

Type Based XML Projection

VLDB 2006, Seoul, Korea

Véronique Benzaken*

Giuseppe Castagna[◇]

Dario Colazzo*

Kim Nguyễn*

* : Équipe Bases de Données, LRI, Université Paris-Sud 11, Orsay, France

◇ : Équipe Langage, DI, École Normale Supérieure, Paris, France

- ① Introduction
- ② Notations
- ③ Algorithm
- ④ Formal results
- ⑤ Experiments
- ⑥ Conclusion

- 1 Introduction
- 2 Notations
- 3 Algorithm
- 4 Formal results
- 5 Experiments
- 6 Conclusion

Main memory query (XQuery) processing

Main memory query (XQuery) processing

Pros :

- Lightweight system

Main memory query (XQuery) processing

Pros :

- Lightweight system
- Easy to program and integrate into an existing architecture

Main memory query (XQuery) processing

Pros :

- Lightweight system
- Easy to program and integrate into an existing architecture

Cons :

- Unlike DBMS solution, whole document in memory (DOM → lots of metadata)

Main memory query (XQuery) processing

Pros :

- Lightweight system
- Easy to program and integrate into an existing architecture

Cons :

- Unlike DBMS solution, whole document in memory (DOM → lots of metadata)

We can use **pruning** to reduce the amount of data needed in main memory!

We can use **pruning** to reduce the amount of data needed in main memory!

Exemple :

```
/descendant-or-self::title/child::text()
```

We can use **pruning** to reduce the amount of data needed in main memory!

Exemple :

```
/descendant-or-self::title/child::text()
```

```
<bib>  
  <book year="1994">  
    <title>TCP/IP Illustrated</title>  
    <author><last>Stevens</last><first>W.</first></author>  
    <publisher>Addison-Wesley</publisher>  
    <price>65.95</price>  
  </book>  
  <book year="1992">  
    <title>Advanced Programming in the Unix environment</title>  
    <author><last>Stevens</last><first>W.</first></author>  
    <publisher>Addison-Wesley</publisher>  
    <price>65.95</price>  
  </book>  
  ...  
</bib>
```

Projecting XML Documents, Amélie Marian and Jérôme Siméon, VLDB 2003.

Queries + document based optimizations, it has some drawbacks :

- Does not take into account backward axis.
- Performances degrade in the presence of //

Indeed, with a query like :

```
/descendant-or-self::title/child::text()
```

one has to iterate through the whole document.

Solution : Queries + **type** based optimisations

By using a DTD and a given set of queries, we can statically infer a *projector* for the set of queries and use it to prune the document.

Main advantages :

- Pruning is efficient
- Take into account // and backward axis
- Pruning is sound : executing the query on the projected document gives the same result as on the original.
- Pruning is precise and even exact for a large class of documents

- 1 Introduction
- 2 Notations**
- 3 Algorithm
- 4 Formal results
- 5 Experiments
- 6 Conclusion

Definition (DTD)

A DTD is a pair (X, E) where X is a distinguished name (the root) and E is a set of productions of the form

$$\{X_1 \rightarrow R_1, \dots, X_n \rightarrow R_n\}$$

where each R_i is of the form :

$$a_i[\text{Regex}] \text{ or } \text{String}$$

where a_i is a unique *tag* name and *Regex* a regular expression of X_i . $\text{Names}(E)$ is the set of names occurring in E .

Definition (Projector)

For a given DTD (X, E) a type projector for (X, E) is a set of names \mathcal{P} such that :

- 1 $\mathcal{P} \subseteq \text{Names}(E)$
- 2 $\forall X_i \in \mathcal{P}$ there is at least one derivation from X to X_i in E

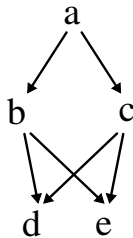
Definition (Pruning)

The pruning (or projection) of a document D valid w.r.t a DTD (X, E) with the projector \mathcal{P} is a document D' where every node not generated by a name in \mathcal{P} is erased (replaced by the empty sequence).

Definition (Pruning)

The *pruning* (or *projection*) of a document D valid w.r.t a DTD (X, E) with the projector \mathcal{P} is a document D' where every node not generated by a name in \mathcal{P} is erased (replaced by the empty sequence).

```
<!ELEMENT a    (b,c)>  
<!ELEMENT b    (d,e)>  
<!ELEMENT c    (d,e)>  
<!ELEMENT d    (#PCDATA)>  
<!ELEMENT e    (#PCDATA)>
```



`/descendant-or-self::c`

/descendant-or-self::c

$X_a \rightarrow \mathbf{a}[X_b, X_c]$

$X_b \rightarrow \mathbf{b}[X_d, X_e]$ $X_c \rightarrow \mathbf{c}[X_d, X_e]$

$X_d \rightarrow \mathbf{String}$ $X_e \rightarrow \mathbf{String}$

<a>

<d>LotsOfData</d>

<e>SuperLongString<e/>

<c>

<d>bar</d>

<e>foo<e/>

</c>

/descendant-or-self::c

$X_a \rightarrow \mathbf{a}[X_b, X_c]$

$X_b \rightarrow \mathbf{b}[X_d, X_e]$ $X_c \rightarrow \mathbf{c}[X_d, X_e]$

$X_d \rightarrow \mathbf{String}$ $X_e \rightarrow \mathbf{String}$

$\mathcal{P} = \{X_a, X_c, X_d, X_e\}$

```
<a>
  <b>
    <d>LotsOfData</d>
    <e>SuperLongString<e/>
  </b>
  <c>
    <d>bar</d>
    <e>foo<e/>
  </c>
</a>
```

```
<a>
  <b>
    <d>LotsOfData</d>
    <e>Superlongstring<e/>
  </b>
  <c>
    <d>foo</d>
    <e>bar<e/>
  </c>
</a>
```

/descendant-or-self::c

$X_a \rightarrow a[X_b, X_c]$

$X_b \rightarrow b[X_d, X_e]$ $X_c \rightarrow c[X_d, X_e]$

$X_d \rightarrow \text{String}$ $X_e \rightarrow \text{String}$

$\mathcal{P} = \{X_a, X_c, X_d, X_e\}$

```
<a>
  <b>
    <d>LotsOfData</d>
    <e>SuperLongString<e/>
  </b>
  <c>
    <d>bar</d>
    <e>foo<e/>
  </c>
</a>
```

```
<a>
  <b>
    <d>LotsOfData</d>
    <e>Superlongstring<e/>
  </b>
  <c>
    <d>foo</d>
    <e>bar<e/>
  </c>
</a>
```

Definition (XPath (1.0))

Q	::=	$Axis :: Test$	$Expr$::=	Q
		$Axis :: Test[Expr]$			$Expr \text{ op } Expr$
		Q/Q			$f(Expr, \dots, Expr)$
					$Base$

$Axis$::= self | child | descendant | parent | ancestor | ...

$Test$::= tag | node() | text()

Definition (Simple path)

Q	::=	$Axis :: Test$	$Cond$::=	$SPath$
		$Axis :: Test[Cond]$			$Cond \text{ or } Cond$
		Q/Q			

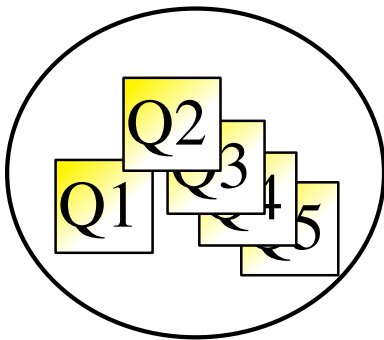
$SPath$::= $Axis :: Test$
| $Axis :: Test/SPath$

- 1 Introduction
- 2 Notations
- 3 Algorithm**
- 4 Formal results
- 5 Experiments
- 6 Conclusion

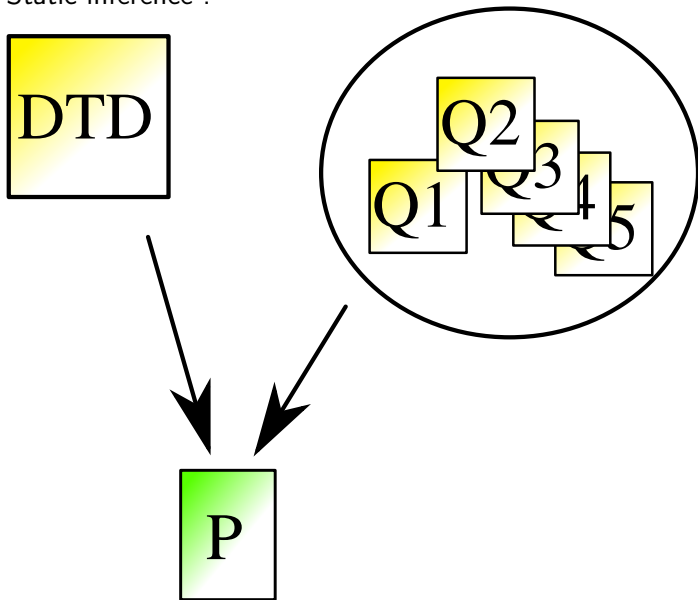
Static inference :



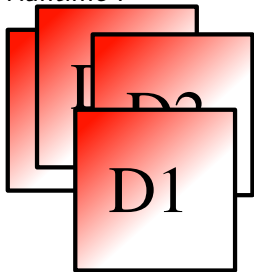
Static inference :



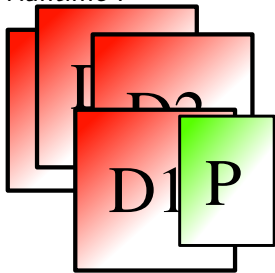
Static inference :



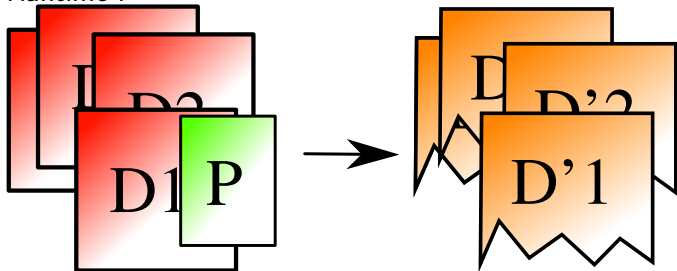
Runtime :



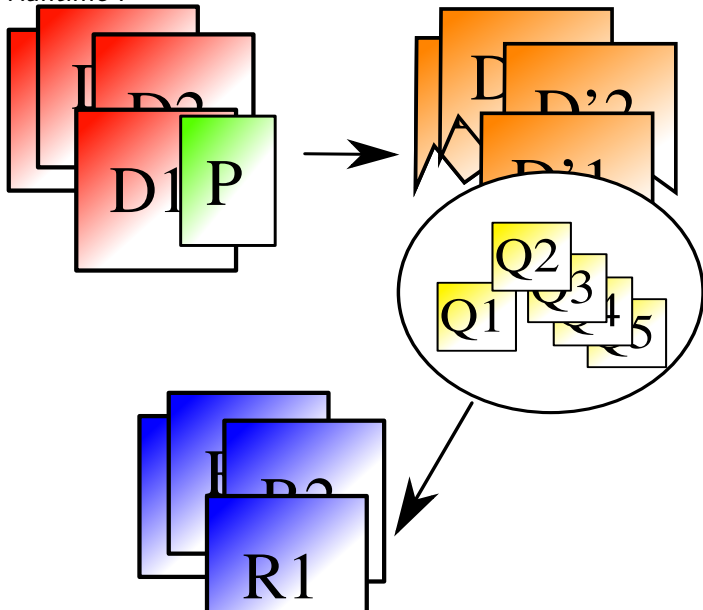
Runtime :

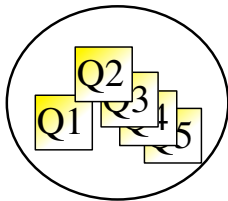


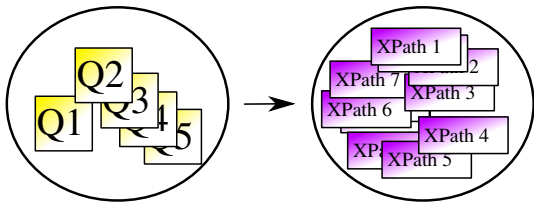
Runtime :

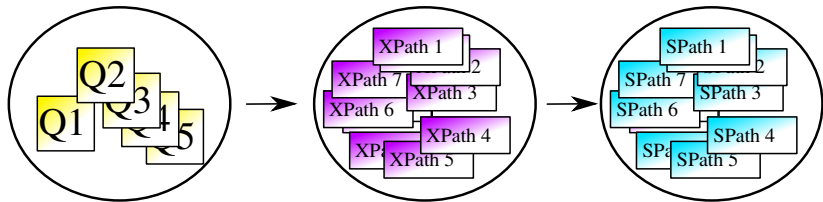


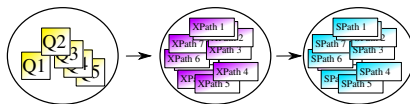
Runtime :



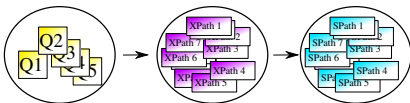






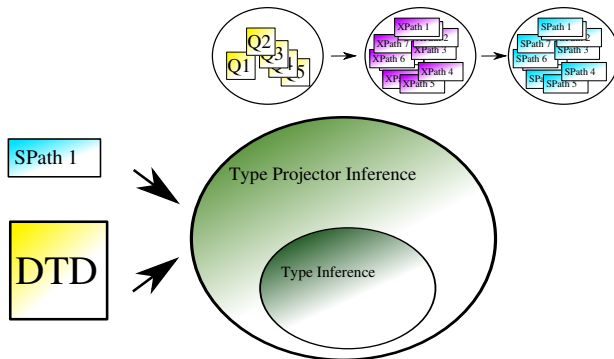


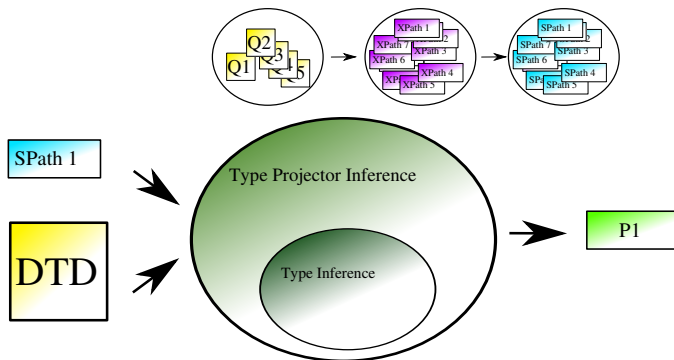
SPath 1

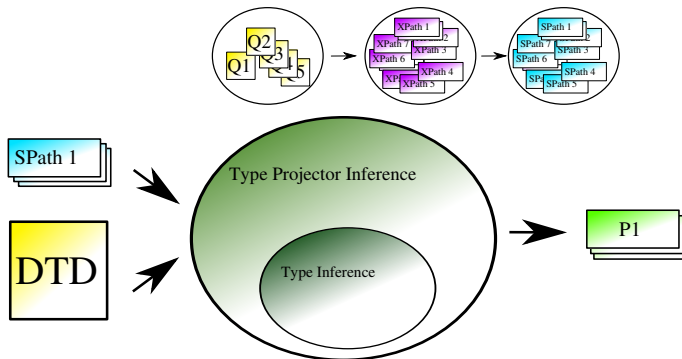


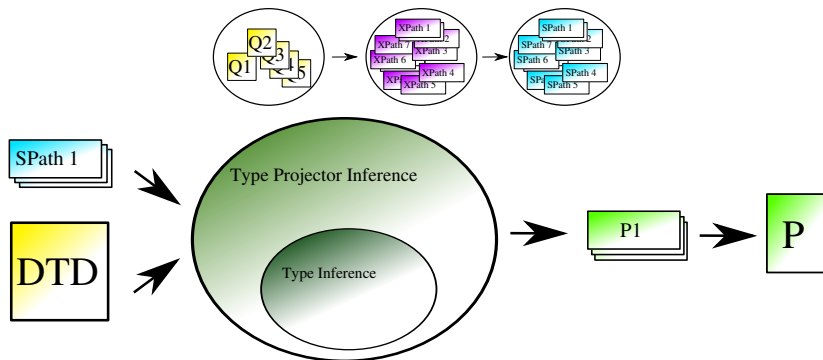
SPath 1

DTD









$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query**.

$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query**.

`/self::a/child::b/child::d`

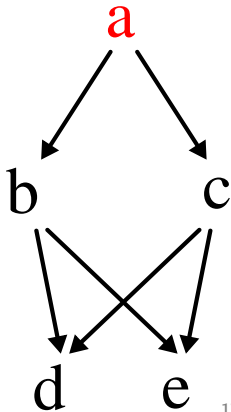
$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query.**

`/self::a/child::b/child::d`

$$\textcircled{1} \text{ typeinf}(\text{DTD}, \{X_a\}, \text{self}::\text{a}) = \{X_a\}$$



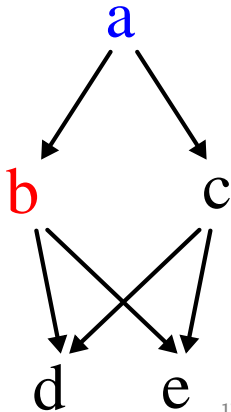
$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query.**

`/self::a/child::b/child::d`

- 1 $\text{typeinf}(\text{DTD}, \{X_a\}, \text{self}::\text{a}) = \{X_a\}$
- 2 $\text{typeinf}(\text{DTD}, \{X_a\}, \text{child}::\text{b}) = \{X_b\}$



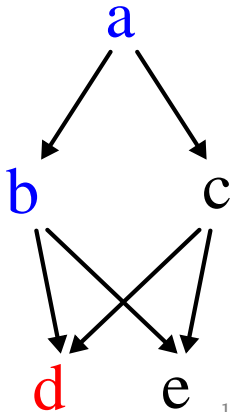
$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query.**

`/self::a/child::b/child::d`

- 1 $\text{typeinf}(\text{DTD}, \{X_a\}, \text{self}::\text{a}) = \{X_a\}$
- 2 $\text{typeinf}(\text{DTD}, \{X_a\}, \text{child}::\text{b}) = \{X_b\}$
- 3 $\text{typeinf}(\text{DTD}, \{X_b\}, \text{child}::\text{d}) = \{X_d\}$



$$\text{typeinf}(\text{DTD}, \{X\}, \text{path}) = T$$

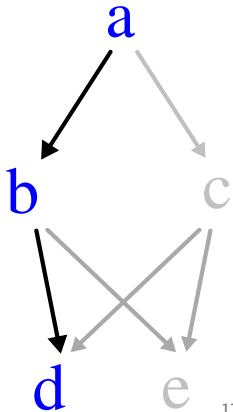
Type T is the set of names of types of nodes in the result.

$T = \emptyset \rightarrow$ **empty query.**

`/self::a/child::b/child::d`

- ❶ $\text{typeinf}(\text{DTD}, \{X_a\}, \text{self}::\text{a}) = \{X_a\}$
- ❷ $\text{typeinf}(\text{DTD}, \{X_a\}, \text{child}::\text{b}) = \{X_b\}$
- ❸ $\text{typeinf}(\text{DTD}, \{X_b\}, \text{child}::\text{d}) = \{X_d\}$

$$\mathcal{P} = \{X_a, X_b, X_d\}$$

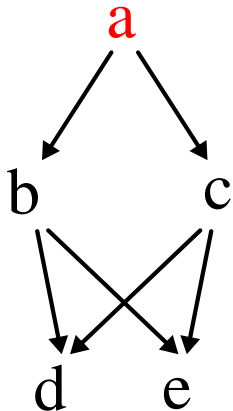


```
/**/self::b/child::d
```



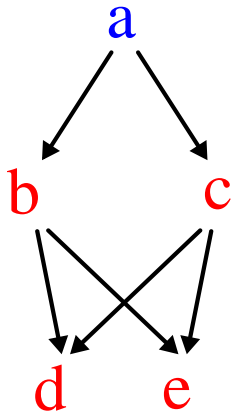
```
/**/self::b/child::d
```

❶ `typeinf(DTD, {Xa}, /**)`



```
/*/self::b/child::d
```

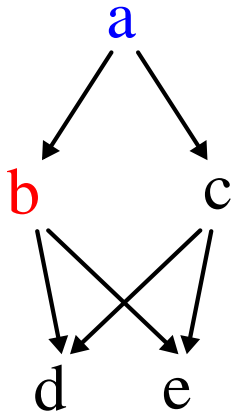
❶ `typeinf(DTD, {Xa}, /*/)`



```
/**/self::b/child::d
```

❶ `typeinf(DTD, {Xa}, /**)`

❶ `typeinf(DTD, {Xb}, self::b/child::d) = {Xd}`

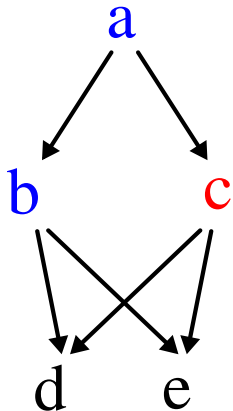


`/*/self::b/child::d`

① `typeinf(DTD, {Xa}, /*)`

① `typeinf(DTD, {Xb}, self::b/child::d) = {Xd}`

② `typeinf(DTD, {Xc}, self::b/child::d) = ∅`



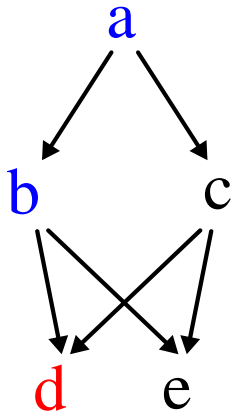
`/*/self::b/child::d`

① `typeinf(DTD, {Xa}, /*)`

① `typeinf(DTD, {Xb}, self::b/child::d) = {Xd}`

② `typeinf(DTD, {Xc}, self::b/child::d) = ∅`

③ `typeinf(DTD, {Xd}, self::b/child::d) = ∅`



`/*/self::b/child::d`

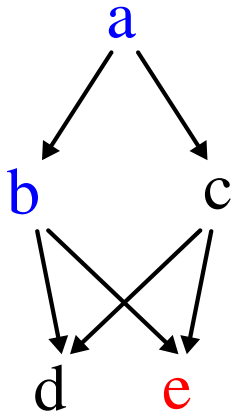
① `typeinf(DTD, {Xa}, /*)`

① `typeinf(DTD, {Xb}, self::b/child::d) = {Xd}`

② `typeinf(DTD, {Xc}, self::b/child::d) = ∅`

③ `typeinf(DTD, {Xd}, self::b/child::d) = ∅`

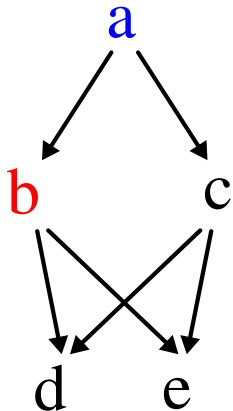
④ `typeinf(DTD, {Xe}, self::b/child::d) = ∅`



```
/**/self::b/child::d
```

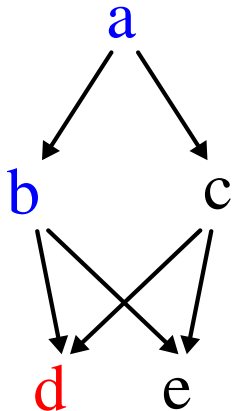
❶ `typeinf(DTD, {Xa}, /**)`

❷ `typeinf(DTD, {Xb}, self::b) = {Xb}`



```
/**/self::b/child::d
```

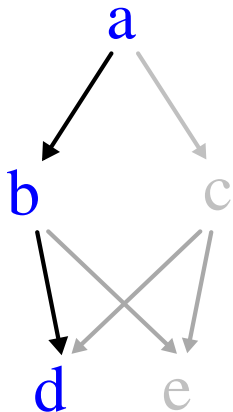
- 1 `typeinf(DTD, {Xa}, /**)`
- 2 `typeinf(DTD, {Xb}, self::b) = {Xb}`
- 3 `typeinf(DTD, {Xb}, child::d) = {Xd}`




```
/**/self::b/child::d
```

- 1 `typeinf(DTD, {Xa}, /**)`
- 2 `typeinf(DTD, {Xb}, self::b) = {Xb}`
- 3 `typeinf(DTD, {Xb}, child::d) = {Xd}`

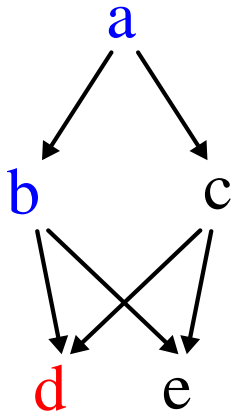
$$\mathcal{P} = \{X_a, X_b, X_d\}$$



```
/self::a/child::b/child::d/parent::node()/child::d
```

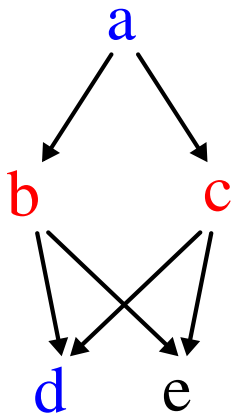
`/self::a/child::b/child::d/parent::node()/child::d`

- 1 `typeinf(DTD, {Xa}, self::a) = {Xa}`
- 2 `typeinf(DTD, {Xa}, child::b) = {Xb}`
- 3 `typeinf(DTD, {Xb}, child::d) = {Xd}`



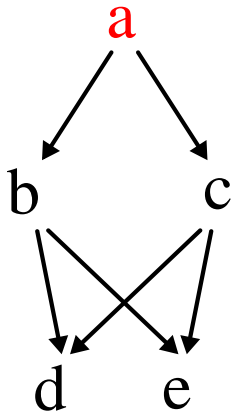
`/self::a/child::b/child::d/parent::node()/child::d`

- 1 `typeinf(DTD, {Xa}, self::a) = {Xa}`
- 2 `typeinf(DTD, {Xa}, child::b) = {Xb}`
- 3 `typeinf(DTD, {Xb}, child::d) = {Xd}`
- 4 `typeinf(DTD, {Xd}, parent::node()) = {Xb, Xc}`



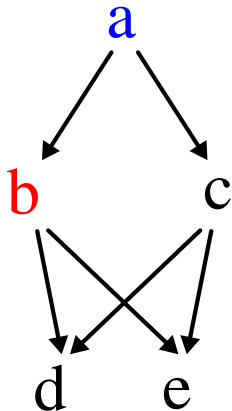
`/self::a/child::b/child::d/parent::node()/child::d`

① `typeinf(DTD, {Xa}, {Xa}, self::a) = {Xa}`



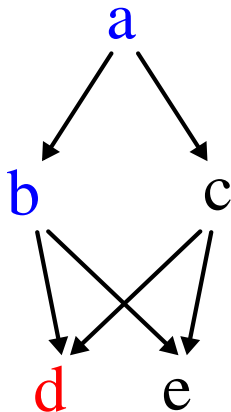
`/self::a/child::b/child::d/parent::node()/child::d`

- 1 `typeinf(DTD, {Xa}, {Xa}, self::a) = {Xa}`
- 2 `typeinf(DTD, {Xa}, {Xa}, child::b) = {Xb}`



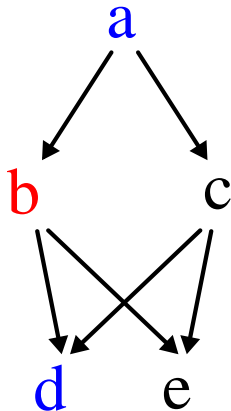
`/self::a/child::b/child::d/parent::node()/child::d`

- 1 `typeinf(DTD, {Xa}, {Xa}, self::a) = {Xa}`
- 2 `typeinf(DTD, {Xa}, {Xa}, child::b) = {Xb}`
- 3 `typeinf(DTD, {Xb}, {Xa, Xb}, child::d) = {Xd}`



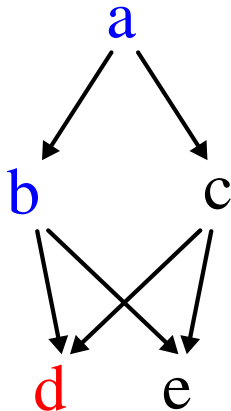
`/self::a/child::b/child::d/parent::node()/child::d`

- 1 `typeinf(DTD, { X_a }, { X_a }, self::a) = { X_a }`
- 2 `typeinf(DTD, { X_a }, { X_a }, child::b) = { X_b }`
- 3 `typeinf(DTD, { X_b }, { X_a, X_b }, child::d) = { X_d }`
- 4 `typeinf(DTD, { X_d }, { X_a, X_b, X_d }, parent::node()) = { X_b }`



`/self::a/child::b/child::d/parent::node()/child::d`

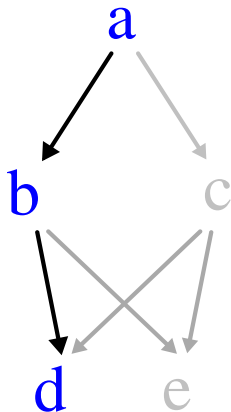
- 1 `typeinf(DTD, { X_a }, { X_a }, self::a) = { X_a }`
- 2 `typeinf(DTD, { X_a }, { X_a }, child::b) = { X_b }`
- 3 `typeinf(DTD, { X_b }, { X_a, X_b }, child::d) = { X_d }`
- 4 `typeinf(DTD, { X_d }, { X_a, X_b, X_d }, parent::node()) = { X_b }`
- 5 `typeinf(DTD, { X_b }, { X_a, X_b, X_d }, child::d) = { X_d }`



`/self::a/child::b/child::d/parent::node()/child::d`

- 1 $\text{typeinf}(\text{DTD}, \{X_a\}, \{X_a\}, \text{self}::\text{a}) = \{X_a\}$
- 2 $\text{typeinf}(\text{DTD}, \{X_a\}, \{X_a\}, \text{child}::\text{b}) = \{X_b\}$
- 3 $\text{typeinf}(\text{DTD}, \{X_b\}, \{X_a, X_b\}, \text{child}::\text{d}) = \{X_d\}$
- 4 $\text{typeinf}(\text{DTD}, \{X_d\}, \{X_a, X_b, X_d\}, \text{parent}::\text{node}()) = \{X_b\}$
- 5 $\text{typeinf}(\text{DTD}, \{X_b\}, \{X_a, X_b, X_d\}, \text{child}::\text{d}) = \{X_d\}$

$$\mathcal{P} = \{X_a, X_b, X_d\}$$



- 1 Introduction
- 2 Notations
- 3 Algorithm
- 4 Formal results**
- 5 Experiments
- 6 Conclusion

Theorem (Soundness)

Let D be a document valid w.r.t a DTD (X, E) and p a path. Let \mathcal{P} the type projector deduced from p . Let D' be the projection of D with \mathcal{P} .

$$\text{eval}(p, D) = \text{eval}(p, D')$$

Pruning is precise . . .

Pruning is precise ...

Theorem (Completeness)

$\mathcal{P} = \{X_1, \dots, X_n\}$ the type projector associated with a path p and a DTD (X, E) . Let $\mathcal{P}' = \mathcal{P} \setminus \{X_i\}$. There exists D a document and its projection D' such that :

$$\text{eval}(p, D) \neq \text{eval}(p, D')$$

Completeness holds with :

- some restrictions on the DTD (*star-guarded*, non-recursive, ...)

Pruning is precise ...

Theorem (Completeness)

$\mathcal{P} = \{X_1, \dots, X_n\}$ the type projector associated with a path p and a DTD (X, E) . Let $\mathcal{P}' = \mathcal{P} \setminus \{X_i\}$. There exists D a document and its projection D' such that :

$$\text{eval}(p, D) \neq \text{eval}(p, D')$$

Completeness holds with :

- some restrictions on the DTD (*star-guarded*, non-recursive, ...)
- some restrictions on the path (no upward axis in predicates, ...)

Pruning is precise ...

Theorem (Completeness)

$\mathcal{P} = \{X_1, \dots, X_n\}$ the type projector associated with a path p and a DTD (X, E) . Let $\mathcal{P}' = \mathcal{P} \setminus \{X_i\}$. There exists D a document and its projection D' such that :

$$\text{eval}(p, D) \neq \text{eval}(p, D')$$

Completeness holds with :

- some restrictions on the DTD (*star-guarded*, non-recursive, ...)
- some restrictions on the path (no upward axis in predicates, ...)

Pruning is precise ...

Theorem (Completeness)

$\mathcal{P} = \{X_1, \dots, X_n\}$ the type projector associated with a path p and a DTD (X, E) . Let $\mathcal{P}' = \mathcal{P} \setminus \{X_i\}$. There exists D a document and its projection D' such that :

$$\text{eval}(p, D) \neq \text{eval}(p, D')$$

Completeness holds with :

- some restrictions on the DTD (*star-guarded*, non-recursive, ...)
- some restrictions on the path (no upward axis in predicates, ...)

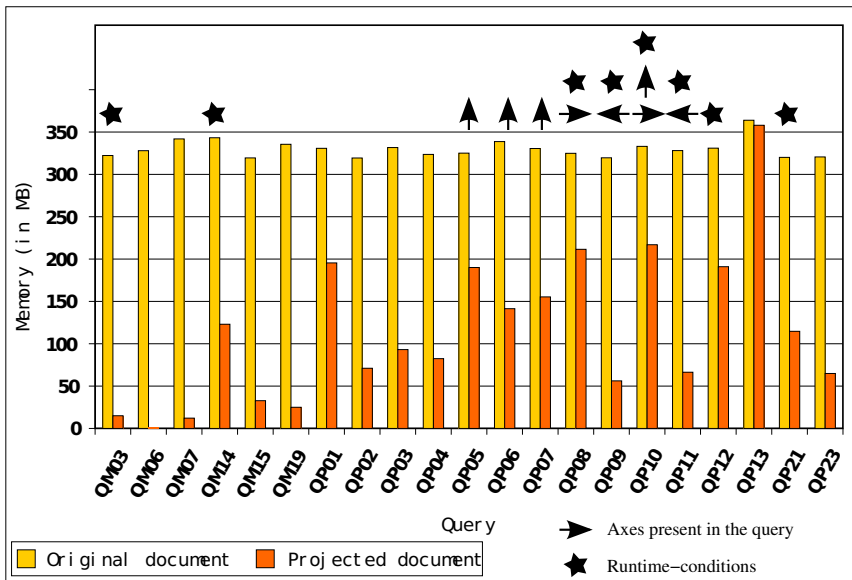
- 1 Introduction
- 2 Notations
- 3 Algorithm
- 4 Formal results
- 5 Experiments**
- 6 Conclusion

Protocol :

- Linux desktop, 512MB Ram, 3Ghz x86 CPU (and no swap).
- Implementation in OCaml.
- Validity of the result checked with the Galax query engine.
- Pruning tested on XMark and XPathMark benchmarks.

Protocol :

- Linux desktop, 512MB Ram, 3Ghz x86 CPU (and no swap).
 - Implementation in OCaml.
 - Validity of the result checked with the Galax query engine.
 - Pruning tested on XMark and XPathMark benchmarks.
- Pruning is **one pass bufferless** traversal of the document.
 - In practice, computing the type projector is fast.



Memory used to process a 56 MB document with Galax.

	QM03	QM06	QM07	QM14	QM15	QM19
	★			★		
Original Size (MB)	930	2048	1100	202	2048	964
Pruned Size(MB)	25	5,3	42	139	24	24
Memory Usage (MB)	374	90	380	512	245	512
% of original size	2.5	0.3	3.4	69.6	1.15	2.5
Gain in Speed (\times faster)	17.8	110.1	28.2	3.9	62.6	7.5

Qualitative : with one exception, type projectors are always equal or more efficient.

Qualitative : with one exception, type projectors are always equal or more efficient.

Performances : pruning is linear in time (in the size of the document) and constant in memory. *It can be done while validating the document.*

Qualitative : with one exception, type projectors are always equal or more efficient.

Performances : pruning is linear in time (in the size of the document) and constant in memory. *It can be done while validating the document.*

Features : handles backward as well as following/preceding axes.

- 1 Introduction
- 2 Notations
- 3 Algorithm
- 4 Formal results
- 5 Experiments
- 6 Conclusion**

- Formal foundations ensure validity of the pruning and it's efficiency.

- Formal foundations ensure validity of the pruning and it's efficiency.
- Handle any XQuery query, either directly or by approximating it.

- Formal foundations ensure validity of the pruning and it's efficiency.
- Handle any XQuery query, either directly or by approximating it.
- Whereas approximations are made for runtime conditions, the technique still preserve a *high degree of precision*.

- Formal foundations ensure validity of the pruning and it's efficiency.
- Handle any XQuery query, either directly or by approximating it.
- Whereas approximations are made for runtime conditions, the technique still preserve a *high degree of precision*.
- *No additional cost at runtime.*

Many directions :

- adapt our approach to work on *untyped* documents by using data-guides or path summaries.

Many directions :

- adapt our approach to work on *untyped* documents by using data-guides or path summaries.
- extend the formalism to handle XML-Schema rather than DTDs.

Many directions :

- adapt our approach to work on *untyped* documents by using data-guides or path summaries.
- extend the formalism to handle XML-Schema rather than DTDs.
- integration with classical databases techniques.

Many directions :

- adapt our approach to work on *untyped* documents by using data-guides or path summaries.
- extend the formalism to handle XML-Schema rather than DTDs.
- integration with classical databases techniques.
- integration with a query engine.