

XML and Databases

Lecture 2
Memory Representations for XML: Space vs Access Speed

Sebastian Maneth
NICTA and UNSW

CSE@UNSW -- Semester 1, 2009

Reminder

You can freely choose to program your assignments in

- C / C++, or
- Java

However, your code **must compile with gcc / g++, javac**, as installed on CSE linux systems!

Assignment 1 is due Monday 23:59, 25th of March!
Submit your code using

```
% gl ve cs4317 ass1 filename.cpp
```

```
% gl ve cs4317 ass1 filename.java
```

Lecture 2

XML into Memory

Problem with DOM

- Uses massive amounts of memory.
- Even if application touches only a single element node, the **DOM API** has to maintain a data structure that represents the **whole XML input document**.

Example

XML size		DOM process size	
81M	x 2	164M	Text only, with one embracing element
52M	x 13.1	680M	Treebank, deep tree structure with short texts
....			

Usually: more than **10-times** blow up!!

To remedy the memory hunger of DOM ...

Preprocess (i.e., filter) the input XML document to reduce its overall size.

- Use an XPath/XSLT processor to preselect *interesting* document regions.
- CAVE: *no updates* on the input XML document are possible
- CAVE: make sure the XPath/XSLT processor is *not* implemented on top of DOM!

- Use a **completely different** approach to XML processing (→ **SAX**)

"design your own XML data structure
and fill in with what you need..."

To remedy the memory hunger of DOM ...

Preprocess (i.e., filter) the input XML document to reduce its overall size.

- Use an XPath/XSLT processor to preselect *interesting* document regions.
- CAVE: *no updates* on the input XML document are possible
- CAVE: make sure the XPath/XSLT processor is *not* implemented on top of DOM!

- Use a **completely different** approach to XML processing (→ **SAX**)

"design your own XML data structure
and fill in with what you need..."

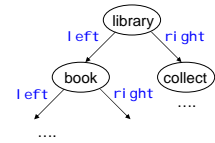
Outline

1. Tree Pointer Structures
2. Binary Tree Encodings
3. Minimal Unique DAGs
4. How to use SAX

1. Tree pointer structures

1. Consider *binary trees*

```
Type Node {
  label : String,
  left  : Node,
  right : Node
}
```

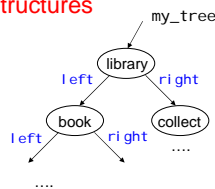


How much memory for *n*-node binary tree?

1. Tree pointer structures

1. Consider *binary trees*

```
Type Node {
  label : String,
  left  : Node,
  right : Node
}
```



How much memory for *n*-node binary tree?

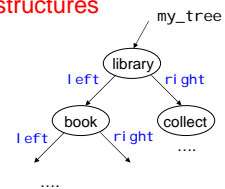
cadr1: |l|l|b|r|a|r|y|0| my_tree: cadr1 tadr1 tadr2
cadr2: |b|o|o|k|0| tadr1: cadr2 tadr3 tadr4
...

length(label_1) + 1 3 * length(pointer) * n
+ length(label_2) + 1
+ ... + length(label_n) + 1

1. Tree pointer structures

1. Consider *binary trees*

```
Type Node {
  label : String,
  left  : Node,
  right : Node
}
```



How much memory for *n*-node binary tree?

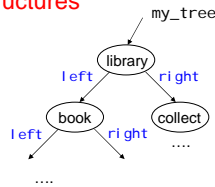
cadr1: |l|l|b|r|a|r|y|0| my_tree: cadr1 tadr1 tadr2
cadr2: |b|o|o|k|0| tadr1: cadr2 tadr3 tadr4
...

length(label_1) + 1 3 * length(pointer) * n
+ length(label_2) + 1
+ ... + length(label_n) + 1 typical: 4 bytes

Tree pointer structures

1. Consider *binary trees*

```
Type Node {
  label : String,
  left  : Node,
  right : Node
}
```



How much memory for *n*-node binary tree?

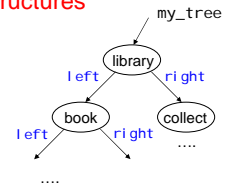
→ Whatever is needed for the labels
PLUS 12 bytes per node.

length(label_1) + 1 3 * length(pointer) * n
+ length(label_2) + 1
+ ... + length(label_n) + 1 typical: 4 bytes

Tree pointer structures

1. Consider *binary trees*

```
Type Node {
  label : String,
  left  : Node,
  right : Node
}
```



How much memory for *n*-node binary tree?

→ Whatever is needed for the labels
PLUS 12 bytes per node.

Can easily be optimized:
E.g., store each distinct
string only once!

length(label_1) + 1 3 * length(pointer) * n
+ length(label_2) + 1
+ ... + length(label_n) + 1 typical: 4 bytes

13

Tree pointer structures

1. Consider *binary trees*

```

Type Node {
  label : String,
  left  : Node,
  right : Node
}

```

Serialization to **XML**

```

<library><book>< ... > ... </book></library>
  ^         ^         ^         ^         ^

```

#characters per node: $5 + 2 * \text{Length}(\text{label})$

→ E.g., one node w. 4-character ASCII label: **13 bytes** (assuming UTF-8!)

14

Tree pointer structures

1. Consider *binary trees*

```

Type Node {
  label : String Byte,
  left  : Node,
  right : Node
}

```

Often #distinct node labels is small, *100. → Fits in *one Byte*
Then, only **9 bytes per node**.

→ **MEM**(n-node binary tree pointer struc, *256 labels)
= **SIZE**(n-node binary tree in XML, average label length=2)

#characters per node: $5 + 2 * \text{Length}(\text{label})$

→ One node w. 2-character ASCII label: **9 bytes** (assuming UTF-8!)

15

Tree pointer structures

Nice
Following pointers is fast!
→ *much higher access speed!*
(than on *doc seen as string..*)

E.g.
at root, get right-child.

Often #distinct node labels is small, *100. → Fits in *one Byte*
Then, only **9 bytes per node**.

→ **MEM**(n-node binary tree pointer struc, *256 labels)
= **SIZE**(n-node binary tree in XML, average label length=2)

#characters per node: $5 + 2 * \text{Length}(\text{label})$

→ One node w. 2-character ASCII label: **9 bytes** (assuming UTF-8!)

16

Tree pointer structures

1. Consider *binary trees*

Plain no attributes, no text nodes, ...

Question

Using a (top-down) pointer structure, as the one above,
how can you implement a **DOM interface**?

Node	nodeName	: DOMString	
	parentNode	: Node	
	firstChild	: Node	leftmost child
	nextSibling	: Node	returns NULL for root elem
	childNodes	: NodeList	

17

Tree pointer structures

1. Consider *binary trees*

Plain no attributes, no text nodes, ...

Question

Using a (top-down) pointer structure, as the one above,
how can you implement a **DOM interface**?

Node	nodeName	: DOMString	
	parentNode	: Node	
	firstChild	: Node	leftmost child
	nextSibling	: Node	returns NULL for root elem
	childNodes	: NodeList	

→ **At run-time** a node is represented as a pointer, PLUS a **stack of pointers of all its ancestors**.

(Node, [parent(Node)::parent(parent(Node))::...::root-node])

18

Tree pointer structures

Access speed of **parentNode** should be approx same, as in a native DOM.
→ What about access speed of **nextSibling**?

What is the *run-time size* of our "binary DOM-tree" data structure? (WC/average)

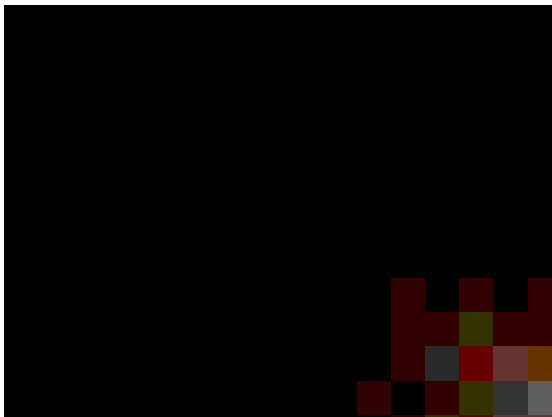
Question

Using a (top-down) pointer structure, as the one above,
how can you implement a **DOM interface**?

Node	nodeName	: DOMString	
	parentNode	: Node	
	firstChild	: Node	leftmost child
	nextSibling	: Node	returns NULL for root elem
	childNodes	: NodeList	

→ **At run-time** a node is represented as a pointer, PLUS a **stack of pointers of all its ancestors**.

(Node, [parent(Node)::parent(parent(Node))::...::root-node])



19

```
interface Node { // NodeType
const unsigned short ELEMENT_NODE = 1;
const unsigned short ATTRIBUTE_NODE = 2;
const unsigned short TEXT_NODE = 3;
const unsigned short CDATA_SECTION_NODE = 4;
const unsigned short ENTITY_REFERENCE_NODE = 5;
const unsigned short ENTITY_NODE = 6;
const unsigned short PROCESSING_INSTRUCTION_NODE = 7;
const unsigned short COMMENT_NODE = 8;
const unsigned short DOCUMENT_NODE = 9;
const unsigned short DOCUMENT_TYPE_NODE = 10;
const unsigned short DOCUMENT_FRAGMENT_NODE = 11;
const unsigned short NOTATION_NODE = 12;
readonly attribute DOMString nodeName;
attribute DOMString nodeValue; // raises(DOMException) on setting
// raises(DOMException) on retrieval
readonly attribute unsigned short nodeType;
readonly attribute Node parentNode;
readonly attribute NodeList childNodes;
readonly attribute Node firstChild;
readonly attribute Node lastChild;
readonly attribute Node previousSibling;
readonly attribute Node nextSibling;
readonly attribute NamedNodeMap attributes;
readonly attribute Document ownerDocument;
Node insertBefore(in Node newChild, in Node refChild) raises(DOMException);
Node replaceChild(in Node newChild, in Node oldChild) raises(DOMException);
Node removeChild(in Node oldChild) raises(DOMException);
Node appendChild(in Node newChild) raises(DOMException);
boolean hasChildNodes(); Node cloneNode(in boolean deep); }
```

20

Tree pointer structures

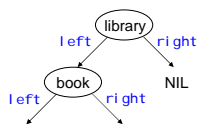
To slash memory hunger (of, e.g., DOM...)

LESSON 1

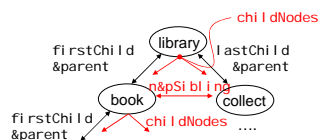
→ Avoid all backward pointers (build them online, dynamically)

21

binary trees



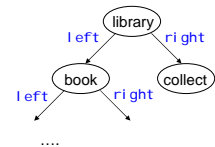
DOM



Tree pointer structures

1. Consider *binary trees*

```
Type Node {
label : String,
left : Node,
right : Node
}
```



22

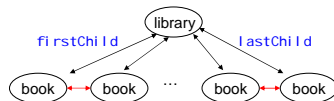
How much memory for n-node binary tree?

How to add *attributes and text nodes*?

→ e.g., "into the label" ...

Tree pointer structures

2. Consider *unranked trees*



unranked = no a priori bound on #children of a node.

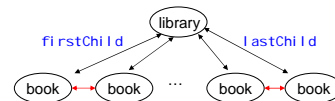
Tree structure of XML: **unranked trees! (not binary)**

```
Type Node {
label : String,
children : List[Node]
}
```

23

Tree pointer structures

2. Consider *unranked trees*



unranked = no a priori bound on #children of a node.

Tree structure of XML: **unranked trees! (not binary)**

```
Type Node {
label : String,
children : List[Node]
}
```

→ How much memory for `List[Node]` of n nodes?

24

25

Tree pointer structures

2. Consider **unranked** trees

unranked = no a priori bound on #children of a node.

Tree structure of XML: **unranked trees!** Typically

```

Type Node {
  label : String,
  children : List[Node]
}

```

→ How much memory for `List[Node]` of n nodes?

26

Tree pointer structures

2. Consider **unranked** trees

unranked = no a priori bound on #children of a node.

Tree structure of XML: **unranked trees!** Typically

```

Type Node {
  label : String,
  children : List[Node]
}

```

→ How much memory for `List[Node]` of n nodes? $2 \cdot n$ pointers

27

Tree pointer structures

2. Consider **unranked** trees

→ In this way, a node of a *binary tree* needs **5 pointers** ☹ (plus label info/pointer..)

unranked = no a priori bound on #children of a node.

Tree structure of XML: **unranked trees!** Typically

```

Type Node {
  label : String,
  children : List[Node]
}

```

→ How much memory for `List[Node]` of n nodes? $2 \cdot n$ pointers

28

Tree pointer structures

2. Consider **unranked** trees

→ In this way, a node of a *binary tree* needs **5 pointers** ☹ (plus label info/pointer..)

More efficient possibilities:

(1) Use **arrays**. Store #children (e.g., in label). n pointers + $(\log d)$ Bits

(2) Encode tree as binary tree.

Typically

```

Type Node {
  label : String,
  children : List[Node]
}

```

→ How much memory for `List[Node]` of n nodes? $2 \cdot n$ pointers

29

Tree pointer structures

2. Consider **unranked** trees

→ In this way, a node of a *binary tree* needs **5 pointers** ☹ (plus label info/pointer..)

More efficient possibilities:

(1) Use **arrays**. Store #children (e.g., in label). n pointers + $(\log d)$ Bits

(2) → Encode tree as binary tree. ←

Typically

```

Type Node {
  label : String,
  children : List[Node]
}

```

→ How much memory for `List[Node]` of n nodes? $2 \cdot n$ pointers

30

2. Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

31

2. Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

32

2. Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

33

2. Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

34

2. Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

35

Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

Popular encoding: **"firstChild/nextSibling" encoding.**

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

36

Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

"firstChild/nextSibling" encoding

n-node unranked tree $\xrightarrow{\text{encoding}}$ n-node binary tree

decoding

Questions

- Time overhead for simulating **firstChild** access, on the **binary encoding**?
- Can you think of other binary tree encodings?
- How to simulate preceding-sibling?

37

Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

n-node unranked tree

“firstChild/nextSibling” encoding

→

←

n-node binary tree

decoding

Good Property of the firstChild/nextSibling encoding:

→ XML types (e.g., DTD, XML Schema, Relax NG) are preserved when going from unranked to binary (and vice versa).

38

Binary Tree Encodings

Any unranked tree can be encoded as a binary tree.

n-node unranked tree

“firstChild/nextSibling” encoding

→

←

n-node binary tree

decoding

Good Property of the firstChild/nextSibling encoding:

→ XML types (e.g., DTD, XML Schema, Relax NG) are preserved when going from unranked to binary (and vice versa).

LESSON 2 ... against memory hunger ...

→ Use binary trees instead of unranked trees. (... or use efficient arrays)

+ Fast child-n access

- Expensive to update (insert/delete)

39

Tree Pointer Structures

Question

Give a datatype for binary trees which stores only non-NIL pointers.

Then, n-node tree: <n pointers

```

Type Node {
  label : String,
  left  : Node,
  right : Node
}

```

40

3. Minimal Unique DAGs

L1: no backward pointers

L2: use binary trees or efficient arrays

Type Node {

label : Byte,

left : Node,

right : Node

}

binary tree

2 pointers per node

Can we do with even less pointers?

41

3. Minimal Unique DAGs

L1: no backward pointers

L2: use binary trees or efficient arrays

Type Node {

label : Byte,

left : Node,

right : Node

}

binary tree

2 pointers per node

Can we do with even less pointers?

n-node tree: 2n pointers

left

c

right

a

a

a

a

a

a

a

a

a

a

a

a

42

3. Minimal Unique DAGs

L1: no backward pointers

L2: use binary trees or efficient arrays

Type Node {

label : Byte,

left : Node,

right : Node

}

binary tree

2 pointers per node

Can we do with even less pointers?

n-node tree: 2n pointers

YES!

→ Directed Acyclic Graph DAG

left

c

right

a

a

a

a

a

a

a

a

a

a

a

a

share identical subtrees

left

c

right

a

a

a

a

a

a

a

a

a

a

a

a

7

3. Minimal Unique DAGs

```
Type Node {
  label : Byte,
  left : Node,
  right : Node
}
```

binary tree
2 pointers per node

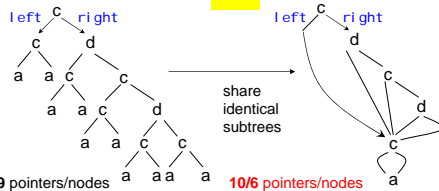
"18 pointers" really means
"18 non-NIL pointers"

Can we do with even less pointers?

n-node tree: 2n pointers

YES!

→ Directed Acyclic Graph DAG



3. Minimal Unique DAGs

A DAG representation of a tree has always

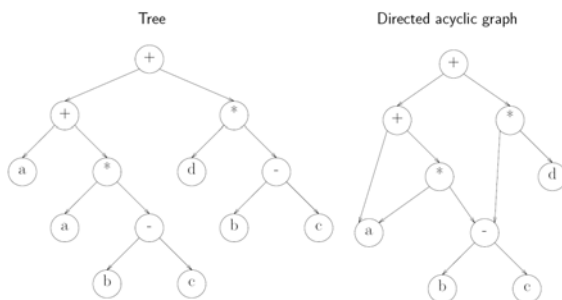
→ Less than or equal #nodes than the tree

→ Less than or equal #pointers than the tree.

3. Minimal Unique DAGs

Local optimizations

Consider the expression: $a + a * (b - c) + (b - c) * d$



3. Minimal Unique DAGs

Local optimizations

Common subexpressions (CSE)

- portion of expressions
- repeated multiple times
- computes same value
- can reuse previously computed value

Directed acyclic graph (DAG)

- program representation
- nodes can have multiple parents
- no cycles allowed
- exposes common subexpressions

Building a DAG for an expression

- maintain hash table for leafs, expressions
- unique name for each node — its *value number*
- reuse nodes found in hash table

3. Minimal Unique DAGs

(minimal) DAGs have many applications!

→ CSE (Common Subexpression Elimination)
for efficient evaluation of expressions
(do "term graph" rewriting, instead of term rewriting)

→ Model checking with BDDs
Binary Decision Diagrams
for efficient evaluation of logic formulas

→ Efficient XML query evaluation

3. Minimal Unique DAGs

(minimal) DAGs have many applications!

→ CSE (Common Subexpression Elimination)
for efficient evaluation of expressions
(do "term graph" rewriting, instead of term rewriting)

→ Model checking with BDDs
Binary Decision Diagrams
for efficient evaluation of logic formulas

→ Efficient XML query evaluation

Btw, inside of a DAG, you have "referential completeness"
→ structural equality = equality of pointers ☺

49

3. Minimal Unique DAGs

→ Every tree has a minimal, unique DAG!

→ The DAG is at most *exponentially* smaller than the tree.

→ Building the minimal unique DAG is easy!
Can be done in (amortized) *linear time*.

50

3. Minimal Unique DAGs

→ Every tree has a minimal, unique DAG!

→ The DAG is at most exponentially smaller than the tree.

→ Building the minimal unique DAG is easy!
Can be done in (amortized) *linear time*.

How?

51

3. Minimal Unique DAGs

→ Every tree has a minimal, unique DAG!

→ The DAG is at most exponentially smaller than the tree.

→ Building the minimal unique DAG is easy!
Can be done in (amortized) *linear time*.

How?

(even while parsing)
→ Build a *hash table* of all subtrees seen so far

(we don't want to compare many trees, node by node, later on..)

Question Give a simple hash function that works for the tree above.

52

3. Minimal Unique DAGs

Hash Table HT

1	a
2	c(a,a), c(c(a,a),a)
3	a(b,c)
4	.
5	.
6	.

Hash function *f*

a hash "bucket"

We want

→ *f* distributes trees uniformly into buckets
→ test if a tree *T* is in HT, time $O(\text{size}(T))$

Question Give a simple hash function that works for the tree above.

53

Minimal Unique DAGs

Example "Parse & DAGify"

1: startElement(c)

hash	content
------	---------

54

Minimal Unique DAGs

1: bib [2, 3, 4, 5]	<bib>
2: book [6, 7]	<book>
3: article [8, 9]	<author></author>
4: book [10, 11]	<title></title>
5: article [12, 13]	</book>
6: author	<article>
7: title	<author></author>
8: author	<title></title>
9: title	</article>
10: price	<book>
11: title	<price></price>
12: price	<title></title>
13: title	</book>

Minimal Unique DAGs

- ```

1: bib [2, 3, 4, 5]
2: book [6, 7]
3: article [8, 9]
4: book [10, 11]
5: article [12, 13]
6: author
7: title
8: author
9: title
10: price
11: title
12: price
13: title

```

```
<bi b>
 <book>
 <author></author>
 <title></title>
 </book>
 <article>
 <author></author>
 <title></title>
 </article>
 <book>
 <price></price>
 <title></title>
 </book>
 <article>
 <price></price>
 <title></title>
 </article>
</book>
```

## Minimal Unique DAGs

- ```

1: bib [2, 3, 4, 5]
2: book [6, 7]
3: article [8, 9]
4: book [10, 11]
5: article [12, 13]
6: author
7: title
8: author
9: title
10: price
11: title
12: price
13: title

```

```
<bi b>
<book>
  <author></author>
  <ti tle></ti tle>
</book>
<arti cl e>
  <author></author>
  <ti tle></ti tle>
</arti cl e>
<book>
  <pr i ce></pr i ce>
  <ti tle></ti tle>
</book>
<arti cl e>
  <pr i ce></pr i ce>
  <ti tle></ti tle>
</arti cl e>
</book>
```

Minimal Unique DAGs

-
- 1: bib [2, 3, 4, 5]
 2: book [6, 7]
 3: article [8, 9]
 4: book [10, 11]
 5: article [12, 13]
 6: author
 7: title
 8: author
 9: title
 10: price
 11: title
 12: price
 13: title

```
<bi b>
</book>
<author></author>
<ti tle></ti tle>
</book>
<arti cle>
<author></author>
<ti tle></ti tle>
</arti cle>
</book>
<price></pri ce>
<ti tle></ti tle>
</book>
<arti cle>
<price></pri ce>
<ti tle></ti tle>
</arti cle>
</book>
```

Minimal Unique DAGs

- 1: bib [2, 3, 4, 5] 6
2: book [6, 7]
3: article [8, 9] 7
4: book [10, 11]
5: article [12, 13]
6: author
7: title
8: author
9: title
10: price
11: title
12: price
13: title

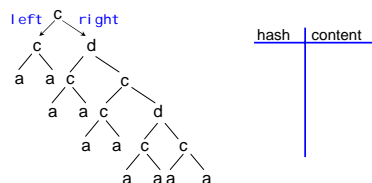
```
<bi b>
<book>
<author></author>
<ti tle></ti tle>
</book>
<article>
<author></author>
<ti tle></ti tle>
</article>
<book>
<pr ice></pr ice>
<ti tle></ti tle>
</book>
<article>
<pr ice></pr ice>
<ti tle></ti tle>
</article>
</book>
```

minimal unique DAG
8 nodes (vs 13 nodes in the original tree)

Minimal Unique DAGs

Example “Parse & DAGify”

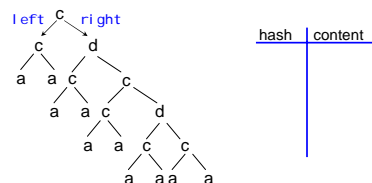
- ```
1: startElement(c)
2: startElement(c)
```



## Minimal Unique DAGs

### Example “Parse & DAGify”

- ```
1: startElement(c)
2: startElement(c)
3: startElement(a)
```



61

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a) →

1. p1=hashT.find(a)
2. if(p1==NULL) { p1=new("a-node", NULL, NULL)
hashT.insert(p1) }

hash	content

62

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a) →

1. p1=hashT.find(a)
2. if(p1==NULL) { p1=new("a-node", NULL, NULL)
hashT.insert(p1) }

hash	content
1	p1
2	

Memory location p1 is a DAG with root node a, and child1-pointer=NULL child2-pointer=NULL

63

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a) →

1. p1=hashT.find(a)
2. if(p1==NULL) { p1=new("a-node", NULL, NULL)
hashT.insert(p1) }

→ must store children lists: [[], [p1]]

hash	content
1	
2	p1

Memory location p1 is a DAG with root node a, and child1-pointer=NULL child2-pointer=NULL

64

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a)
- 5: startElement(a)

hash	content
1	
2	p1

65

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a)
- 5: startElement(a)
- 6: endElement(a) →

1. p2=hashT.find(a)
2. if(p2==NULL) { p2=new("a-node", NULL, NULL)
hashT.insert(p2) }

hash	content
1	
2	p1

66

Minimal Unique DAGs

Example "Parse & DAGify"

- 1: startElement(c)
- 2: startElement(c)
- 3: startElement(a)
- 4: endElement(a)
- 5: startElement(a)
- 6: endElement(a) →

1. p2=hashT.find(a)=p1
2. if(p2==NULL) { p2=new("a-node", NULL, NULL)
hashT.insert(p2) }

hash	content
1	
2	p1

67

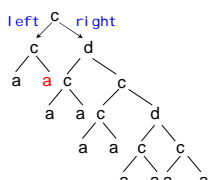
Minimal Unique DAGs

Example "Parse & DAGify"

- startElement(c)
- startElement(c)
- startElement(a)
- endElement(a)
- startElement(a)
- endElement(a)

- $p_2 = \text{hashT.find}(a) = p_1$
- $\text{if}(p_2 == \text{NULL}) \{ p_2 = \text{new}(\text{"a-node"}, \text{NULL}, \text{NULL}); \text{hashT.insert}(p_2) \}$

→ store children lists: [[], [p1, p1]]



hash	content
1	
2	p1

68

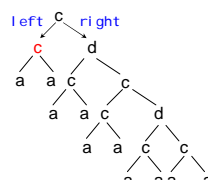
Minimal Unique DAGs

Example "Parse & DAGify"

- startElement(c)
- startElement(c)
- startElement(a)
- endElement(a)
- startElement(a)
- endElement(a)
- endElement(c)

- $p_2 = \text{hashT.find}(a) = p_1$
- $\text{if}(p_2 == \text{NULL}) \{ p_2 = \text{new}(\text{"a-node"}, \text{NULL}, \text{NULL}); \text{hashT.insert}(p_2) \}$

→ store children lists: [[], [p1, p1]]



hash	content
1	
2	p1

69

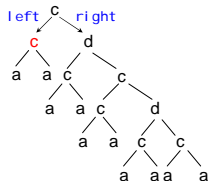
Minimal Unique DAGs

Example "Parse & DAGify"

- startElement(c)
- startElement(c)
- startElement(a)
- endElement(a)
- startElement(a)
- endElement(a)
- endElement(c)

- $p = \text{hashT.find}(a)$
- $\text{if}(p == \text{NULL}) \{ p = \text{new}(\text{"c-node"}, p_1, p_1); \text{hashT.insert}(p) \}$

→ store children lists: [[], [p1, p1]] → use children list!!



hash	content
1	
2	p1
...	
20201	p

Memory location **p** is a DAG with root node **c**, and child1-pointer=p1 child2-pointer=p1

70

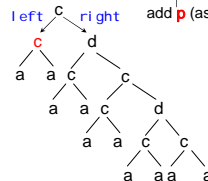
Minimal Unique DAGs

Example "Parse & DAGify"

- startElement(c)
- startElement(c)
- startElement(a)
- endElement(a)
- startElement(a)
- endElement(a)
- endElement(c)

- $p = \text{hashT.find}(a)$
- $\text{if}(p == \text{NULL}) \{ p = \text{new}(\text{"c-node"}, p_1, p_1); \text{hashT.insert}(p) \}$

→ store children lists: [[], [~~p1~~, p1]] ← now update!!



hash	content
1	
2	p1
...	
20201	p

Memory location **p** is a DAG with root node **c**, and child1-pointer=p1 child2-pointer=p1

71

Minimal Unique DAGs

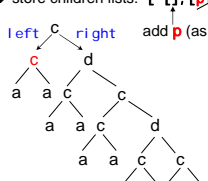
Example "Parse & DAGify"

New children lists: [[p]]

children of root node (so far, one)

- $p = \text{hashT.find}(a)$
- $\text{if}(p == \text{NULL}) \{ p = \text{new}(\text{"c-node"}, p_1, p_1); \text{hashT.insert}(p) \}$

→ store children lists: [[], [~~p1~~, p1]] ← now update!!



hash	content
1	
2	p1
...	
20201	p

Memory location **p** is a DAG with root node **c**, and child1-pointer=p1 child2-pointer=p1

72

Minimal Unique DAGs

Example "Parse & DAGify"

New children lists: [[p]]

children of root node (so far, one)

- $p = \text{hashT.find}(a)$
- $\text{if}(p == \text{NULL}) \{ p = \text{new}(\text{"c-node"}, p_1, p_1); \text{hashT.insert}(p) \}$

hash	content
1	
2	p1
...	
20201	p

→ Assume · 100 element names

Example hash function:

```
(#elementName
+ 100 * #elementName(1st child)
+ 100 * 100 * #elementName(2nd child)
+ 100^3 * #elementName(1st child of 1st ch.)
+ ... ) MOD sizeOf(hashT)
```

73

Minimal Unique DAGs

→ DOM interfact to the DAG?

parentNode / p&nsi bl i ng as before

→ Updates can be expensive (copying!)

How to attach attribute & text nodes to the DAG?

74

Minimal Unique DAGs

→ DOM interfact to the DAG?

parentNode / p&nsi bl i ng as before

→ Updates can be expensive (copying!)

How to attach attribute & text nodes to the DAG?

→Store them seperately in a table.

Index by e.g., Node number (in doc-order)
or number of atr/text nodes

Store index in each DAG node / or compute it online. (pre-traversal)

75

Minimal Unique DAGs

What about *unranked*, vs binary DAGs?

firstChild lastChild

More precisely,

What about size of minimal-unique-unranked-DAG(Tree)
vs size of minimal-unique-binary-DAG(fCnS-enc(Tree))

firstChild/nextSibling

76

Minimal Unique DAGs

size of minimal-unique-unranked-DAG(Tree)
vs size of minimal-unique-binary-DAG(fCnS-enc(Tree))

Questions

Give a tree for which first is smaller than the second.

Give a tree for which the second is smaller than the first.

77

Unranked vs Binary Trees

18/19 3/4

Min. DAG

13/9

Min. DAG

78

Unranked vs Binary Trees

18/19 3/4

Min. DAG

Can it be vica versa? (min bin. DAG is smaller)

YES!!

→ Has **18 edges**

→ DAG of bin.encoding only **12 edges**

DAG compression is **sensible** to rank/unrankedness!

79

YES: the binary trees become very “regular” (deep, to the right)”



80

1

81

The Figure on the next slide is from the “EXI Measurement Note”
-- new version of the note came out 25 July 2007...!)

82

The diagram illustrates the E3I architecture, showing the flow from document preprocessing to results. It is divided into two main horizontal paths: the top path for XML and the bottom path for XSL documents.

Top Path (XML):

- Text corpus preprocessing** leads to **XML Parsing**.
- XML Parsing** leads to **Deserialization**.
- Deserialization** leads to **Encoding**.
- Encoding** leads to **Event store**.
- Event store** leads to **Contributor encoders**.
- Contributor encoders** leads to **Results**.

Bottom Path (XSL):

- Rapid document preprocessing** leads to **Data binding**.
- Data binding** leads to **Serialization**.
- Serialization** leads to **Decoding**.
- Decoding** leads to **Event store**.
- Event store** leads to **Contributor encoders**.
- Contributor encoders** leads to **Results**.

Intermediate Components and Data Flow:

- XML Path Intermediate Steps:**
 - Build test corpus, XML content, use-case coverage, etc.** (cloud shape) feeds into **XML, XSL content, Text Encoding** (green box).
 - XML, XSL content, Text Encoding** feeds into **Parse XML, API** (green box).
 - Parse XML, API** feeds into **SAX or typed event interfaces** (green box).
 - SAX or typed event interfaces** feeds into **In-memory data structures, full or block** (blue box).
 - In-memory data structures, full or block** feeds into **Compression Algorithms** (blue box).
 - Compression Algorithms** (Grammar-based, Statistical, Other options?) feeds into **E3I encoding of XML** (cloud shape).
 - E3I encoding of XML** leads to **binary blob**.
- XSL Path Intermediate Steps:**
 - Comparison** (Confirms proper hand-off, conversion of original document (B6)) feeds into **Parse checks and format, valid outputs** (green box).
 - Parse checks and format, valid outputs** feeds into **XML, Document, Text Encoding** (green box).
 - XML, Document, Text Encoding** feeds into **E3I XML, API** (green box).
 - E3I XML, API** feeds into **SAX or typed event interfaces** (green box).
 - SAX or typed event interfaces** feeds into **In-memory data structures, full or block** (blue box).
 - In-memory data structures, full or block** feeds into **Decompression Algorithms** (blue box).
 - Decompression Algorithms** (Lookup tables, Computational, Other options?) feeds into **E3I encoding of XML** (cloud shape).
 - E3I encoding of XML** leads to **binary blob**.

Annotations:

- Between **XML, XSL content, Text Encoding** and **Parse XML, API**: *one or more namespaces*
- Between **Parse XML, API** and **SAX or typed event interfaces**: *also handles streaming from, to disk or network*
- Between **SAX or typed event interfaces** and **In-memory data structures, full or block**: *Measurement endpoint, in CPU*
- Between **In-memory data structures, full or block** and **Compression Algorithms**: *asymmetric performance*

83

Instead, we need a *more flexible parser* that gives us the freedom of what exactly to store, and how.

84

You never need to write a parser again!

How to use SAX

Remember one of the promises of XML...

You never need to write a parser again!

... but, of course if you want to build up your own (e.g. memory-efficient) data structure, you need to "talk" to the parser.

You want to tell the parser:

Give me low level access to the data:

- Bracket by bracket,
 - text-node by text-node.
- In "document order".

How to use SAX

Remember one of the promises of XML...

You never need to write a parser again!

... but, of course if you want to build up your own (e.g. memory-efficient) data structure, you need to "talk" to the parser.

You want to tell the parser:

Give me low level access to the data:

- Bracket by bracket,
 - text-node by text-node.
- In "document order".

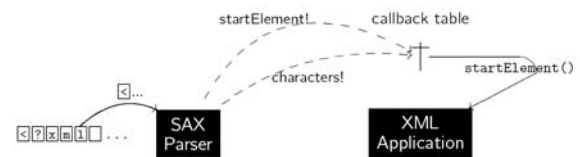
→ SAX

SAX—Simple API for XML

- **SAX⁷ (Simple API for XML)** is, unlike DOM, *not* a W3C standard, but has been developed jointly by members of the XML-DEV mailing list (ca. 1998).
- SAX processors use **constant space**, regardless of the XML input document size.
 - ▶ Communication between the SAX processor and the backend XML application does *not* involve an intermediate tree data structure.
 - ▶ Instead, the **SAX parser sends events** to the application whenever a certain piece of XML text has been recognized (*i.e.*, parsed).
 - ▶ The **backend acts on/ignores events** by populating a **callback function table**.

⁷<http://www.saxproject.org/>

Sketch of SAX's mode of operations



- A SAX processor reads its input document **sequentially** and **once** only.
- No memory of what the parser has seen so far is retained while parsing. As soon as a **significant bit of XML text** has been recognized, an **event** is sent.
- The application is able to act on events **in parallel** with the parsing progress.

SAX Events

- To meet the constant memory space requirement, SAX reports **fine-grained parsing events** for a document:

Event	... reported when seen	Parameters sent
<i>startDocument</i>	<code><?xml...?></code> ⁸	
<i>endDocument</i>	(EOF)	
<i>startElement</i>	<code><t a₁=v₁ ... a_n=v_n></code>	<i>t</i> , (a ₁ , v ₁), ..., (a _n , v _n)
<i>endElement</i>	<code></t></code>	<i>t</i>
<i>characters</i>	text content	Unicode buffer ptr, length
<i>comment</i>	<code><!--c--></code>	<i>c</i>
<i>processingInstruction</i>	<code><?t pi?></code>	<i>t</i> , <i>pi</i>

⁸N.B.: Event *startDocument* is sent even if the optional XML text declaration should be missing.

SAX Events

```

1 <?xml encoding="utf-8"?> *1
2 <bubbles> *2
3 <!-- Dilbert looks stunned --> *3
4 <bubble speaker="phb" to="dilbert"> *4
5   Tell the truth, but do it in your usual engineering way
6   so that no one understands you. *5
7 </bubble> *6
8 </bubbles> *7 *8
    
```

Event ⁹ ¹⁰	Parameters sent
*1 <i>startDocument</i>	
*2 <i>startElement</i>	<i>t</i> = "bubbles"
*3 <i>comment</i>	<i>c</i> = "Dilbert looks stunned."
*4 <i>startElement</i>	<i>t</i> = "bubble", ("speaker", "phb"), ("to", "dilbert")
*5 <i>characters</i>	<i>buf</i> = "Tell the...understands you.", <i>len</i> = 99
*6 <i>endElement</i>	<i>t</i> = "bubble"
*7 <i>endElement</i>	<i>t</i> = "bubbles"
*8 <i>endDocument</i>	

⁹Events are reported in **document reading order** *1, *2, ..., *8.

¹⁰N.B.: Some events suppressed (white space).

SAX Callbacks

- To provide an efficient and tight **coupling** between the SAX **frontend** and the application **backend**, the SAX API employs **function callbacks**.¹¹
 - Before parsing starts, the application **registers function references** in a table in which each event has its own slot:

Event	Callback
startElement	?
endElement	?

→ SAX.register(startElement, startElement ());
SAX.register(endElement, endElement ());
 - The application alone decides on the implementation of the functions it registers with the SAX parser.
 - Reporting an event** e_i then amounts to call the function (with parameters) registered in the appropriate table slot.

¹¹Much like in event-based GUI libraries.

Marc H. Schohl (DBIS, Uni KN) XML and Databases Winter 2005/06 87

SAX Callbacks

Java SAX API

In Java, populating the callback table is done via implementation of the SAX `ContentHandler` interface: a `ContentHandler` object represents the callback table, its methods (e.g., `public void endDocument ()`) represent the table slots.

Example: Reimplement `content.cc` shown earlier for DOM (find all XML text nodes and print their content) using SAX (pseudo code):

```

content (File f)
// register the callback,
// we ignore all other events
SAXregister (characters, printText);
SAXparse (f);
return;

printText ((Unicode) buf, Int len)
Int i;
foreach i ∈ 1...len do
    print (buf[i]);
return;

```

Marc H. Schohl (DBIS, Uni KN) XML and Databases Winter 2005/06 88

SAX and XML Trees

SAX and the XML Tree Structure

- Looking closer, the **order** of SAX events reported for a document is determined by a **preorder traversal** of its document tree¹²:

Sample XML document

```

1 <?xml version="1.0" />
2 <foo>
3   <bar> foo+4 </bar>+5
4   <!--sample-->+6
5   <baz>
6     <baz> baz+9 </baz>+10
7     <baz>+11 baz+12 </baz>+13
8   </baz>+14
9 </foo>+15 +16

```

N.B.: An `Elem [Doc]` node is associated with two SAX events, namely `startElement` and `endElement` [`startDocument`, `endDocument`].

¹²Sequences of sibling `Char` nodes have been collapsed into a single `Text` node.

Marc H. Schohl (DBIS, Uni KN) XML and Databases Winter 2005/06 89

94

Deadline: 6th April

For **Assignment 2**, you only need to register `startElement` and `endElement`.

In that way, you automatically receive only element nodes..

Of course you can use SAX for other things than building up a data structure.

E.g.

- answer path queries while parsing (on a "stream") (low memory consumption!)

95

END

Lecture 2