

Capturing Functional Dependencies and Relational Schemas in RDFS

Motivation

- ▶ Enormous growth in Semantic Web data
- ▶ Most Semantic Web data originates from relational databases
- ▶ Normal forms for relational databases are important
- ▶ **GOAL: Is there a Normal Form for Semantic Web data?**

Semantic Web

- ▶ Data is stored in RDF graphs extended with RDFS
- ▶ Focus on *DL-Lite_{RDFS}* with disjointness

Relational Databases

- ▶ Data is stored in n-ary relations
- ▶ Functional dependencies (FDs) ensure *consistency*
- ▶ Normal forms avoid *redundancy*

Direct Mapping of Relational Data to RDFS Graph Data

- ▶ each database instance corresponds to a model of the ontology and vice versa.

Example

Table: *loc*

room	house	addr
Entrance	White House	1600_PA_Av
Oval Office	White House	1600_PA_Av

1. Translate schema information

- ▶ Concept disjointness:

$locs \sqsubseteq \neg room$ $room \sqsubseteq \neg house$
 $room \sqsubseteq \neg addr$ $house \sqsubseteq \neg addr$

- ▶ Role typing:

$\exists room \sqsubseteq has\ room$ $\exists has\ room \sqsubseteq room$
 $\exists house \sqsubseteq has\ house$ $\exists has\ house \sqsubseteq house$
 $\exists addr \sqsubseteq has\ addr$ $\exists has\ addr \sqsubseteq addr$

- ▶ Functionality and reification assertions:

(*func* *has room*)

(*func* *has house*)

(*func* *has addr*)

loc?(*has room*, *has house*, *has addr*)

2. Translate data

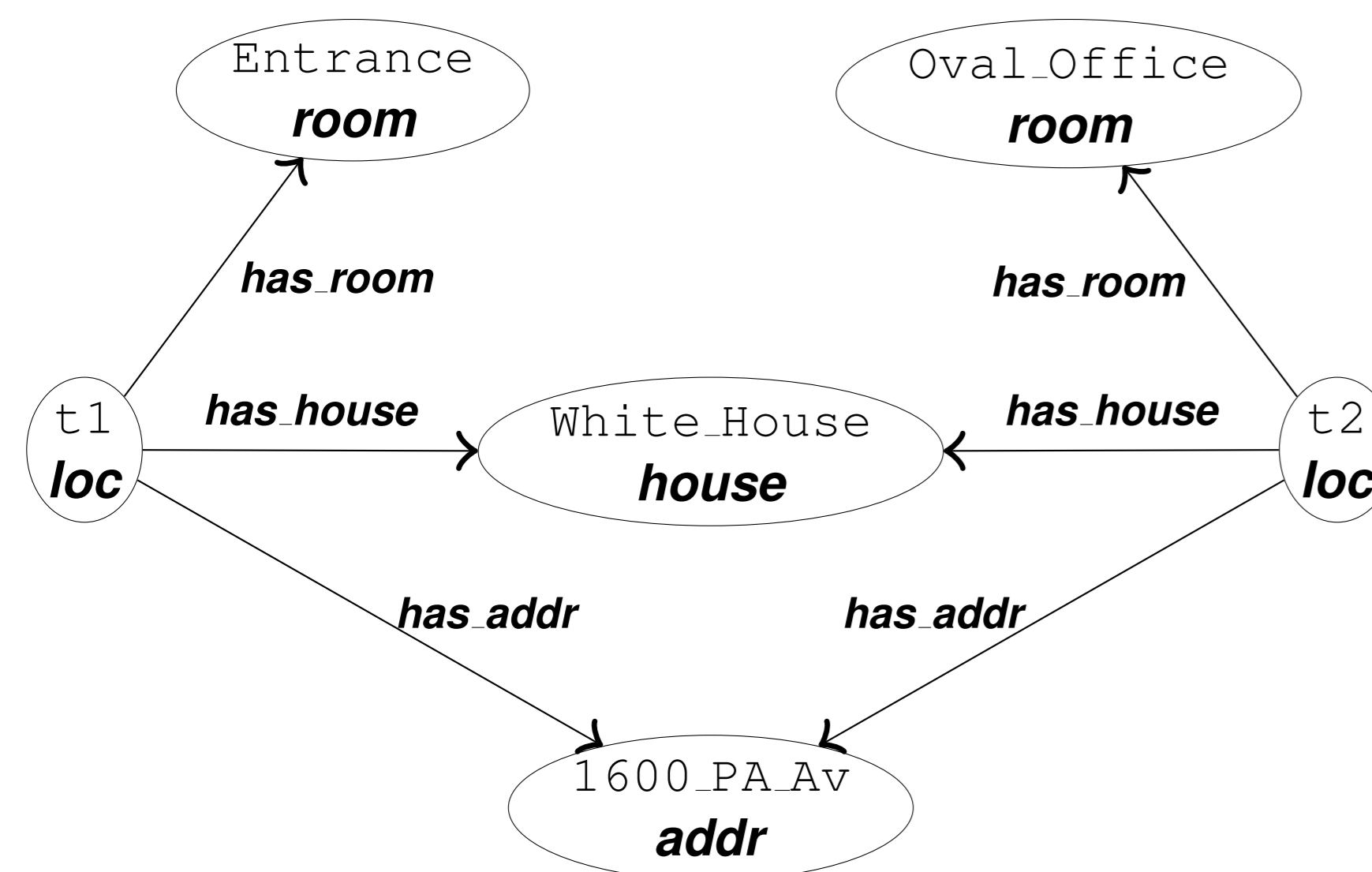


Figure 1: Translated relational table

Capturing Functional Dependencies

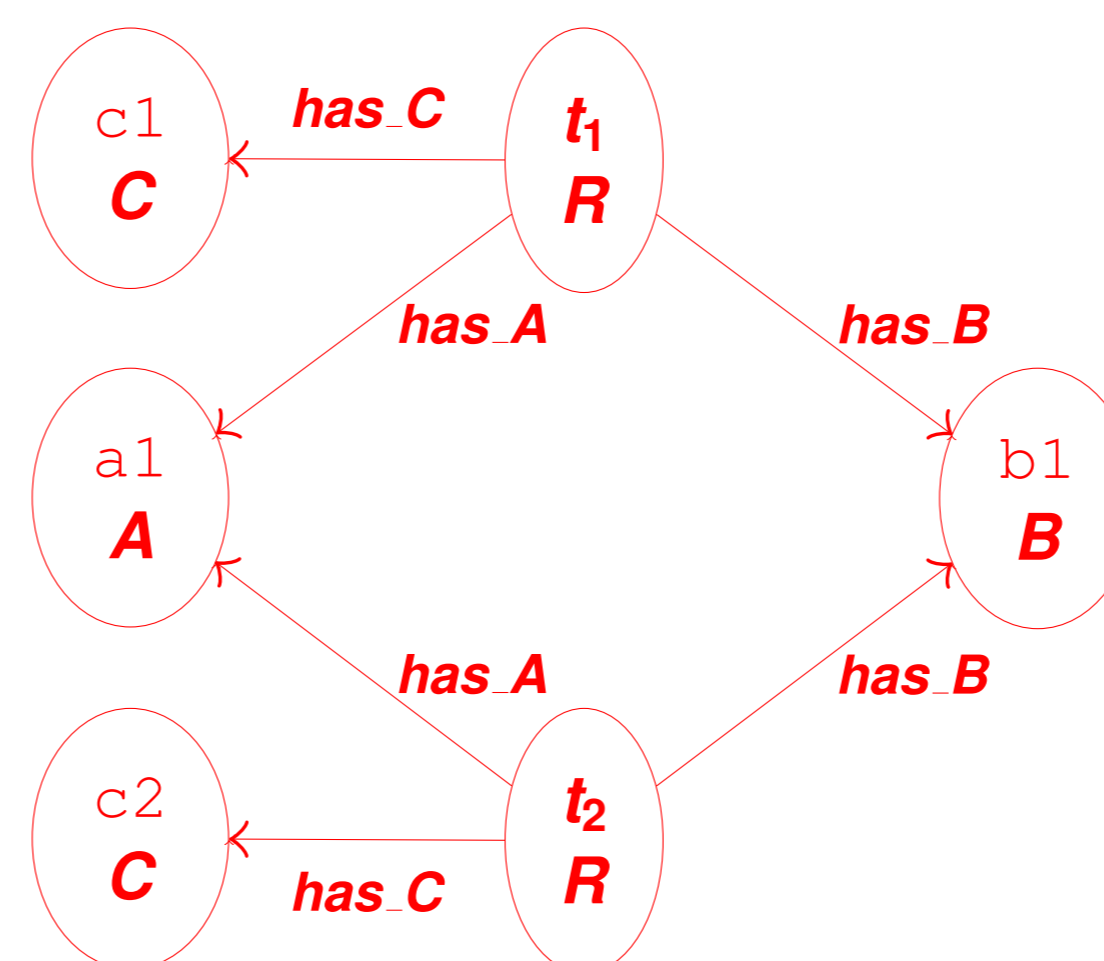
- ▶ the translation preserves satisfaction of dependencies

- ▶ Use path-based identification constraints (pids) *[Calvanese et al., KR 2008]
- ▶ There are FDs for which the above property does not hold for any set of pids.

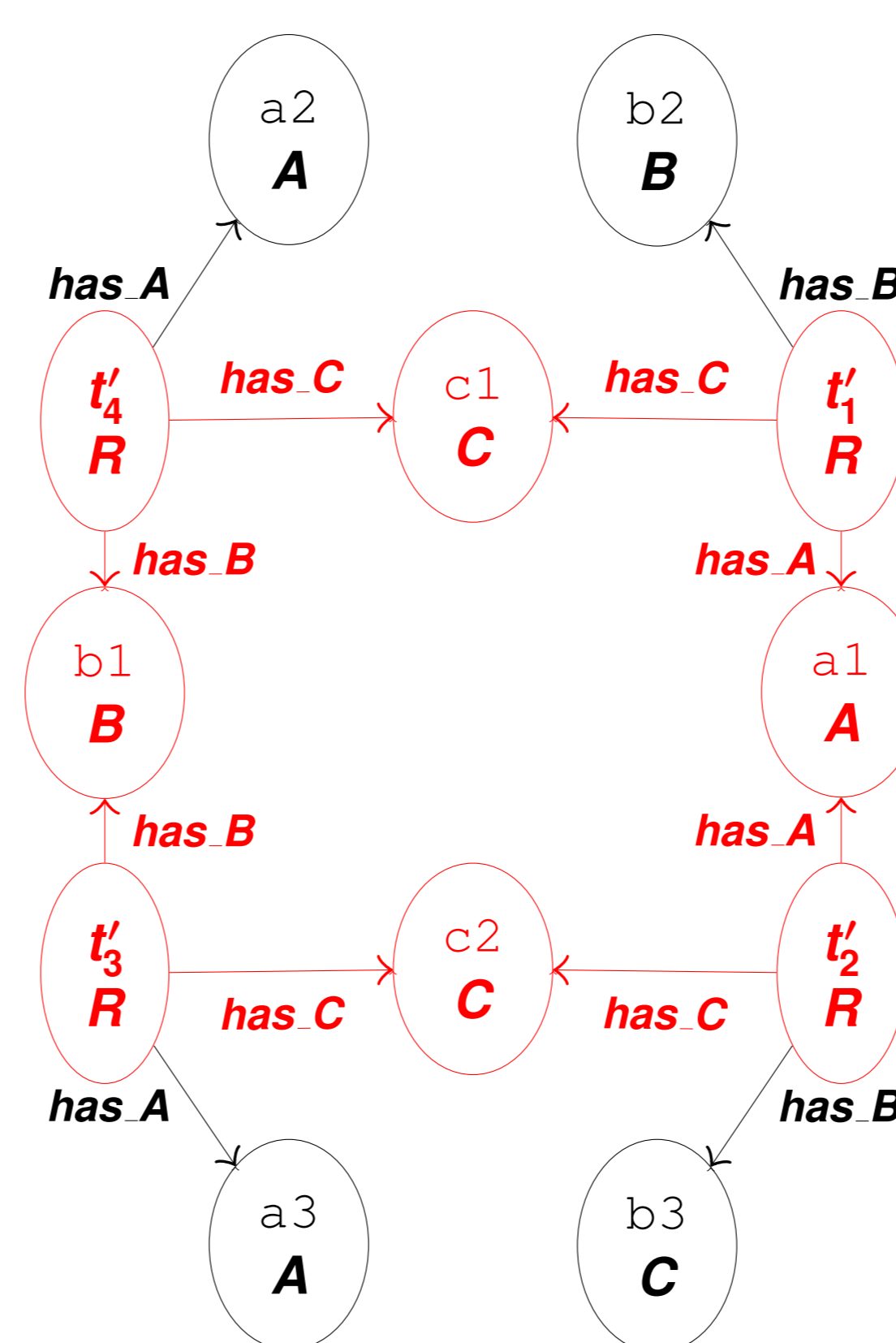
Example

- ▶ Consider the FD $\{A, B\} \rightarrow C$ and its translation to a pid $C?(has\ C^- \cdot has\ A, has\ C^- \cdot has\ B)$.
- ▶ Translation *does not* preserve satisfaction of dependencies.

	A	B	C
t_1	a1	b1	c1
t_2	a1	b1	c2



	A	B	C
t'_1	a1	b2	c1
t'_2	a2	b1	c1
t'_3	a3	b1	c2
t'_4	a4	b3	c3



FD and pid violation in red.

Normal Form for Semantic Web Data

- ▶ Boyce-Codd Normal Form (BCNF) in the Semantic Web with **tree-based identification constraints**.

- ▶ BCNF avoids update anomalies originating from redundancies induced by FDs.

Example

- ▶ $addr?(has\ addr^- \cdot has\ house)$ leads to a redundancy.
- ▶ Avoided by storing *addr* directly to the *house*.

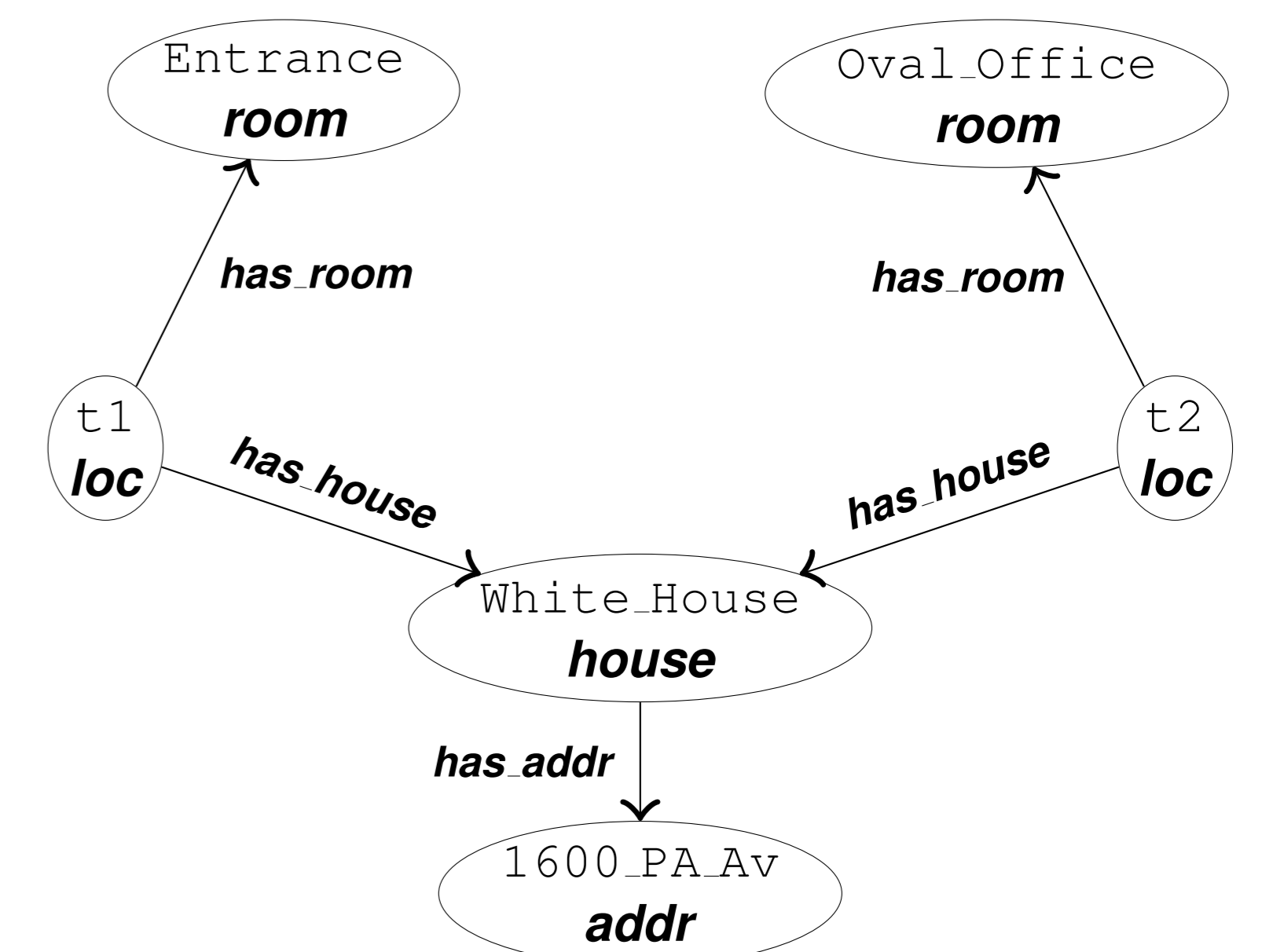


Figure 2: An RDF graph in RNF

Definition

An RDF graph is in *RDF Normal Form* (RNF) if and only if all tids can be substituted by tids with depth 1.

Theorem

A relational database is in BCNF iff the RDF graph with tids (translated with the direct mapping) is in RNF.

Theorem

Deciding whether an RDF graph is in RNF is feasible in polynomial time.

Future Work

- ▶ Extend the definition of RNF to more expressive DLs and investigate in these DLs the recognizability of RNF.
- ▶ Consider relaxations of our RNF definition.
- ▶ Investigate the complexity of reasoning with tids.
- ▶ Compare RNF with normal forms of other non-relational data formats (such as XML Normal Form).

Tree-based identification constraints (tids)

Definition

A tid is an expression constructed using the following grammar:

$$\tau ::= \epsilon \mid \sigma \cdot \tau \mid (\tau, \tau)$$

where σ is an ordinary role or a test role of the form $C?$, and C is a concept.

Example

The tid $C?(has\ C^- \cdot (has\ A, has\ B))$ captures the FD $\{A, B\} \rightarrow C$.

Theorem

Tree-based identification constraints are able to express FDs over RDF graphs.