

# Discovering XSD Keys from XML Data

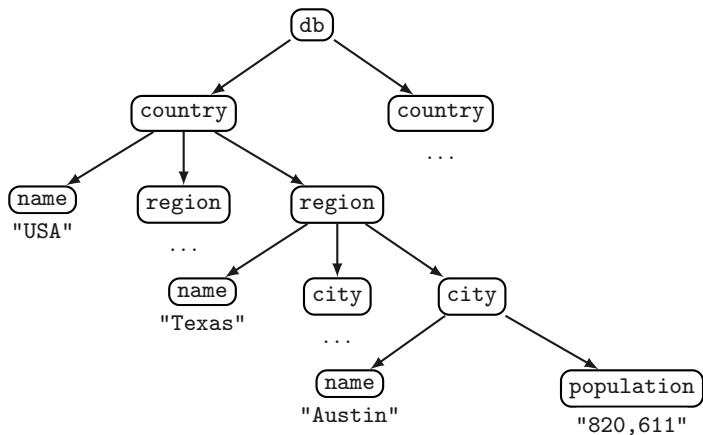
Marcelo Arenas

PUC Chile & U. of Oxford

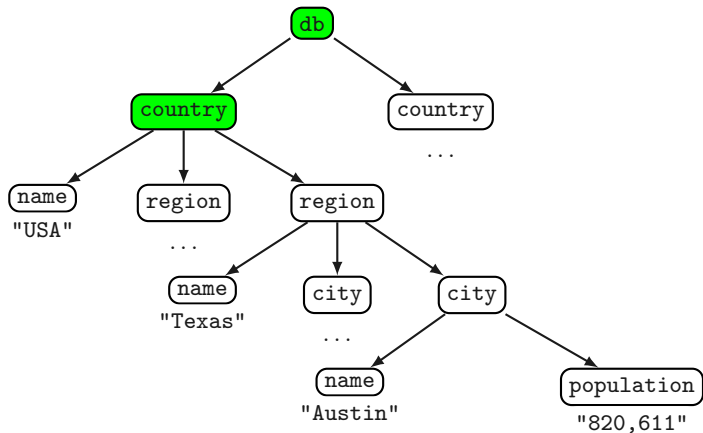
Joint work with Jonny Daenen, Frank Neven, Martin Ugarte, Jan Van den Bussche and Stijn Vansummeren

Beihang University, August 2013

## XML trees: An example

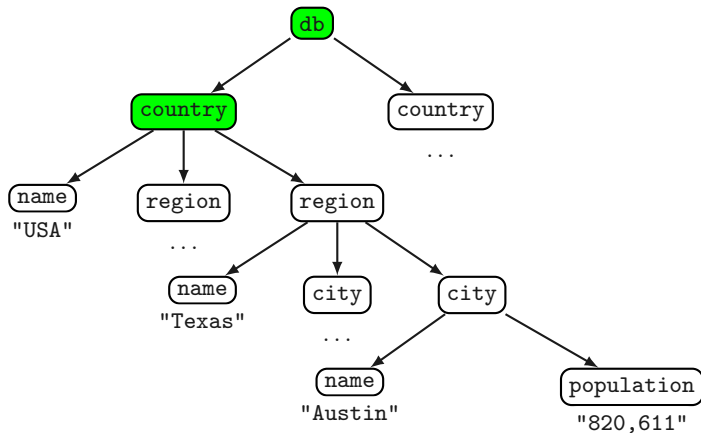


## XML schema (XSD): An example



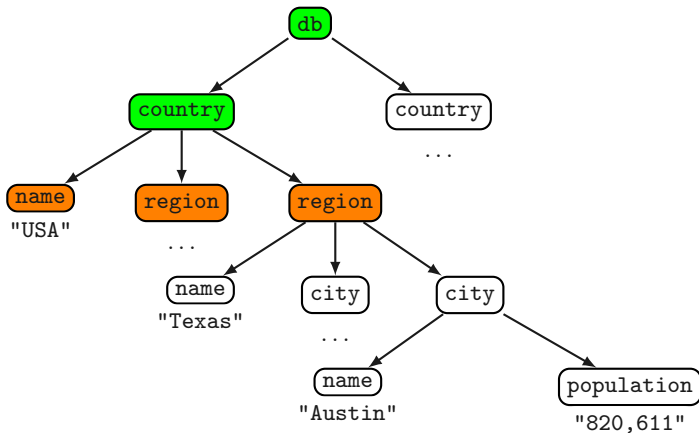
# XML schema (XSD): An example

`db/country`  $\mapsto$  `name region*`



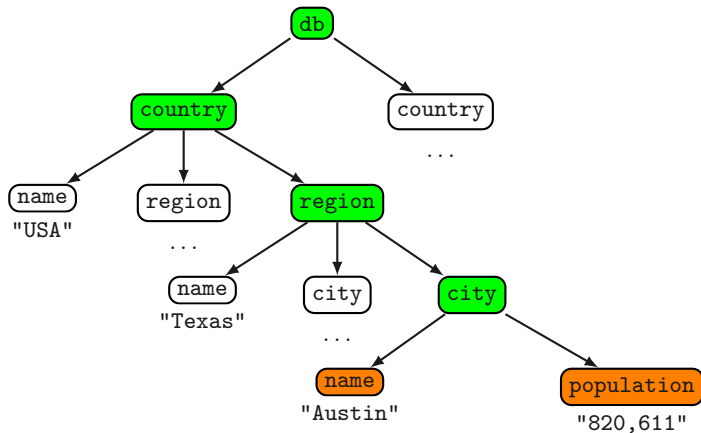
# XML schema (XSD): An example

`db/country`  $\mapsto$  `name region*`



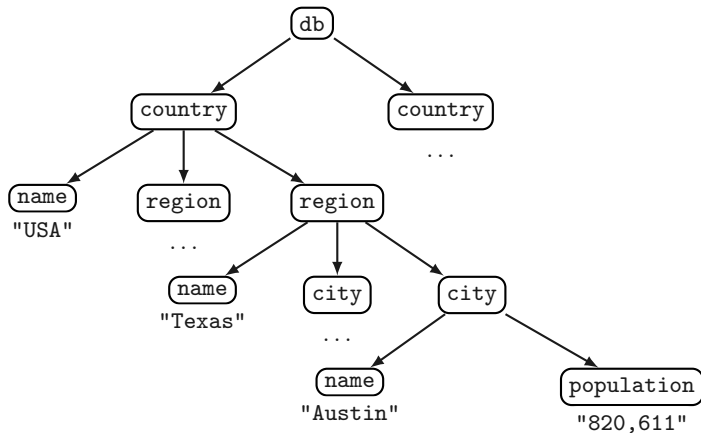
# XML schema (XSD): An example

db/country/region/city ↦ name population



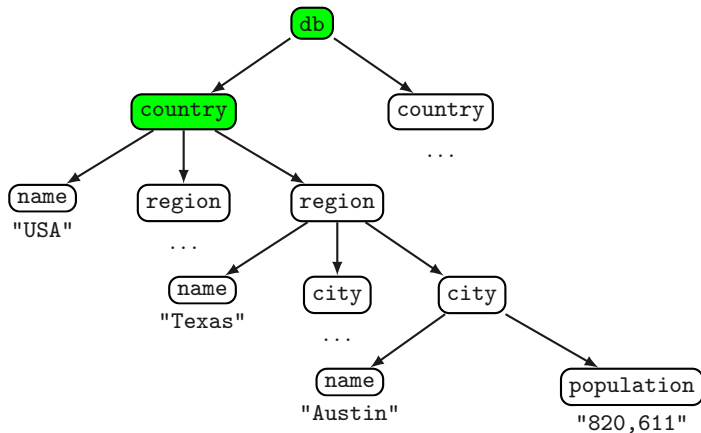
# XSD Keys: An example

(db/country, ../city, (./name))



# XSD Keys: An example

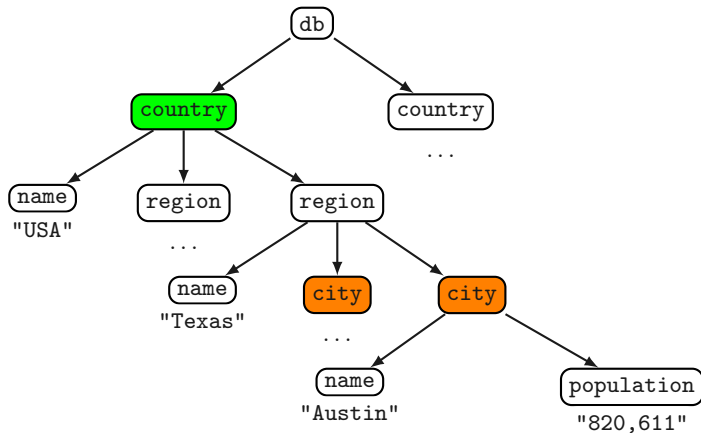
`(db/country, ../city, (./name))`





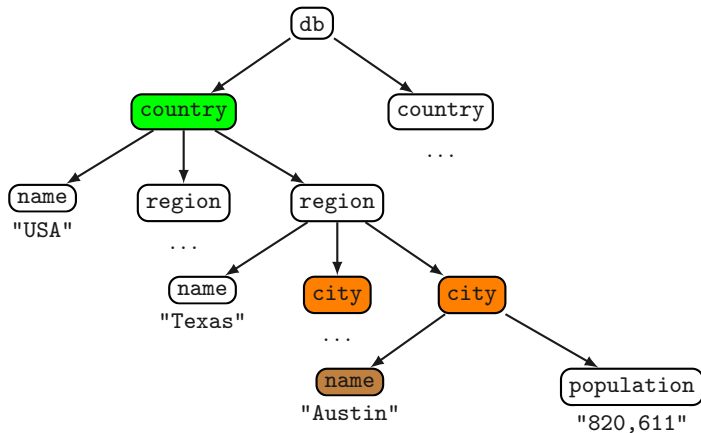
# XSD Keys: An example

(db/country, **./city**, (./name))



# XSD Keys: An example

(db/country, ../city, (./name))



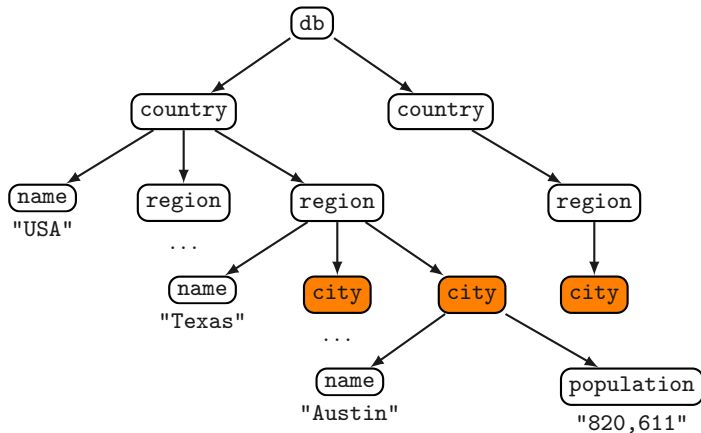
# XSD key mining

Support of a key in an XML tree: Number of nodes captured by the key

# XSD key mining

Support of a key in an XML tree: Number of nodes captured by the key

`(db/country, ../city, (./name))`



## XSD key mining (cont'd)

Given an XSD  $X$ , an XML tree  $t$  satisfying  $X$  and a threshold  $N$ , we want to find *all* keys  $\phi$  such that:

1.  $t$  satisfies  $\phi$
2. the support of  $\phi$  in  $t$  is greater than  $N$
3.  $\phi$  meets some quality requirements (which are defined only with respect to  $X$ )

# XSD key mining (cont'd)

Important issues:

- ▶ Collections of XML trees are not considered, as a collection can always be combined into a single XML tree by introducing a common root
- ▶ The standard definition of XML keys (given by the W3C) is considered in this work
- ▶ The quality requirements play a fundamental role in this work

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message



# XML trees

Given: Finite set  $\Sigma$  of element names and an infinite set **Data** of data elements.

- ▶  $\Sigma$ , **Data** are assumed to be disjoint

# XML trees

Given: Finite set  $\Sigma$  of element names and an infinite set **Data** of data elements.

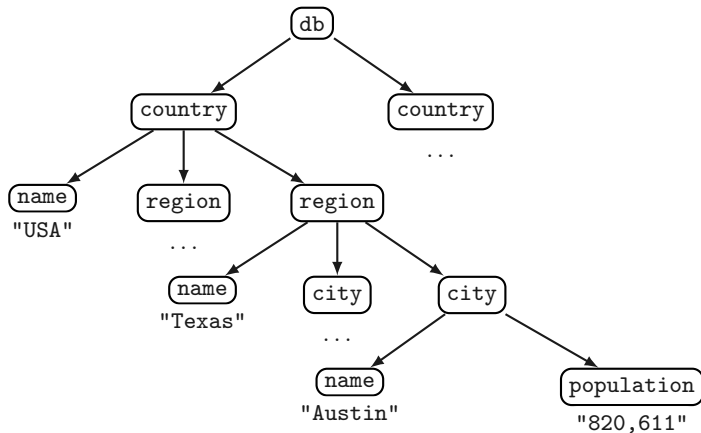
- ▶  $\Sigma$ , **Data** are assumed to be disjoint

XML tree (or just tree):  $(t, \text{lab}_t)$

- ▶  $t$ : Finite, unranked and sibling-ordered tree
- ▶  $\text{lab}_t$ : node labeling function with range  $(\Sigma \cup \mathbf{Data})$ , which labels non-leaf nodes with  $\Sigma$

Mixed content is not allowed: When a node is labeled with **Data**, then it is the only child of its parent.

## XML trees: An example



# XSD: Definition

An XSD defines the structure of a collection of XML trees.

Formally:  $X = (A, \lambda)$

- ▶  $A = (Q, \Sigma \cup \{\text{data}\}, q_0, \delta)$  is a DFA without final states
- ▶  $\lambda$  assigns to each state in  $Q$  a (one-unambiguous) regular expression over  $(\Sigma \cup \{\text{data}\})$

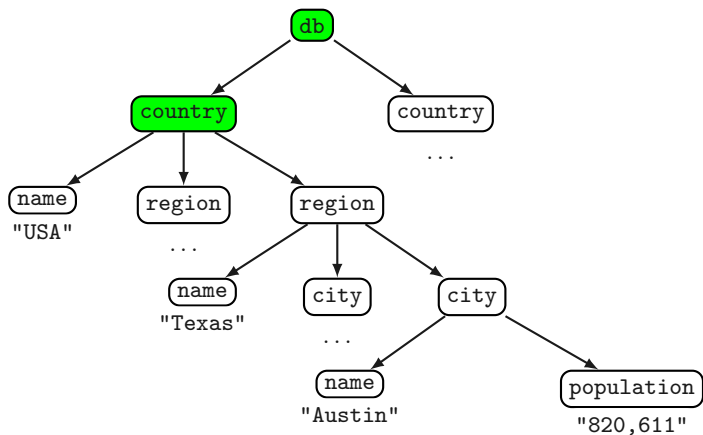
# XSD: Definition

An XSD defines the structure of a collection of XML trees.

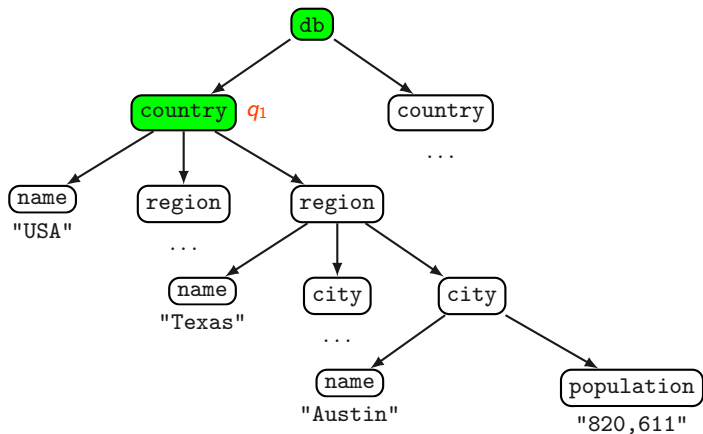
Formally:  $X = (A, \lambda)$

- ▶  $A = (Q, \Sigma \cup \{\text{data}\}, q_0, \delta)$  is a DFA without final states
- ▶  $\lambda$  assigns to each state in  $Q$  a (one-unambiguous) regular expression over  $(\Sigma \cup \{\text{data}\})$ 
  - ▶ There exists a polynomial-time algorithm that, given  $q \in Q$ , generates a DFA accepting  $\lambda(q)$

# Conforming to an XSD

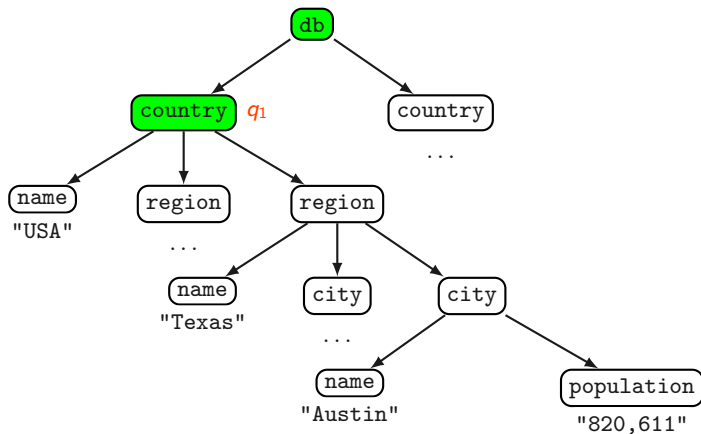


# Conforming to an XSD



# Conforming to an XSD

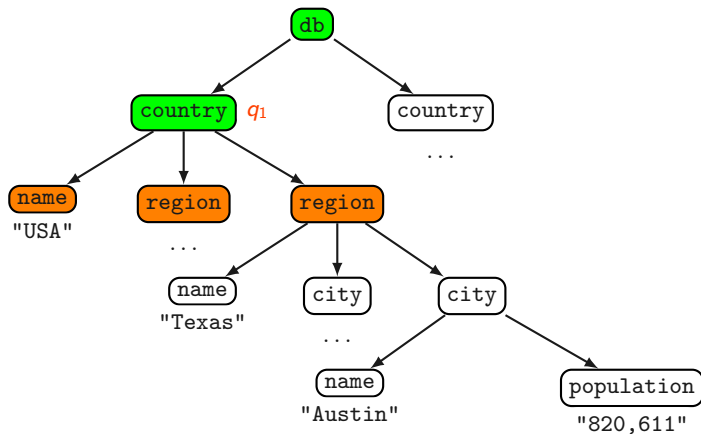
$$\lambda(q_1) = \text{name region}^*$$



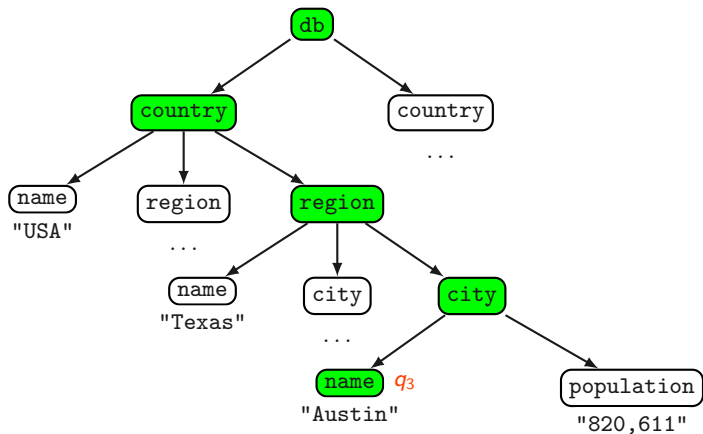


# Conforming to an XSD

$$\lambda(q_1) = \text{name region}^*$$

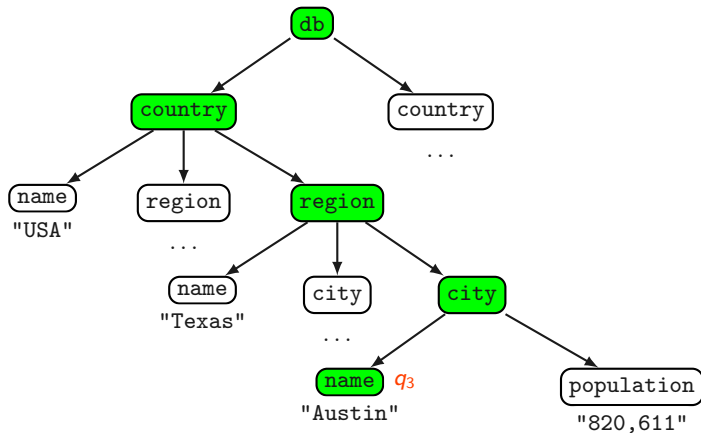


# Conforming to an XSD



# Conforming to an XSD

$\lambda(q_3) = \text{data}$



# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

## Defining XSD keys: Contexts

Let  $X = (A, \lambda)$  be an XSD, where  $Q$  is the set of states of  $A$

A context  $c = (q, a)$  is a pair in  $Q \times \Sigma$

## Defining XSD keys: Contexts

Let  $X = (A, \lambda)$  be an XSD, where  $Q$  is the set of states of  $A$

A context  $c = (q, a)$  is a pair in  $Q \times \Sigma$

Given an XML tree  $t$ ,  $\text{CNodes}_t(c)$  is the set of nodes  $v$  in  $t$  such that

1.  $\text{lab}_t(v) = a$
2.  $A$  halts in state  $q$  on input the string formed by the labels on the path from  $t$ 's root to  $v$

# Defining XSD keys: Selector expressions

A selector expression is of one of the following three forms:

- ▶ . (the dot symbol)
- ▶  $./a_1/\cdots/a_k$  (starting with the child axis), where  $k \geq 1$  and for every  $i \in \{1, \dots, k\}$ , either  $a_i \in \Sigma$  or  $a_i = *$
- ▶  $././a_1/\cdots/a_k$  (starting with the descendant axis), where  $k \geq 1$  and for every  $i \in \{1, \dots, k\}$ , either  $a_i \in \Sigma$  or  $a_i = *$

A union of selector expressions is of the form  $(\tau_1 | \cdots | \tau_m)$ , where  $m \geq 1$  and each  $\tau_i$  is a selector expression.

$\mathcal{SE}$  denotes the class of selector expressions and  $\mathcal{DSE}$  the class of union of selector expressions.

# Selector expressions: Semantics

Let  $w = w_1 \dots w_n$  be a string, where each  $w_i \in \Sigma$ .

- ▶  $w$  is said to match  $./a_1/\dots/a_k$  if  $n = k$  and for every  $i \in \{1, \dots, k\}$ :  $w_i = a_i$  or  $a_i = *$
- ▶  $w$  is said to match  $./a_1/\dots/a_k$  if a suffix of  $w$  matches  $./a_1/\dots/a_k$



## Selector expressions: Semantics (cont'd)

Let  $t$  be an XML tree,  $v$  a node in  $t$  and  $\tau$  a selector expression.

$\tau(t, v)$ : Nodes in  $t$  that are reachable from  $v$  by following  $\tau$ .

## Selector expressions: Semantics (cont'd)

Let  $t$  be an XML tree,  $v$  a node in  $t$  and  $\tau$  a selector expression.

$\tau(t, v)$ : Nodes in  $t$  that are reachable from  $v$  by following  $\tau$ .

- ▶  $\tau(t, v) = \{v\}$  if  $\tau = .$
- ▶ Otherwise,  $\tau(t, v)$  contains every node  $v'$  such that
  - ▶  $v'$  is a descendant of  $v$  in  $t$ , and
  - ▶ the path of labels from  $v$  to  $v'$  (but excluding the label of  $v$ ) matches  $\tau$

## Selector expressions: Semantics (cont'd)

Let  $t$  be an XML tree,  $v$  a node in  $t$  and  $\tau$  a selector expression.

$\tau(t, v)$ : Nodes in  $t$  that are reachable from  $v$  by following  $\tau$ .

- ▶  $\tau(t, v) = \{v\}$  if  $\tau = .$
- ▶ Otherwise,  $\tau(t, v)$  contains every node  $v'$  such that
  - ▶  $v'$  is a descendant of  $v$  in  $t$ , and
  - ▶ the path of labels from  $v$  to  $v'$  (but excluding the label of  $v$ ) matches  $\tau$

If  $\tau' = (\tau_1 | \dots | \tau_m)$  is a disjunction of selector expressions:

$$\tau'(t, v) = \tau_1(t, v) \cup \dots \cup \tau_m(t, v)$$

## XSD keys: Definition

### Definition

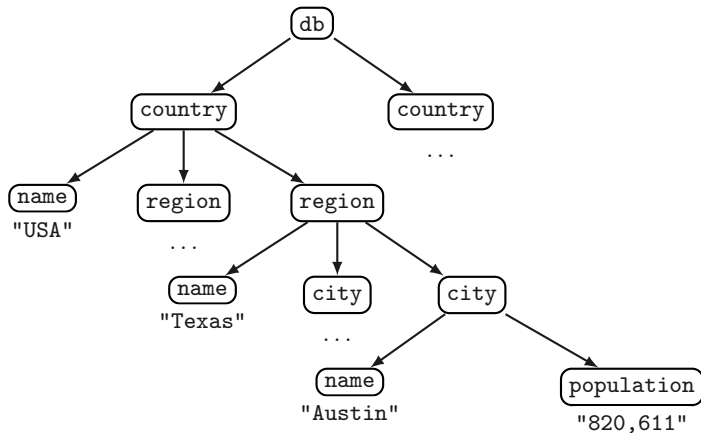
An XSD key, defined w.r.t. an XSD  $X$ , is a tuple  $\phi = (c, \tau, P)$ :

- ▶  $c$  is a context in  $X$ ,
- ▶  $\tau \in \mathcal{DSE}$ , and
- ▶  $P = (p_1, \dots, p_k)$ , where  $k \geq 1$  and each  $p_i \in \mathcal{DSE}$

Notation:  $\tau$  is the *target* path and each  $p_i$  is a *key* path

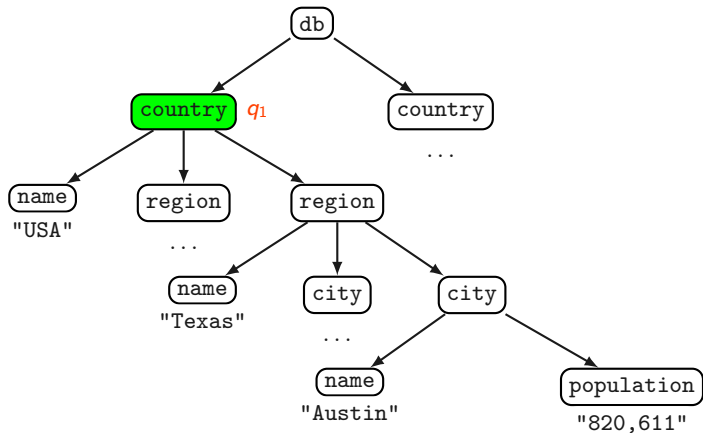
# XSD Keys: Some examples

`((country, q1), .//city, (./name))`



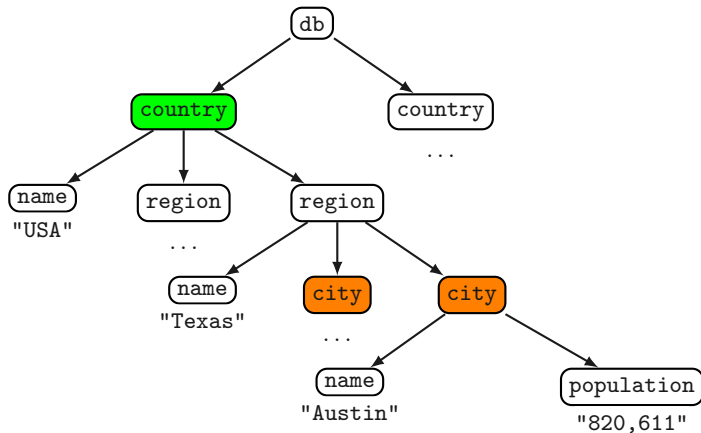
# XSD Keys: Some examples

`((country, q1), .//city, (./name))`



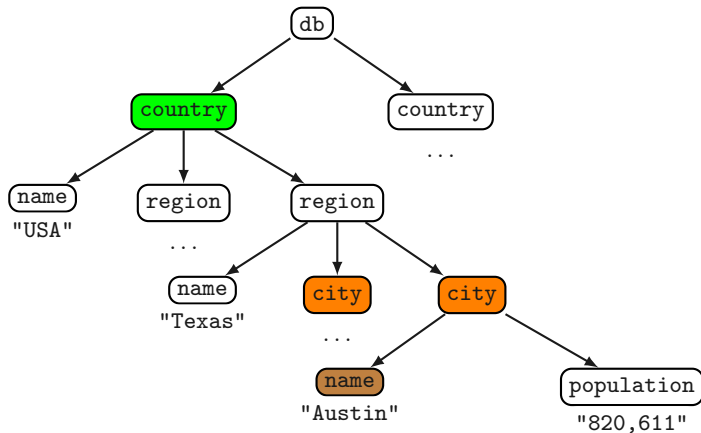
# XSD Keys: Some examples

`((country, q1), .//city, (./name))`



# XSD Keys: Some examples

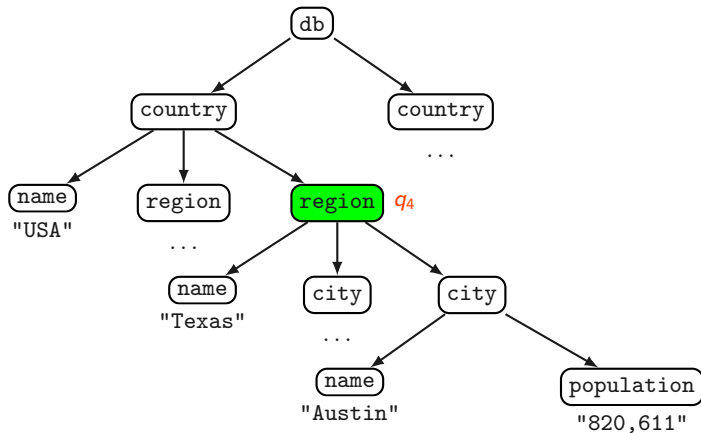
`((country, q1), .//city, (./name))`





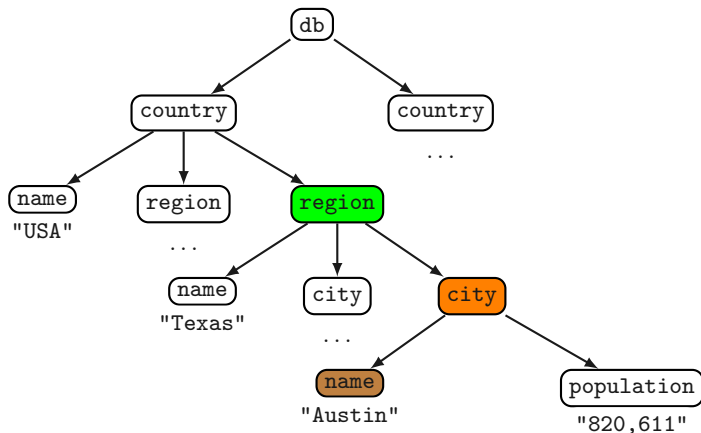
# XSD Keys: Some examples

`((region, q4), ../city, (./name))`



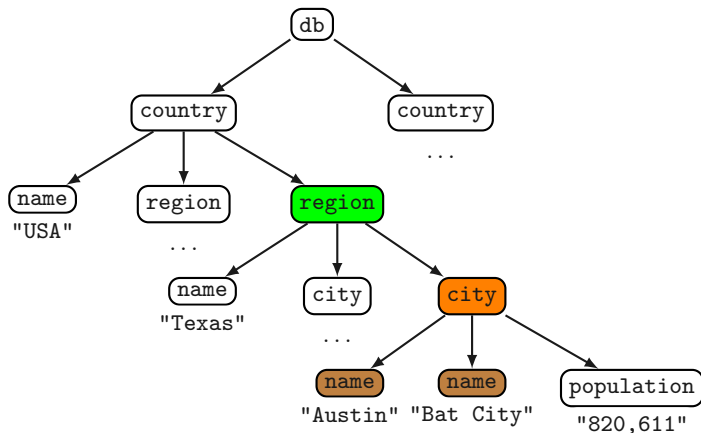
# Satisfaction of XSD keys: Structural requirements

`((region, q4), ../city, (./name))`



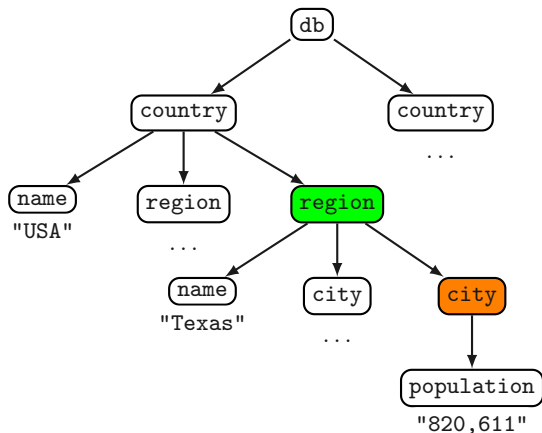
# Satisfaction of XSD keys: Structural requirements

`((region, q4), ../city, (./name))`



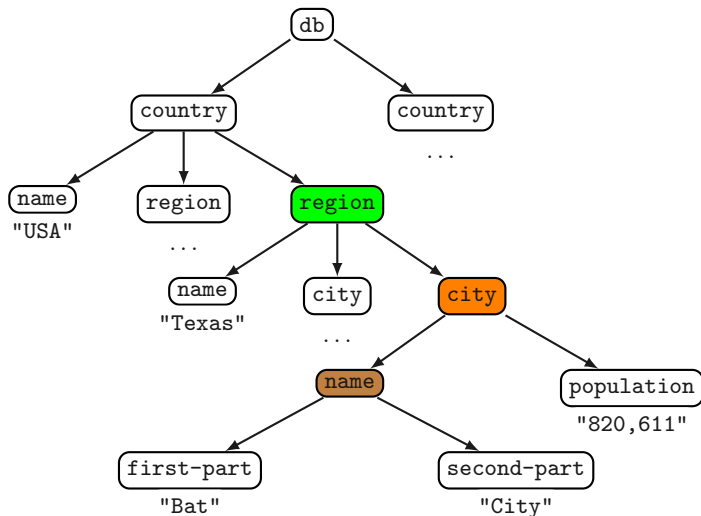
# Satisfaction of XSD keys: Structural requirements

$((\text{region}, q_4), \text{.//city}, (\text{./name}))$



# Satisfaction of XSD keys: Structural requirements

`((region, q4), ../city, (./name))`



# Satisfaction of XSD keys: Formalizing the structural requirements

$v$  is a **Data**-node in an XML tree  $t$  if  $v$  has a single child  $u$ , and  $\text{lab}_t(u) \in \mathbf{Data}$ .

- ▶ Then  $\text{value}_t(v) = \text{lab}_t(u)$

# Satisfaction of XSD keys: Formalizing the structural requirements

$v$  is a **Data**-node in an XML tree  $t$  if  $v$  has a single child  $u$ , and  $\text{lab}_t(u) \in \mathbf{Data}$ .

- ▶ Then  $\text{value}_t(v) = \text{lab}_t(u)$

## Definition

A key  $\phi = (c, \tau, (p_1, \dots, p_k))$  qualifies in an XML tree  $t$  if for every  $v \in \text{CNodes}_t(c)$ , every  $u \in \tau(t, v)$  and every  $p_i$ :

$p_i(t, u)$  is a singleton containing a **Data**-node

# Satisfaction of XSD keys: Complete definition

Let  $u$  be a node in an XML tree  $t$ , and  $P = (p_1, \dots, p_k)$  a sequence of disjunctions of selector expressions.

Assume that  $p_i(t, u) = \{u_i\}$  with  $u_i$  a **Data**-node, for every  $i \in \{1, \dots, k\}$ .

- ▶ Then  $\text{record}_P(t, u) = [\text{value}_t(u_1), \dots, \text{value}_t(u_k)]$



# Satisfaction of XSD keys: Complete definition

Let  $u$  be a node in an XML tree  $t$ , and  $P = (p_1, \dots, p_k)$  a sequence of disjunctions of selector expressions.

Assume that  $p_i(t, u) = \{u_i\}$  with  $u_i$  a **Data**-node, for every  $i \in \{1, \dots, k\}$ .

- ▶ Then  $\text{record}_P(t, u) = [\text{value}_t(u_1), \dots, \text{value}_t(u_k)]$

## Definition

An XML tree  $t$  satisfies a key  $\phi = (c, \tau, P)$ , denoted by  $t \models \phi$ , if:

1.  $\phi$  qualifies in  $t$
2. for every  $v \in \text{CNodes}_t(c)$  and every pair  $u_1, u_2$  of distinct nodes in  $\tau(t, v)$ , it holds that  $\text{record}_P(t, u_1) \neq \text{record}_P(t, u_2)$

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

# XSD key mining problem: Definition

First, we need to define the support of a key in an XML tree.

Assume that  $\phi = (c, \tau, P)$ . Given an XML tree  $t$ :

$$\text{TNodes}_t(\phi) = \bigcup_{v \in \text{CNodes}_t(c)} \tau(t, v)$$

Then  $\text{supp}(\phi, t) = |\text{TNodes}_t(\phi)|$

- ▶  $\text{supp}(c, \tau, t)$  is also used to denote  $\text{supp}(\phi, t)$

## XSD key mining problem: Definition (cont'd)

Second, we need to introduce a fundamental quality requirement:

### Definition

A key  $\phi$  is consistent w.r.t. an XSD  $X$  if  $\phi$  qualifies in every XML tree conforming to  $X$

# XSD key mining problem: Definition (cont'd)

Second, we need to introduce a fundamental quality requirement:

## Definition

A key  $\phi$  is consistent w.r.t. an XSD  $X$  if  $\phi$  qualifies in every XML tree conforming to  $X$

Now we can formally define the XSD mining problem:

## Definition

Given an XSD  $X$ , and XML tree  $t$  conforming to  $X$  and a threshold  $N$ , the XSD key mining problem consists of finding all keys  $\phi$  such that:

- ▶  $t \models \phi$
- ▶  $\text{supp}(\phi, t) > N$
- ▶  $\phi$  is consistent w.r.t.  $X$

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

# First problem to solve: The consistency problem

$\text{CONSISTENCY}(\mathcal{DSE})$  is the problem of verifying, given an XSD  $X$  and a key  $\phi$ , whether  $\phi$  is consistent w.r.t.  $X$ .

- ▶  $\text{CONSISTENCY}(\mathcal{SE})$  is defined analogously, but assuming that target and key paths are selector expressions

To determine the complexity of these problems, we study a closely related problem.

# Cardinality of selector expression results

Notation:

- ▶  $\mathcal{L}(X)$ : set of trees  $t$  conforming to XSD  $X$
- ▶  $\tau(t)$ : defined as  $\tau(t, r)$ , where  $r$  is the root of tree  $t$



# Cardinality of selector expression results

Notation:

- ▶  $\mathcal{L}(X)$ : set of trees  $t$  conforming to XSD  $X$
- ▶  $\tau(t)$ : defined as  $\tau(t, r)$ , where  $r$  is the root of tree  $t$

## Definition

Given  $\bullet \in \{<, =, >\}$  and  $k \geq 0$ :

$$\forall_{\text{tree}}^{\bullet k, \mathcal{SE}} = \{(X, \rho) \mid X \text{ is an XSD,} \\ \rho \in \mathcal{SE} \text{ and for every } t \in \mathcal{L}(X) : |\rho(t)| \bullet k\}$$

$\forall_{\text{tree}}^{\bullet k, \mathcal{DSE}}$  is defined analogously.

# Cardinality of selector expression results and the consistency problem

## Proposition

- ▶  $\forall_{\text{tree}}^{\leq 1, \mathcal{SE}}$  is polynomially equivalent to  $\text{CONSISTENCY}(\mathcal{SE})$
- ▶  $\forall_{\text{tree}}^{\leq 1, \mathcal{DSE}}$  is polynomially equivalent to  $\text{CONSISTENCY}(\mathcal{DSE})$

We study the more general family of problems  $\forall_{\text{tree}}^{\bullet k, \mathcal{SE}}$

# A useful detour: The string case

Notation:

- ▶  $\mathcal{L}(A)$ : set of strings  $s$  accepted by DFA  $A$
- ▶  $\tau(s)$ : defined as before but considering string  $s$  as a tree

## Definition

Given  $\bullet \in \{<, =, >\}$  and  $k \geq 0$ :

$$\forall_{\text{string}}^{\bullet, k, \mathcal{SE}} = \{(A, p) \mid A \text{ is a DFA, } p \in \mathcal{SE} \text{ and for every } s \in \mathcal{L}(A) : |p(s)| \bullet k\}$$

The complexity of  $\forall_{\text{string}}^{<2, \mathcal{SE}}$

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

The complexity of  $\forall_{\text{string}}^{<2, \mathcal{SE}}$

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

The complexity of  $\forall_{\text{string}}^{<2, \mathcal{SE}}$

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

## Proposition

$\forall_{\text{string}}^{<2, \mathcal{SE}}$  is in PTIME

# The complexity of $\forall_{\text{string}}^{<2, \mathcal{SE}}$

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

## Proposition

$\forall_{\text{string}}^{<2, \mathcal{SE}}$  is in PTIME

**Proof idea:** Assume given a DFA  $A$  and selector expression  $p$ .

Construct an NFA  $B_p$  accepting every string  $s$  such that  $|p(s)| \geq 2$

Verify whether  $\mathcal{L}(A \times B_p) = \emptyset$

- ▶  $(A, p) \in \forall_{\text{string}}^{<2, \mathcal{SE}}$  if and only if  $\mathcal{L}(A \times B_p) = \emptyset$



The complexity of  $\forall_{\text{string}}^{>0, \mathcal{SE}}$ : A more difficult problem

### Proposition

$\forall_{\text{string}}^{>0, \mathcal{SE}}$  is PSPACE-complete



# The complexity of $\forall_{\text{string}}^{=k, \mathcal{SE}}$ : A more surprising result

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

- ▶ We conclude that  $\forall_{\text{string}}^{=1, \mathcal{SE}}$  is in PSPACE

# The complexity of $\forall_{\text{string}}^{=k, \mathcal{SE}}$ : A more surprising result

We have that  $\forall_{\text{string}}^{=1, \mathcal{SE}} = \forall_{\text{string}}^{<2, \mathcal{SE}} \cap \forall_{\text{string}}^{>0, \mathcal{SE}}$

- ▶ We conclude that  $\forall_{\text{string}}^{=1, \mathcal{SE}}$  is in PSPACE

But a lot more can be proved:

## Theorem

$\forall_{\text{string}}^{=1, \mathcal{SE}}$  is in PTIME

# $\forall_{\text{string}}^{=1, \mathcal{SE}}$ is in PTIME: Proof idea

Assume given a DFA  $A$  and selector expression  $p$ .

The following polynomial-time algorithm verifies whether  $(A, p) \in \forall_{\text{string}}^{=1, \mathcal{SE}}$ :

1. Check whether  $(A, p) \in \forall_{\text{string}}^{<2, \mathcal{SE}}$ . If this condition holds, then go to step 2. Otherwise return **no**.
2. Construct an NFA  $C_p$  such that:
  - ▶  $s \in \mathcal{L}(C_p)$  iff  $|p(s)| \geq 1$
  - ▶ the number of accepting runs of  $C_p$  on  $s$  is equal to  $|p(s)|$
3. Compute  $A \times C_p$

## $\forall_{\text{string}}^{=1, \mathcal{SE}}$ is in PTIME: Proof idea (cont'd)

4. Check whether  $\mathcal{L}(A) \subseteq \mathcal{L}(A \times C_p)$ . If this condition holds, then return **yes**, otherwise return **no**
  - ▶ This is equivalent to verifying whether  $\mathcal{L}(A) \subseteq \mathcal{L}(C_p)$

# $\forall_{\text{string}}^{=1, \mathcal{SE}}$ is in PTIME: Proof idea (cont'd)

4. Check whether  $\mathcal{L}(A) \subseteq \mathcal{L}(A \times C_p)$ . If this condition holds, then return **yes**, otherwise return **no**
  - ▶ This is equivalent to verifying whether  $\mathcal{L}(A) \subseteq \mathcal{L}(C_p)$

Key observations:

- ▶ Containment problem for unambiguous finite automata (UFAs) can be solved in polynomial time [SH85].
  - ▶ An NFA  $D$  is unambiguous if for every  $s \in \mathcal{L}(D)$ , there is only one accepting run of  $D$  on  $s$
- ▶  $A$  and  $A \times C_p$  are UFAs
  - ▶ Given that  $A$  is a DFA,  $(A, p) \in \forall_{\text{string}}^{<2, \mathcal{SE}}$  and the number of accepting runs of  $C_p$  on a string  $s$  is equal to  $|p(s)|$



# String case: Summary of results

| $\mathcal{P}$       | $\forall_{\text{string}}^{>k, \mathcal{P}}$ | $\forall_{\text{string}}^{<k, \mathcal{P}}$ | $\forall_{\text{string}}^{=k, \mathcal{P}} (k \geq 1)$ |
|---------------------|---|---|--|
| $\mathcal{RE}$      | PSPACE-complete                             | in PTIME                                    | PSPACE-complete  |
| $\mathcal{DSE}$     | PSPACE-complete                             | in PTIME                                    | CONP-complete  |
| $\mathcal{SE}$      | PSPACE-complete                             | in PTIME                                    | in PTIME   |
| $\mathcal{SE}^*$    | in PTIME                                    | in PTIME                                    | in PTIME   |
| $\mathcal{SE}^{//}$ | in PTIME                                    | in PTIME                                    | in PTIME   |

Notation:

$\mathcal{RE}$ : regular expressions

$\mathcal{SE}^*$ : selector expressions not using \*

$\mathcal{SE}^{//}$ : selector expressions not using //

# Tree case: Summary of results

| $\mathcal{P}$       | $\forall_{\text{tree}}^{>k, \mathcal{P}}$ | $\forall_{\text{tree}}^{<k, \mathcal{P}}$ | $\forall_{\text{tree}}^{=k, \mathcal{P}} \ (k \geq 1)$ |
|---------------------|---|---|--|
| $\mathcal{RE}$      | EXPTIME-complete                          | in PTIME                                  | in EXPTIME<br>PSPACE-hard                              |
| $\mathcal{DSE}$     | EXPTIME-complete                          | in PTIME                                  | in EXPTIME<br>coNP-hard                                |
| $\mathcal{SE}$      | EXPTIME-complete                          | in PTIME                                  | in PTIME   |
| $\mathcal{SE}^*$    | in EXPTIME                                | in PTIME                                  | in PTIME   |
| $\mathcal{SE}^{//}$ | in PTIME                                  | in PTIME                                  | in PTIME   |

# The consistency problem: Main results

## Theorem

- ▶  $\text{CONSISTENCY}(\mathcal{SE})$  is in PTIME
- ▶  $\text{CONSISTENCY}(\mathcal{DSE})$  is CONP-hard and in EXPTIME



# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

# An XSD key mining algorithm

Input: XSD  $X$ , a tree  $t$  conforming to  $X$  and a threshold  $N$

Algorithm:

```
for all  $c \in \text{ContextMiner}_{t,X}$  do  
  for all  $\tau \in \text{TargetPathMiner}_{t,X,N}(c)$  do  
     $S := \text{OneKeyPathMiner}_{t,X}(c, \tau)$   
     $\mathcal{P} := \text{MinimalKeyPathSetMiner}_{t,X}(c, \tau, S)$   
    for each  $P \in \mathcal{P}$  return  $(c, \tau, P)$ 
```

# An XSD key mining algorithm

Input: XSD  $X$ , a tree  $t$  conforming to  $X$  and a threshold  $N$

Algorithm:

```
for all  $c \in \text{ContextMiner}_{t,X}$  do  
  for all  $\tau \in \text{TargetPathMiner}_{t,X,N}(c)$  do  
     $S := \text{OneKeyPathMiner}_{t,X}(c, \tau)$   
     $\mathcal{P} := \text{MinimalKeyPathSetMiner}_{t,X}(c, \tau, S)$   
    for each  $P \in \mathcal{P}$  return  $(c, \tau, P)$ 
```

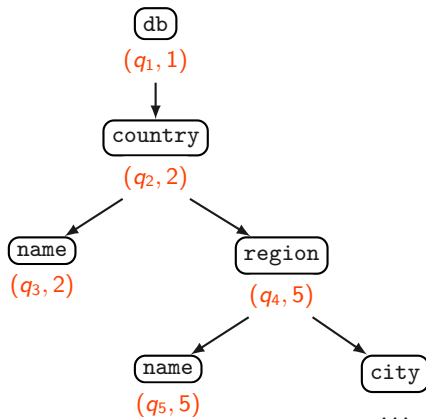
Restrictions:

- ▶ We focus on selector expressions (we disregard the union operator)
- ▶ We consider target and key paths up to a given length  $k_{\max}$  (which can be at most the depth of  $t$ )

# ContextMiner<sub>t,X</sub>

It returns a list of possible contexts based on  $X$  and  $t$ .

Prefix tree  $PT(t)$  is used in this procedure: Obtained by collapsing all nodes with the same ancestor string.



# TargetPathMiner<sub>t, X, N</sub>(c)

Given a context node  $c$ , it returns a list of target paths with support greater than  $N$  in  $t$ .

- ▶ It follows a framework of levelwise search [MT97]

Consider the following partial order on selector expressions:

$$\tau_1 \preceq \tau_2 \text{ iff for every tree } t: \tau_2(t) \subseteq \tau_1(t)$$

If  $\tau_1 \preceq \tau_2$ , then  $\tau_2$  is more specific than  $\tau_1$

# Deciding $\preceq$

Consider the following “one-step specialization relation” on selector expressions:

- ▶  $\tau \prec_1 \tau'$  if  $\tau'$  is obtained from  $\tau$  by one of the following operations:
  - ▶ if  $\tau$  starts with the descendant axis, replace it by the child axis
  - ▶ if  $\tau$  starts with the descendant axis, insert a wildcard step right after it
  - ▶ replacing a wildcard with an element name

# Deciding $\preceq$

Consider the following “one-step specialization relation” on selector expressions:

- ▶  $\tau \prec_1 \tau'$  if  $\tau'$  is obtained from  $\tau$  by one of the following operations:
  - ▶ if  $\tau$  starts with the descendant axis, replace it by the child axis
  - ▶ if  $\tau$  starts with the descendant axis, insert a wildcard step right after it
  - ▶ replacing a wildcard with an element name

We can decide whether  $\tau \preceq \tau'$  by using relation  $\prec_1$ :

## Proposition

$\preceq$  is equal to the reflexive and transitive closure of  $\prec_1$ .

# TargetPathMiner<sub>t,X,N</sub>(c) (cont'd)

Let  $U$  be the set of all selector expressions of length at most  $k_{\max}$ .

TargetPathMiner<sub>t,X,N</sub>(c) uses partial order  $\preceq$ :

$C_0 := \{./\}$

$i := 0$

**while**  $C_i \neq \emptyset$  **do**

$F_i := \{\tau \in C_i \mid \text{supp}(c, \tau, t) > N\}$

$C_{i+1} := \{\tau \in U \mid \forall \tau' \in U : \tau' \prec \tau \Rightarrow \tau' \in \bigcup_{j \leq i} F_j\} \setminus \bigcup_{j \leq i} C_j$

$i := i + 1$

**return**  $\bigcup_{j \leq i} F_j$



# TargetPathMiner<sub>t,X,N</sub>(c) (cont'd)

Let  $U$  be the set of all selector expressions of length at most  $k_{\max}$ .

TargetPathMiner<sub>t,X,N</sub>(c) uses partial order  $\preceq$ :

$C_0 := \{./\}$

$i := 0$

**while**  $C_i \neq \emptyset$  **do**

$F_i := \{\tau \in C_i \mid \text{supp}(c, \tau, t) > N\}$

$C_{i+1} := \{\tau \in U \mid \forall \tau' \in U : \tau' \prec \tau \Rightarrow \tau' \in \bigcup_{j \leq i} F_j\} \setminus \bigcup_{j \leq i} C_j$

$i := i + 1$

**return**  $\bigcup_{j \leq i} F_j$

Some remarks:

# TargetPathMiner<sub>t,X,N</sub>(c) (cont'd)

Let  $U$  be the set of all selector expressions of length at most  $k_{\max}$ .

TargetPathMiner<sub>t,X,N</sub>(c) uses partial order  $\preceq$ :

$C_0 := \{./\}$

$i := 0$

**while**  $C_i \neq \emptyset$  **do**

$F_i := \{\tau \in C_i \mid \text{supp}(c, \tau, t) > N\}$

$C_{i+1} := \{\tau \in U \mid \forall \tau' \in U : \tau' \prec \tau \Rightarrow \tau' \in \bigcup_{j \leq i} F_j\} \setminus \bigcup_{j \leq i} C_j$

$i := i + 1$

**return**  $\bigcup_{j \leq i} F_j$

Some remarks:

- ▶ Prefix tree  $PT(t)$  is used when checking  $\text{supp}(c, \tau, t) > N$

# TargetPathMiner<sub>t,X,N</sub>(c) (cont'd)

Let  $U$  be the set of all selector expressions of length at most  $k_{\max}$ .

TargetPathMiner<sub>t,X,N</sub>(c) uses partial order  $\preceq$ :

$C_0 := \{./\}$

$i := 0$

**while**  $C_i \neq \emptyset$  **do**

$F_i := \{\tau \in C_i \mid \text{supp}(c, \tau, t) > N\}$

$C_{i+1} := \{\tau \in U \mid \forall \tau' \in U : \tau' \prec \tau \Rightarrow \tau' \in \bigcup_{j \leq i} F_j\} \setminus \bigcup_{j \leq i} C_j$

$i := i + 1$

**return**  $\bigcup_{j \leq i} F_j$

Some remarks:

- ▶ Prefix tree  $PT(t)$  is used when checking  $\text{supp}(c, \tau, t) > N$
- ▶ For more specific results:  $\bigcup_{j \leq i} F_j$  is replaced by  $F_i$

# TargetPathMiner<sub>t,X,N</sub>(c) (cont'd)

Let  $U$  be the set of all selector expressions of length at most  $k_{\max}$ .

TargetPathMiner<sub>t,X,N</sub>(c) uses partial order  $\preceq$ :

$C_0 := \{./\}$

$i := 0$

**while**  $C_i \neq \emptyset$  **do**

$F_i := \{\tau \in C_i \mid \text{supp}(c, \tau, t) > N\}$

$C_{i+1} := \{\tau \in U \mid \forall \tau' \in U : \tau' \prec \tau \Rightarrow \tau' \in \bigcup_{j \leq i} F_j\} \setminus \bigcup_{j \leq i} C_j$

$i := i + 1$

**return**  $\bigcup_{j \leq i} F_j$

Some remarks:

- ▶ Prefix tree  $PT(t)$  is used when checking  $\text{supp}(c, \tau, t) > N$
- ▶ For more specific results:  $\bigcup_{j \leq i} F_j$  is replaced by  $F_i$
- ▶ A more efficient implementation uses  $\prec_1$  instead of  $\preceq$

# OneKeyPathMiner <sub>$t, X$</sub> ( $c, \tau$ )

It returns the set of all key paths  $p$  of length at most  $k_{\max}$  for which  $(c, \tau, (p))$  is consistent with  $X$

It uses the consistency algorithm presented in the previous slides.

- ▶ To reduce the number of consistency tests, it only considers the candidates  $p$  for which  $(c, \tau, (p))$  qualifies in  $t$

# MinimalKeyPathSetMiner<sub>t, X</sub>(c, τ, S)

It returns a set  $\mathcal{P}$  of minimal subsets  $P$  of  $S$  for which  
 $t \models (c, \tau, P)$

- ▶ We capitalize on existing relational techniques for mining functional dependencies

Assume that  $S = \{p_1, \dots, p_k\}$ , and let  $R$  be a relation such that:

- ▶ The schema of  $R$  is  $(CID, TID, p_1, \dots, p_k)$
- ▶  $(v, u, a_1, \dots, a_k) \in R$  iff  $v \in CNodes_t(c)$ ,  $u \in \tau(t, v)$  and  $record_S(t, u) = [a_1, \dots, a_k]$

## MinimalKeyPathSetMiner<sub>t,X</sub>(c, τ, S) (cont'd)

For every  $P = \{p_{i_1}, \dots, p_{i_n}\}$ , with  $1 \leq i_1 < \dots < i_n \leq k$ , we have that:

$$t \models (c, \tau, (p_{i_1}, \dots, p_{i_n}))$$

iff

$$\text{CID}, p_{i_1}, \dots, p_{i_n} \rightarrow \text{TID holds in } R$$

Thus, we can now plug in any existing functional dependency discovery algorithm.

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message



# Experimental evaluation: Setting

We use a corpus of 90 XML trees and associated XSDs from the University of Amsterdam XML Web Collection.

- ▶ Maximal and average number of elements occurring in trees is 91K and 5K, respectively
- ▶ Maximal and average number of elements occurring in XSDs is 532 and 52, respectively

We search for keys with target path length at most 4 and key path length at most 2.

## Experimental evaluation: A few words

- ▶ In  $\text{ContextMiner}_{t,X}$ , the use of prefix trees allows to reduce the number of contexts: **52% do not need to be considered**
- ▶ In  $\text{TargetPathMiner}_{t,X,N}(c)$ , the framework of levelwise search allows a significant reduction in the number of target paths considered:

|               |                      |
|---------------|----------------------|
| possible TPs  | $2.4 \times 10^{11}$ |
| candidate TPs | $6.7 \times 10^6$    |
| supported TPs | $8.4 \times 10^4$    |

## Experimental evaluation: A few words (cont'd)

- ▶ In  $\text{OneKeyPathMiner}_{t,X}(c, \tau)$ , start by checking consistency on trees allows a significant reduction on the number of expensive consistency tests:

|                        |       |
|------------------------|-------|
| candidate KPs          | 48144 |
| consistent KPs on tree | 484   |
| consistent KPs on XSD  | 288   |

- ▶ Resulting keys from  $\text{MinimalKeyPathSetMiner}_{t,X}(c, \tau, S)$ :

| consistency test             | no  |     | yes |     |
|------------------------------|-----|-----|-----|-----|
| support                      | 10  | 100 | 10  | 100 |
| derived keys                 | 107 | 54  | 43  | 16  |
| trees with keys              | 27  | 15  | 19  | 10  |
| average nr. of keys per tree | 4   | 3.6 | 2.2 | 1.6 |
| max nr. of keys per tree     | 23  | 23  | 9   | 4   |

# Outline

- ▶ A bit of notation: XML trees and XSDs
- ▶ XSD keys: Formal definition
- ▶ XSD key mining problem: Formal definition
- ▶ The consistency problem: Complexity
- ▶ An XSD key mining algorithm
- ▶ Experimental evaluation: A few words
- ▶ Take-home message

# Take-home message

- ▶ We consider the problem of discovering XSD keys from a given XML tree that conforms to a given XSD.
- ▶ We consider the semantics for XSD keys proposed by the W3C, and we incorporate in our discovering algorithm some quality criteria like consistency, which was shown to be decidable in polynomial time.

**Thank you!**

# Bibliography

- [MT97] H. Mannila, H. Toivonen: Levelwise Search and Borders of Theories in Knowledge Discovery. *Data Min. Knowl. Discov.* 1(3): 241–258, 1997
- [S90] H. Seidl: Deciding Equivalence of Finite Tree Automata. *SIAM J. Comput.* 19(3):424–437, 1990
- [SH85] R. E. Stearns, H. B. Hunt III: On the Equivalence and Containment Problems for Unambiguous Regular Expressions, Regular Grammars and Finite Automata. *SIAM J. Comput.* 14(3):598–611, 1985

# Backup slides



# The complexity of $\forall_{\text{string}}^{>0, \mathcal{SE}}$ : A more difficult problem

We show how to reduce  $\overline{\text{CNF-SAT}}$  to  $\forall_{\text{string}}^{>0, \mathcal{SE}}$ .

Let  $\varphi$  be a propositional formula in CNF.

- ▶ We define a DFA  $A$  and a selector expression  $p$  such that  $\varphi$  is not satisfiable if and only if  $(A, p) \in \forall_{\text{string}}^{>0, \mathcal{SE}}$

Assume that  $\varphi = C_1 \wedge C_2 \wedge \dots \wedge C_m$  and that  $\varphi$  mentions variables  $x_1, \dots, x_n$ .

- ▶ For each  $x_i$ : strings 01, 10 are used to encode values *false* and *true*, respectively

# The complexity of $\forall_{\text{string}}^{>0, \mathcal{SE}}$ : A more difficult problem (cont'd)

A string  $w$  is accepted by  $A$  iff  $w = s_1 \# s_2 \# \cdots \# s_m$ , where

- ▶  $s_j$  is a string of length  $2n$  representing a satisfying assignment for clause  $C_j$

## Example

If  $\varphi = (x_1 \vee \neg x_2 \vee x_3) \wedge (x_1 \vee \neg x_3 \vee \neg x_4)$ , then the following strings are accepted by  $A$ :

10010101#10010101

01011001#10100101

# The complexity of $\forall_{\text{string}}^{>0, \mathcal{SE}}$ : A more difficult problem

Assignments  $s_i, s_j$  ( $i \neq j$ ) can be inconsistent.

- ▶ Inconsistent bits will be at distance  $2n$ . In particular, one will find a 0 followed by  $2n$  symbols, and then followed by a 1

Thus, we consider the following selector expression:

$$p = .//0 / \underbrace{* / \dots / *}_{2n \text{ times}} / 1$$



# The complexity of $\forall_{\text{string}}^{=k, \mathcal{SE}}$ : The general case

## Theorem

$\forall_{\text{string}}^{=k, \mathcal{SE}}$  is in PTIME (for every  $k \geq 0$ )

# The complexity of $\forall_{\text{string}}^{=k, \mathcal{SE}}$ : The general case

## Theorem

$\forall_{\text{string}}^{=k, \mathcal{SE}}$  is in PTIME (for every  $k \geq 0$ )

**Proof idea:**  $\forall_{\text{string}}^{=0, \mathcal{SE}} = \forall_{\text{string}}^{<1, \mathcal{SE}}$ , so let's consider the case  $k \geq 1$

Polynomial-time algorithm to verify whether  $(A, p) \in \forall_{\text{string}}^{=k, \mathcal{SE}}$ :

1. Check whether  $(A, p) \in \forall_{\text{string}}^{<k+1, \mathcal{SE}}$ . If this condition holds, then go to step 2. Otherwise return **no**.
2. Construct an NFA  $B_p^k$  such that:
  - ▶  $s \in \mathcal{L}(B_p^k)$  iff  $|p(s)| \geq k$
  - ▶ if  $s \in \mathcal{L}(B_p^k)$ , then the number of accepting runs of  $B_p^k$  on  $s$  is

$$\frac{|p(s)|!}{(|p(s)| - k)!}$$

# The complexity of $\forall_{\text{string}}^{=k, \mathcal{SE}}$ : The general case (cont'd)

3. Compute  $A \times B_p^k$
4. Check whether  $\mathcal{L}(A) \subseteq \mathcal{L}(A \times B_p^k)$ . If this condition holds, then return **yes**, otherwise return **no**

Key observations:

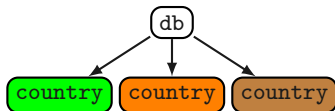
- ▶ Containment problem for  $c$ -unambiguous finite automata ( $c$ -UFAs) can be solved in polynomial time [SH85].
- ▶  $A \times B_p$  is a  $k!$ -UFA:
  - ▶ If  $(A, p) \in \forall_{\text{string}}^{<k+1, \mathcal{SE}}$  and  $s \in \mathcal{L}(A \times B_p^k)$ , then  $|p(s)| = k$  and the number of accepting runs of  $A \times B_p^k$  on  $s$  is bounded by:

$$\frac{|p(s)|!}{(|p(s)| - k)!} = \frac{k!}{(k - k)!} = k!$$

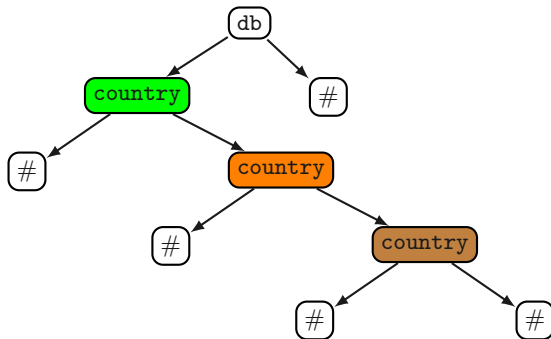


# A polynomial-time algorithm for $\forall_{tree}^{=1, \mathcal{SE}}$ : Binary representations

Tree  $t$ :



Binary representation  $fcns(t)$  of  $t$ :



# Binary representations: Main results

We use binary top-down tree automata (BTA) to process binary trees.

- ▶ There exists a deterministic BTA  $A_{\#}$  such that for every binary tree  $t'$ :  $t' \in \mathcal{L}(A_{\#})$  iff there exists a tree  $t$  for which  $t' = \text{fcns}(t)$

The following problems can be solved in polynomial-time:

- ▶ Given an XSD  $X$ , construct a deterministic BTA  $A_X$  such that for every tree  $t$ :  $\text{fcns}(t) \in \mathcal{L}(A_X)$  if and only if  $t \in \mathcal{L}(X)$
- ▶ Given a selector expression  $p$ , construct a BTA  $B_p$  such that for every tree  $t$ :
  - ▶  $\text{fcns}(t) \in \mathcal{L}(B_p)$  if and only if  $|p(t)| \geq 1$
  - ▶ the number of accepting runs of  $B_p$  on  $\text{fcns}(t)$  is equal to  $|p(t)|$



## A polynomial-time algorithm for $\forall_{\text{tree}}^{=1, \mathcal{SE}}$ : Putting all together

Assume given an XSD  $X$  and selector expression  $p$ .

The following polynomial-time algorithm verifies whether  $(X, p) \in \forall_{\text{tree}}^{=1, \mathcal{SE}}$ :

1. Check whether  $(X, p) \in \forall_{\text{tree}}^{<2, \mathcal{SE}}$ . If this condition holds, then go to step 2. Otherwise return **no**.
2. Construct a BTA  $B_p$  such that for every tree  $t$ :
  - ▶  $\text{fcns}(t) \in \mathcal{L}(B_p)$  iff  $|p(t)| \geq 1$
  - ▶ the number of accepting runs of  $B_p$  on  $\text{fcns}(t)$  is equal to  $|p(t)|$
3. Compute  $A_{\#} \times A_X$  and  $A_{\#} \times A_X \times B_p$

## A polynomial-time algorithm for $\forall_{\text{tree}}^{=1, \mathcal{SE}}$ : Putting all together

4. Check whether  $\mathcal{L}(A_{\#} \times A_X) \subseteq \mathcal{L}(A_{\#} \times A_X \times B_p)$ . If this condition holds, then return **yes**, otherwise return **no**

Key observations:

- ▶ Containment problem for unambiguous BTAs can be solved in polynomial time [S90].
- ▶  $A_{\#} \times A_X$  and  $A_{\#} \times A_X \times B_p$  are unambiguous BTAs
  - ▶ Given that  $A_{\#}, A_X$  are deterministic,  $(X, \rho) \in \forall_{\text{string}}^{<2, \mathcal{SE}}$  and the number of accepting runs of  $B_p$  on a tree  $t' = \text{fcns}(t)$  is equal to  $|\rho(t)|$

