

# Contrôle de version incertain dans l'édition collaborative ouverte de documents arborescents

M. Lamine BA, Talel Abdessalem & Pierre Senellart



BDA'13 – 22-25 October 2013 (Nantes, France)

<http://dbweb.enst.fr/>

# Versioning on the Web is Uncertain (I)



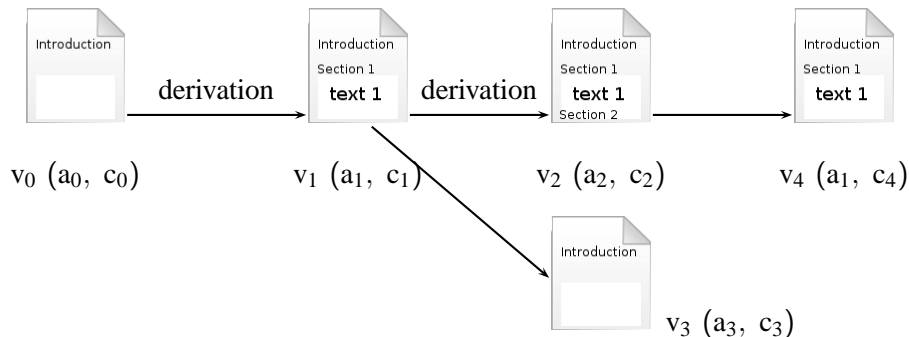
- ▶ Large-scale, open and collaborative editing platforms, e.g., wikipedia
  - ▶ **Unreliable contributors**, Novice vs. Experts, etc.
  - ▶ Malicious edits, Vandalism acts, **Contradictions**

## Versioning on the Web is Uncertain (II)

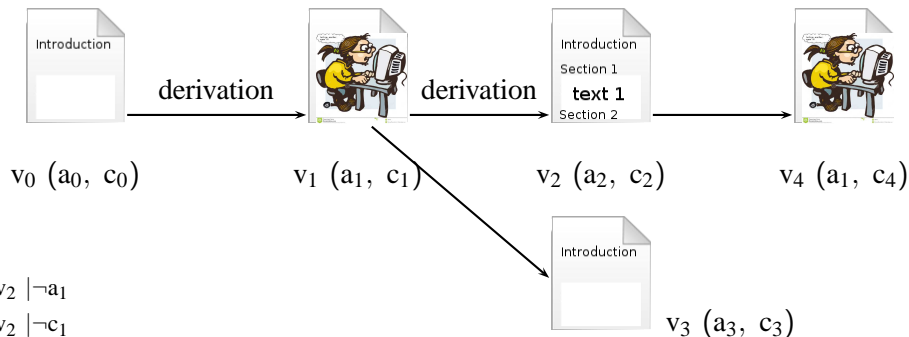
**Version control** is used in large-scale web collaborative editing platforms in order to **integrate** contributions from different sources and to support **fixing errors** with the possibility to query previous **data versions**

- ▶ No notion of more **relevant** versions or contributions which will fit to **user-preferences**, but just the concept of last **valid** revisions
- ▶ **Deterministic** version control models [Kerstin et al.(2009), Al Koc et al.(2011)] in the literature

# Versioning on the Web is Uncertain (III)



# Versioning on the Web is Uncertain (III)

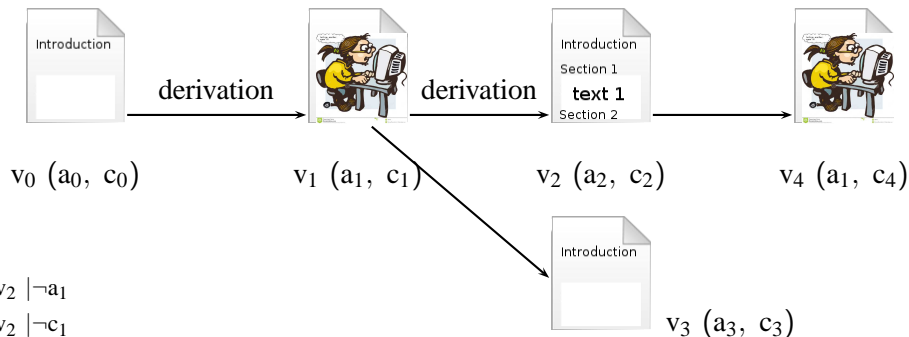


$Q_1 : v_2 \mid \neg a_1$

$Q_2 : v_2 \mid \neg c_1$

$Q_3 : \text{all } v_i \mid \text{Pr}(v_i) \neq 0$

# Versioning on the Web is Uncertain (III)



$Q_1 : v_2 \mid \neg a_1$

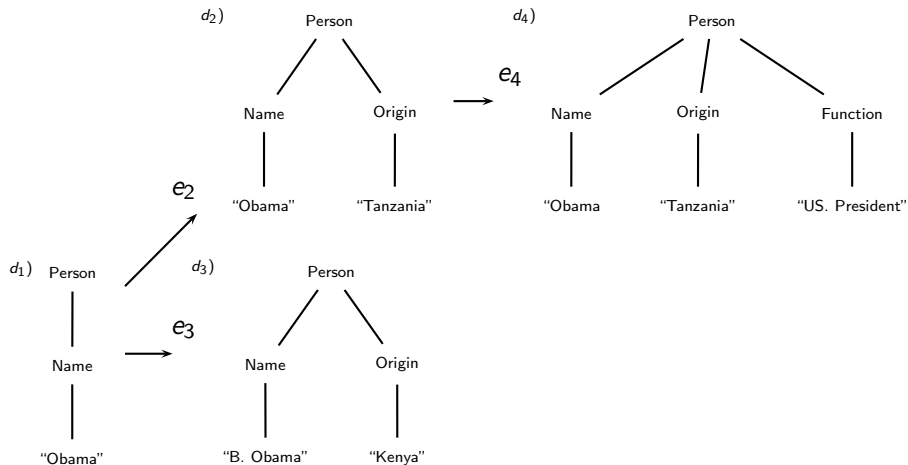
$Q_2 : v_2 \mid \neg c_1$

$Q_3 : \text{all } v_i \mid \text{Pr}(v_i) \neq 0$

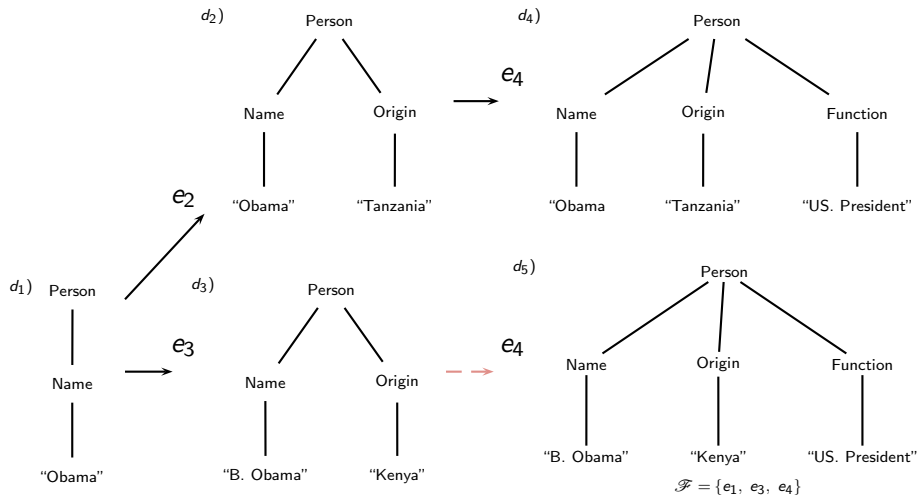


**Uncertain version control in Open Collaborative Contexts**

# Versioning on the Web is Uncertain (IV)

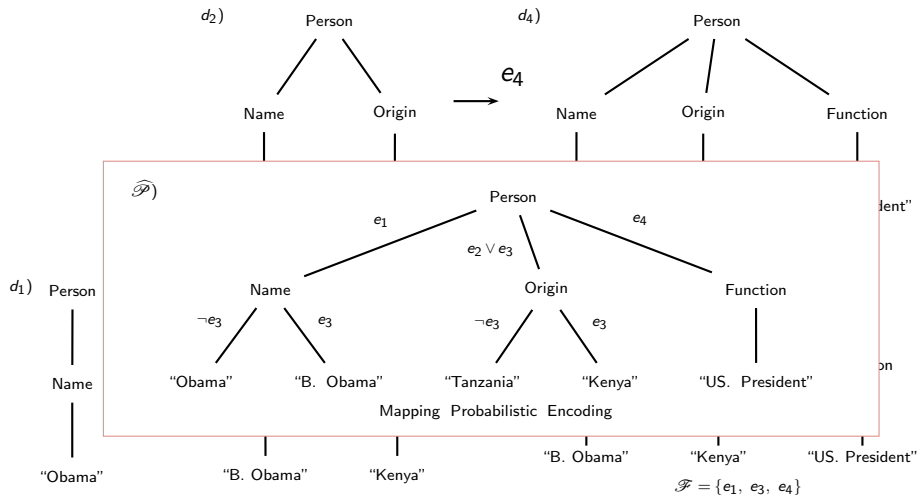


# Versioning on the Web is Uncertain (IV)





## Versioning on the Web is Uncertain (IV)



# Plan

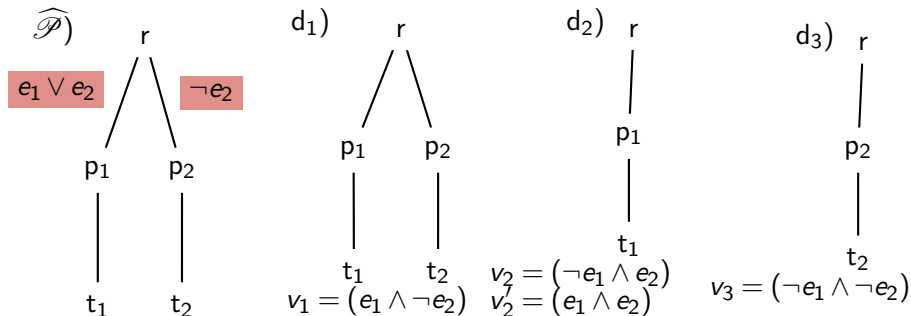
Uncertain Tree-Structured Data Model

Uncertain Multi-Version XML Document

Conclusion and Further work

# Uncertain Tree-Structured Data

Probabilistic XML [Kimelfeld & Senellart.(2013)]



$$\Pr(e_1) = 0.2; \Pr(e_2) = 0.8$$

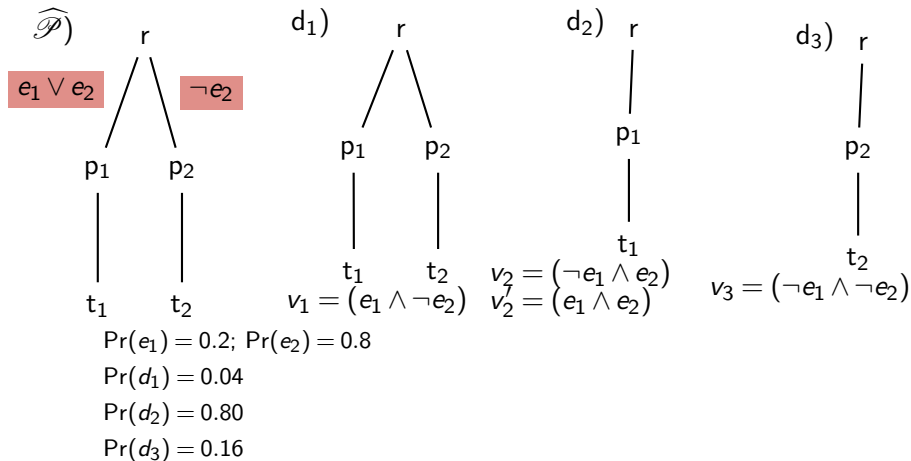
$$\Pr(d_1) = \Pr(e_1) \times \Pr(\neg e_2)$$

$$\Pr(d_2) = (\Pr(\neg e_1) \times \Pr(e_2)) + (\Pr(e_1) \times \Pr(e_2))$$

$$\Pr(d_3) = \Pr(\neg e_1) \times \Pr(\neg e_2)$$

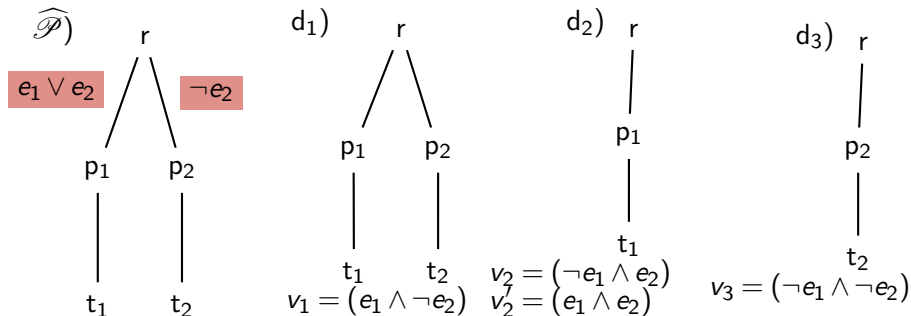
## Uncertain Tree-Structured Data

Probabilistic XML [Kimelfeld &amp; Senellart.(2013)]



# Uncertain Tree-Structured Data

Probabilistic XML [Kimelfeld & Senellart.(2013)]



►  $\widehat{\mathcal{P}} = \langle \mathcal{T}, C(E), \text{fie}, \text{Pr} \rangle$

►  $[[\widehat{\mathcal{P}}]] = \langle D, \text{Pr} \rangle$ ;  $\sum \{\text{Pr}(d) \mid d \in D\} = 1$

# Plan

Uncertain Tree-Structured Data Model

Uncertain Multi-Version XML Document

- Uncertain Version Control Model

- Probabilistic XML Encoding

- Updating Uncertain Multi-version XML documents

- Performance Analysis

Conclusion and Further work

# Uncertain Multi-Version XML Document

## Uncertain Version Control Model

- ▶ **Probability space** (PSV) over Uncertain versions of XML documents
- ▶ **Random** derivation graph (DG) over the document versions produced

**Intuition:** states in DG are complex variables  $e_i$ 's called **version control events** based on simpler variables  $b_1 \dots b_m$  managing uncertainty in data

$\mathcal{D} \langle \mathcal{G}, \Omega \rangle$  defines a multi-version XML document with uncertain data

- ▶  $\mathcal{G}$  is DG over a set of versioning events  $\mathcal{V} \cup e_0$  with  $\mathcal{V} = \{e_1 \dots e_n\}$
- ▶  $\Omega : 2^{\mathcal{V} \setminus \{e_0\}} \rightarrow \mathcal{D}$  a mapping computing the PSV according to sets of valid events

# Uncertain Multi-Version XML Document

## Possible versions and Probabilities

### Possible versions

- ▶ Version control events come with **edit scripts** updating content
- ▶  $\forall i, \forall \mathcal{F} \subseteq 2^{\mathcal{V} \setminus \{e_i\}}$ , the possible version  $\Omega(\{e_i\} \cup \mathcal{F}) = [\Omega(\mathcal{F})]^{\Delta_i}$

### Probability of possible versions

- ▶ Assume a prior probability distribution over simple variables  $b_1 \dots b_m$
- ▶ The probability of a given possible version  $\Omega(\mathcal{F})$  is the probability of  $\bigvee_{\mathcal{F} \subseteq \mathcal{V}, \Omega(\mathcal{F})} \mathcal{F}$



# Uncertain Multi-Version XML Document

## Probabilistic XML Encoding

- ▶ **Compact** representation of all possible versions (PSV) in  $\Omega$  mapping

**Intuition:** Represent PSV using **propositional formulas** of simple variables  $b_1 \dots b_m$  attached to nodes in a global tree  $\mathcal{T}$  containing all provided data

A probabilistic XML encoding of  $\langle \mathcal{G}, \Omega \rangle$  is a couple  $\langle \mathcal{G}, \widehat{\mathcal{P}} \rangle$  with

- ▶  $\widehat{\mathcal{P}} = \langle \mathcal{T}, C(\mathcal{V}), fie, Pr \rangle$  is a probabilistic XML document

**Thm1:**  $[[\widehat{\mathcal{P}}]]$  defines the same probability distribution over  $\mathcal{D}$  as  $\Omega$ , i.e.,

$$\langle \mathcal{G}, [[\widehat{\mathcal{P}}]] \rangle = \langle \mathcal{G}, \Omega \rangle$$

# Uncertain Multi-Version XML Document

## Updating Uncertain Versions (I)

- **Uncertain** editions in  $\Delta$  over nodes of a possible tree version

An update is an **uncertain** version control event defined based on a triple  $\langle e_i, e_j, \Delta \rangle$  ( $e_i \in \mathcal{G}$  and  $e_j \notin \mathcal{G}$ )

$\text{updOP}(\langle e_i, e_j, \Delta \rangle)[\langle \mathcal{G}, \Omega \rangle]$

- $\mathcal{G} := \mathcal{G} \cup (\{e_j\}, \{(e_i, e_j)\})$
- Extension of  $\Omega$  to a  $\Omega'$  mapping with for each  $\mathcal{F} \in 2^{(\mathcal{V} \setminus \{e_0\}) \cup \{e_j\}}$   
 $\Omega'(\mathcal{F}) = [\Omega(\mathcal{F} \setminus \{e_j\})]^\Delta$  if  $e_j \in \mathcal{F}$  and  $\Omega'(\mathcal{F}) = \Omega(\mathcal{F})$  otherwise

# Uncertain Multi-Version XML Document

## Updating Uncertain Versions (II)

$\text{updPrXML}(\langle e_i, e_j, \Delta \rangle)[\langle \widehat{\mathcal{P}}, \Omega \rangle]$

- $\mathcal{G} := \mathcal{G} \cup (\{e_j\}, \{(e_i, e_j)\})$
- For all  $\text{ins}(x, i)$  in  $\Delta$ 
  - ▶  $\text{Set}(x, \text{fie}(x) \vee (e_j))$  if  $x$  already in  $\widehat{\mathcal{P}}$
  - ▶ Insert  $x$  in  $\widehat{\mathcal{P}}$  and  $\text{Set}(x, (e_j))$  otherwise
- For all  $\text{del}(x)$  in  $\Delta$ 
  - ▶  $\text{Set}(x, \text{fie}(x) \wedge \neg(e_j))$

**Thm2:**  $\text{updOP}(\langle e_i, e_j, \Delta \rangle) \equiv \text{updPrXML}(\langle e_i, e_j, \Delta \rangle)$

**Thm3:**  $\text{updPrXML}$  runs in  $O(1)$  while  $\widehat{\mathcal{P}}$  grows linearly according to  $|\Delta|$

# Uncertain Multi-version XML documents

## Performance Analysis (Metrics, datasets and baseline)

- ▶ Estimation of the commit time and checkout cost of the model

### Baseline Systems

- ☞ Versioning tools SubVersion and Git
- Use of their Java implementations based on the APIs SvnKit and JGit

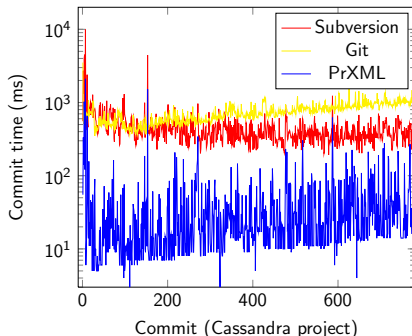
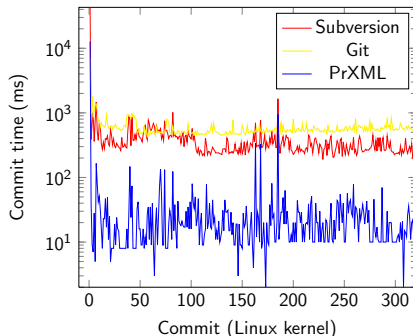
### Real Datasets

- History of commits over two large file systems (shared tree-structured data)
- ☞ Linux kernel development
  - ☞ Cassandra project

- ▶ Implementation of our system (PrXML) in Java
- ▶ Measures are obtained with all accesses in RAM Disk

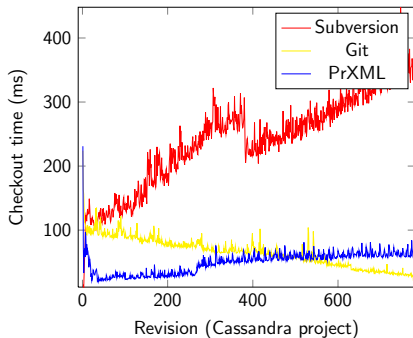
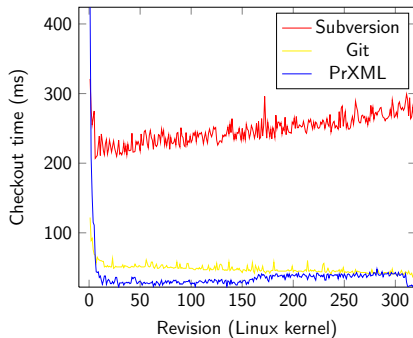
# Evaluation of the model

## Performance Analysis (Commit time)



# Evaluation of the model

## Performance Analysis (Checkout time)



# Plan

Uncertain Tree-Structured Data Model

Uncertain Multi-Version XML Document

Conclusion and Further work

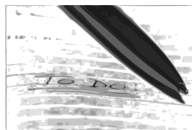
## Conclusion and Further work (i)

- ▶ Design of a probabilistic version control approach for uncertain tree-structured documents
  - ▶ Both logical description and an efficient probabilistic compact XML encoding
  - ▶ Set-up of the most used version control operation, i.e., update and with an efficient mapping algorithms
  - ▶ Both theoretical and practical complexity analysis of the proposed model
- ▶ **Extension** of our model in [Ba et al.(DChanges, 2013)] with the **merging operation** over uncertain versions



## Conclusion and Further work (ii)








- ▶ Support of more complex versioning operations such as copying, renaming etc.
- ▶ Study the impact of introducing some constraints over the order of nodes in uncertain version control



# MERCI!



# References

-  Kerstin Altmanninger, Martina Seidl, and Manuel Wimmer, *A survey on model versioning approaches*, IJWIS 5 (2009).
-  M. Lamine Ba, Talel Abdesslem, and Pierre Senellart, *Towards a version control model with uncertain data*, PIKM, 2011.
-  \_\_\_\_\_, *Merging uncertain multi-version XML documents*, Proc. DChanges (Florence, Italy), sep 2013.
-  \_\_\_\_\_, *Uncertain version control in open collaborative editing of tree-structured documents*, Proc. DocEng (Florence, Italy), sep 2013.
-  Evgeny Kharlamov, Werner Nutt, and Pierre Senellart, *Updating Probabilistic XML*, Updates in XML, 2010.
-  Benny Kimelfeld and Pierre Senellart, *Probabilistic XML: Models and complexity*, Advances in Probabilistic Databases for Uncertain Information Management (Zongmin Ma and Li Yan, eds.), Springer-Verlag, 2013.
-  Al Koc and Abdullah Uz Tansel, *A survey of version control systems*, ICEME, 2011.