

Containment of Shape Expression Schemas for RDF

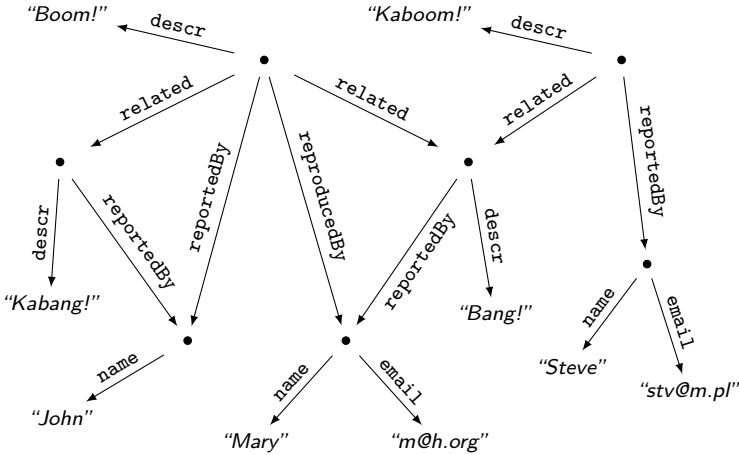
Sławek Staworko
University of Lille & INRIA LINKS
(joint work with P. Wiecek, University of Wrocław)

Principles of Database Systems 2019

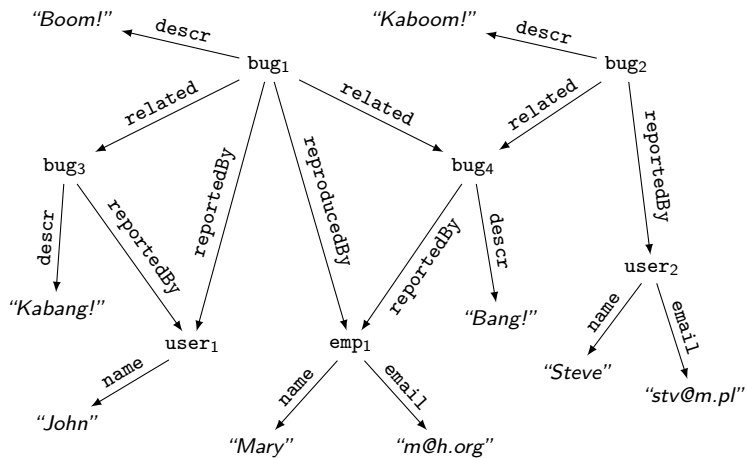
Amsterdam, Netherlands

2 July 2019

What is RDF and does it need schemas?



What is RDF and does it need schemas?



What is RDF and does it need schemas? (cont'd.)

Originally, free-range RDF

- ▶ The *driving* technology of Web 3.0
- ▶ “*Just publish your data so others can access it!*”
- ▶ Intentionally schema-free and ontology oriented (RDF Schema)

Nowadays, industrial-strength RDF

- ▶ Produced and consumed by applications (data exchange format)
- ▶ Often obtained from exporting data from relational databases (e.g., R2RML)
- ▶ Follows a strict structure

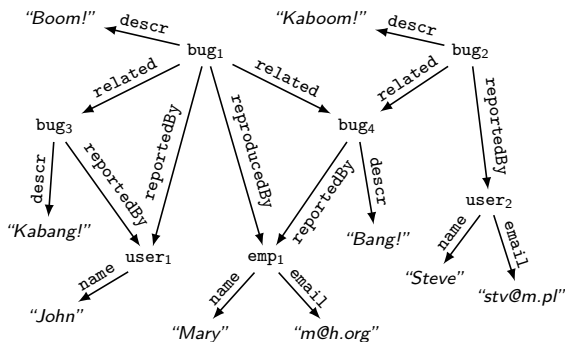
What are schemas for?

- ▶ Provide a semantic insight into data
- ▶ Capture the structure of the graph (summary)
- ▶ Enable validation i.e., checking data conformance

Shape Expression Schema (ShEx)

Syntax

ShEx is a set of rules of the form $Type \rightarrow RegExp(Predicate \times Type)$



```
Bug → descr :: str,  
      reportedBy :: User,  
      reproducedBy :: Employee?,  
      related :: Bug*
```

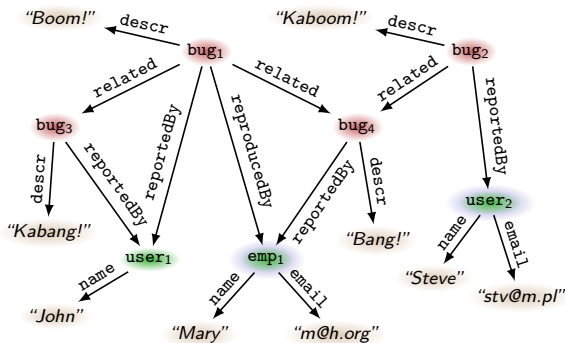
```
User → name :: str,  
      email :: str?
```

```
Employee → name :: str,  
          email :: str
```

Shape Expression Schema (ShEx)

Syntax

ShEx is a set of rules of the form $Type \rightarrow RegExp(Predicate \times Type)$



```
Bug → descr :: str,  
      reportedBy :: User,  
      reproducedBy :: Employee?,  
      related :: Bug*
```

```
User → name :: str,  
      email :: str?
```

```
Employee → name :: str,  
          email :: str
```

Semantics

Graph satisfies a schema if every node has **at least one** type

Background information

Shape Expressions Schemas (ShEx)

- ▶ Inspired by XML Schema and reminiscent of (tree) automata
- ▶ Based on regular expressions under commutative closure
membership **NP-c** [Kopczynski&To'10]; containment **coNEXP-c** [Haase&Hofman'16]
- ▶ Envisioned as a potential XSLT-like transformation engine for RDF

ShEx vs SHACL

- ▶ ShEx is a schema language with a growing base of users and a host of applications
- ▶ SHACL is Shape Constraint Language (e.g., path constraints)
- ▶ significant overlap (upcoming paper) but also differences (recursion, negation etc.)
- ▶ comparable validation complexity (**NP-complete**)
- ▶ both have been developed under the tutelage of W3C
- ▶ SHACL ended up a W3C Recommendation (yay!), ShEx a W3C Community Group Project

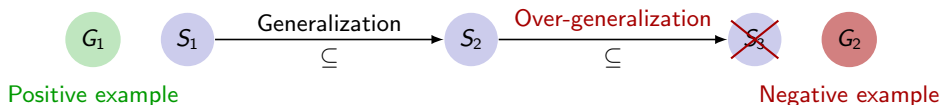
Containment problem

Containment $S_1 \subseteq S_2$

Does every graph that satisfies S_1 also satisfies S_2 ?

Motivation

- ▶ Fundamental problem (static analysis: query optimization, schema minimization etc.)
- ▶ Inference of ShEx (work in progress)



The challenge

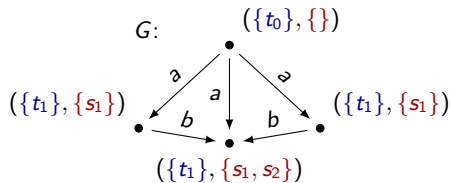
- ▶ Commutative (unordered) REs = Presburger Arithmetic (PA)
- ▶ $\text{MSO}_G \subsetneq \text{ShEx} \subseteq \text{MSO}_G + \text{PA}$
- ▶ MSO_G with very little arithmetic is undecidable [Elgot&Rabin'66]

Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$

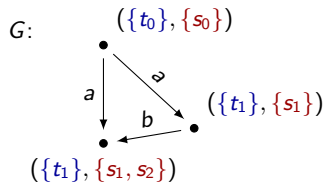


Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$

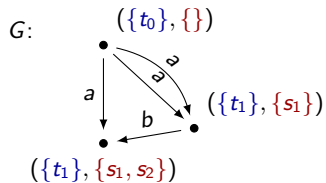


Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$

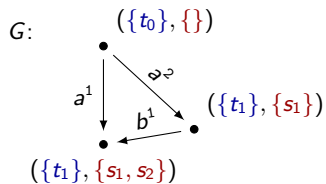


Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$

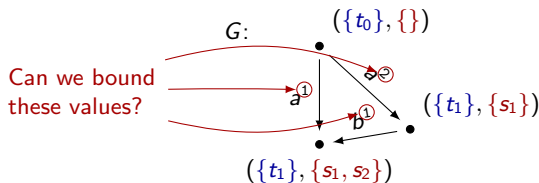


Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$



Containment of ShEx is in $\text{co2NEXP}^{\text{NP}}$

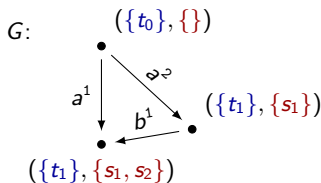
- ▶ The counter-example is a graph with at most exponential number of nodes, one node per (A, B) -kind
- ▶ A PA formula that describes the **multiplicities**
- ▶ PA enjoys an upper bound $O(|\varphi|^{3|\bar{x}|^k})$ on minimal solutions [Weispfenning'90]
- ▶ **Double exponential** upper bound on the size of a counter-example

Decidability of Containment

$$S_1 : \begin{array}{l} t_0 \rightarrow a :: t_1^* \\ t_1 \rightarrow b :: t_1^? \end{array}$$

$\not\subseteq$

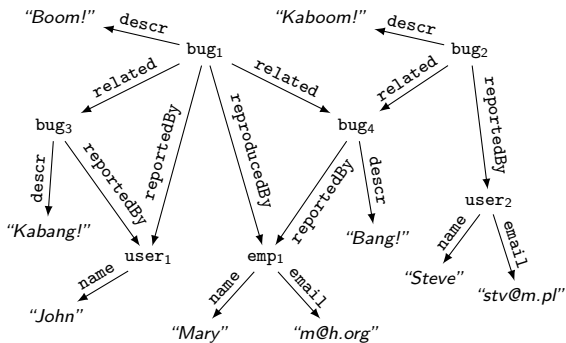
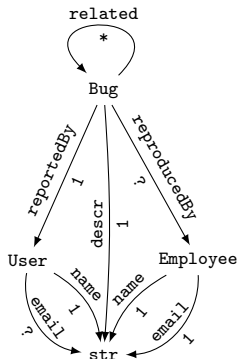
$$S_2 : \begin{array}{l} s_0 \rightarrow a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ s_1 \rightarrow b :: s_2^? \quad s_2 \rightarrow \epsilon \end{array}$$



Containment of ShEx is in $\text{co2NEXP}^{\text{NP}}$ and coNEXP-hard

- ▶ The counter-example is a graph with at most exponential number of nodes, one node per (A, B) -kind
- ▶ A PA formula that describes the **multiplicities**
- ▶ PA enjoys an upper bound $O(|\varphi|^{3|\bar{x}|^k})$ on minimal solutions [Weispfenning'90]
- ▶ **Double exponential** upper bound on the size of a counter-example
- ▶ Containment of commutative REs has recently been shown to be **coNEXP-hard** [Haase&Hofman'16]

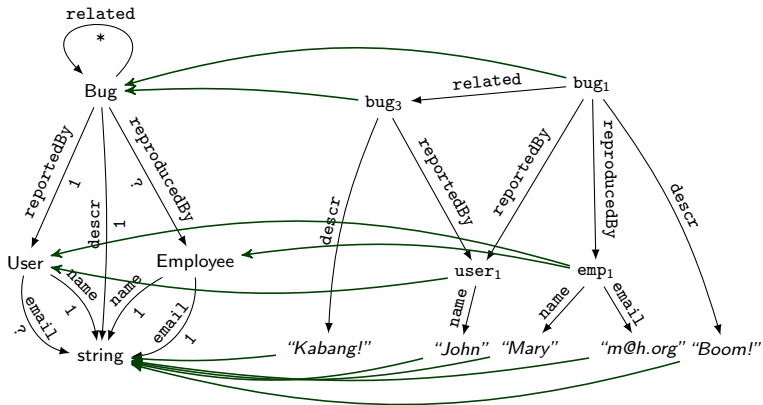
- ▶ **no disjunction** ($a :: t_1 \mid b :: t_2$) and **no grouping** ($(a :: t_1, b :: t_2)^*$)
- ▶ **Shape Graphs** – an equivalent graphical representation



Bug \rightarrow descr :: str, reportedBy :: User, reproducedBy :: Employee[?], related :: Bug^{*}
 User \rightarrow name :: str, email :: str[?]
 Employee \rightarrow name :: str, email :: str

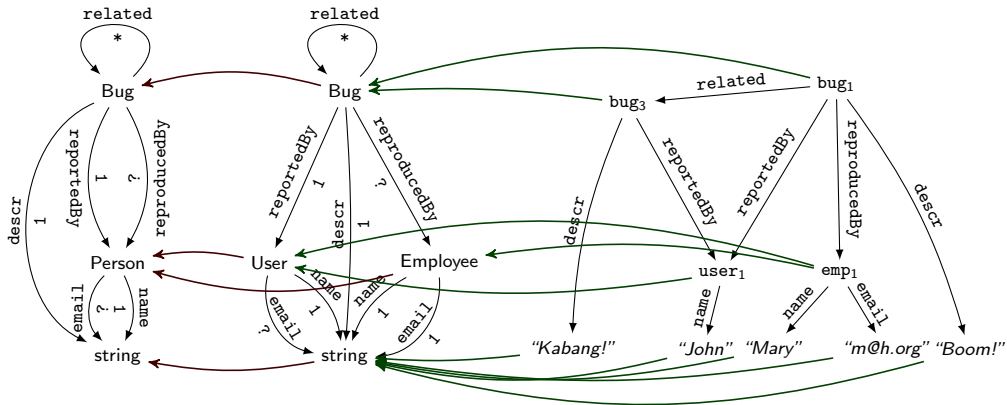
Embeddings

- ▶ Generalized simulations (graph morphism with occurrence constraints)
- ▶ Capture semantics of ShEx_0 by means of structural comparison



Embeddings

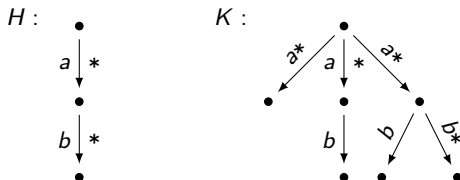
- ▶ Generalized simulations (graph morphism with occurrence constraints)
- ▶ Capture semantics of ShEx₀ by means of structural comparison
- ▶ Embeddings generalize naturally to pairs of shape graphs



Properties of embeddings

Embedding and containment

- ▶ Embedding implies containment
- ▶ In general, the converse does not hold



H cannot be embedded into K ($b :: t^*$ is equivalent to $\epsilon \mid b :: t \mid b :: t^+$)

Theorem

Constructing embeddings is

- ▶ in **PTIME** if only 1, ?, *, + are used
- ▶ **NP-complete** if arbitrary occurrence constraints are allowed $a :: t^{[n;m]}$

When does containment implies embedding ?

Determinism

- ▶ DetShEx_0 every type uses each predicate symbol at most once
- ▶ DetShEx_0^- no + are allowed and ? must be dominated by *

Characterizing graph

For any $H \in \text{DetShEx}_0^-$ there is a polynomially-sized graph G characterizing H under containment i.e.,

$$\forall K \in \text{DetShEx}_0^-. G \text{ satisfies } K \Rightarrow H \subseteq K.$$

Theorem

Containment for DetShEx_0^- is in **PTIME**

Theorem

Containment for DetShEx_0 is **coNP-hard**

Two equivalent ShEx₀ schemas and their shape graphs

H :

Bug \rightarrow descr :: str, reportedBy :: User, reproducedBy :: Employee[?],

related :: Bug^{*}

User \rightarrow name :: str, email :: str[?]

Employee \rightarrow name :: str, email :: str

K :

User₁ \rightarrow name :: str

User₂ \rightarrow name :: str, email :: str

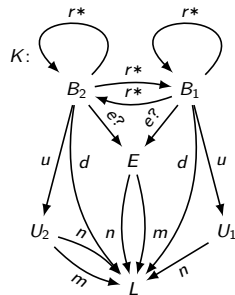
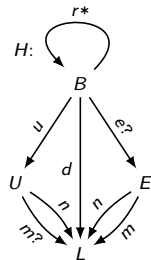
Bug₁ \rightarrow descr :: str, reportedBy :: User₁, reproducedBy :: Employee[?],

related :: Bug₁^{*}, related :: Bug₂^{*}

Bug₂ \rightarrow descr :: str, reportedBy :: User₂, reproducedBy :: Employee[?],

related :: Bug₁^{*}, related :: Bug₂^{*}

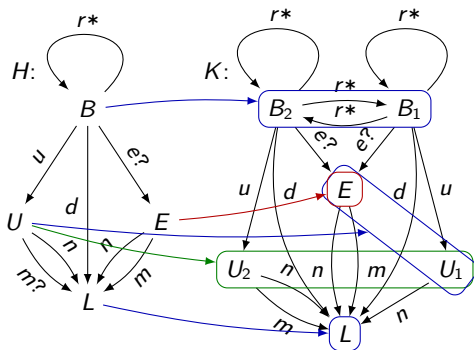
Employee \rightarrow name :: str, email :: str



Coverings

Generalization of embeddings

A type t is **covered** by a set of types $S = \{s_1, \dots, s_k\}$ iff any node satisfying t also satisfies one of the types in S

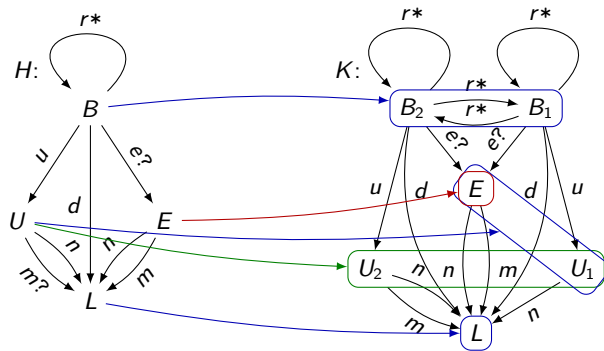


Lemma (Constructing covering)

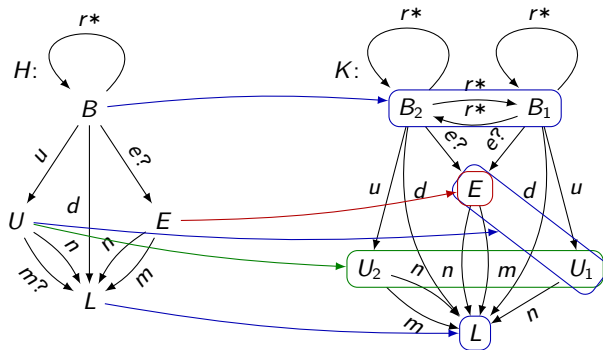
Covering is the maximum relation $R \subseteq \text{Types}(H) \times \mathcal{P}(\text{Types}(K))$ such that

$$\forall (t, S) \in R. \text{def}(t) \xrightarrow{\text{Unfold}}_R \{\text{def}(s) \mid s \in S\}.$$

Unfolding



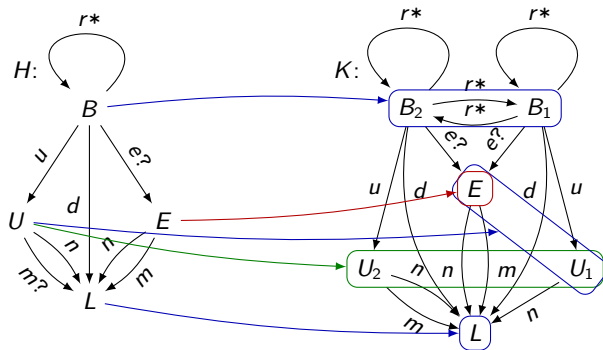
Unfolding



Unfolding U into $\{U_1, U_2\}$

$$U \rightarrow n :: L, m :: L^? \equiv n :: L, (\epsilon \mid m :: L) \equiv (n :: L) \mid (n :: L, m :: L) \leftarrow U_1 \mid U_2$$

Unfolding



Unfolding B into $\{B_1, B_2\}$

$$\begin{aligned}
 B &\rightarrow r :: B^*, u :: U, d :: L, e :: E^? \\
 &\equiv (r :: B^*, u :: U_1, d :: L, e :: E^?) \mid (r :: B^*, u :: U_2, d :: L, e :: E^?) \\
 &\equiv (r :: B_1^*, r :: B_2^*, u :: U_1, d :: L, e :: E^?) \mid (r :: B_1^*, r :: B_2^*, u :: U_2, d :: L, e :: E^?) \\
 &\leftarrow B_1 \mid B_2
 \end{aligned}$$

Complexity of ShEx_0

Theorem

Containment for ShEx_0 is in **EXP**

- ▶ Covering is a relation of exponential size
- ▶ Covering can be obtained with an iterative refinement process (starting with maximal relation and remove at least one element at each iteration until stabilization)
- ▶ At each step unfoldings are constructed and each unfolding is a tree whose size is bounded exponentially

Theorem

Containment for ShEx_0 is **EXP-complete**

- ▶ Reduction from containment for binary tree automata

Conclusions and future work

Summary of results

- ▶ Containment for ShEx is decidable
- ▶ There is a (arguably practical) class DetShEx_0^- with tractable containment
- ▶ ShEx is very different from tree automata and requires novel techniques

ShEx	DetShEx	ShEx ₀	DetShEx ₀	DetShEx ₀ ⁻
coNEXP-h and co2EXP ^{NP}	co2EXP	EXP-c	coNP-h	PTIME

Further work

- ▶ Since ShEx₀ still can capture (limited) disjunction, can the lower bounds be adapted to ShEx₀ with disjunction?
- ▶ How many of our results transfer to SHACL and at what cost?
- ▶ What is the precise impact of determinism on complexity of containment?

Questions