# Valid Query Answers for XML

Slawomir Staworko[1]     Jan Chomicki[2]

[1]University at Buffalo and INRIA Futurs

[2]University at Buffalo and Warsaw University

June 25-29, 2007

## Querying Invalid XML

- Integration of XML documents
- Slight differences between schemas (e.g. different cardinality constraints)
- Legacy XML databases
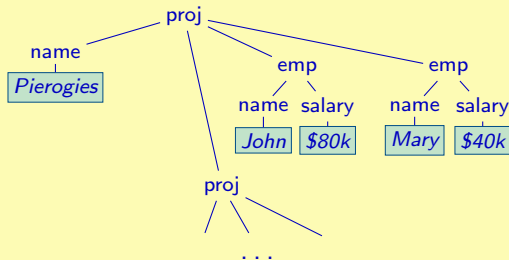- Database updates

# Invalid XML documents

## Querying Invalid XML

- ▶ Integration of XML documents
- ▶ Slight differences between schemas (e.g. different cardinality constraints)
- ▶ Legacy XML databases
- ▶ Database updates

## Document Type Definition ($D_0$)

```
proj   → (name, emp, proj*, emp*)
emp    → (name, salary)
name → #PCDATA
salary → #PCDATA
```

## Query: get salaries of all employees that are not managers

//proj/name/emp/following_sibling::emp/salary
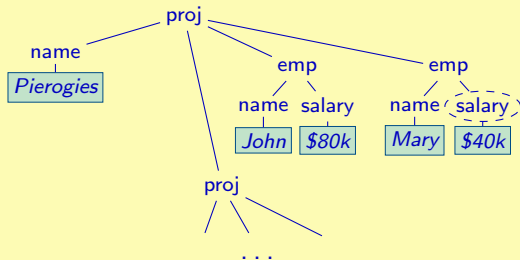
# Invalid XML documents

## Querying Invalid XML

- ▶ Integration of XML documents
- ▶ Slight differences between schemas (e.g. different cardinality constraints)
- ▶ Legacy XML databases
- ▶ Database updates

## Document Type Definition ($D_0$)

proj   → (name, emp, proj*, emp*)
emp   → (name, salary)
name → #PCDATA
salary → #PCDATA

## Document with errors



*Query*: get salaries of all employees that are not managers

//proj/name/emp/following_sibling::emp/salary

## Querying Invalid XML

- ▶ Integration of XML documents
- ▶ Slight differences between schemas (e.g. different cardinality constraints)
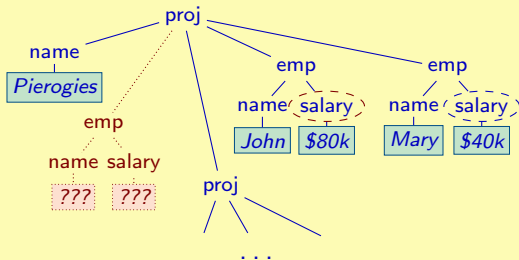- ▶ Legacy XML databases
- ▶ Database updates

## Document Type Definition ($D_0$)

proj  → (name, emp, proj*, emp*)
emp   → (name, salary)
name → #PCDATA
salary → #PCDATA

## Document with errors



*Query*: get salaries of all employees that are not managers

//proj/name/emp/following_sibling::emp/salary

# Invalid XML documents

## Querying Invalid XML

- ▶ Integration of XML documents
- ▶ Slight differences between schemas (e.g. different cardinality constraints)
- ▶ Legacy XML databases
- ▶ Database updates

## Document Type Definition ($D_0$)

proj   → (name, emp, proj*, emp*)
emp    → (name, salary)
name → #PCDATA
salary → #PCDATA

## Document with errors



**Query**: get salaries of all employees that are not managers

//proj/name/emp/following_sibling::emp/salary

# Core XPath

## Core XPath Queries

- text values but no attributes
- all standard axes
- subexpressions, value tests, and joins:
  `//*[A/B]`, `//*[B/text()=C/text()]`
- negation and disjunction:
  `//*[not A/B]`, `//A | //B`
- no functions
  `count(/A/*)`



## Tree Reachability Fact: $(x, Q, y)$

$$\boxed{\text{node } x} \xrightarrow{\text{query } Q} \boxed{\text{object } y}$$

Basic facts use only `/*` and `following_sibling::`

$$(N_0, /*, N_1), (N_1, /*, N_2), ...$$

Other facts are inferred with implications:
$$(X, Q/P, Y) \Leftarrow (X, Q, Z) \wedge (Z, P, Y).$$

$$(N_0, /*/*, N_2).$$

## Query Answers

- given query $Q$ and document $T$ with the root node $r$
- find all tree facts that hold in $T$
- $x$ in an answer to $Q$ in $T$ iff the tree fact $(r, Q, x)$ holds in $T$

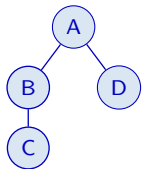# Query Evaluation for Positive Core XPath (no negation)

## Bottom-up approach

- computing tree facts for query $Q$
- tree facts for $T_1, \ldots, T_m$ computed before
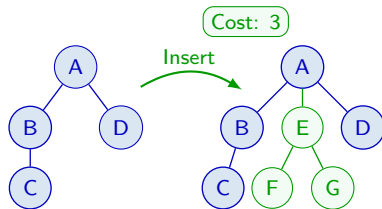- including inferred facts (involving subqueries of $Q$)

## Tree



## Algorithm

I start with $\varnothing$

II for subtree $T_i$ $(i = 1, \ldots, m)$

    1 add all facts of the subtree (obtained by recursion)

    2 add $(N_0, /*, N_i)$

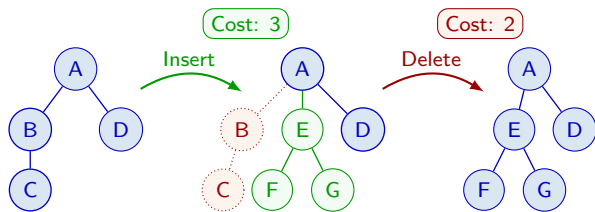    3 if $i > 1$ add $(N_{i-1}, \texttt{following\_sibling::*}, N_i)$

- Inserting a subtree

Cost: 3

Insert

Cost: 2

Delete

## Editing operations

- ▶ Inserting a subtree
- ▶ Deleting a subtree

## Editing operations

▶ Inserting a subtree
▶ Deleting a subtree
▶ Modifying a node's label

## Distance between documents

$dist(T, S)$ is the minimal cost of transforming $T$ into $S$

## Distance to a DTD

$dist(T, D)$ is the minimal cost of repairing $T$ w.r.t $D$ i.e.,

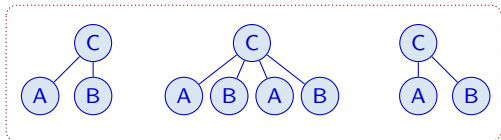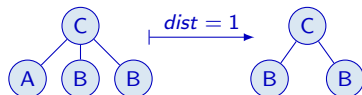$$\min\{dist(T, S) | S \text{ valid w.r.t } D\}$$

## DTD

$C \rightarrow (A,B)*$
$A \rightarrow$ EMPTY
$B \rightarrow$ EMPTY

## Repair

$T'$ is a repair of $T$ w.r.t $D$ iff

$$dist(T', T) = dist(T, D)$$



There can be an exponential number of repairs

# Valid Query Answers

## Valid Query Answers

$x$ is a valid answer to query $Q$ in $T$ w.r.t. $D$ iff
$x$ is an answer to $Q$ in every repair of $T$ w.r.t. $D$.

## XML Document



## Document Type Definition ($D_0$)

proj $\rightarrow$ (name, emp, proj*, emp*)
emp $\rightarrow$ (name, salary)
name $\rightarrow$ #PCDATA
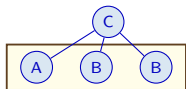salary $\rightarrow$ #PCDATA

## Queries and Valid Answers

```
//proj[emp[1]/salary='$90K']/name/text()      → {Pierogies}
//proj[name='Pierogies']/emp[1]/salary/text() → {$90K}
//proj[name='Pierogies']/emp[1]/name/text()   → ∅
//proj[name='Pierogies']/emp[1]/salary        → ∅
```

## DTD
C → (A,B)*
A → EMPTY
B → EMPTY

Read A

$Q_0$   $Q_1$

Read B
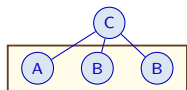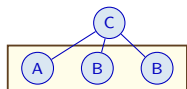
A        B        B

# Trace graph



DTD

C → (A,B)*
A → EMPTY
B → EMPTY

## DTD

$C \rightarrow (A,B)*$
$A \rightarrow EMPTY$
$B \rightarrow EMPTY$

Read A

Del    Ins A    Del

$Q_0$    $Q_1$

Ins B

Read B

A    B    B

$Q_0$    $Q_0$    $Q_0$    $Q_0$

$Q_1$    $Q_1$    $Q_1$    $Q_1$

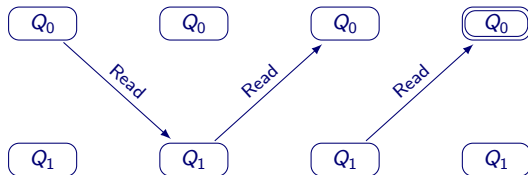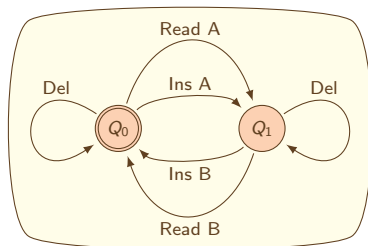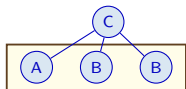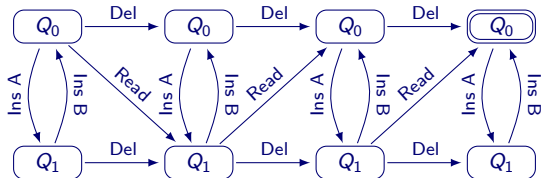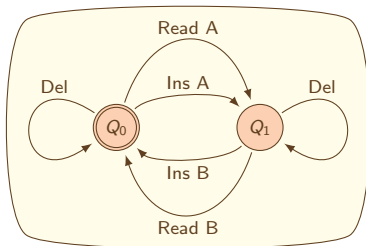# Trace graph



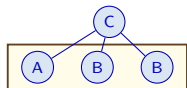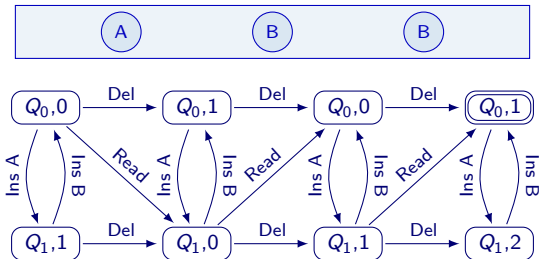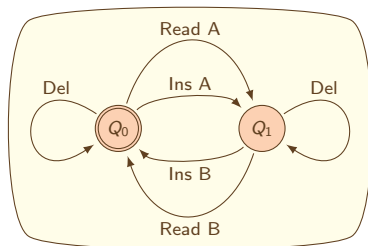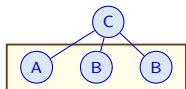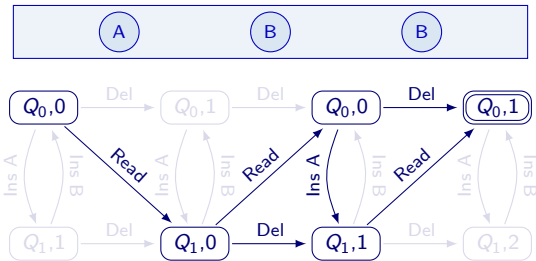**DTD**

$C \rightarrow (A,B)*$
$A \rightarrow EMPTY$
$B \rightarrow EMPTY$

# Trace graph

## DTD

C → (A,B)*
A → EMPTY
B → EMPTY

# Trace graph



DTD

C → (A,B)*
A → EMPTY
B → EMPTY

Compact representation of all repairs



Repairing Paths:

- (Read, Read, Del)
- (Read, Read, Ins A, Read)
- (Read, Del, Read)

# Computing Valid Query Answers

## Certain Tree Facts

Tree facts present in every repair of a given tree

## Bottom-up approach

Precomputed values:
- certain facts for all children
- certain facts common for every minimal tree satisfying DTD

## Obtain certain facts by

Intersection of the sets of all repairing paths



## For every repairing path, construct:

set of facts collected "so far":
- start with $\varnothing$
- **Read** adds certain facts of the corresponding child
- **Del** adds nothing
- **Ins A** adds certain facts common for minimal trees labeled with A

## Problem

Possibly an exponential number of paths

## Solution: Eager Intersection

Intersect all sets of certain facts for paths sharing the same last operation adding facts (Read/Ins).

## Computing VQA

| Queries | Combined-complexity | Data-complexity |
|---|---|---|
| Descending axes | PTIME | PTIME |
| + Ascending axes | co-NP-complete | ??? |
| + Sliding axes | co-NP-complete | ??? |
| + Negation/Disjunction | co-NP-complete | ??? |
| + Joins | co-NP-complete | co-NP-complete |

## Compared algorithms

| | |
|---|---|
| PARSE | base line |
| VALIDATE | regular automata |
| DIST | distance computation |

## Data generation

1. random valid document
2. removing and adding random nodes
3. invalidity ratio

   $$dist(T, D)/|T| \simeq 0.1\%$$

4. small height (8-10)

## Compared algorithms

| PARSE | base line |
|---|---|
| VALIDATE | regular automata |
| DIST | distance computation |

## Data generation

1. random valid document
2. removing and adding random nodes
3. invalidity ratio

   $dist(T, D)/|T| \simeq 0.1\%$

4. small height (8-10)

## Compared algorithms
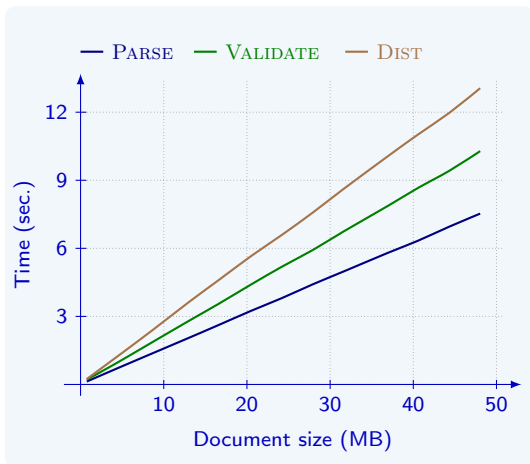
PARSE       base line
VALIDATE    regular automata
DIST        distance computation

## Data generation

1. random valid document

2. removing and adding
   random nodes

3. invalidity ratio

   $dist(T, D)/|T| \simeq 0.1\%$

4. small height (8-10)



Chart legend: — VALIDATE   — DIST

Y-axis: Time (sec.)
X-axis: DTD size $|D|$

## Compared algorithms

QA    base line
VQA  eager intersection

## Data generation

- invalidity ratio

    $dist(T, D)/|T| \simeq 0.1\%$

- small height (8-10)

## Implemented Queries

- `//, /, following_sibling::`
- `name()=A,text()='`*str*`'`
- QA works in linear time

## Compared algorithms

QA    base line
VQA  eager intersection

## Data generation

- invalidity ratio

  $dist(T, D)/|T| \simeq 0.1\%$

- small height (8-10)

## Implemented Queries

- `//, /, following_sibling::`
- `name()=A,text()='`*str*`'`
- QA works in linear time

## Compared algorithms
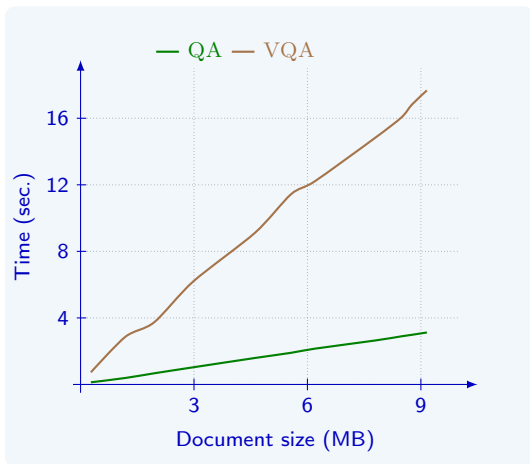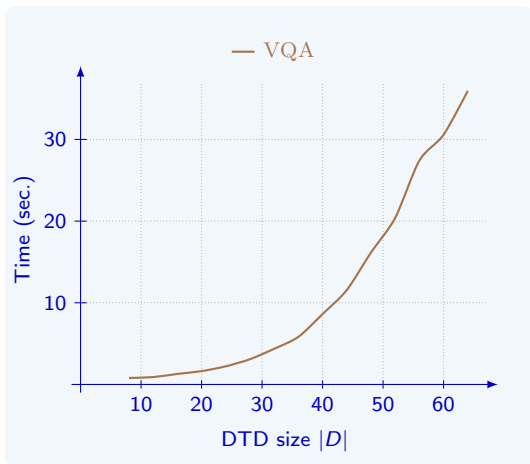
QA    base line
VQA  eager intersection

## Data generation

- invalidity ratio

    $dist(T, D)/|T| \simeq 0.1\%$

- small height (8-10)

## Implemented Queries

- `//, /, following_sibling::`
- `name()=A,text()='str'`
- QA works in linear time

## Conclusions

- ▶ Framework for querying documents with validity violations of local nature (missing or superfluous nodes)
- ▶ Efficient algorithm for computing valid answers to a class of XPath queries

## Future Work

- ▶ Valid answers by query rewriting
- ▶ In-depth analysis of data complexity
- ▶ Other tree operations: inner node deletion/insertion, node move, . . . [1]
- ▶ Semantic inconsistencies: keys, functional dependencies,. . . [2]

Marcelo Arenas

Leo Bertossi

Wenfei Fan

Jerzy Marcinkowski

Slawomir Staworko

Jef Wijsen

P. Bille.
A Survey on Tree Edit Distance and Related Problems.
*Theoretical Computer Science*, 337(1-3):217–239, 2003.

S. Flesca, F. Furfaro, S. Greco, and E. Zumpano.
Querying and Repairing Inconsistent XML Data.
In *Web Information Systems Engineering*, pages 175–188. Springer, LNCS 3806, 2005.

S. Staworko and J. Chomicki.
Validity-Sensitive Querying of XML Databases.
In *EDBT Workshops (dataX)*, pages 164–177. Springer, LNCS 4254, 2006.