Shape Expression Schemas for RDF Semantics, Complexity, and Inference

Sławek Staworko University of Lille & INRIA LINKS

GT-ALGA 2019

Paris, France

10 October 2019

# What is RDF and does it need schemas?



# What is RDF and does it need schemas?



# RDF at its conception and RDF now

#### What is RDF?

set of triples (subject, predicate, object)

### Originally, free-range RDF

- The driving technology of Web 3.0 [Tim Berners-Lee]
- "Just publish your data so others can access it!"
- Intentionally schema-free and ontology oriented

# 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

# Outline

#### Uses for schemas

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

### If RDF that does not come with schema?

Infer the schema

### If the RDF does not satisfy its schema?

Repair the RDF

# Initial research [ICDT'15]

- Define the semantics of ShEx
- Complexity of validation
- Expressive power

Joint work with joint work with I. Boneva, J. Labra Gayo, S. Hym, E. G. Prud'hommeaux, and H. Solbrig

#### Current research

- Containment problem [PODS'19]
- Grammatical inference (ongoing)

Joint work with P. Wieczorek, A. Lemay, and B. Groz

#### Future research

# Previously existing schema formalisms for RDF RDF Schema (RDFS) [W3C]:

- lightweight ontology language (types and type inclusion relations)
- range and domain constraints for properties (predicate types)
- virtually no power to constrain the structure of the graph



# Previously existing schema formalisms for RDF (cont'd.)

### OWL + CWA + UNA [Sirin, RR'10]

- Potentially confusing nonstandard semantics
- Potentially high complexity of validation
- SPARQL (SPIN) [Bolleman et al., SWAT4LS'12]
  - Very powerful and expressive
  - High complexity

Resource Shapes [IBM, Ryman et al., LDOW'13]

Extends RDFS with simple cardinality constraints on the outbound neighborhood of a node

### What does exactly RDF validation mean

Verification the typing is given (with rdf:type) and its correctness is to be verified (this meaning is employed by the above schema proposals)

Model checking no typing is given and the goal is to construct a valid typing (this is a more general problem that we are interested in)

# Shape Expressions Schemas

# Shape Expression Schema (ShEx)

#### Syntax

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



# Shape Expression Schema (ShEx)

#### Syntax

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



### 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 (e.g., Wikidata)
- 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

# Flavorful ShEx

```
<Emp> {
    name xsd:string | (first-name xsd:string, last-name xsd:string),
    email xsd:string,
    department @<Dept>
}

CDept> {
    name xsd:string,
    reportedBy @<User>,
    reportedOn xsd:dateTime,
    (manager @<Emp>, appointedOn xsd:dateTime)?,
    employees @<Emp>*
}
```

### Bags of symbols (unordered words)

Bag (multiset) is a function mapping a symbol to the number of its occurrences.

 $w_0 = \{ | a, a, a, c, c \}$  represents the function  $w_0(a) = 3$ ,  $w_0(b) = 0$ , and  $w_0(c) = 2$ .

The collection of outgoing labels is a bag:

$$out-lab_G(n) = \{ | a | (n, a, m) \in E_G \}$$



Bag union: [a, c, c]  $\bigcup$  [a, b] = [a, a, b, c, c] (concatenation of unordered words).

# Regular Bag Expressions (RBEs)

Language of regular expressions for defining bags (unordered concatenation , )

 $E ::= \epsilon \mid a \mid E^* \mid (E'' \mid E) \mid (E'', E)$ 

with natural macros: E? := ( $\epsilon \mid E$ ) and E<sup>+</sup> := (E, E<sup>\*</sup>)

#### Examples

- ( $a^*$ ,  $b^+$ , c, c) arbitrary number of a's, positive number of b's, and two c's
- ► (a, b)\* the same number of a's and b's
- $(a, b, c)^*$  the same number of a's, b's, and c's.

#### RBEs are equivalent to

- 1. Presburger arythmetic (PA),
- 2. Parikh images of context-free languages,
- 3. semilinear sets.

#### Computational properties

- Membership  $w \in E$  is NP-complete,
- Emptiness  $E_1 \cap E_2 = \emptyset$  is coNP-complete.

# $RBE_0$ and $ShEx_0 = ShEx(RBE_0)$

# RBE<sub>0</sub>

- ▶ RBEs using only symbols with multiplicities  $\{0, 1, *, +, ?\}$  and , operator only
- ▶ can be canonized  $a, a^? \equiv a^{[1,2]}, b^+, b^+ \equiv b^{[2,\infty]}$ , etc.
- ▶ the canonical form is  $a^{[n,n']}$ ,  $b^{[m,m']}$ , ...
- ▶ Presburger formulas: conjunctions of atoms #a < n and #a > n
- captures IBM's Resource Shapes

#### Computational properties: simple arithmetic

A lightweight class enjoying tractability of a number of problems:

- membership
- containment
- ▶ intersection (also with RBE<sub>1</sub>)

Also learnable from positive examples [DBPL'13]



 $G_0$ :

$$\begin{array}{ll} S_0: & t_0 \rightarrow (a::t_1)^+, \ b::t_2 \\ & t_1 \rightarrow (a::t_1 \mid b::t_2)^* \\ & t_2 \rightarrow b::t_2 \mid c::t_1 \end{array}$$





n

A single-type typing is a function  $\lambda : V \to \Gamma$ .

 $\lambda$  is valid if every node *n* satisfies its type definition i.e.,

$$out-lab-type_G^{\lambda}(n) = \{ a :: \lambda(m) \mid (n, a, m) \in E \} \in \delta(\lambda(n)).$$

$$S_0: \quad t_0 \to (a :: t_1)^+, \ b :: t_2$$
$$t_1 \to (a :: t_1 \mid b :: t_2)^* \qquad G_0:$$
$$t_2 \to b :: t_2 \mid c :: t_1$$
$$out-lab-type_{C}^{\lambda_0}(n_0) = \{|a :: t_1, a :: t_1, b :: t_2|\}$$



A single-type typing is a function  $\lambda: V \to \Gamma$ .

 $\lambda$  is valid if every node *n* satisfies its type definition i.e.,

$$out-lab-type_G^{\lambda}(n) = \{ a :: \lambda(m) \mid (n, a, m) \in E \} \in \delta(\lambda(n)).$$

A valid single-type typing of  $G_0$  w.r.t.  $S_0$ 

$$\lambda_0(n_0) = t_0, \qquad \lambda_0(n_1) = t_1, \qquad \lambda_0(n_2) = t_2, \qquad \lambda_0(n_3) = t_1, \qquad \lambda_0(n_4) = t_2.$$

# Intractability of single-type validation

#### Validation problem

Checking if there exists a valid typing of given input graph w.r.t. a given input schema.

### Sources of complexity

- 1. guessing a typing
- 2. checking that it is valid (RBE membership is already NP-complete)

# Intractability of single-type validation

#### Validation problem

Checking if there exists a valid typing of given input graph w.r.t. a given input schema.

### Sources of complexity

- 1. guessing a typing
- 2. checking that it is valid (RBE membership is already NP-complete)

#### Theorem

Single-type validation is NP-complete (even if only RBE<sub>0</sub> are used).

Reduction from graph 3-colorability:

$$t_r 
ightarrow \_ :: t_b^*, \_ :: t_g^*$$
  $t_g 
ightarrow \_ :: t_r^*, \_ :: t_b^*$   $t_b 
ightarrow \_ :: t_g^*, \_ :: t_r^*$ 

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$G_{1}: n_{0} \xrightarrow{b} n_{1} \xrightarrow{c} n_{2}$$

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$G_{1}: n_{0} \xrightarrow{a} n_{1} \xrightarrow{c} n_{2}$$

$$t_{0}$$

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$G_{1}: n_{0} \xrightarrow[t_{0}]{a} n_{1} \xrightarrow[t_{1}]{c} n_{2}$$

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$G_{1}: n_{0} \xrightarrow[t_{0}]{} n_{1} \xrightarrow[t_{1}]{} n_{2}$$

$$t_{0} \xrightarrow[t_{1}]{} t_{2} \xrightarrow[t_{3}]{} t_{3}$$

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

 $\lambda_1(n_0) = \{t_0\}, \qquad \lambda_1(n_1) = \{t_1, t_2\}, \qquad \lambda_1(n_2) = \{t_3\}.$ 

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

$$S_{1}: \quad t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$\lambda_{1}(n_{0}) = \{t_{0}\}, \qquad \lambda_{1}(n_{1}) = \{t_{1}, t_{2}\}, \qquad \lambda_{1}(n_{2}) = \{t_{3}\}.$$

1-

 $\lambda$  is valid if every node satisfies every of its associated types.

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$\lambda_{1}(n_{0}) = \{t_{0}\}, \qquad \lambda_{1}(n_{1}) = \{t_{1}, t_{2}\}, \qquad \lambda_{1}(n_{2}) = \{t_{3}\}.$$

.

 $\lambda$  is valid if every node satisfies every of its associated types.

When defining that a node n satisfies a type t...

▶ we inspect the outbound neighborhood *out-lab-node*<sub>G</sub>(n) = {(a, m) | (n, a, m) \in E}

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.



 $\lambda$  is valid if every node satisfies every of its associated types.

When defining that a node n satisfies a type t...

- ▶ we inspect the outbound neighborhood out-lab-node<sub>G</sub>(n) = {(a, m) | (n, a, m) \in E}
- > a node m may assume any of its assigned types  $\lambda(m)$ , one per edge incoming from n

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

$$S_{1}: t_{0} \rightarrow a :: t_{1}$$

$$t_{1} \rightarrow b :: t_{2}, c :: t_{3}$$

$$t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3}$$

$$t_{3} \rightarrow \epsilon$$

$$\lambda_{1}(n_{0}) = \{t_{0}\}, \qquad \lambda_{1}(n_{1}) = \{t_{1}, t_{2}\}, \qquad \lambda_{1}(n_{2}) = \{t_{3}\}.$$

1.

 $\lambda$  is valid if every node satisfies every of its associated types.

When defining that a node *n* satisfies a type *t*...

- ▶ we inspect the outbound neighborhood out-lab-node<sub>G</sub>(n) = {(a, m) | (n, a, m) \in E}
- ▶ a node *m* may assume any of its assigned types  $\lambda(m)$ , one per edge incoming from *n*
- $(n, a, m) \in E$  yields the choice  $|_{t \in \lambda(m)} a :: t$

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

$$S_{1}: t_{0} \rightarrow a :: t_{1} \qquad \qquad b \\ t_{1} \rightarrow b :: t_{2}, c :: t_{3} \qquad \qquad G_{1}: n_{0} \xrightarrow{a \ n_{1}} c \ n_{2} \\ t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3} \qquad \qquad f_{1} \ t_{2} \rightarrow t_{3} \\ t_{3} \rightarrow \epsilon \qquad \qquad Out Type(n_{1}, \lambda_{1}) = (b :: t_{1} \mid b :: t_{2}), c :: t_{3} \\ \lambda_{1}(n_{0}) = \{t_{0}\}, \qquad \lambda_{1}(n_{1}) = \{t_{1}, t_{2}\}, \qquad \lambda_{1}(n_{2}) = \{t_{3}\}.$$

,

 $\lambda$  is valid if every node satisfies every of its associated types.

When defining that a node n satisfies a type t...

- ▶ we inspect the outbound neighborhood out-lab-node<sub>G</sub>(n) = {(a, m) | (n, a, m) \in E}
- ▶ a node *m* may assume any of its assigned types  $\lambda(m)$ , one per edge incoming from *n*
- $(n, a, m) \in E$  yields the choice  $|_{t \in \lambda(m)} a :: t$

• 
$$OutType(n, \lambda) = \bigcirc_{(n,a,m)\in E} (|_{t\in\lambda(m)} a :: t)$$

A multi-type typing is a function  $\lambda: V \to 2^{\Gamma}$  that assign to every node a set of types.

$$S_{1}: t_{0} \rightarrow a :: t_{1} \qquad \qquad b \\ t_{1} \rightarrow b :: t_{2}, c :: t_{3} \qquad \qquad G_{1}: n_{0} \xrightarrow{a \rightarrow n_{1}} c \rightarrow n_{2} \\ t_{2} \rightarrow (b :: t_{2})^{?}, c :: t_{3} \qquad \qquad t_{0} \xrightarrow{t_{1}} t_{2} \xrightarrow{t_{3}} t_{3} \\ t_{3} \rightarrow \epsilon \qquad \qquad Out Type(n_{1}, \lambda_{1}) = (b :: t_{1} \mid b :: t_{2}), c :: t_{3} \\ \lambda_{1}(n_{0}) = \{t_{0}\}, \qquad \lambda_{1}(n_{1}) = \{t_{1}, t_{2}\}, \qquad \lambda_{1}(n_{2}) = \{t_{3}\}.$$

,

 $\lambda$  is valid if every node satisfies every of its associated types.

*n* satisfies *t* w.r.t.  $\lambda$  if  $OutType(n, \lambda) \cap \delta(t) \neq \emptyset$ , where  $OutType(n, \lambda) = \bigcirc_{(n,a,m)\in E} (|_{t\in\lambda(m)} a:: t)$ 

### Refinement algorithm

The set of all valid multi-type typings of G w.r.t. S is a semi-lattice.

- 1. Start with the universal typing  $\lambda(n) := \Gamma$
- 2. Iteratively refine it  $\lambda := Refine(\lambda)$

 $[Refine(\lambda)](n) = \{t \in \lambda(n) \mid OutType(n, \lambda) \cap \delta(t) \neq \emptyset\}.$ 

- 3. Until a fix-point is reached
- 4. The graph satisfies the schema iff the fix-point  $\lambda$  is valid ... ... and then  $\lambda$  is also the maximal valid multi-type typing.

# Satisfiability of RBEs

*OutType* yields expressions of the form (RBE<sub>1</sub>)

$$(a_1 \mid \cdots \mid a_k), \ \cdots, \ (z_1 \mid \ldots \mid z_m)$$

The essential decision problem for a class  ${\mathcal C}$  of RBEs used in the ShEx schema is

 $\mathsf{INTER}_1(\mathcal{C}) = \{ (E_0, E) \in \mathsf{RBE}_1 \times \mathcal{C} \mid E_0 \cap E \neq \emptyset \}.$ 

#### Lemma

Tractability of INTER<sub>1</sub> is a necessary and sufficient condition for tractability of multi-type validation.

### Corollary

Multi-type validation is NP-complete.

Theorem Multi-type validation for  $ShEx_0 = ShEx(RBE_0)$  is in PTIME

# Determinism

#### Determinism of shape expressions

Given the type (of a node) and the label of an outgoing edge, the expression specifies the type that the target node must satisfy.

 $\begin{array}{ccc} a :: t_1, \ b :: t_2^*, \ a :: t_1, \ c :: t_2 & (a :: t_1, \ b :: t_2) \mid (a :: t_3, \ c :: t_4) & a :: t_1, \ b :: t_2^*, \ a :: t_3 \\ \\ & \text{deterministic} & \text{not deterministic} & \text{not deterministic} \end{array}$ 

### Determinism

#### Determinism of shape expressions

Given the type (of a node) and the label of an outgoing edge, the expression specifies the type that the target node must satisfy.

 $a :: t_1, b :: t_2^*, a :: t_1, c :: t_2$  $(a :: t_1, b :: t_2) | (a :: t_3, c :: t_4)$  $a :: t_1, b :: t_2^*, a :: t_3$ deterministicnot deterministicnot deterministic

#### Lemma

For schemas using only deterministic shape expressions, tractability of membership is a sufficient and necessary condition for tractability of multi-type validation

### Proof sketch

- Knowing the label a of an outgoing edge determines the type t<sub>a</sub> for the target node
- $OutType(n, \lambda) = \bigcirc_{(n,a,m)\in E} (|_{t\in\lambda(m)} a :: t) \text{ becomes } \bigcirc_{(n,a,m)\in E} (a :: t_a)$
- $\blacktriangleright \bigcirc_{(n,a,m)\in E} (a :: t_a) \text{ defines a singleton } \{w\} \text{ with } w = \{a :: t_a \mid (n,a,m)\}$
- $OutType(n, \lambda) \cap \delta(t) \neq \emptyset \equiv w \in \delta(t).$
# Single-occurrence REBs (SORBEs)

SORBE allows a symbol to be used at most once in an expression but also allows  $a^{[n,m]}$ 

Theorem Membership for SORBE is in PTIME :)

 $a:: t_1, b:: t_2^*, a:: t_1$   $(a:: t_1, b:: t_2) | (a:: t_3, c:: t_4)$   $(a:: t_1, b:: t_2)^*, c:: t_3$ 

deterministic not deterministic deterministic but yet not single-occurrence single-occurence single-occurrence

#### Theorem

Multi-type validation for deterministic shape expressions using SORBE is in PTIME. :)

and

## Expressive power of ShEx





	FO <sub>G</sub>	$ShEx_m$	ShEx <sub>s</sub>	∃MSO <sub>G</sub>
$L: G_{\diamond} \in L, \ G_{<} \in L$	1	1	1	✓
$L: G_{\diamond} \not\in L, \ G_{<} \in L$	1	Х	1	1
$L: G_{\diamond} \in L, \ G_{<} \not\in L$	1	Х	X	✓

### Quick recap

### Complexity

	$RBE_0$	RBE	SORBE	SORBE det.
multi-type	PTIME	NP-complete		PTIME
single-type		NP-c	omplete	

### Expressive power

- automata-like formalism
- incomparable with FO and MSO (unless we forbid \* over expressions)
- captured with MSO+PA
- incomparable with NR and HR graph grammars
- closed under intersection but not under union or negation
- single-type semantics is more expressive than multi-type semantics

# Containment of ShEx

### 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

- RBEs = Presburger Arithmetic (PA)
- $\blacktriangleright \mathsf{MSO}_{\mathsf{G}} \not\supseteq \mathsf{ShEx} \subseteq \mathsf{MSO}_{\mathsf{G}} + \mathsf{PA}$
- MSO<sub>G</sub> with very little arithmetic becomes undecidable [Elgot&Rabin'66]

$$\begin{array}{ccc} S_1: & t_0 \to a :: t_1^* & \\ & t_1 \to b :: t_1^? & \end{array} \qquad \qquad \begin{array}{ccc} S_2: & s_0 \to a :: s_1 \mid (a :: s_1, a :: s_2)^* \\ & s_1 \to b :: s_2^2 & s_2 \to \epsilon \end{array}$$





# Containment of ShEx is in co2NEXP^{NP}

- ▶ The counter-example is a graph with at most exponential number of nodes, one node per (A, B)-kind
- There is a PA formula  $\varphi$  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 (binary) size of the values of multiplicities
- Validation of graphs with multiplicities remains in NP



# Containment of ShEx is in co2NEXP<sup>NP</sup> and coNEXP-hard

- The counter-example is a graph with at most exponential number of nodes, one node per (A, B)-kind
- $\blacktriangleright$  There is a PA formula  $\varphi$  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 (binary) size of the values of multiplicities
- Validation of graphs with multiplicities remains in NP
- Containment of commutative REs recently shown to be coNEXP-hard [Haase&Hofman'16]

# $\mathsf{ShEx}_0$

• no disjunction  $(a :: t_1 | b :: t_2)$  and no grouping  $(a :: t_1, b :: t_2)^*$ 





$$\begin{split} &\text{Bug} \to \texttt{descr}::\texttt{str}, \; \texttt{reportedBy}::\texttt{User}, \; \texttt{reproducedBy}::\texttt{Employee}^?, \; \texttt{related}::\texttt{Bug}^*\\ &\text{User} \to \texttt{name}::\texttt{str}, \; \texttt{email}::\texttt{str}^?\\ &\text{Employee} \to \texttt{name}::\texttt{str}, \; \texttt{email}::\texttt{str} \end{split}$$

### Embeddings

- Graph morphism with occurrence constraints, closely related to graph simulations
- Capture semantics of ShEx<sub>0</sub> by means of structural comparison



### Embeddings

- Graph morphism with occurrence constraints, closely related to graph simulations
- Capture semantics of ShEx<sub>0</sub> by means of structural comparison
- 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*<sup>[*n*;*m*]</sup>

# When does containment implies embedding?

### Determinism

- DetShEx<sub>0</sub> every type uses each predicate symbol at most once
- DetShEx<sub>0</sub><sup>-</sup> 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* satisfies  $K \Rightarrow H \subseteq K$ .

Theorem Containment for  $DetShEx_0^-$  is in PTIME

Theorem Containment for DetShEx<sub>0</sub> is coNP-hard



# Two equivalent $ShEx_0$ schemas and their shape graphs

```
H:
Bug \rightarrow descr :: str, reportedBy :: User, reproducedBy :: Employee<sup>?</sup>,
             related :: Bug*
User \rightarrow name :: str. email :: str?
\texttt{Employee} \rightarrow \texttt{name} :: \texttt{str}, \texttt{email} :: \texttt{str}
K:
\texttt{User}_1 \rightarrow \texttt{name} :: \texttt{str}
User_2 \rightarrow name :: str, email :: str
Bug_1 \rightarrow descr :: str, reportedBy :: User_1, reproducedBy :: Employee<sup>?</sup>,
              related :: Bug<sup>*</sup><sub>1</sub>, related :: Bug<sup>*</sup><sub>2</sub>
\operatorname{Bug}_2 \to \operatorname{descr} :: \operatorname{str}, \operatorname{reportedBy} :: \operatorname{User}_2, \operatorname{reproducedBy} :: \operatorname{Employee}^?
              related :: Bug<sub>1</sub><sup>*</sup>, related :: Bug<sub>2</sub><sup>*</sup>
\texttt{Employee} \rightarrow \texttt{name} :: \texttt{str}, \texttt{ email} :: \texttt{str}
```



### Coverings

#### Generalization of embeddings

A type t is covered by a set of types  $S = \{s_1, \ldots, s_k\}$  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

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

# 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  $\{B_1, B_2\}$ 

$$B \to r :: B^*, \ u :: U, \ d :: L, \ e :: E^?$$
  

$$\equiv (r :: B^*, \ u :: U_1, \ d :: L, \ e :: E^?) | (r :: B^*, \ u :: U_2, \ d :: L, \ e :: E^?)$$
  

$$\equiv (r :: B_1^*, \ r :: B_2^*, \ u :: U_1, \ d :: L, \ e :: E^?) | (r :: B_1^*, \ r :: B_2^*, \ u :: U_2, \ d :: L, \ e :: E^?)$$
  

$$\leftarrow B_1 | B_2$$

Sławek S. (ULille & INRIA LINKS)

# 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<sub>0</sub> is EXP-complete

Reduction from containment for binary tree automata

- 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 <sup>NP</sup>	co2EXP	EXP-c	coNP-h	PTIME

# Inference of ShEx

### Constructing ShEx from an RDF Graph

What for?

- RDF originally schema-free but schemas are useful
- Free-range RDF is dirty
- ▶ Industrial-strength RDF is relatively clean and exhibits regular structure

Relational database  $\xrightarrow{\text{R2RDF}} \text{RDF}$ 

# Constructing ShEx from an RDF Graph

What for?

- RDF originally schema-free but schemas are useful
- Free-range RDF is dirty
- Industrial-strength RDF is relatively clean and exhibits regular structure

Relational database  $\xrightarrow{\text{R2RDF}}$  RDF

#### What do we want?

For a given RDF Graph G construct a ShEx schema S that captures the structure of G:

Soundness G satisfies S

Succinctness S is small enough for a human user to consume

Specificity S is not overly general (as the universal schema is)

#### Implicit node similarity

Two nodes should have the same type if they are similar in some way. Two possible criteria:

- content the outbound neighborhood is similar
- context the inbound neighborhood is similar

# Inference and Fitting

### Fitting

For a input G construct  $\subseteq$ -minimal schema S that G satisfies.

#### Inference

An inference algorithm A is an algorithm that is

- sound i.e., it returns a schema that the input graph satisfies;
- complete i.e., it can return any goal schema provided that the input graph is sufficiently informative (or rich enough).

Both approaches are parameterized by a class  ${\mathcal C}$  of goal ShEx schemas.

### Fitting for $ShEx_0$ is trivial and potentially too verbose

An RDF graph interpreted as a shape graph is its own  $ShEx_0$  fitting



#### Fitting for $ShEx_0$ is trivial and potentially too verbose

An RDF graph interpreted as a shape graph is its own ShEx<sub>0</sub> fitting



#### Fitting for ShEx<sub>0</sub> is trivial and potentially too verbose

An RDF graph interpreted as a shape graph is its own  $\mathsf{ShEx}_0$  fitting



#### Fitting for ShEx<sub>0</sub> is trivial and potentially too verbose

An RDF graph interpreted as a shape graph is its own  $\mathsf{ShEx}_0$  fitting



#### Bisimulation can identify redundant/repetitive information

But it won't generalize the result (it won't introduce \* or recursion).

### Bad news for Inference: Limit Point

#### Limit point

A class of languages C has the *limit point* property iff C contains an ascending chain of languages  $L_1 \subsetneq L_2 \subsetneq \ldots$  whose limit point  $L_{\infty} = \bigcup_i L_i$  also belongs to C.

#### Folklore result: Limit point precludes inference

No family with limit point property has an inference algorithm.

$$L_0: \bullet \qquad L_1: \bullet \xrightarrow{a?} \bullet \qquad L_2: \bullet \xrightarrow{a?} \bullet \xrightarrow{a?} \bullet \qquad \dots \qquad L_{\infty}: \bullet \bigcirc a?$$

#### Lemma

For any  $M \subseteq \{1,?,*,+\}$  with at least two elements the class  $ShEx_0(M)$  has the limit point property.

### We compromise

- Find a suitable subclass that allows inference and remains relatively practical.
- Motivation: first draft of a ShEx schema for an architect to work on.

### We compromise

- Find a suitable subclass that allows inference and remains relatively practical.
- Motivation: first draft of a ShEx schema for an architect to work on.

### Our results

	ShEx <sub>0</sub>	$SingShEx_0$	$DetShEx_0$	Typed graph
Fitting	trivial	undefined (?)	exponential	NP-hard
Inference	unfeasible	$SingShEx_0(1, *)$	$DetShEx_0^-$	$ShEx_0$ (?)

# Embeddings (recall)

- Generalize simulations
- Capture semantics of ShEx by means of structural comparison



# Singular Shape Expression Schemas
## Inference of SingShEx<sub>0</sub>: Generalization



## Inference of SingShEx<sub>0</sub>: Generalization



## Inference of SingShEx $_0$ : Generalization



## Inference of SingShEx<sub>0</sub>: Generalization



1. Put \* on every edge of the input graph

## Inference of SingShEx<sub>0</sub>: Reduction



- 1. Put \* on every edge of the input graph
- 2. Construct the autoembedding

## Inference of SingShEx<sub>0</sub>: Reduction



- 1. Put \* on every edge of the input graph
- 2. Construct the autoembedding
- 3. Remove any dominated nodes

G\*



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m



- remove n and all of its outgoing edges
- redirect its incoming edges to m

## Inference of SingShEx<sub>0</sub>: Specialization



## Inference of SingShEx<sub>0</sub>: Specialization



• Construct the embedding of the input graph G into the reduct of  $G^*$ 

## Inference of SingShEx<sub>0</sub>: Specialization



- Construct the embedding of the input graph G into the reduct of  $G^*$
- Use the embedding to replace \*'s with fitter multiplicities

## Inference of SingShEx $_0$ : Specialization



• Construct the embedding of the input graph G into the reduct of  $G^*$ 

Use the embedding to replace \*'s with fitter multiplicities

## Singular Shape Expression Schemas

#### Definition (SingShEx<sub>0</sub>)

#### A shape graph H is singular if

- 1.  $H^*$  is reduced i.e., there are no two types  $t_1$  and  $t_2$  in  $H^*$  that  $t_1$  can be embedded into  $t_2$ ,
- 2.  $H^*$  has no two edges with the same label and the same source and target nodes.

#### Singularity as a restriction is

- ▶ stronger than forbidding any two types  $t_1$  and  $t_2$  such that  $t_1 \subseteq t_2$ .
- weaker than forbidding any two types with comparable signatures (sets of outgoing edge labels).



## Deterministic Shape Expression Schemas

## Deterministic Shape Expression Schemas

#### Definition (DetShEx<sub>0</sub>)

A shape graph H is deterministic if for every label a every node has at most one outgoing edges labeled with a.



## Determinizing Shape Graphs



#### Determinization of a shape graph

- gives the fitting of the input graph (the ⊆-minimal DetShEx<sub>0</sub> schema that validates the input graph)
- might produce exponential output

## Quotient Determinization of Shape Graphs

#### Avoid explosion by refusing to duplicate types

- ▶ if  $n_1$  is fused with  $n_2$  and  $n_2$  is fused with  $n_3$ , then fuse  $n_1$ ,  $n_2$ , and  $n_3$  together.
- produces linear output



Quotient determinization is an inference algorithm for DetShEx<sub>0</sub>

Example



## Conclusions and Future Work

## Summary

#### What is ShEx?

- Automata-like formalism for graphs but unlike any automata we have seen.
- Strong connections to graph bisimulation.

#### Complexity of Validation

	ShEx <sub>0</sub>	ShEx	DetShEx (SORBE)		
multi-type	PTIME	NP-complete	PTIME		
single-type	NP-complete				

#### Complexity of Containment

ShEx	DetShEx(SORBE)	ShEx <sub>0</sub>	$DetShEx_0$	$DetShEx_0^-$
coNEXP-h and co2EXP <sup>NP</sup>	co2EXP	EXP-c	coNP-h	PTIME

#### Inference and Fitting

	ShEx <sub>0</sub>	$SingShEx_0$	$DetShEx_0$	Typed graph
Fitting	trivial	undefined (?)	exponential	NP-hard
Inference	unfeasible	$SingShEx_0(1, *)$	$DetShEx_0^-$	ShEx <sub>0</sub> (?)

## Future work: Repairing RDF



## Questions

# Appendix

## Formal Definition of Grammatical Inference of Shape Expression Schemas

#### Definition

A class of shape graphs C is *learnable in polynomial time and data from* a class of graphs G iff there exists a polynomial inference algorithm A such that the following two conditions are satisfied:

Soundness For every input graph  $G \in \mathcal{G}$  the inference algorithm returns a graph schema A(G) = H such that  $H \in C$  and  $G \in L(H)$ .

Completeness For every graph schema  $H \in C$  there exists a polynomially-sized *characteristic* graph  $G_H \in L(H)$  such that for any G that extends  $G_H$  consistently with H we have  $A(G') \equiv H$ .

#### When does G extend G'?

3 alternative definitions

- 1. G is a disjoint union of G' and some G''
- 2. G is obtained from G' by adding new nodes and new edges
- 3. as 2. but no node of G may loose a type as a result of adding outgoing edges to it.