# XML and Databases

**Lecture 8**
*Streaming Evaluation: how much memory do you need*?

Sebastian Maneth
NICTA and UNSW

*CSE@UNSW  --  Semester 1, 2009*

---

Small XPath Quiz

Can you give an expression that returns the  last / first  occurrence
of each distinct price element?

```
<b>
<price>3</price>
<price>1</price>
<price>3</price>
<price>1</price>
<price>3</price>
<price>4</price>
<price>1</price>
<price>7</price>
</b>
```

Should return

```
<price>3</price>
<price>4</price>
<price>1</price>
<price>7</price>
```

Should return

```
<price>3</price>
<price>1</price>
<price>4</price>
<price>7</price>
```

2

---

Small XPath Quiz

Can you give an expression that returns the  last / first  occurrence
of each distinct price element?

```
<b>
<price>3.0</price>
<price>1</price>
<price>3.00</price>
<price>1</price>
<price>3</price>
<price>4</price>
<price>1.000</price>
<price>7</price>
</b>
```

Should return

```
<price>3</price>
<price>4</price>
<price>1.000</price>
<price>7</price>
```

Should return

```
<price>3.0</price>
<price>1</price>
<price>4</price>
<price>7</price>
```

What if we mean  *number-distinctness*  (not strings)?

3

---

0.  Recall

→   Evaluation of  Simple Paths  //a/b/c
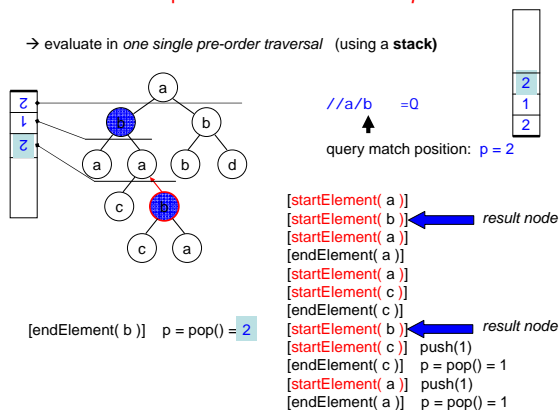→   Arbitrary Queries over  //, /, *

## Outline

1.  Automaton Approach

2.  Parallel Evaluation of Multiple Queries

3.  Sizes of Automata

4.  How to deal with Filters

5.  Existing Systems for Streaming XPath Evaluation
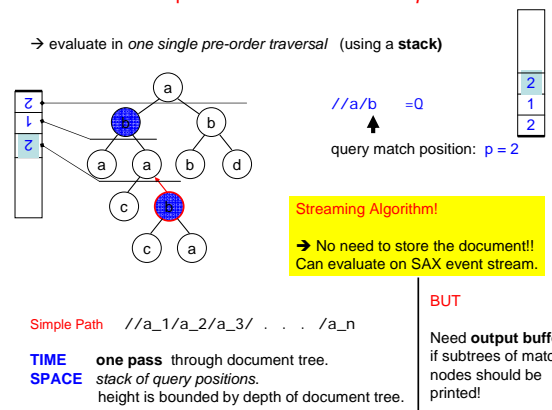
4

---

## Recall:  Top-Down Evaluation of *Simple Paths*

→ evaluate in *one single pre-order traversal*   (using a **stack)**

//a/b   =Q

query match position: p = 2

```
[startElement( a )]
[startElement( b )]        ← result node
[startElement( a )]
[endElement( a )]
[startElement( a )]
[startElement( c )]
[endElement( c )]
[startElement( b )]        ← result node
[startElement( c )]   push(1)
[endElement( c )]     p = pop() = 1
[startElement( a )]   push(1)
[endElement( a )]     p = pop() = 1
```

[endElement( b )]   p = pop() = 2

5

---

## Recall:  Top-Down Evaluation of *Simple Paths*

→ evaluate in *one single pre-order traversal*   (using a **stack)**

//a/b   =Q

query match position: p = 2

Streaming Algorithm!

➔ No need to store the document!!
Can evaluate on SAX event stream.

BUT

Need **output buffers**,
if subtrees of match
nodes should be
printed!

Simple Path   //a_1/a_2/a_3/ . . . /a_n

**TIME**     **one pass**  through document tree.
**SPACE**   *stack of query positions.*
             height is bounded by depth of document tree.

6

---

1

**Slide 7**

Recall: Top-Down Evaluation of *Simple Paths*

→ evaluate in *one single pre-order traversal* (using a **stack**)

If we print **node-IDs**, then no output buffers are needed!

//a/b  =Q

query match position: p = 2

→ True Streaming, with memory need proportional to height.

Streaming Algorithm!

→ No need to store the document!! Can evaluate on SAX event stream.

Simple Path  //a_1/a_2/a_3/ . . . /a_n

**TIME**  **one pass** through document tree.
**SPACE**  *stack of query positions.* height is bounded by depth of document tree.

BUT

Need **output buffers**, if subtrees of match nodes should be printed!

7

**Slide 8**

Recall: Top-Down Evaluation of *Simple Paths*

→ evaluate in *one single pre-order traversal* (using a **stack**)

If we print **node-IDs**, then no output buffers are needed!

→ any good implementation of this algorithm should work for documents with *depth up to a couple of millions*, and **NO restriction on document size**!

//a/b  =Q

query match position: p = 2

Streaming Algorithm!

→ No need to store the document!! Can evaluate on SAX event stream.

Simple Path  //a_1/a_2/a_3/ . . . /a_n

**TIME**  **one pass** through document tree.
**SPACE**  *stack of query positions.* height is bounded by depth of document tree.

1 Byte is enough for small queries!
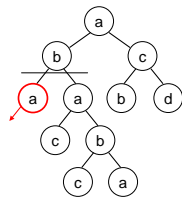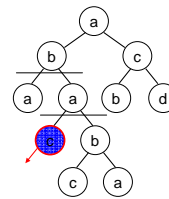
8

**Slide 9**

Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with /, //, *

multiple //'s

//a/b**//**c

query match position: p = 3

no match stay in p=3!

…
[startElement( a )]  push(3)
[endElement( a )]  p = pop() = 3

9

**Slide 10**

Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with /, //, *

multiple //'s

//a/b**//**c

query match position: p = 3

no match stay in p=3!

…
[startElement( a )]  push(3)
[endElement( a )]  p = pop() = 3
[startElement( a )]  push(3)
[startElement( c )]  push(3)

Result node!
Mark it, and stay in p=3.

10
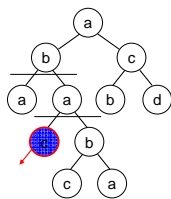
**Slide 11**

Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with /, //, *

multiple //'s

//a/b**//**c

query match position: p = 3

no match stay in p=3!

…
[startElement( a )]  push(3)
[endElement( a )]  p = pop() = 3
[startElement( a )]  push(3)
[startElement( c )]  push(3)

Result node!
Mark it, and stay in p=3.

*Output Node-ID*  *Start copying to Output Buffer*

11

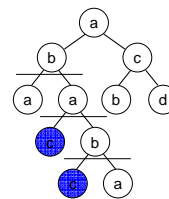**Slide 12**

Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with /, //, *

multiple //'s

//a/b**//**c

query match position: p = 3

no match stay in p=3!

…
[startElement( a )]  push(3)
[endElement( a )]  p = pop() = 3
[startElement( a )]  push(3)
[startElement( c )]  push(3)
[endElement( c )]  p = pop() = 3
[startElement( b )]  push(3)
[startElement( c )]  push(3)
…

12

2

## Slide 13

### Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with `/, //, *`

multiple `//`'s

`//a/b//c`

query match position: p = 3

| 3 |
| 3 |
| 3 |
| 2 |

Stay at position 3,
for the *complete subtree*!

Never go back to pos. 1 or pos. 2!

13

## Slide 14

### Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with `/, //, *`

multiple `//`'s

`//a/b//c`

query match position: p = 3

| 3 |
| 3 |
| 3 |
| 2 |

Optimizations (for Output Buffers)

(1) If *inside a matched subtree*, record position (or range within buffer), instead of creating a new output buffer.

(2) If *subtree is finished* (we are not inside a match), then we can write its buffer out and can start with empty buffer again.
[ Worst Case:
  root node selected. size of doc. Needed. ]

14

## Slide 15

### Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with `/, //, *`

multiple `//`'s

`//a/b//c`

query match position: p = 3

| 3 |
| 3 |
| 3 |
| 2 |

`//a/b//c/d/*/e//f/g//h`

➔ **Same as before**

jump back within `/`-sequence.
AT MOST to the beginning of the last `//`.

Use KMP within `/`-sequence.

For `*`'s: build several KMP-tables.

15

## Slide 16

### Arbitrary *Slash+Slashslash*

→ evaluate in *one single pre-order traversal* (using a **stack**)

Arbitrary queries with `/, //, *`

multiple `//`'s

`//a/b//c`

query match position: p = 3

| 3 |
| 3 |
| 3 |
| 2 |

`//a/b//c/d/*/e//f/g//h`

*Query Problem* is solved!

Leave optimizations of

`>cat file.xml` [1.2.7, 1.3, 1.3.1.1, …]

To OS/UNIX hackers.. ☺

If Node-IDs are printed, then no output buffers are needed.

Then:
Memory proportional to height.
Should run for arbitrary large docs!

16

## Slide 17

### 1. Automaton Approach

`//a/b//c/d/*/e//f/g//h`

➔ **Same as before**

jump back within `/`-sequence.
AT MOST to the beginning of the last `//`.

Use KMP within `/`-sequence.

For `*`'s: build several KMP-tables.

**Recall**

Deterministic Automaton runs in

→ *linear time* and
→ *constant space*

(plus *stack of states*, if we run on paths of a tree)

17

## Slide 18

### 1. Automaton Approach

¬a,¬b
*not* a
and *not* b

`//a/b//c/d/*/e//f/g//h`

➔ **Same as before**

jump back within `/`-sequence.
AT MOST to the beginning of the last `//`.

Use KMP within `/`-sequence.

For `*`'s: build several KMP-tables.

**Recall**

Deterministic Automaton runs in

→ *linear time* and
→ *constant space*

(plus *stack of states*, if we run on paths of a tree)

18

3

# 1. Automaton Approach

a
b
c
d
*
e
f
g
h

a
b
c
d
¬a,¬b
¬c,¬d

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

➔ **Same as before**

jump back within ∕-sequence.
AT MOST to the beginning of the last **//**.

Use KMP within ∕-sequence.

For *'s:  build several KMP-tables.

**Recall**

Deterministic Automaton  runs in

➔ *linear time*  and
➔ *constant space*

(plus *stack of states*, if we run
on paths of a tree)

19

---

# 1. Automaton Approach

a
b
c
d
X
*
e
f
g
h

a
b
c
d
¬a,¬b
¬c,¬d

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

**Problem**
If it is NOT an e here, then what to do??

E.g.,
a b c d c d

We should be in state **X**!

20

---

# 1. Automaton Approach

a
b
c
d
X *=c
e
f
g
h

a
b
c
d
¬a,¬b
¬c,¬d

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

**Problem**
If it is NOT an e here, then what to do??

➔ Need to know what the * was!!

E.g.,
a b c d c d

We should be in state **X**!

21

---

# 1. Automaton Approach

a
b
c
d
X *=c
e
f
g
h

a
b
c
d
¬a,¬b
¬c,¬d
c
¬c,¬d,¬e

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

*=**?**  Which other letters need to be considered?

c d x y
≠c
≠e

22

---

a
b
c
d
X *=c
e
f
g
h
¬c
e
d
c
¬c,¬d
¬c,¬d,¬e

a
b
¬a,¬b

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

*=**?**  Which other letters need to be considered?

c d x y
≠c
≠e

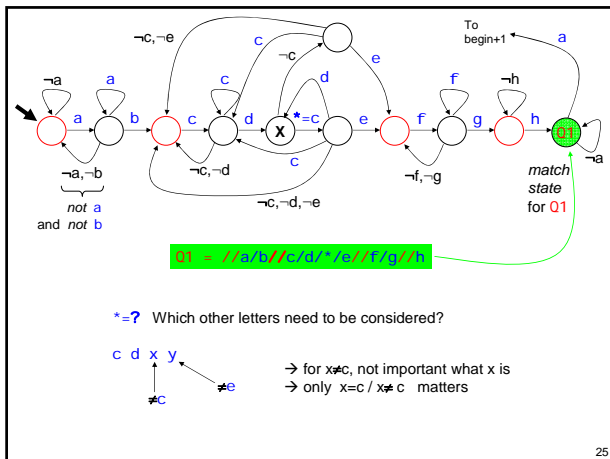➔ for x≠c, not important what x is
➔ only  x=c / x≠ c   matters

23

---

"splitting" – can be at most
#different symbols many

a
b
c
d
X *=c
e
f
g
h
¬c
e
d
c
¬c,¬d
¬c,¬d,¬e

a
b
¬a,¬b

*not* a
and *not* b

//a/b**//**c/d/*/e//f/g//h

*=**?**  Which other letters need to be considered?

c d x y
≠c
≠e

➔ for x≠c, not important what x is
➔ only  x=c / x≠ c   matters

24

4

**Slide 25**

¬c,¬e
¬a  a  c  ¬c  c  d  e  f  ¬h
To begin+1  a
a  b  c  X *=c  e  f  g  h  Q1
¬a,¬b  ¬c,¬d  ¬f,¬g  ¬a
*not* a and *not* b
¬c,¬d,¬e
*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

*=**?** Which other letters need to be considered?

c  d  x  y
≠c        ≠e
→ for x≠c, not important what x is
→ only  x=c / x≠ c   matters

25

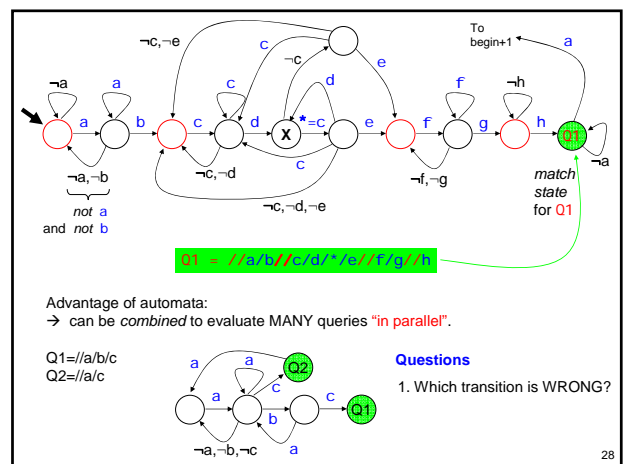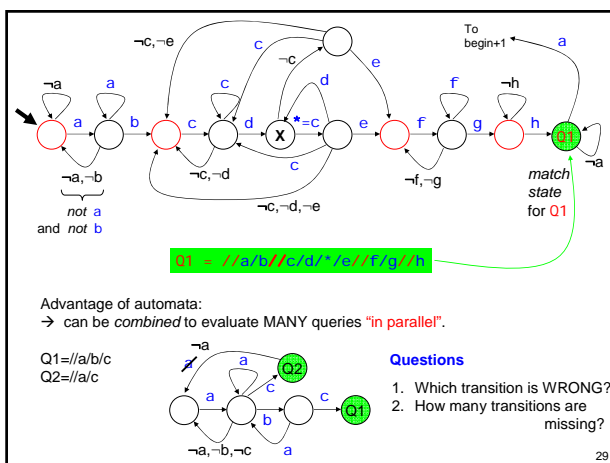**Slide 26**

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

26

**Slide 27**

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

a  a  Q2  c  c  Q1
a  b
¬a,¬b,¬c  a

27

**Slide 28**

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

**Questions**
1. Which transition is WRONG?

¬a,¬b,¬c

28

**Slide 29**

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

¬a,¬b,¬c

29

**Slide 30**

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

¬a

¬a,¬b,¬c

30

5

**Slide 31**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  a  Q2  a  a  c  c  Q1  a  b  ¬a,¬b,¬c  a  ¬a,¬c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

31

---

**Slide 32**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  a  Q2  a  a  c  c  Q1  a  b  ¬a,¬b,¬c  a  ¬a,¬c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

32

---

**Slide 33**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  ¬a  a  Q2  a  a  c  c  Q1  a  b  ¬a,¬b,¬c  a  ¬a,¬c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

33

---

**Slide 34**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  ¬a  a  Q2  a  a  c  c  Q1  ¬a  b  ¬a,¬b,¬c  a  ¬a,¬c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?

34

---

**Slide 35**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  ¬a  a  Q2  a  a  c  c  Q1  ¬a  b  ¬a,¬b,¬c  a  ¬a,¬c

**Questions**
1. Which transition is WRONG?
2. How many transitions are missing?
→ 5

35

---

**Slide 36**

To begin+1

¬c,¬e  c  ¬c  c  d  e  f  ¬h  a

¬a  a

a b c d X *=c e f g h Q1  ¬a

¬a,¬b

*not* a and *not* b

¬c,¬d  ¬c,¬d,¬e  c  ¬f,¬g

*match state* for Q1

Q1 = //a/b//c/d/*/e//f/g//h

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

¬a  ¬a  a  Q2  a  a  c  c  Q1  ¬a  b  ¬a,¬b,¬c  a  ¬a,¬c

ONE look-up per node!

Combined automaton:
SIZE ≤ SIZE(A1) x SIZE(A2)

36

**Question**

What is SIZE(A1) wrt size of Q1?

Take
(1) SIZE(A) = #states
(2) SIZE(A) = #transitions

Advantage of automata:
→ can be *combined* to evaluate MANY queries "in parallel".

Q1=//a/b/c
Q2=//a/c

ONE look-up per node!

Combined automaton:
SIZE ≤ SIZE(A1) x SIZE(A2)

*not a* and *not b*

*match state* for Q1

37

---

## 3. The Size of the DFA

$$//a/*/*/*/b$$

Size of DFA = exponential in *'s

(not a real concern)



**NFA**

**DFA** (fragment, and without back edges)

---

## 3. The Size of the DFA

**Theorem** [GMOS'02] The number of states in the DFA for one linear XPath expression P is at most:

$$k + |P| \ k \ s^m$$

k = number of //
s = size of the alphabet (number of tags)
m = max number of * between two consecutive //

---

How to deal with filters?

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



When we meet the c-nodes (in pre order traversal) we do not know yet if the **filter** will evaluate to true.

40

---

How to deal with filters?

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



When we meet the c-nodes (in pre order traversal) we do not know yet if the **filter** will evaluate to true.

→ We have to use *buffers,* as before.

However, now buffers may be **deleted** without being used.

**Question**

If we output node ID's, then how much memory is needed in the worst case for queries with filters?

Must be stored in memory

41

---

How to deal with filters?

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



→ Size of largest documents that can be streamed in this way depends on - #filters,
   - sizes of (pre) selected trees,
   - quality of (1), (2), etc..

Optimizations

(1) Store potential match trees as DAGs
(2) Release potential match trees as early as possible!

Must be stored in memory

42

---

## Slide 43

How to deal with filters?

//a[./d/e]/b//c



➔ Size of largest documents that can be streamed in this way depends on -  #filters,
- sizes of (pre) selected trees,
- quality of (1), (2), etc..

➔ Release potential match trees *as early as possible*!

Find earliest point at which we know the filter is true.

Must be stored in memory

43

## Slide 44

How to deal with filters?

//a[./d/e]/b//c



➔ Size of largest documents that can be streamed in this way depends on -  #filters,
- sizes of (pre) selected trees,
- quality of (1), (2), etc..

➔ Release potential match trees *as early as possible*!

Find earliest point at which we know the filter is true.

No need to store. Stream! ☺

44

## Slide 45

How to deal with filters?

//a[./d/e]/b//c



➔ Size of largest documents that can be streamed in this way depends on -  #filters,
- sizes of (pre) selected trees,
- quality of (1), (2), etc..

Find earliest point at which we know the filter is true.

Harder for *Boolean combinations*:

[not(./d/e) and (./c/d or //b/c)]

**Question**  where is the earliest point for this filter?

No need to store. Stream! ☺

45

## Slide 46

How to deal with filters?

//a[./d/e]/b//c



➔ Size of largest documents that can be streamed in this way depends on -  #filters,
- sizes of (pre) selected trees,
- quality of (1), (2), etc..

Find earliest point at which we know the filter is true.

Harder for *Boolean combinations*:

[not(./d/e) and (./c/d or //b/c)]

**Question**  where is the earliest point for this filter?
➔ and now?

No need to store. Stream! ☺

46

## Slide 47

How to deal with filters?

//a[./d/e]/b//c



➔ Size of largest documents that can be streamed in this way depends on -  #filters,
- sizes of (pre) selected trees,
- quality of (1), (2), etc..

We can also construct automata for filter expressions!

Use a *push-down* for potential candidates.

Push-Down Automaton
can probably be designed so that it pops/outputs candidates *as early as possible*.

47

## Slide 48

How to deal with filters?

//a[./d/e]/b//c

Another Idea

Use **2-pass algorithm**: first (bottom-up) phase to mark subtrees with filter information.
Second (top-down) phase to determine match nodes.

Why is this interesting?

→ Fast main memory evaluation
→ Use disk as intermediate store  (stream twice)

48

## 5. Streaming XPath Algorithms

- XFilter and YFilter   [Altinel and Franklin 00] [Diao et al 02]
- X-scan  [Ives, Levy, and Weld 00]
- XMLTK  [Avila-Campillo et al 02]
- XTrie  [Chan et al 02]
- SPEX  [Olteanu, Kiesling, and Bry 03]
- Lazy DFAs  [Green et al 03]
- The XPush Machine  [Gupta and Suciu 03]
- XSQ  [Peng and Chawathe 03]
- TurboXPath [Josifovski, Fontoura, and Barta 04]
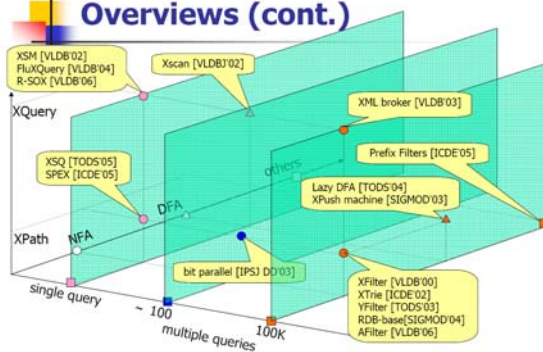- …

49

## 5. Streaming XPath Algorithms

Some following slides are by T. Amagasa and M Onizuka (Japan)
See http://www.dasfaa07.ait.ac.th/DASFAA2007_tutorial3_1.pdf

Most of the following slides are by Dan Suciu (the above slides are
Actually also based on Suciu's slides ☺ )
See
http://www.cs.washington.edu/homes/suciu/talk-spire2002.ppt

50

## XFilter (cont.)
## NFA, view class: //tag

Decomposing XPath Query

/a/b//c



Query ID, Position, Relative Position, Level

| | a | b | c |
|---|---|---|---|
| Query ID | Q1 | Q1 | Q1 |
| Position | 1 | 2 | 3 |
| Relative Position | 0 | 1 | -1 |
| Level | 1 | 0 | -1 |

Every Path Node represents a state in the FSM

---

## XFilter (cont.)
## NFA, view class: //tag

node-test hash table

Q1= /a/b//c
Q1-1  Q1-2  Q1-3

Q2= //b/*/c/d
Q2-1  Q2-2  Q2-3

Q3= /*/a/c//d
Q3-1  Q3-2  Q3-3



---

## YFilter
## NFA, view class: XP{/,//,*}

- prefix sharing
- Predicates are processed by labels

Q1=/a/b
Q2=/a/c
Q3=/a/b/c
Q4=/a/b/c
Q5=/a/*/c
Q6=/a//c
Q7=/a/*/*/c
Q8=/a/b/c

{Q1}
{Q3, Q8}
{Q2}
{Q4}
{Q6}
{Q5}
{Q7}

(a) XPath queries    (b) A corresponding NFA (YFilter)

---

## Shared data structure

- Sharing identical structures among query trees
  What to share? node-test, simple path, branch, etc.

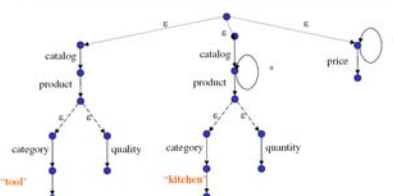| What to share? | View class | Algorithms |
|---|---|---|
| node-test | //tag | XFilter [VLDB'00] |
| simple sub-path | //tag1/.../tagN | XTrie [ICDE'02] |
| simple path | XP{/,//,*} | YFilter [TODS'03], Lazy DFA [TODS'04], Prefix Filters [ICDE'05], AFilter [VLDB'06] |
| branch | XP{[],/,//,*} | XPush machine [SIGMOD'03] |
| ... | ... | ... |

---

## XPath Processing with FA
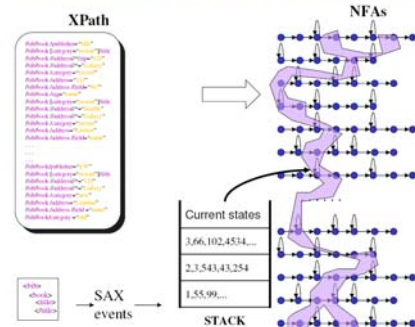## -- From XPath (XP{[],/,//,*) to NFA --

/catalog/product[category="tools"]/quantity
/catalog//product[category="kitchen"]/quality
//price



---

## NFA-based XPE Processing

XPath                                NFAs

Current states

3,66,102,4534,...

2,3,543,43,254

1,55,99,...

SAX events

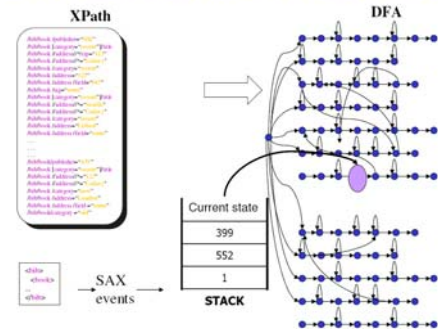STACK

## Basic NFA Evaluation

Properties:
  ☺ Space = linear
  ☹ Throughput = decreases linearly

Systems:
- XFilter [Altinel&Franklin'99], YFilter.
- XTrie [Chan et al.'02]

---



Duality  -> XML databases -> XML streams

## DFA-based XPE Processing

---

## Basic DFA Evaluation

Properties:
  ☺ Throughput = constant !
  ☹ Space = GOOD QUESTION

System:
- XML Toolkit [University of Washington]
  **http://xmltk.sourceforge.net**

---

## The Size of the DFA

**Theorem** [GMOS'02] The number of states in the DFA for one linear XPath expression P is at most:
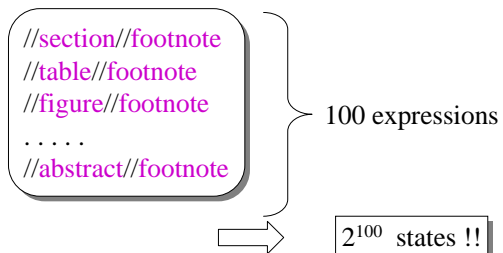
$$k+|P|\ k\ s^m$$

k = number of  //
s = size of the alphabet (number of tags)
m = max number of * between two consecutive //

---

## Size of DFA:
## Multiple Expressions

//section//footnote
//table//footnote
//figure//footnote
. . . . .
//abstract//footnote

100 expressions

$2^{100}$  states !!

There is a theorem here too, but it's not useful…

---

## Solution:
## Compute the DFA Lazily

- Also used in text searching
- But will it work for $10^6$ XPath expressions ?
- YES !
- For XPath it is *provably* effective, for two reasons:
  - XML data is not very deep
  - The nesting structure in XML data tends to be predictable

## Lazy DFA
## DFA, view class: XP{/,//,*}

### Features
- Sharing the process of / and //, * and tag
- DFA-based
- Compute DFA lazily (on demand)
- # of DFA states
  - Independent from # of XPath exprs.
  - Depends on DataGuide size (schema)

### Issue
- Predicates: XPush machine [SIGMOD'03]

---

## Lazy DFA and "Simple" DTDs

- Document Type Definition (DTD)
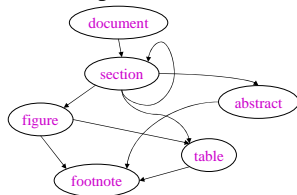  - Part of the XML standard
  - Will be replaced by XML Schema
- Example DTD:

```
<!ELEMENT document (section*)>
<!ELEMENT section ((section|abstract|table|figure)*)>
<!ELEMENT figure (table?,footnote*)>
. . . . .
```

**Definition** A DTD is simple if all cycles are loops

---

## Lazy DFA and "Simple" DTDs

Simple DTD:



XPath expressions

```
//section//footnote
//table//footnote
//figure//footnote
//abstract//footnote
```

⟹ Eager DFA "remembers" $2^4$ sets
Lazy DFA "remembers" only 4 sets

---

## Lazy DFA and "Simple" DTDs

**Theorem** [GMOS'02] If the XML data has a "simple" DTD, then lazy DFA has at most:

$$1+D(1+n)^d$$

states.

n = max depths of XPath expressions
D = size of the "unfolded" DTD
d = max depths of self-loops in the DTD

**Fact of life**: "Data-like" XML has simple DTDs

---

## Lazy DFA and Data Guides

- "Non-simple" DTDs are useless for the lazy DFA
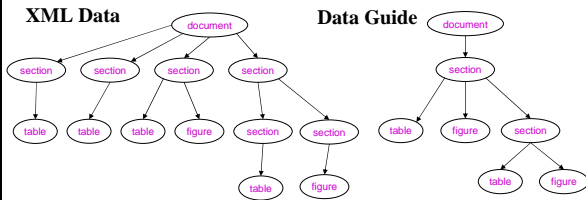- "Everything may contain everything"

```
<!ELEMENT document (section*)>
<!ELEMENT section ((section|table|figure|abstract|footnote)*)>
<!ELEMENT table   ((section|table|figure|abstract|footnote)*)>
<!ELEMENT figure  ((section|table|figure|abstract|footnote)*)>
<!ELEMENT abstract ((section|table|figure|abstract|footnote)*)>
```

Fact of life: "Text"-like XML has non-simple DTDs

---

## Lazy DFA and Data Guides

**Definition** [Goldman&Widom'97]
The data guide for an XML data instance is the Trie of all its root-to-leaf paths

---

## Lazy DFA and Data Guides

**XML Data**

document
- section
- section
- section → table, table, table, figure
- section → section → table, section → figure

**Data Guide**

document
- section → table, figure, section → table, figure

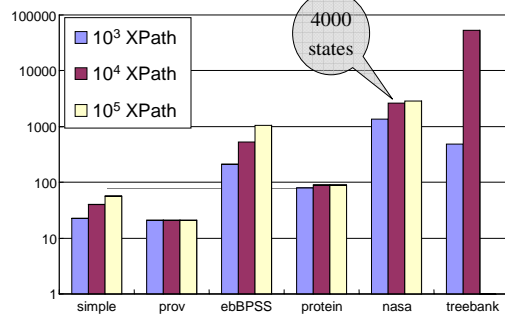Fact of life: real XML data has "small" data guide [Liefke&S.'00]

## Lazy DFA and "Simple" DTDs

**Theorem** [GMOS'02] If the XML data has a data guide with G nodes, then the number of states in the lazy DFA is at most:
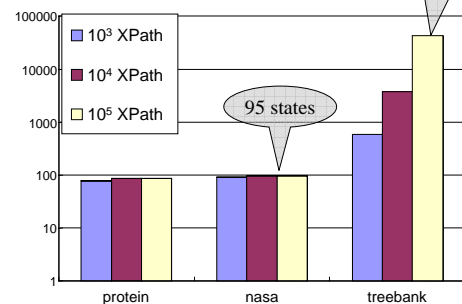
$$1+G$$

G = number of nodes in the data guide
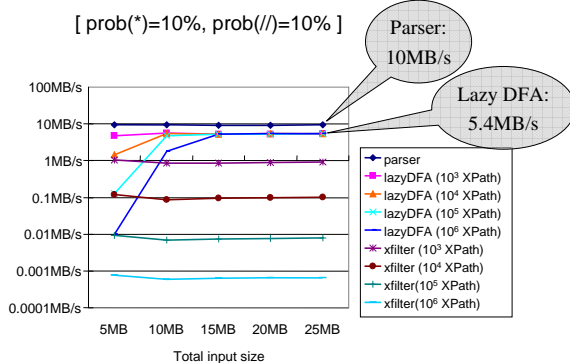
---

Number of Lazy DFA States - SYNTHETIC Data

- $10^3$ XPath
- $10^4$ XPath
- $10^5$ XPath

4000 states

(simple, prov, ebBPSS, protein, nasa, treebank)

---

Number of Lazy DFA States - REAL Data

40000 states
G = 350000

95 states

- $10^3$ XPath
- $10^4$ XPath
- $10^5$ XPath

(protein, nasa, treebank)

---

Throughput for $10^3$, $10^4$, $10^5$, $10^6$ XPath expressions

[ prob(*)=10%, prob(//)=10% ]

Parser: 10MB/s

Lazy DFA: 5.4MB/s

- parser
- lazyDFA ($10^3$ XPath)
- lazyDFA ($10^4$ XPath)
- lazyDFA ($10^5$ XPath)
- lazyDFA ($10^6$ XPath)
- xfilter ($10^3$ XPath)
- xfilter ($10^4$ XPath)
- xfilter($10^5$ XPath)
- xfilter($10^6$ XPath)

Total input size

---

END
Lecture 9

78