\text{XPath}(/, //, [], *, *)$ in the presence of DTDs can be solved in **EXPTIME**.

→ Automata can also be Tested for Finiteness!
 Is $p \subseteq_0 q$, for all trees but finitely many exceptions?
 ↑ solvable!

Emptiness test for automata

Proof Idea construct automaton for **all possible counter example trees**. Test if this automaton accepts any tree.

65

3. XPath Containment Test

Chase Technique -- 1979 relational DB's to check query containment in the presence of *integrity constraints*.

Example

DTD $E =$

| | | |
|------|---|---------|
| root | → | a* |
| a | → | b* c* |
| b | → | d+C+ |
| c | → | b?c? |

("the chase" extends the relational homomorphism technique)

$p = /a/b//d$
 $q = /a//c$ Is p contained in q for E-conform documents?

First Possibility: use tree automata

- Construct automata A_p, A_q, A_E
- Construct B_q for the complement of A_q (=not q)
- Intersect B_q with A_p with A_E (gives automaton A)
- Check if A accepts any tree.

66

3. XPath Containment Test

Chase Technique -- 1979 relational DB's to check query containment in the presence of *integrity constraints*.

Example
 DTD $E =$

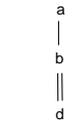
| | | |
|------|---|---------|
| root | → | a* |
| a | → | b* c* |
| b | → | d+c+ |
| c | → | b?c? |

 ("the chase" extends the relational homomorphism technique)

$p = /a/b//d$
 $q = /a//c$ Is p contained in q for E-conform documents?

Each b-element has a d-child and a c-child
 → *constraints*

c1: $b \rightarrow d$
 c2: $b \rightarrow c$



p's pattern tree

67

3. XPath Containment Test

Chase Technique -- 1979 relational DB's to check query containment in the presence of *integrity constraints*.

Example
 DTD $E =$

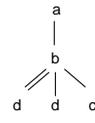
| | | |
|------|---|---------|
| root | → | a* |
| a | → | b* c* |
| b | → | d+c+ |
| c | → | b?c? |

 ("the chase" extends the relational homomorphism technique)

$p = /a/b//d$
 $q = /a//c$ Is p contained in q for E-conform documents?

Each b-element has a d-child and a c-child
 → *constraints*

c1: $b \rightarrow d$
 c2: $b \rightarrow c$



p's pattern tree after *chasing* with c1,c2

68

3. XPath Containment Test

Chase Technique -- 1979 relational DB's to check query containment in the presence of *integrity constraints*.

Example
 DTD $E =$

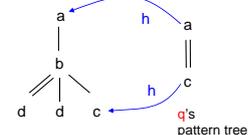
| | | |
|------|---|---------|
| root | → | a* |
| a | → | b* c* |
| b | → | d+c+ |
| c | → | b?c? |

 ("the chase" extends the relational homomorphism technique)

$p = /a/b//d$
 $q = /a//c$ Is p contained in q for E-conform documents?

Each b-element has a d-child and a c-child
 → *constraints*

c1: $b \rightarrow d$
 c2: $b \rightarrow c$



p is contained in q in the presence of the DTD E

p's pattern tree after *chasing* with c1,c2

69

END
 Lecture 9

70