

## XML and Databases

Lecture 9  
Properties of XPath

Sebastian Maneth  
NICTA and UNSW

CSE@UNSW -- Semester 1, 2010

## 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

3

## 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

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

### Recall

child::\* all child nodes that are **elements**  
 child::comment() all child nodes that are **comments**  
 child::processing-instruction() all child nodes that are **proc. instr.'s**  
 child::node() all child nodes that are **element/comments/PI's**

(only way to get to an *attribute*, is via the *attribute-axis*)

### Question

Which axes can bring you from an attribute-node back to an element-node?

```
<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()]]`

7

## 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))` ☺

8

## A Note on Equality Test

### Question

Can you give an XPath 1.0 filter expression for checking whether the **node set** of **p1** is equal to the **node set** of **p2**?

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

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

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

9

## 1. XPath Equivalence

**p1, p2** XPath (1.0) Expressions

`(p1 = 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

<code>/a//*/b</code>	■	<code>/a//*/b</code>
<code>//a/b/c/.../..</code>	■	<code>//a[. b/c/]</code>
<code>//a[b   c]</code>	■	<code>//a[b]   //a[c]</code>
<code>//*[a = /b]</code>	■	<code>/..</code>

NOT equivalent: `child::*/parent::*/self::*`  
→ show a counter example!

10

## 1. XPath Equivalence

EBNF for XPath's that we want to consider now:

```
path ::= path | path / path | path / path [qualif] | axis :: nodetest | ⊥
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** (will be evaluated *relative* to a given context node).

(Note: This is **Core XPath** w/o negation, but with `=` and `==` operators)

11

## 1. XPath Equivalence

**p1, p2** XPath's

**p** arbitrary XPath

**q** arbitrary qualifier

**Rel→Abs** If **p1** ■ **p2**, then **/p1** ■ **/p2**.

**Adjunct** If **p1** ■ **p2** and **p** is a relative, then **p1/p** ■ **p2/p**.  
 If **p1** ■ **p2** and **p1, p2** relative, then **p/p1** ■ **p/p2**.  
 If **p1** ■ **p2**, then **p1[q]** ■ **p2[q]** and **p[p1]** ■ **p[p2]**.

**Qualifier Flattening** **p[p1/p2]** ■ **p[p1[p2]]**

**ancestor-or-self::n** ■ **ancestor::n** | **self::n**  
**descendant-or-self::n** ■ **descendant::n** | **self::n**

**p[p1 = /p2]** ■ **p[p1[self::node() = /p2]]**  
**p[p1 == /p2]** ■ **p[p1[self::node() == /p2]]**

12

## 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:

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

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

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

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

(same holds for  $a = \text{parent}$ )

13

## 2. No Looking Back

Dual

backward

forward

parent	child
ancestor	descendant
ancestor-or-self	ancestor-or-self
preceding	following
preceding-sibling	following-sibling

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

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

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

14

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

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

any "m[s]-node"  
in the tree

but, via dual axis, must  
reach context node

E.g.  $ax = \text{ancestor}$

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

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

15

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

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

any "m[s]-node"  
in the tree

but, via dual axis, must  
reach context node

E.g.  $ax = \text{preceding-sibling}$

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

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

16

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

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

any "m[s]-node"  
in the tree

but, via dual axis, must  
reach context node

E.g.  $ax = \text{preceding-sibling}$

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

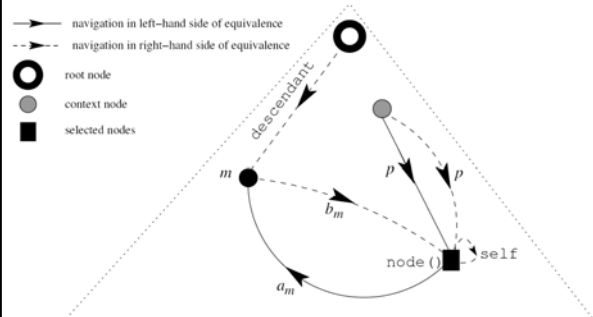
"any m-node from which the context node can be reached via `following-sibling`, must be a `preceding-sibling` of the context node."

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

17

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

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



18

**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**  $p[p1/p2] = p[p1[p2]]$

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

19

E.g.

$/descendant::price/preceding::name$

is rewritten via Rule #2a into:

$/descendant::name[following::price == /descendant::price]$

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

20

E.g.

$/descendant::price/preceding::name$

is rewritten via Rule #2a into:

$/descendant::name[following::price == /descendant::price]$

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

$/descendant::name[following::price]$

a tautology

Not needed, in this example.

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

21

E.g.

$/descendant::journal[child::title]/descendant::price/preceding::name$

becomes

$/descendant::name[following::price == /descendant::journal[child::title]/descendant::price]$

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

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

22

```
path ::= path | path / path | path / path [qualif] | axis :: nodetest | ⊥ .
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() .
```

- (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?

23

## 2. No Looking Back

Dual

backward

forward



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

$p/forw/back \rightarrow p\_new$

$p/forw[back] \rightarrow p\_new$

24

## 2. No Looking Back

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

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

25

## 2. No Looking Back

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

$$\begin{aligned} \text{descendant}::n/\text{parent}::m &\equiv \text{descendant-or-self}::m[\text{child}::n] & (3) \\ \text{child}::n/\text{parent}::m &\equiv \text{self}::m[\text{child}::n] & (4) \end{aligned}$$

26

## 2. No Looking Back

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

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

27

## 2. No Looking Back

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

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

28

## 2. No Looking Back

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

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

29

## 2. No Looking Back

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

$$\begin{aligned} \text{descendant}::n/\text{parent}::m &\equiv \text{descendant-or-self}::m[\text{child}::n] & (3) \\ \text{child}::n/\text{parent}::m &\equiv \text{self}::m[\text{child}::n] & (4) \\ p/\text{self}::n/\text{parent}::m &\equiv p[\text{self}::n]/\text{parent}::m & (5) \\ p/\text{following-sibling}::n/\text{parent}::m &\equiv p[\text{following-sibling}::n]/\text{parent}::m & (6) \\ p/\text{following}::n/\text{parent}::m &\equiv p/\text{following}::m[\text{child}::n] & (7) \\ &\quad | p/\text{ancestor-or-self}::*[ \text{following-sibling}::n ] \\ &\quad \quad / \text{parent}::m \\ \\ \text{descendant}::n [ \text{parent}::m ] &\equiv \text{descendant-or-self}::m/\text{child}::n & (8) \\ \text{child}::n [ \text{parent}::m ] &\equiv \text{self}::m/\text{child}::n & (9) \\ p/\text{self}::n [ \text{parent}::m ] &\equiv p [ \text{parent}::m ] / \text{self}::n & (10) \\ p/\text{following-sibling}::n [ \text{parent}::m ] &\equiv p [ \text{parent}::m ] / \text{following-sibling}::n & (11) \\ p/\text{following}::n [ \text{parent}::m ] &\equiv p/\text{following}::m/\text{child}::n & (12) \\ &\quad | p/\text{ancestor-or-self}::*[ \text{parent}::m ] \\ &\quad \quad / \text{following-sibling}::n \end{aligned}$$

30

## 2. No Looking Back

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

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

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

31

## 2. No Looking Back

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

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

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

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

32

## 2. No Looking Back

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

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

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

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

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

33

## 2. No Looking Back

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

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

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

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

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

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

34

## 2. No Looking Back

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

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

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

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

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

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

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

35

## 2. No Looking Back

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

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

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

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

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

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

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

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

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

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

$$\quad \quad /\text{ancestor}::m$$

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

36

## 2. No Looking Back

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

$$\begin{aligned}
 p/\text{descendant}::n/\text{ancestor}::m &\equiv p[\text{descendant}::n]/\text{ancestor}::m & (13) \\
 &\quad | p/\text{descendant-or-self}::m[\text{descendant}::n] \\
 / \text{descendant}::n/\text{ancestor}::m &\equiv / \text{descendant-or-self}::m[\text{descendant}::n] & (13a) \\
 p/\text{child}::n/\text{ancestor}::m &\equiv p[\text{child}::n]/\text{ancestor-or-self}::m & (14) \\
 p/\text{self}::n/\text{ancestor}::m &\equiv p[\text{self}::n]/\text{ancestor}::m & (15) \\
 p/\text{following-sibling}::n/\text{ancestor}::m &\equiv p[\text{following-sibling}::n]/\text{ancestor}::m & (16) \\
 p/\text{following}::n/\text{ancestor}::m &\equiv p/\text{following}::m[\text{descendant}::n] & (17) \\
 &\quad | p/\text{ancestor-or-self}::* \\
 &\quad \quad [\text{following-sibling}::*/\text{descendant-or-self}::n] \\
 &\quad \quad / \text{ancestor}::m
 \end{aligned}$$

Similar rules for **ancestor** in filters.

E.g., what is the forward query for: `/**[ancestor::a]`

37

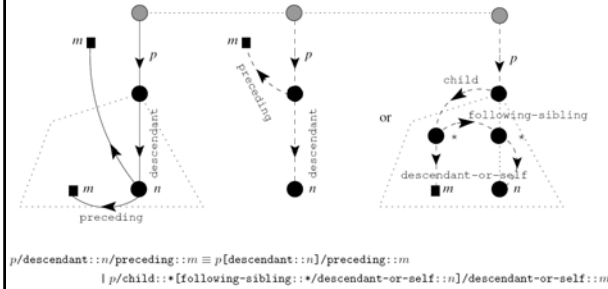
## 2. No Looking Back

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

$$\begin{aligned}
 p/\text{descendant}::n/\text{preceding}::m &\equiv p[\text{descendant}::n]/\text{preceding}::m & (33) \\
 &\quad | p/\text{child}::* \\
 &\quad \quad [\text{following-sibling}::*/\text{descendant-or-self}::n] \\
 &\quad \quad / \text{descendant-or-self}::m \\
 / \text{descendant}::n/\text{preceding}::m &\equiv / \text{descendant}::m[\text{following}::n] & (33a) \\
 p/\text{child}::n/\text{preceding}::m &\equiv p[\text{child}::n]/\text{preceding}::m & (34) \\
 &\quad | p/\text{child}::*[\text{following-sibling}::n] \\
 &\quad \quad / \text{descendant-or-self}::m \\
 p/\text{self}::n/\text{preceding}::m &\equiv p[\text{self}::n]/\text{preceding}::m & (35) \\
 p/\text{following-sibling}::n/\text{preceding}::m &\equiv p[\text{following-sibling}::n]/\text{preceding}::m & (36) \\
 &\quad | p/\text{following-sibling}::*[\text{following-sibling}::n] \\
 &\quad \quad / \text{descendant-or-self}::m \\
 &\quad \quad | p[\text{following-sibling}::n]/\text{descendant-or-self}::m \\
 p/\text{following}::n/\text{preceding}::m &\equiv p[\text{following}::n]/\text{preceding}::m & (37) \\
 &\quad | p/\text{following}::m[\text{following}::n] \\
 &\quad \quad | p[\text{following}::n]/\text{descendant-or-self}::m
 \end{aligned}$$

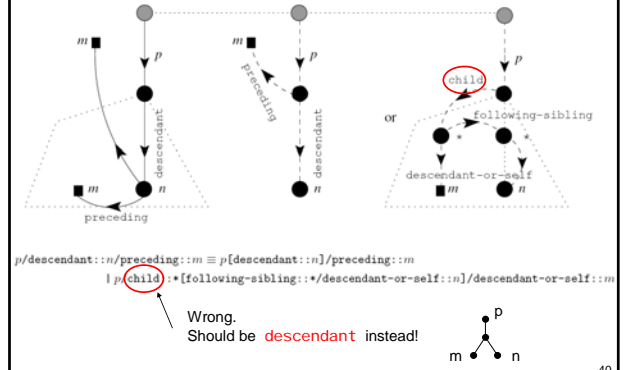
38

### Rule 33



39

### Rule 33



40

## 2. No Looking Back

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

is rewritten via Rule #2a into:

`/descendant::name[fol owi ng::pri ce=]/descendant::pri ce`

Now, let us use Rule (33a)

`/descendant::n/preceding::m → /descendant::m[fol owi ng::n]`

We obtain

`/descendant::name[fol owi ng::pri ce]`

41

**p**

`/descendant::journal [ch i l d::t i t l e]/descendant::pri ce/preceding::name`

becomes

`/descendant::name[fol owi ng::pri ce=] /descendant::journal [ch i l d::t i t l e]/descendant::pri ce`

Rule (33a)

`/descendant::n/preceding::m → /descendant::m[fol owi ng::n]`  
doesn't work because descendant is absolute here.

Rule (33):

`p/descendant::n/preceding::m → p[descendant::n]/preceding::m`  
`| p/ch i l d::*[fol owi ng-si b l i ng::*/descendant-or-sel f::n] /descendant-or-sel f::m`

We obtain

`p[descendant::pri ce]/preceding::name`  
`| p/ch i l d::*[fol owi ng-si b l i ng::*/descendant-or-sel f::pri ce] /descendant-or-sel f::name`

42

**p**

```
/descendant::journal[chiId::title]/descendant::price/preceding::name
```

becomes

```
/descendant::name[following::price=
/descendant::journal[chiId::title]/descendant::price]
```

Rule (33a)  
 /descendant::n/preceding::m → /descendant::m[following::n]  
 doesn't work because descendant is absolute here.

Rule (33):  
**p**/descendant::n/preceding::m → p[descendant::n]/preceding::m  
 | p/child::\*[following-sibling::\*]/descendant-or-self::n  
 /descendant-or-self::m

→ Rule (33a) with **n** = journal[chiId::title][descendant::price]

```
p[descendant::price]/preceding::name
| p/child::*[following-sibling::*]/descendant-or-self::price
/descendant-or-self::name
```

43

**p**

```
/descendant::journal[chiId::title]/descendant::price/preceding::name
```

becomes

```
/descendant::name[following::price=
/descendant::journal[chiId::title]/descendant::price]
```

Rule (33a)  
 /descendant::n/preceding::m → /descendant::m[following::n]  
 doesn't work because descendant is absolute here.

```
/descendant::name[following::journal[chiId::title][descendant::price]]
| p/child::*[following-sibling::*]/descendant-or-self::price
/descendant-or-self::name
```

→ Rule (33a) with **n** = journal[chiId::title][descendant::price]

```
p[descendant::price]/preceding::name
| p/child::*[following-sibling::*]/descendant-or-self::price
/descendant-or-self::name
```

44

**p**

```
/descendant::journal[chiId::title]/descendant::price/preceding::name
```

becomes

```
/descendant::name[following::price=
/descendant::journal[chiId::title]/descendant::price]
```

Rule (33a)  
 /descendant::n/preceding::m → /descendant::m[following::n]  
 doesn't work because descendant is absolute here. seems it does work! ☺

```
/descendant::name[following::journal[chiId::title][descendant::price]]
| p/child::*[following-sibling::*]/descendant-or-self::price
/descendant-or-self::name
```

**p**[p1/p2]  
 = p[p1[p2]]

What about this one:

```
/descendant::name[following::journal[chiId::title]/descendant::price]
```

45

**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**

- Why rewriting takes exponential time?
- 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?

46

**Theorem**  
 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**  
 (moreover: *no joins* are introduced when computing **u**)

**More Questions**

- Give an example of Core backward XPath with **negation**, for which there is no forward XPath query.
- Give an example of Core backward XPath with **data values**, for which there is no forward XPath query.
- Give an example of a Core backward XPath with **counting**, for which there is no forward XPath query.

47

### 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 the 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

48



### 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?)

49

### 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  $p$ ?

50

### 3. XPath Containment Test

Given two XPath expressions  $p, q$

Sometimes we want to test containment wrt a given DTD:

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

Want to check if  $p \subseteq_0 q$ .

NO!

a  
|  
b  
|  
d

Boolean!

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

root  $\rightarrow a^*$   
 $a \rightarrow b^* \mid c^*$   
 $b \rightarrow d+c+$   
 $c \rightarrow b?c?$

51

### 2. XPath Containment Test

from:

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

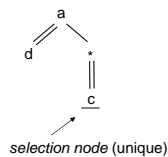
PSPACE	$XP(/, //, *, *)$ [21] $XP(/, //, *, *)$ (see [19]) $XP(/, //, //)$ [2], with fixed bounded SXICs [9] $XP(/, //) + DTDs$ [22] $XP(/, //) + DTDs$ [22]
coNP	$XP(/, //, //, *)$ [19] $XP(/, //, //, *, *)$ [22] $XP(/, //) + DTDs$ [22] $XP(/, //) + DTDs$ [22]
$\Pi_2^P$	$XP(/, //, //, //) +$ existential 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	$XP(/, //, //, //, *)$ and $XP(/, //, //, //)$ 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] $XP(/, //, //, //, //) +$ node-set equality + simple DTDs [22] $XP(/, //, //, //, //) +$ existential variables with inequality [22]

52

#### Pattern trees

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

Note: child order has no meaning in  
pattern trees!

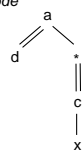


Test  $\subseteq_1$  (node set inclusion) using  $\subseteq_0$  (Boolean inclusion)

→ 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.



53

### 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

54

### 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  $/a/b/*c \mid /*a/b/*c \mid /*a/b/*c \mid /*a/b/*c \mid /*a/b/*c$ ...

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

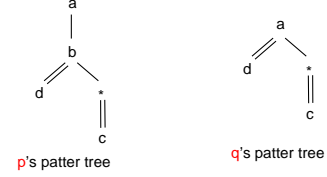
- every  $*$  by some new symbol "z"
  - every  $//$  by  $z/, z/z/, z/z/z/, \dots z/z/..z/$
- N+1 many z's
- N = length of longest  $*/./$  chain in  $q$

55

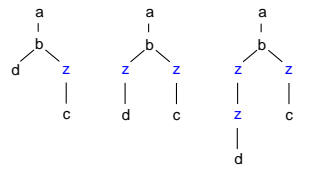
### 3. XPath Containment Test

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.

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

- (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(/, //, [])

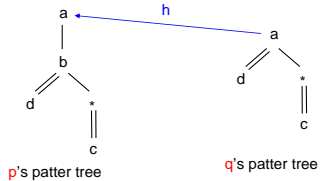
**Theorem**

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

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

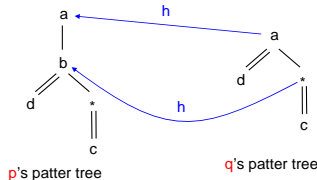


- (1) root of  $Q$  is mapped to root of  $P$
- (2) if  $(u,v)$  is child-edge of  $Q$  then  $(h(u),h(v))$  is child-edge of  $P$
- (3) if  $(u,v)$  is descendant-edge of  $Q$ , then  $h(v)$  is a "below"  $h(u)$  in  $P$
- (4) if  $u$  is labeled by "e" (not  $*$ ), then  $h(u)$  is also labeled by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

58

### 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

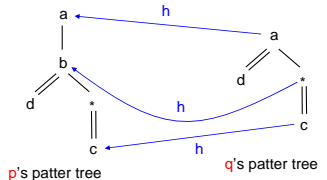


- (1) root of  $Q$  is mapped to root of  $P$
- (2) if  $(u,v)$  is child-edge of  $Q$  then  $(h(u),h(v))$  is child-edge of  $P$
- (3) if  $(u,v)$  is descendant-edge of  $Q$ , then  $h(v)$  is a "below"  $h(u)$  in  $P$
- (4) if  $u$  is labeled by "e" (not  $*$ ), then  $h(u)$  is also labeled by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

59

### 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

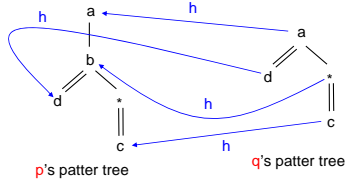


- (1) root of  $Q$  is mapped to root of  $P$
- (2) if  $(u,v)$  is child-edge of  $Q$  then  $(h(u),h(v))$  is child-edge of  $P$
- (3) if  $(u,v)$  is descendant-edge of  $Q$ , then  $h(v)$  is a "below"  $h(u)$  in  $P$
- (4) if  $u$  is labeled by "e" (not  $*$ ), then  $h(u)$  is also labeled by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

60

### 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".

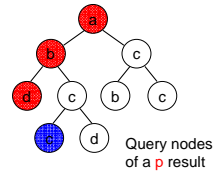
61

### 3. XPath Containment Test

"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  $p$ ?

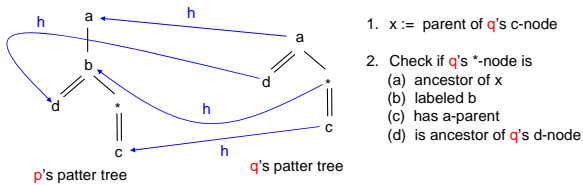


→ With homomorphism technique:

Use a result node of  $q$  together with run-time info on pattern nodes. Enables to search "inside", only on paths between pattern nodes.

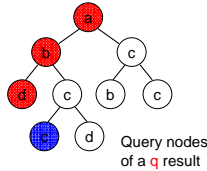
62

### 3. XPath Containment Test



1.  $x :=$  parent of  $q$ 's c-node
2. Check if  $q$ 's \*-node is
  - (a) ancestor of  $x$
  - (b) labeled  $b$
  - (c) has a-parent
  - (d) is ancestor of  $q$ 's d-node

Cave: we will have to try all homomorphisms ...



→ With homomorphism technique:

Use a result node of  $q$  together with run-time info on pattern nodes. Enables to search "inside", only on paths between pattern nodes.

63

### 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  $P$ .

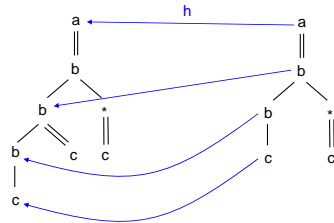
**Cave** If we add the star (\*) then homomorphism need not exist!

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

64

### 3. XPath Containment Test

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



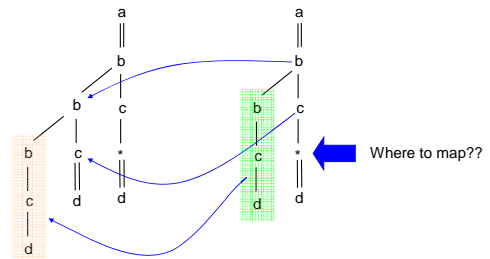
IS there a homomorphism??

**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  $P$  ☹

65

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



**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  $P$  ☹

66

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

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 ☹

67

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

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

This is a Java implementation of the theoretical work of  
Gerome Miklau and Dan Suciu (Containment and Equivalence  
for XPath with DTDs, ACM SIGMOD 2004) and  
Andreas Podelski (in a Journal of XPath, 2005, 2007)

**Instructions:**  
Enter two XPath expressions in the abbreviated syntax and press the  
button.  
For instance:  
if  $p = /a[b]$  and  $p' = /a[*]$   
the algorithm will detect that p is a subset of p'.  
Or if  $p = /a[*]b$  and  $p' = /a[*]b$   
the algorithm will detect that p is equal to p'  
because the subset equation holds in both directions.

[Download the Java Source Code](#)  
[Download Khaled's bachelor thesis \(in German\)](#)

If there is no application on the right side please contact  
our system administrator: k.haj-yalaya at gmx.de or k.haj-yalaya at uni-leipzig.de

$p = a/*//a$   
 $q = a/*//a$

Clearly, p is equivalent to q.  
(containment holds in both directions)

But, no homomorphisms exist.

**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 ☹

69

### 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.

70

### 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.

71

### 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**.

Emptiness test for automata

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

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

72

### 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<sub>p</sub>, A<sub>q</sub>, A<sub>E</sub>
- Construct B<sub>q</sub> for the complement of A<sub>q</sub> (=not q)
- Intersect B<sub>q</sub> with A<sub>p</sub> with A<sub>E</sub> (gives automaton A)
- Check if A accepts any tree.

73

### 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→d  
c2: b→c

a  
|  
b  
||  
d

p's pattern tree

74

### 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→d  
c2: b→c

a  
|  
b  
||  
d d c

p's pattern tree after chasing with c1,c2

75

### 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→d  
c2: b→c

a  
|  
b  
||  
d d c

h h

q's pattern tree

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

p's pattern tree after chasing with c1,c2

76

END  
Lecture 8

77