

XML and Databases

Lecture 3
XML into RDBMS

Sebastian Maneth
NICTA and UNSW

CSE@UNSW -- Semester 1, 2010

2

Lecture 3

XML Into RDBMS

Memory Representations

Facts

- DOM is easy to use, but memory heavy.
in-memory size usually **5-10 times larger** than size of original file.
- SAX is very flexible.
Using arrays or binary trees w/o backward pointers,
in-memory size is **approx. same** as size of original file.
- Can be further improved using DAGs/sharing-Graphs
& coding/compression for data values.

Memory Representations

Facts

- DOM is easy to use, but memory heavy.
in-memory size usually **5-10 times larger** than size of original file.
- SAX is very flexible.
Using arrays or binary trees w/o backward pointers,
in-memory size is **approx. same** as size of original file.
- Can be further improved using DAGs/sharing-Graphs
& coding/compression for data values.

TODAY

- ➔ How can we map XML into a **relational DB**?

Memory Representations

Facts

- DOM is easy to use, but memory heavy.
in-memory size usually **5-10 times larger** than size of original file.
- SAX is very flexible.
Using arrays or binary trees w/o backward pointers,
in-memory size is **approx. same** as size of original file.
- Can be further improved using DAGs/sharing-Graphs
& coding/compression for data values.

TODAY

- ➔ How can we map XML into a **relational DB**?

... but first: **Memory efficient tree traversals** using e.g. DOM?

Traversals and Pre/Post-Encoding

1. Pre-Order Traversal (recursively)
2. Post-Order
3. Pre-Order (iteratively)
4. Into RDBMS with **Pre/Post-encoding**

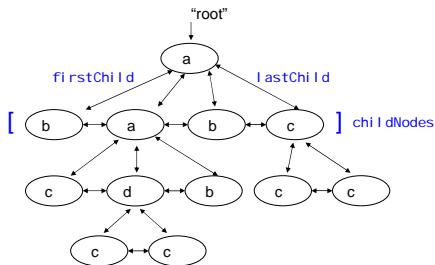
TODAY

- ➔ How can we map XML into a **relational DB**?

... but first: **Memory efficient tree traversals** using e.g. DOM?

Tree Traversals

Start at root node; want to visit every node.

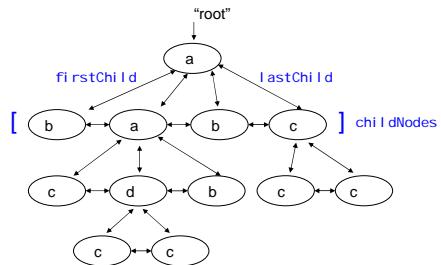


7

Tree Traversals

Start at root node; want to visit every node.

(1) recursively



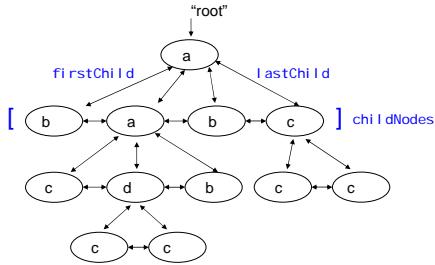
8

Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```



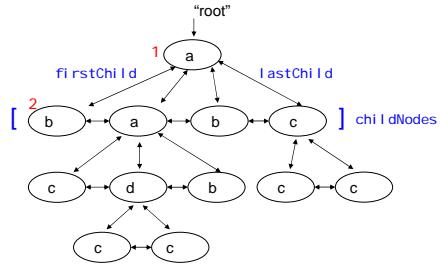
9

Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```



10

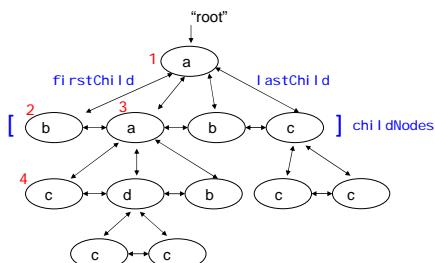
Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```

Tr(1)
pr(1)
Tr(2)
pr(2)
Tr(3)
pr(3)
Tr(4)
pr(4)



11

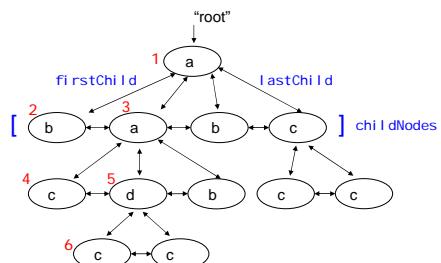
Tree Traversals

Start at root node; want to visit every node.

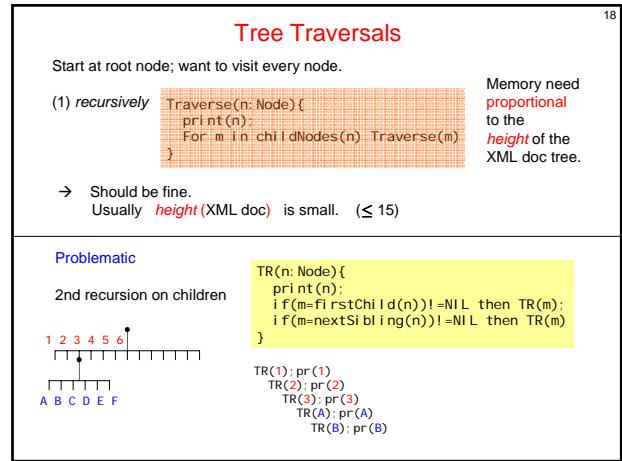
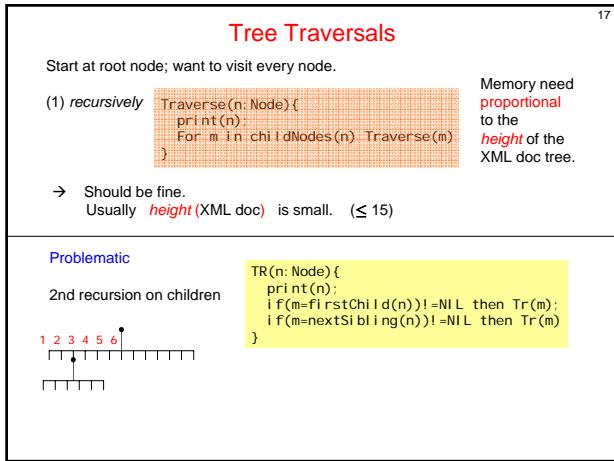
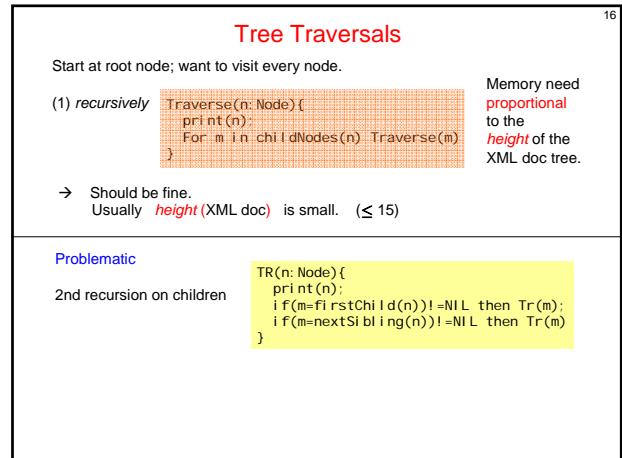
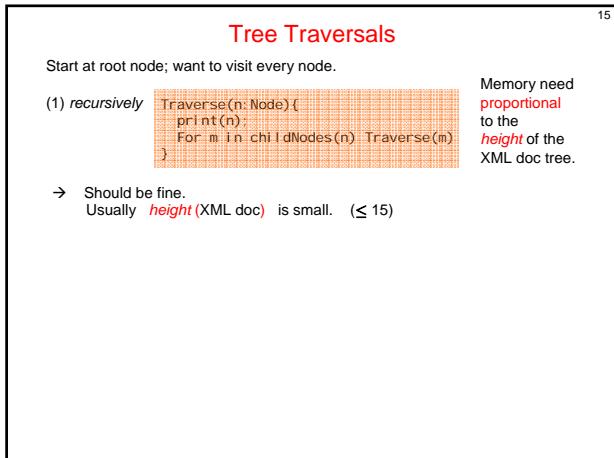
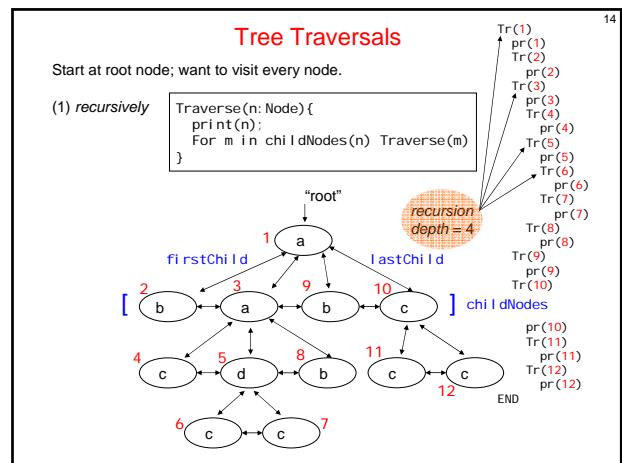
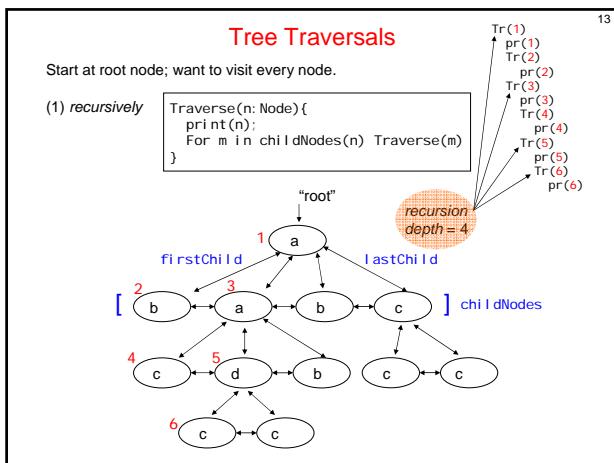
(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```

Tr(1)
pr(1)
Tr(2)
pr(2)
Tr(3)
pr(3)
Tr(4)
pr(4)
Tr(5)
pr(5)
Tr(6)
pr(6)



12



Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```

Memory need proportional to the height of the XML doc tree.

→ Should be fine.
Usually `height` (XML doc) is small. (≤ 15)

Problematic

2nd recursion on children

```
TR(n: Node){  
    print(n);  
    if(m=firstChild(n))!=NIL then Tr(m);  
    if(m=nextSibling(n))!=NIL then Tr(m)  
}
```

```
TR(1): pr(1)  
TR(2): pr(2)  
TR(3): pr(3)  
TR(A): pr(A)  
TR(B): pr(B)
```

Memory need proportional to max. length of `(firstChild | nextSibling)*-path`

Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```

Memory need proportional to the height of the XML doc tree.

→ Should be fine.
Usually `height` (XML doc) is small. (≤ 15)

Problematic

2nd recursion on children

```
TR(n: Node){  
    print(n);  
    if(m=firstChild(n))!=NIL then Tr(m);  
    if(m=nextSibling(n))!=NIL then Tr(m)  
}
```

```
TR(1): pr(1)  
TR(2): pr(2)  
TR(3): pr(3)  
TR(A): pr(A)  
TR(B): pr(B)
```

Memory need proportional to max. length of `(firstChild | nextSibling)*-path`

Can be HUGEEE!! =size(tree)

Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    if(m=firstChild(n))!=NIL then Tr(m);  
    if(m=nextSibling(n))!=NIL then Tr(m)  
}
```

Memory need proportional to max. length of `(firstChild | nextSibling)*-path`

Can be HUGEEE!! =size(tree)

Question
What is the max recursion depth on this tree?

Tree Traversals

Start at root node; want to visit every node.

(1) recursively

```
Traverse(n: Node){  
    print(n);  
    For m in childNodes(n) Traverse(m)  
}
```

Memory need proportional to the height of the XML tree.

→ Should be fine.
Usually `height` (XML doc) is small. (≤ 15)

Problematic

2nd recursion on children

Recall binary tree (firstChild/nextSibling) encoding.

Question
In the binary tree, what corresponds to this number?

```
TR(1): pr(1)  
TR(2): pr(2)  
TR(3): pr(3)  
TR(A): pr(A)  
TR(B): pr(B)
```

max. length of `(firstChild | nextSibling)*-path`

Binary Tree Encoding

Recall "firstChild/nextSibling" encoding.

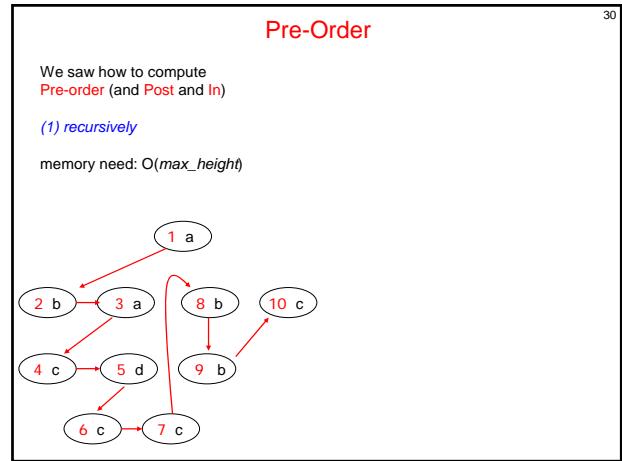
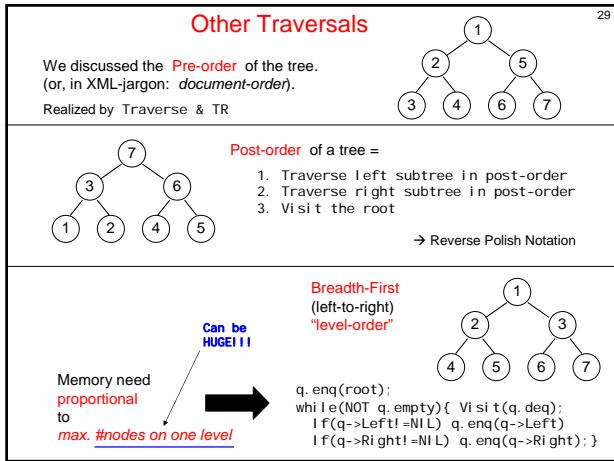
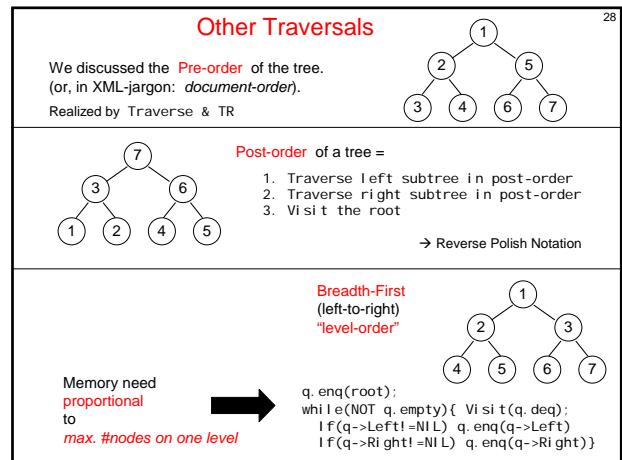
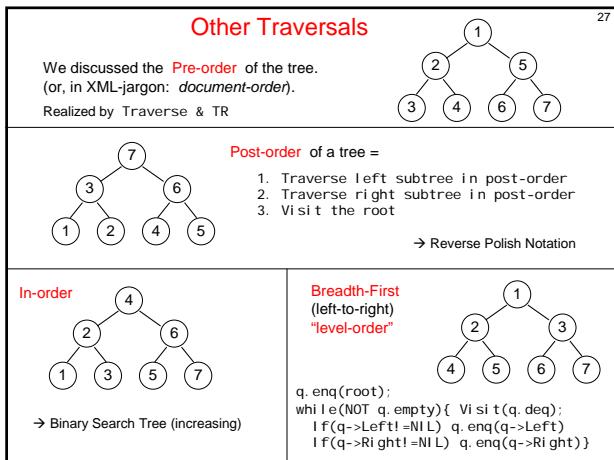
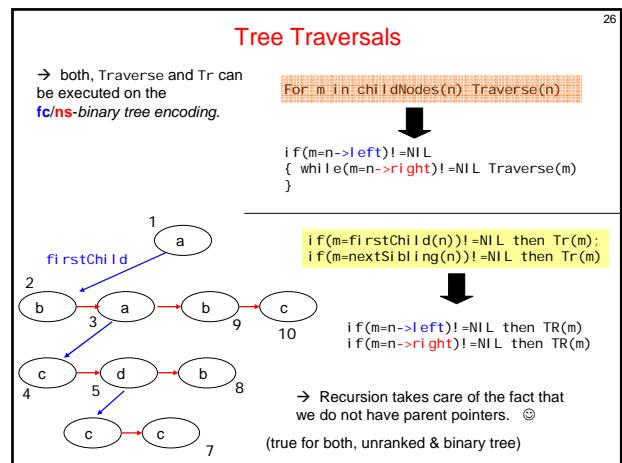
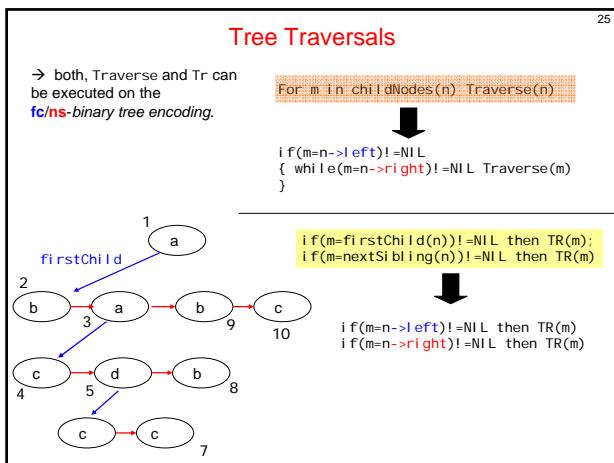
The "firstChild" becomes the `left` pointer
The "nextSibling" becomes the `right` pointer

per Node 2 pointers/IDs + label info

ID	fc: ns: lab
1	(2: -: a)
2	(-: 3: b)
3	(4: 9: a)
4	(-: 5: c)
5	(6: 8: d)
6	(-: 7: c)
7	(-: -: b)
8	(-: -: c)
9	(-: 10: b)
10	(-, -: c)

Tree Traversals

→ both, Traverse and TR can be executed on the `fc/ns-binary tree encoding`.



Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
while(firstChild(n) != NIL)
{   n = firstChild(n);
    pre(++i) = n;
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(2) iteratively

→ Memory need?

$i = 1; n = \text{root}; \text{pre}(i) = n; \text{while}(\text{firstChild}(n) \neq \text{NIL}) \{ n = \text{firstChild}(n); \text{pre}(++i) = n; \}$

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
while(firstChild(n) != NIL)
{   n = firstChild(n);
    pre(++i) = n;
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
    {   n = firstChild(n);
        pre(++i) = n;
    }
    while(nextSibling(n) == NIL)
    {   n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
    }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
    {   n = firstChild(n);
        pre(++i) = n;
    }
    while(nextSibling(n) == NIL)
    {   n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
    }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
    {   n = firstChild(n);
        pre(++i) = n;
    }
    while(nextSibling(n) == NIL)
    {   n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
    }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
        { n = firstChild(n);
        pre(++i) = n;
        }
    while(nextSibling(n) == NIL)
        { n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
        }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(2) iteratively

→ Memory need?

i = 1;
n = root;
pre(i) = n;
repeat {
 while(firstChild(n) != NIL)
 { n = firstChild(n);
 pre(++i) = n;
 }
 while(nextSibling(n) == NIL)
 { n = parent(n);
 n = nextSibling(n);
 pre(++i) = n;
 }
}

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
        { n = firstChild(n);
        pre(++i) = n;
        }
    while(nextSibling(n) == NIL)
        { n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
        }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
        { n = firstChild(n);
        pre(++i) = n;
        }
    while(nextSibling(n) == NIL)
        { n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
        }
}
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(1) recursively

memory need: $O(\max_height)$

```
i = 1;
n = root;
pre(i) = n;
repeat {
    while(firstChild(n) != NIL)
        { n = firstChild(n);
        pre(++i) = n;
        }
    while(nextSibling(n) == NIL)
        { n = parent(n);
        n = nextSibling(n);
        pre(++i) = n;
        }
}
if(n=NIL) then break; Ⓢ
```

Pre-Order

We saw how to compute Pre-order (and Post and In)

(2) iteratively

→ Memory need?

No recursion!
Needs constant memory
(only one pointer)

i = 1;
n = root;
pre(i) = n;
repeat {
 while(firstChild(n) != NIL)
 { n = firstChild(n);
 pre(++i) = n;
 }
 while(nextSibling(n) == NIL)
 { n = parent(n);
 n = nextSibling(n);
 pre(++i) = n;
 }
}
if(n=NIL) then break; Ⓢ

Pre-Order

Question
Given a *binary tree*, (top-down, no parent) how much memory do you need to compute `pre`?

No recursion!
Needs constant memory!
(only one pointer)

```
i=1;
n=root;
pre(i)=n;
repeat {
    while(firstChild(n)!=NIL)
    { n=firstChild(n);
      pre(++i)=n;
    }
    while(nextSibling(n)!=NIL)
    { n=parent(n);
      n=nextSibling(n);
      pre(++i)=n;
    }
    if(n=NIL) then break; ⊗
}
```

Pre-Order

Question
Given a *binary tree*, (top-down, no parent) how much memory do you need to compute `pre`?

No recursion!
Needs constant memory!
(only one pointer)

→ Can you do it w. **constant memory**?

```
i=1;
n=root;
pre(i)=n;
repeat {
    while(firstChild(n)!=NIL)
    { n=firstChild(n);
      pre(++i)=n;
    }
    while(nextSibling(n)!=NIL)
    { n=parent(n); n=parent(n)
      n=nextSibling(n);
      pre(++i)=n;
    }
    if(n=NIL) then break; ⊗
}
```

Pre-Order

Question
Fun (MS ji)
How much memory you need to check for cycles, in a single-linked (pointer) list?

No recursion!
Needs constant memory!
(only one pointer)

```
i=1;
n=root;
pre(i)=n;
repeat {
    while(firstChild(n)!=NIL)
    { n=firstChild(n);
      pre(++i)=n;
    }
    while(nextSibling(n)!=NIL)
    { n=parent(n);
      n=nextSibling(n);
      pre(++i)=n;
    }
    if(n=NIL) then break; ⊗
}
```

Pre-Order

Question
Do you see how to do **Post-** and **In-order** iteratively?

No recursion!
Needs constant memory!
(only one pointer)

```
i=1;
n=root;
pre(i)=n;
repeat {
    while(firstChild(n)!=NIL)
    { n=firstChild(n);
      pre(++i)=n;
    }
    while(nextSibling(n)!=NIL)
    { n=parent(n); n=parent(n)
      n=nextSibling(n);
      pre(++i)=n;
    }
    if(n=NIL) then break; ⊗
}
```

Pre-Order

From `pre()`

we can compute → `PreFollowing(n) = { nodes m with pre(m) > n }`
→ `PrePreceding(n) = { nodes m with pre(m) < n }`

Pre-Order

From `pre()`

we can compute → `PreFollowing(n) = { nodes m with pre(m) > n }`
→ `PrePreceding(n) = { nodes m with pre(m) < n }`

What is this?

Pre-Order

From `pre()`

we can compute

$$\rightarrow \text{PreFollowing}(n) = \{ \text{nodes } m \text{ with } \text{pre}(m) > n \}$$

$$\rightarrow \text{PrePreceding}(n) = \{ \text{nodes } m \text{ with } \text{pre}(m) < n \}$$

What is this?

	<code>pre</code>	<code>folI</code>	<code>prec</code>
1	2-10	-	
2	3-10	1	
3	4-10	1-2	
4	5-10	1-3	
5	6-10	1-4	
6	7-10	1-5	
7	8-10	1-6	
8	9-10	1-7	
9	10	1-8	
10	-	1-9	

Pre-Order

From `pre()`

we can compute

$$\rightarrow \text{PreFollowing}(n) = \{ \text{nodes } m \text{ with } \text{pre}(m) > n \}$$

$$\rightarrow \text{PrePreceding}(n) = \{ \text{nodes } m \text{ with } \text{pre}(m) < n \}$$

...useful???

Not "tree (navigation) complete" ☺

	<code>pre</code>
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	

XML to RDBMS Encoding

Mapping XML to Databases | Introduction

Relational XML processors (2)

Our approach to **relational XQuery processing**:

- The XQuery data model—ordered, unranked trees and ordered item sequences—is, in a sense, alien to a relational database kernel.
- A **relational tree encoding** \mathcal{E} is required to map trees into the relational domain, i.e., tables.

Relational tree encoding \mathcal{E}

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 331

Mapping XML to Databases | Introduction

What makes a good (relational) (XML) tree encoding?

Hard requirements:

- \mathcal{E} is required to reflect **document order** and **node identity**.
 - ▶ Otherwise: cannot enforce XPath semantics, cannot support `<<` and `is`, cannot support node construction.
- \mathcal{E} is required to encode the **XQuery DM node properties**.
 - ▶ Otherwise: cannot support XPath axes, cannot support XPath node tests, cannot support atomization, cannot support validation.
- \mathcal{E} is able to encode any well-formed **schema-less XML fragment** (i.e., \mathcal{E} is "**schema-oblivious**", see below).
 - ▶ Otherwise: cannot process non-validated XML documents, cannot support arbitrary node construction.

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 332

Mapping XML to Databases | Introduction

What makes a good (relational) (XML) tree encoding?

Soft requirements (primarily motivated by performance concerns):

- **Data-bound operations** on trees (potentially delivering/copying lots of nodes) should map into efficient database operations.
 - ▶ `XPath location steps` (12 axes)
- Principal, recurring **operations imposed by the XQuery semantics** should map into efficient database operations.
 - ▶ `Subtree traversal` (atomization, element construction, serialization).

For a relational encoding, "database operations" always mean "table operations" ...

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 333

XML to RDBMS Encoding

`pre()` is not enough

Other possibilities: (1) large (unparsed) text block
 (2) Schema-based encoding
 (3) Adjacency-based encoding

Dead Ends
No good...

Questions Why is (1) a dead end?

Relational Tree Encoding | Dead Ends

Dead end #2: Schema-based encoding

XML address database (excerpt)

```
<person>
  <name><first>John</first><last>Foo</last></name>
  <address><street>13 Main St</street>
    <zip>12345</zip><city>Miami</city>
  </address>
</person>
<person>
  <name><first>Erik</first><last>Bar</last></name>
  <address><street>42 Kings Rd</street>
    <zip>54321</zip><city>New York</city>
  </address>
</person>
```

Schema-based relational encoding: table person

id	first	last	street	zip	city
0	John	Foo	13 Main St	12345	Miami
1	Erik	Bar	42 Kings Rd	54321	New York

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 335

Relational Tree Encoding | Dead Ends

Dead end #2: Schema-based encoding

Irregular hierarchy

```
<a no="0">
  <b><c>X</c></b>
</a>
<a no="1">
  <b><c>Y</c></b>
</a>
<a><b/></a>
<a no="3"/>
```

A relational encoding

id	@no	b	id	b	c
0	0	α	1	α	X
3	1	β	2	α	NULL ^c
5	NULL ^a	γ	4	β	Y
6	3	NULL ^b			

Issues:

- Number of encoding tables depends on nesting depth.
- Empty element c encoded by NULL^c, empty element b encoded by absence of γ (will need outer join on column b).
- NULL^a encodes absence of attribute, NULL^b encodes absence of element.
- Document order/identity of b elements only implicit.

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 337

Relational Tree Encoding | Dead Ends

Dead end #3: Adjacency-based encoding

Adjacency-based encoding of XML fragments

```
<a id="0">
  <b>fo</b>o
  <c>
    <d>b</d><e>ar</e>
  </c>
</a>
```

≡

Resulting relational encoding

id	parent	tag	text	val
0	NULL	a	NULL	NULL
1	0	@id	NULL	"0"
2	0	b	NULL	NULL
3	2	NULL	"fo"	NULL
4	0	NULL	"o"	NULL
5	0	c	NULL	NULL
⋮				

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 338

Relational Tree Encoding | Dead Ends

Dead end #3: Adjacency-based encoding

Pro:

- Since this captures all adjacency, kind, and content information, we can—in principle—**serialize the original XML fragment**.
- Node identity** and **document order** is adequately represented.

Contra:

- The XQuery processing model is not well-supported: subtree traversals require **extra-relational queries (recursion)**.
- This is completely parent-child centric. How to support descendant, ancestor, following, or preceding?

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 339

60

Pre/Post Encoding

→ Add POST order

PRE	POST	lab
1	10	a
2	1	b
3	6	a
4	2	c
5	5	d
6	3	c
7	4	c
8	7	b
9	8	b
10	9	c

CREATE VIEW descendant AS

```
SELECT r1.pre, r2.pre FROM R r1, R r2
  WHERE r1.pre < r2.pre
    AND r1.post > r2.post
```

61

Pre/Post Encoding

→ Add POST order

	PRE	POST	lab
1	10	a	
2	1	b	
3	6	a	
4	2	c	
5	5	d	
6	3	c	
7	4	c	
8	7	b	
9	8	b	
10	9	c	

CREATE VIEW descendant AS
SELECT r1.pre, r2.post FROM R r1, R r2
WHERE r1.pre < r2.pre
AND r1.post > r2.post

"structural join"

XPath Accelerator Encoding | Pre-Order and Post-Order Traversal Ranks

XPath axes in the pre/post plane

Plane partitions ≡ XPath axes, o is arbitrary!

Pre/post plane regions ≡ major XPath axes
The **major XPath axes** descendant, ancestor, following, preceding correspond to rectangular **pre/post plane windows**.

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 352

XPath Accelerator Encoding | Pre-Order and Post-Order Traversal Ranks

XPath Accelerator encoding

XML fragment f and its skeleton tree

```
<a>
  <b><c></b>
  <!--d-->
  <e><f><g/><?h?></f>
    <i>j</i>
  </e>
</a>
```

Pre/post encoding of f: table accel

pre	post	par	kind	tag	text
0	9	NULL	elem	a	NULL
1	1	0	elem	b	NULL
2	0	1	text	NULL	c
3	2	0	com	NULL	d
4	8	0	elem	e	NULL
5	5	4	elem	f	NULL
6	3	5	elem	g	NULL
7	4	5	pi	NULL	h
8	7	4	elem	i	NULL
9	6	8	text	NULL	j

Marc H. Scholl (DBIS, Uni KN) XML and Databases Winter 2005/06 353

64

Pre/Post Encoding

Straightforward how to compute, for a given node,

- descendants
- ancestors
- following
- preceding

Questions

How to do

- lastChild
- parent
- childNodes

Can you find corresponding SQL queries?

65

END
Lecture 3