

Expressive Languages for Querying the Semantic Web

Marcelo Arenas¹

Georg Gottlob²

Andreas Pieris²

¹Department of Computer Science, PUC Chile, Chile

²Department Computer Science, University of Oxford, UK

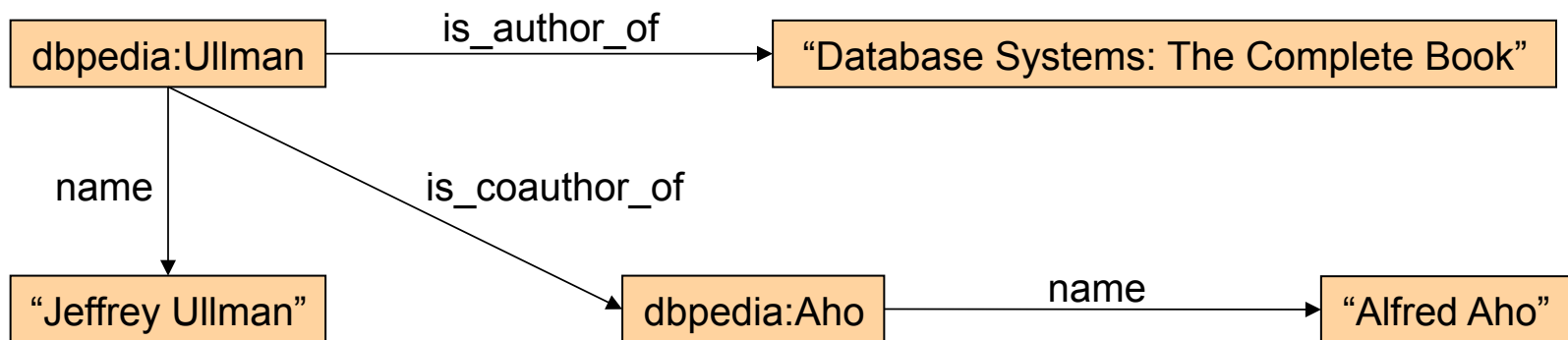
PODS 2014, Snowbird, Utah, USA

Resource Description Framework (RDF)

- **Data model** for representing information about web resources
- **Uniform Resource Identifier** (URI) - http://dbpedia.org/resource/Jeffrey_Ullman
- URIs are organized as **RDF graphs** - (*subject, predicate, object*)

Resource Description Framework (RDF)

- **Data model** for representing information about web resources
- **Uniform Resource Identifier (URI)** - http://dbpedia.org/resource/Jeffrey_Ullman
- URIs are organized as **RDF graphs** - (*subject, predicate, object*)



dbpedia: <<http://dbpedia.org/resource/>>

Resource Description Framework (RDF)

- **Data model** for representing information about web resources
- **Uniform Resource Identifier (URI)** - http://dbpedia.org/resource/Jeffrey_Ullman
- URIs are organized as **RDF graphs** - (*subject, predicate, object*)

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

dbpedia: <<http://dbpedia.org/resource/>>

SPARQL

- Graph-matching query language
- First public working draft in 2004 by W3C
- W3C recommendation in 2008

SPARQL

- Graph-matching query language
- First public working draft in 2004 by W3C
- W3C recommendation in 2008

```
SELECT ?X
WHERE {
    ?Y    is_author_of    ?Z .
    ?Y    name            ?X . }
```

SPARQL

- Graph-matching query language
- First public working draft in 2004 by W3C
- W3C recommendation in 2008

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

algebraic syntax introduced in [Pérez, Arenas & Gutierrez, TODS 2009]

SPARQL

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

SPARQL

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

Answer: "Jeffrey Ullman"

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r_1 , rdf:type, owl:Restriction)

(r_1 , owl:onProperty, is_coauthor_of)

(r_1 , owl:someValuesFrom, owl:Thing)

(r_2 , rdf:type, owl:Restriction)

(r_2 , owl:onProperty, is_author_of)

(r_2 , owl:someValuesFrom, owl:Thing)

(r_1 , rdfs:subClassOf, r_2)

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r_1 , rdf:type, owl:Restriction)

(r_1 , owl:onProperty, is_coauthor_of)

(r_1 , owl:someValuesFrom, owl:Thing)



$r_1 = \{a \mid \text{there exists a URI } b \text{ such that } (a, \text{is_coauthor_of}, b)\}$

(r_2 , rdf:type, owl:Restriction)

(r_2 , owl:onProperty, is_author_of)

(r_2 , owl:someValuesFrom, owl:Thing)



$r_2 = \{a \mid \text{there exists a URI } b \text{ such that } (a, \text{is_author_of}, b)\}$

(r_1 , rdfs:subClassOf, r_2)

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r_1 , rdf:type, owl:Restriction)

(r_1 , owl:onProperty, is_coauthor_of)

(r_1 , owl:someValuesFrom, owl:Thing)

(r_2 , rdf:type, owl:Restriction)

(r_2 , owl:onProperty, is_author_of)

(r_2 , owl:someValuesFrom, owl:Thing)

(r_1 , rdfs:subClassOf, r_2)

each co-author is an author

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r₁, rdf:type, owl:Restriction)

(r₁, owl:onProperty, is_coauthor_of)

(r₁, owl:someValuesFrom, owl:Thing)

(r₂, rdf:type, owl:Restriction)

(r₂, owl:onProperty, is_author_of)

(r₂, owl:someValuesFrom, owl:Thing)

(r₁, rdfs:subClassOf, r₂)

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

Expected answers: Jeffrey Ullman

Alfred Aho

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r₁, rdf:type, owl:Restriction)

(r₁, owl:onProperty, is_coauthor_of)

(r₁, owl:someValuesFrom, owl:Thing)

(r₂, rdf:type, owl:Restriction)

(r₂, owl:onProperty, is_author_of)

(r₂, owl:someValuesFrom, owl:Thing)

(r₁, rdfs:subClassOf, r₂)

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

Expected answers: Jeffrey Ullman ✓

Alfred Aho ✗

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r₁, rdf:type, owl:Restriction)

(r₁, owl:onProperty, is_coauthor_of)

(r₁, owl:someValuesFrom, owl:Thing)

(r₂, rdf:type, owl:Restriction)

(r₂, owl:onProperty, is_author_of)

(r₂, owl:someValuesFrom, owl:Thing)

(r₁, rdfs:subClassOf, r₂)

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

Expected answers: Jeffrey Ullman ✓

Alfred Aho ✗

**we are forced to encode the semantics
of RDFS and OWL in the query**

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r₁, rdf:type, owl:Restriction)

(r₁, owl:onProperty, is_coauthor_of)

(r₁, owl:someValuesFrom, owl:Thing)

(r₂, rdf:type, owl:Restriction)

(r₂, owl:onProperty, is_author_of)

(r₂, owl:someValuesFrom, owl:Thing)

(r₁, rdfs:subClassOf, r₂)

SELECT ?X

(?Y, rdf:type, ?Z) AND

(?Z, rdf:type, owl:Restriction) AND

(?Z, owl:onProperty, is_author_of) AND

(?Z, owl:someValuesFrom, owl:Thing)

AND (?Y, name, ?X)

Jeffrey Ullman

Alfred Aho

RDFS and OWL Vocabularies

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, name, "Jeffrey Ullman")

(dbpedia:Aho, is_coauthor_of, dbpedia:Ullman)

(dbpedia:Aho, name, "Alfred Aho")

(r₁, rdf:type, owl:Restriction)

(r₁, owl:onProperty, is_coauthor_of)

(r₁, owl:someValuesFrom, owl:Thing)

(r₂, rdf:type, owl:Restriction)

(r₂, owl:onProperty, is_author_of)

(r₂, owl:someValuesFrom, owl:Thing)

(r₁, rdfs:subClassOf, r₂)

SELECT ?X

(?Y, rdf:type, ?Z) AND

(?Z, rdf:type, owl:Restriction) AND

(?Z, owl:onProperty, is_author_of) AND

(?Z, owl:someValuesFrom, owl:Thing)

AND (?Y, name, ?X)

Jeffrey Ullman



Alfred Aho



Need for Decoupling

SELECT ?X

(?Y, is_author_of, ?Z) AND (?Y, name, ?X)

vs

SELECT ?X

(?Y, rdf:type, ?Z) AND

(?Z, rdf:type, owl:Restriction) AND

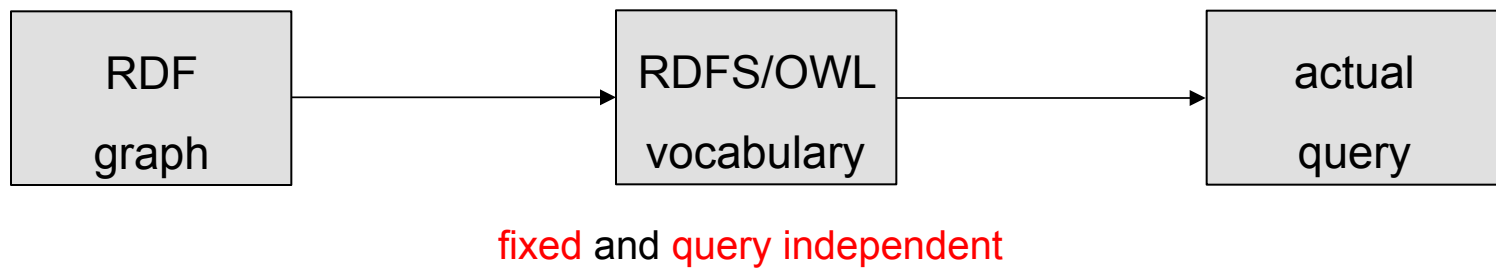
(?Z, owl:onProperty, is_author_of) AND

(?Z, owl:someValuesFrom, owl:Thing)

AND (?Y, name, ?X))

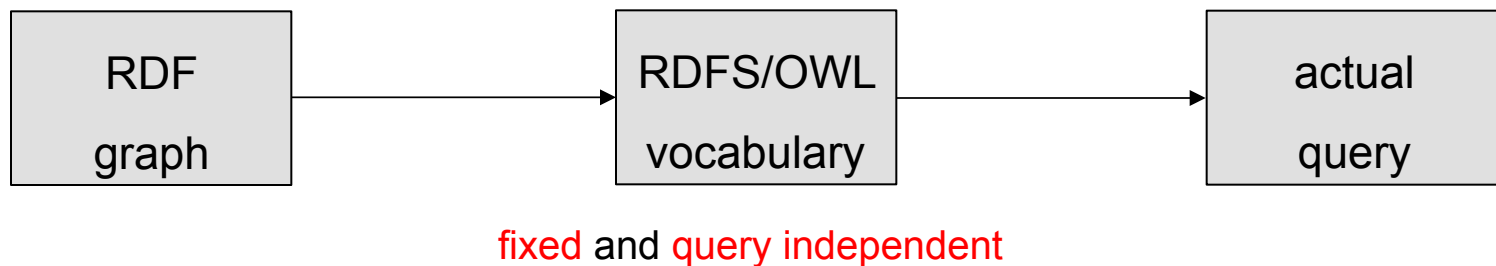
Our Objectives

- **Decouple** the reasoning part and the actual query - simpler queries



Our Objectives

- **Decouple** the reasoning part and the actual query - simpler queries



- **Navigational capabilities** - exploit the graph structure of RDF data
- General form of **recursion** - central feature for graph query languages

The rest of the Talk

- The modular query language **TriQ**
- From SPARQL over OWL 2 QL to **TriQ**
- **TriQ-Lite** - a tractable language
- Is **TriQ-Lite** really necessary?
- Concluding remarks

Triple Query Language (TriQ)

$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$

$$\text{ans}(M, G) = \text{ans}(Q_{\text{SPARQL}}, \text{ans}(Q_{\text{RDFS/OWL}}, G))$$

Triple Query Language (TriQ)

$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$

$$\text{ans}(M, G) = \text{ans}(Q_{\text{SPARQL}}, \text{ans}(Q_{\text{RDFS/OWL}}, G))$$

What is the right **syntax** for $Q_{\text{RDFS/OWL}}$ and Q_{SPARQL} ?

Triple Query Language (TriQ)

$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$

$$\text{ans}(M, G) = \text{ans}(Q_{\text{SPARQL}}, \text{ans}(Q_{\text{RDFS/OWL}}, G))$$

What is the right **syntax** for $Q_{\text{RDFS/OWL}}$ and Q_{SPARQL} ?

- **Datalog[\neg s]** represents every SPARQL query
- **Datalog[\exists, \perp]** is appropriate for ontological reasoning

Datalog[\exists, \neg s, \perp]

Triple Query Language (TriQ)

$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$

$$\text{ans}(M, G) = \text{ans}(Q_{\text{SPARQL}}, \text{ans}(Q_{\text{RDFS/OWL}}, G))$$

What is the right **syntax** for $Q_{\text{RDFS/OWL}}$ and Q_{SPARQL} ?

- **Datalog[\neg s]** represents every SPARQL query
- **Datalog[\exists, \perp]** is appropriate for ontological reasoning


Datalog[\exists, \neg s, \perp]

weakly-guarded Datalog[\exists, \neg s, \perp]

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body



$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z Q(?X, ?Z)$$

Affected positions = ?

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z Q(?X, ?Z)$$

$$\text{Affected positions} = \{T[3], P[1], Q[2]\}$$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(\textcircled{?X}, ?Y), S(?Y, ?Z) \rightarrow \exists ?W \ T(?Y, \textcircled{?X}, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W \ P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z \ Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2]\}$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2], P[2]\}$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2], P[2], Q[1],$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W \ T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W \ P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z \ Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2], P[2], Q[1],$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body

$$P(?X, ?Y), S(\cancel{?Y}, ?Z) \rightarrow \exists ?W \ T(\cancel{?Y}, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W \ P(?W, ?Y)$$


$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z \ Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2], P[2], Q[1]\}$

Weakly-Guarded Datalog[\exists, \neg s]

All body-variables at **affected positions** occur in a positive body-atom

appears in an \exists -position or just in affected positions
in the body



$$P(?X, ?Y), S(?Y, ?Z) \rightarrow \exists ?W T(?Y, ?X, ?W)$$

$$T(?X, ?Y, ?Z) \rightarrow \exists ?W P(?W, ?Y)$$

$$P(?X, ?Y), \neg R(?X) \rightarrow \exists ?Z Q(?X, ?Z)$$

Affected positions = $\{T[3], P[1], Q[2], T[2], P[2], Q[1]\}$

Weakly-Guarded Datalog[$\exists, \neg s, \perp$]

All body-variables at **affected positions** occur in a positive body-atom



appears in an \exists -position or just in affected positions
in the body

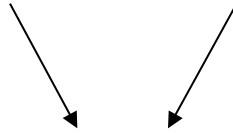
weakly-guarded Datalog[$\exists, \neg s$] + $\Phi(x_1, \dots, x_k) \rightarrow \perp$

=

weakly-guarded Datalog[$\exists, \neg s, \perp$]

Triple Query Language (TriQ)

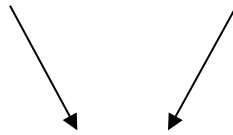
$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$



weakly-guarded Datalog[\exists , \neg s, \perp] queries

Triple Query Language (TriQ)

$$M = [Q_{\text{RDFS/OWL}}, Q_{\text{SPARQL}}]$$



weakly-guarded Datalog $[\exists, \neg, \perp]$ queries

Weakly-guarded Datalog $[\exists, \neg, \perp]$ query: (Π, Λ)

- Π is a weakly-guarded Datalog $[\exists, \neg, \perp]$ program
- Λ is a set of answer rules: $\Phi(x_1, \dots, x_k) \rightarrow \text{answer}(x_1, \dots, x_k)$

Triple Query Language (TriQ): Complexity

- **Theorem:** Query evaluation for TriQ is in **EXPTIME**

Triple Query Language (TriQ): Complexity

- **Theorem:** Query evaluation for TriQ is in **EXPTIME**
- **Theorem** [Gottlob, Rudolph & Šimkus, PODS 2014]: Every query that can be evaluated in EXPTIME **can be expressed** in weakly-guarded Datalog[\exists, \neg] (no order)

Triple Query Language (TriQ): Complexity

- **Theorem:** Query evaluation for TriQ is in **EXPTIME**
- **Theorem** [Gottlob, Rudolph & Šimkus, PODS 2014]: Every query that can be evaluated in EXPTIME **can be expressed** in weakly-guarded Datalog[\exists, \neg] (no order)
- **Corollary:** TriQ and weakly-guarded Datalog[\exists, \neg] **capture** EXPTIME (no order)

From SPARQL to TriQ

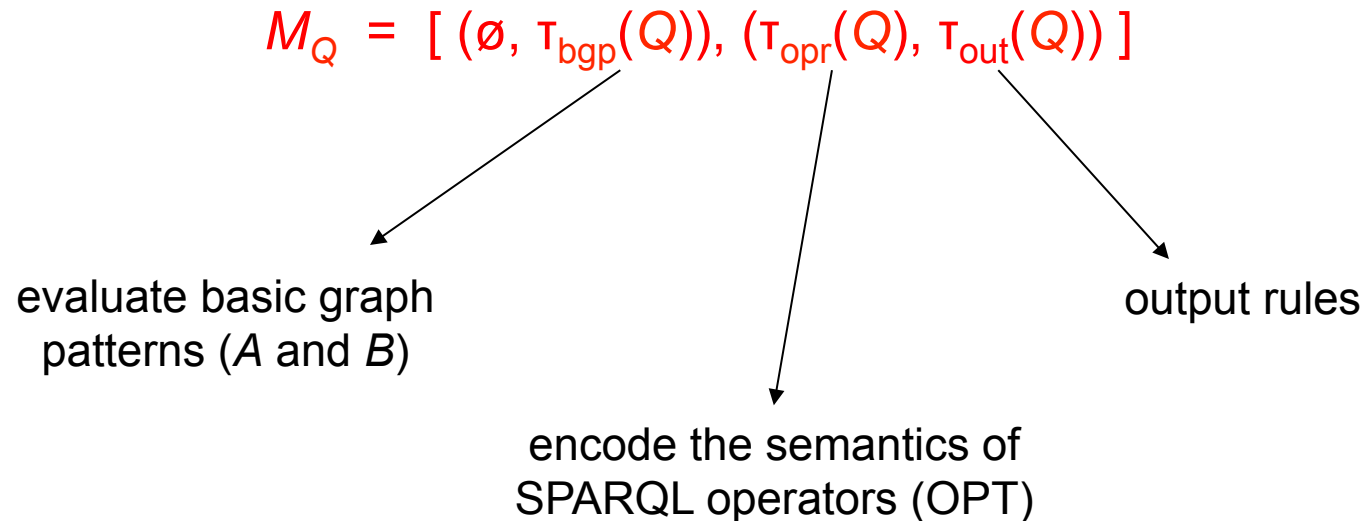
$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

for every object a , we ask for the name and the phone of a , if the phone number of a is available; otherwise, we only ask for the name of a

From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

for every object a , we ask for the name and the phone of a , if the phone number of a is available; otherwise, we only ask for the name of a



From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

The program $\tau_{\text{bgp}}(Q)$:

$\text{triple}(?X, \text{name}, ?Y) \rightarrow \text{query}_A(?X, ?Y)$

$\text{triple}(?X, \text{phone}, ?Y) \rightarrow \text{query}_B(?X, ?Y)$

From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

The program $\tau_{\text{opr}}(Q)$:

$$\text{query}_A(?X, ?Y), \text{query}_B(?X, ?Z) \rightarrow \text{query}_Q(?X, ?Y, ?Z)$$

list of individuals with phone number

$$\text{query}_A(?X, ?Y), \text{query}_B(?X, ?Z) \rightarrow \text{compatible}_Q(?X)$$

From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

The program $\tau_{\text{opr}}(Q)$:

$$\text{query}_A(?X, ?Y), \text{query}_B(?X, ?Z) \rightarrow \text{query}_Q(?X, ?Y, ?Z)$$

list of individuals with phone number

$$\text{query}_A(?X, ?Y), \text{query}_B(?X, ?Z) \rightarrow \text{compatible}_Q(?X)$$

the third argument (phone number) is missing

$$\text{query}_A(?X, ?Y), \neg \text{compatible}_Q(?X) \rightarrow \text{query}_{Q, \{3\}}(?X, ?Y)$$

From SPARQL to TriQ

$Q = (\underbrace{?X, \text{name}, ?Y}_A) \text{ OPT } (\underbrace{?X, \text{phone}, ?Z}_B)$

The program $\tau_{\text{out}}(Q)$:

$query_Q(?X, ?Y, ?Z) \rightarrow answer_Q(?X, ?Y, ?Z)$

$query_{Q,\{3\}}(?X, ?Y) \rightarrow answer_{Q,\{3\}}(?X, ?Y)$

From SPARQL to TriQ

$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$

The program $\tau_{\text{out}}(Q)$:

$query_{Q, \emptyset} (?X, ?Y, ?Z) \rightarrow answer_{Q, \emptyset} (?X, ?Y, ?Z)$

$query_{Q, \{3\}} (?X, ?Y) \rightarrow answer_{Q, \{3\}} (?X, ?Y)$

From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

$$M_Q = [(\emptyset, T_{\text{bgp}}(Q)), (T_{\text{opr}}(Q), T_{\text{out}}(Q))]$$

Given an RDF graph G :

$$\text{evaluation of } Q \text{ over } G = \text{ans}(M_Q, DB(G))$$

From SPARQL to TriQ

$$Q = \underbrace{(?X, \text{name}, ?Y)}_A \text{ OPT } \underbrace{(?X, \text{phone}, ?Z)}_B$$

$$M_Q = [(\emptyset, T_{\text{bgp}}(Q)), (T_{\text{opr}}(Q), T_{\text{out}}(Q))]$$

Given an RDF graph G ,

$$\text{evaluation of } Q \text{ over } G = \text{ans}(M_Q, DB(G))$$

\downarrow
 $\{ \text{triple}(a, b, c) \mid (a, b, c) \text{ belongs to } G \}$

From SPARQL over OWL 2 QL to TriQ

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

- Elements of G that eat something: `SELECT ?X (?X, eats, ?Y)`
- However, the answer is empty due to the **active domain semantics**

From SPARQL over OWL 2 QL to TriQ

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

- Elements of G that eat something: $SELECT ?X (?X, eats, ?Y)$
- However, the answer is empty due to the **active domain semantics**
- We need the query $(?X, rdf:type, \exists eats)$:

$G \models_{OWL2} (dog, rdf:type, \exists eats)$

- This is what is called “the evaluation of Q over G under the OWL 2 direct semantics entailment regime”

From SPARQL over OWL 2 QL to TriQ

$G = \{(\text{dog}, \text{rdf:type}, \text{animal}), (\text{animal}, \text{rdfs:subClassOf}, \exists \text{eats})\}$

$Q = (?X, \text{rdf:type}, \exists \text{eats})$

$M_Q = [(T_{\text{OWL2QL}}, T_{\text{bgp}}(Q)), (\emptyset, T_{\text{out}}(Q))]$

fixed program used to encode
the semantics F_{OWL2}

$query_Q(?X) \rightarrow answer_Q(?X)$

$triple_1(?X, \text{rdf:type}, \exists \text{eats}), dom(?X) \rightarrow query_Q(?X)$

From SPARQL over OWL 2 QL to TriQ

Theorem: Given a SPARQL query Q and an RDF graph G :

The evaluation of Q over G under the OWL 2 direct semantics entailment regime = $ans(M_Q, DB(G))$

- $M_Q = [(T_{OWL2QL}, T_{bgp}(Q)), (T_{opr}(Q), T_{out}(Q))]$ is a **TriQ query**
- T_{OWL2QL} **is fixed**, it does not depend on Q

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

- Elements of G that eat something: `SELECT ?X (?X, eats, ?Y)`
- However, the answer is empty due to the **active domain semantics**
- We need the query `(?X, rdf:type, \exists eats)`

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

$Q = (?X, rdf:type, \exists eats)$

$M_Q = [(T_{OWL2QL}, T_{bqp}(Q)), (\emptyset, T_{out}(Q))]$



$triple_1(?X, rdf:type, \exists eats), dom(?X) \rightarrow query_Q(?X)$

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

$Q = \text{SELECT } ?X \underbrace{(?X, \text{eats}, ?Y)}_P$

$M_Q = [(T_{\text{OWL2QL}}, T_{\text{bgp}}(Q)), (T_{\text{opr}}(Q), T_{\text{out}}(Q))]$

\swarrow
 $\text{triple}_1(?X, \text{eats}, ?Y),$
 $\text{dom}(?X), \text{dom}(?Y) \rightarrow \text{query}_P(?X, ?Y)$

\searrow
 $\text{query}_P(?X, ?Y) \rightarrow \text{query}_Q(?X)$

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

$Q = \text{SELECT } ?X \underbrace{(?X, \text{eats}, ?Y)}_P$

$M_Q = [(T_{\text{OWL2QL}}, T_{\text{bgp}}(Q)), (T_{\text{opr}}(Q), T_{\text{out}}(Q))]$

$\text{triple}_1(?X, \text{eats}, ?Y),$
 $\text{dom}(?X), \text{dom}(?Y) \rightarrow \text{query}_P(?X, ?Y)$

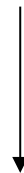
$\text{query}_P(?X, ?Y) \rightarrow \text{query}_Q(?X)$

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

$Q = \text{SELECT } ?X \text{ } (?X, \text{eats}, ?Y)$

$M_Q = [(T_{\text{OWL2QL}}, T_{\text{bgp}}(Q)), (\emptyset, T_{\text{out}}(Q))]$



$\text{triple}_1(?X, \text{eats}, ?Y), \text{dom}(?X) \rightarrow \text{query}_P(?X)$

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

$Q = \text{SELECT } ?X \text{ } (?X, \text{eats}, ?Y)$

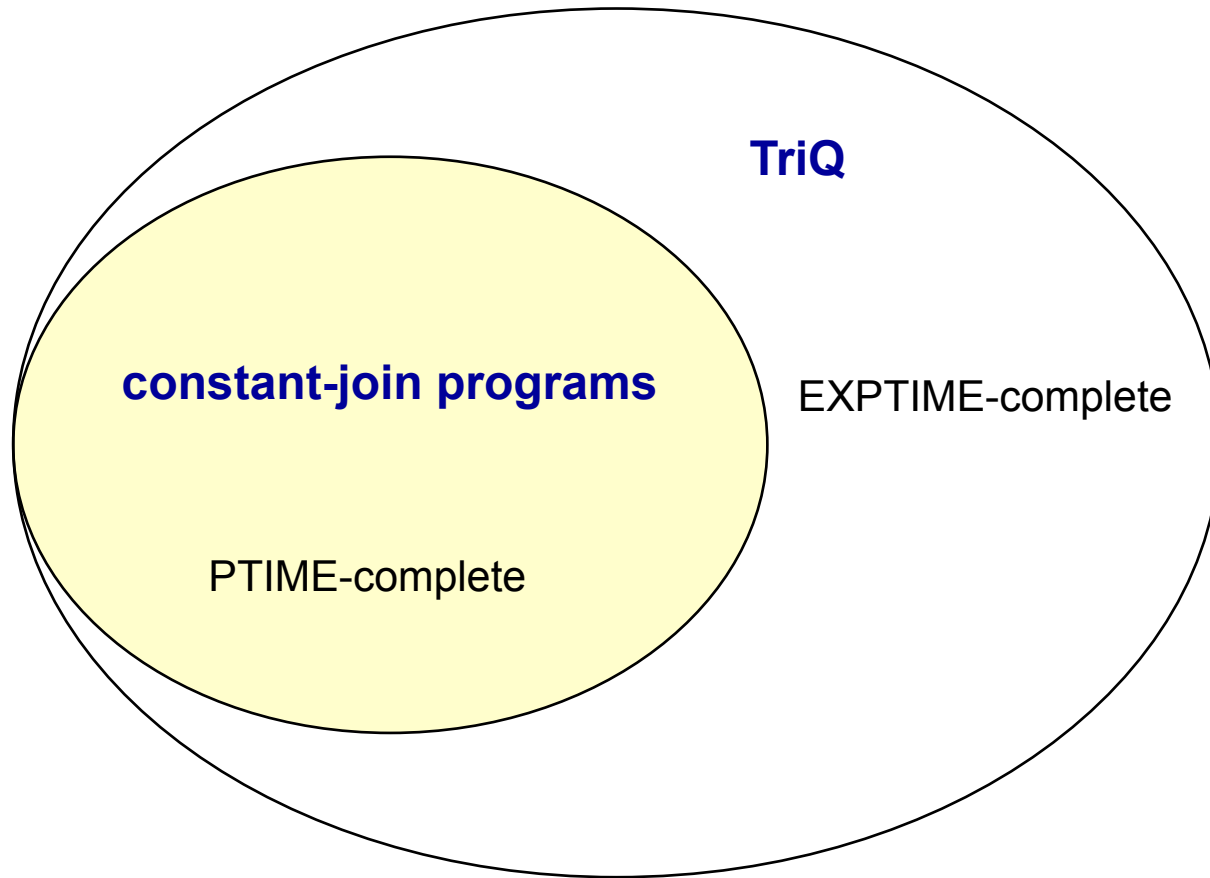
$Q = (?X, \text{eats}, _ :b)$

$M_Q = [(T_{\text{OWL2QL}}, T_{\text{bgp}}(Q)), (\emptyset, T_{\text{out}}(Q))]$



$\text{triple}_1(?X, \text{eats}, ?Y), \text{dom}(?X) \rightarrow \text{query}_P(?X)$

The Language TriQ-Lite



join variables occur only at non-affected positions

The Language TriQ-Lite

Theorem: Every SPARQL query under the entailment regime for OWL 2 QL can be expressed as a TriQ-Lite query (with or without the active domain restriction).

The Language TriQ-Lite

Theorem: Query evaluation for TriQ-Lite is **PTIME-complete**

Proof sketch:

PTIME-membership: enough to consider the evaluation problem for weakly-guarded constant-join Datalog[\exists , \neg s, \perp]

- Let D be an RDF graph and (Π, \wedge) be a weakly-guarded constant-join Datalog[\exists , \neg s, \perp] query
- First step: Constraints are eliminated from Π to generate $ex(\Pi)$
 - Checking for inconsistencies
- Second step: ground chase is computed
 - Atoms $p(a_1, \dots, a_k)$ in $chase(D, ex(\Pi))$ such that a_1, \dots, a_k are URIs
 - ALOGSPACE algorithm

The Language TriQ-Lite

Theorem: Query evaluation for TriQ-Lite is **PTIME-complete**

Proof sketch:

- Third step: negation is eliminated from $\text{ex}(\Pi)$ to generate D_L and $\text{ex}(\Pi)^+$
 - Every atom $\neg p(t_1, \dots, t_k)$ in a rule of $\text{ex}(\Pi)$ is replaced by $\text{cp}(t_1, \dots, t_k)$, where cp stores the complement of p
 - D_L is the extension of D with the atoms $\text{cp}(a_1, \dots, a_k)$ such that a_1, \dots, a_k are URIs, which are computed in each strata of $\text{ex}(\Pi)$ by using the ground chase
- Last step: transform $\text{ex}(\Pi)^+$ into a linear Datalog[\exists] program Π_L
 - Key observation: in every rule of $\text{ex}(\Pi)^+$, every variable that can be assigned non-URI values must occur only in the weak-guard.

The Language TriQ-Lite

Theorem: Query evaluation for TriQ-Lite is **PTIME-complete**

Proof sketch:

- To finish the proof: query evaluation problem for linear Datalog[\exists] is in PTIME in program complexity
 - **Program complexity:** only Λ is fixed, program Π_L and RDF graph D_L depend on Π and D

PTIME-hardness: since **Datalog** is already **PTIME-hard**

Is TriQ-Lite really necessary?

- Is **existential quantification** really necessary?
- TriQ-Lite is **Datalog rewritable** - it seems that it is not
- But, what about our main objective - **need for decoupling**?

Is TriQ-Lite really necessary?

$$M_Q = [(\tau_{\text{OWL2QL}}, \tau_{\text{bgp}}(Q)), (\tau_{\text{opr}}(Q), \tau_{\text{out}}(Q))]$$

$$M_{Q,\Pi} = [(\Pi, \tau_{\text{bgp}}(Q)), (\tau_{\text{opr}}(Q), \tau_{\text{out}}(Q))]$$

Theorem: There exists an RDF graph G and a SPARQL query Q such that, for every Datalog[\neg s, \perp] program Π :

$$\mathit{ans}(M_{Q,\Pi}, DB(G)) \neq \mathit{ans}(M_Q, DB(G))$$

Concluding remarks

1. We introduce the modular query language TriQ
 - TripQ captures EXPTIME (no order)
2. We show that every SPARQL query can be expressed in TriQ
 - Including the OWL 2 direct semantics entailment regime
 - Dropping the active domain restriction
3. We identify the tractable fragment TriQ-Lite of TriQ with the same properties as in 2.
4. We prove that the existential quantification in TriQ-Lite is necessary
 - In the paper, we define and study some other notions of program expressiveness

Thank you!

RDFS and OWL Vocabularies

Can be even worse...

(dbpedia:Ullman, is_author_of, "Database Systems: The Complete Book")

(dbpedia:Ullman, owl:sameAs, yago:Ullman)

(yago:Ullman, name, "Jeffrey Ullman")

⋮

SELECT ?X

((?Y, rdf:type, ?Z) AND

(?Z, rdf:type, owl:Restriction) AND

(?Z, owl:onProperty, is_author_of) AND

(?Z, owl:someValuesFrom, owl:Thing)

AND (?Y, name, ?X))

UNION

((?Y, is_author_of, ?Z) AND (?Y, owl:sameAs, ?W)

AND (?W, name, ?X)))

The Program τ_{OWL2QL}

- Collect the **domain elements**:

triple(?X, ?Y, ?Z) \rightarrow dom(?X), dom(?Y), dom(?Z)

- Store the **elements in the ontology**:

triple(?X, rdf:type, ?Y) \rightarrow type(?X, ?Y)

triple(?X, rdfs:subPropertyOf, ?Y) \rightarrow sp(?X, ?Y)

triple(?X, owl:inverseOf, ?Y) \rightarrow inv(?Y, ?X)

triple(?X, owl:Restriction, ?Y) \rightarrow rest(?X, ?Y)

triple(?X, rdfs:subClassOf, ?Y) \rightarrow sc(?X, ?Y)

triple(?X, owl:DisjointWith, ?Y) \rightarrow disj(?X, ?Y)

triple(?X, ?Y, ?Z) \rightarrow triple₁(?X, ?Y, ?Z)

The Program τ_{OWL2QL}

- Reason about **properties**:

$sp(?X, ?Y), inv(?Z, ?X), inv(?W, ?Y) \rightarrow sp(?Z, ?W)$

$type(?X, owl:ObjectProperty) \rightarrow sp(?X, ?X)$

$sp(?X, ?Y), sp(?Y, ?Z) \rightarrow sp(?X, ?Z)$

- Reason about **classes**:

$sp(?X, ?Y), rest(?Z, ?X), rest(?W, ?Y) \rightarrow sc(?Z, ?W)$

$type(?X, owl:Class) \rightarrow sc(?X, ?X)$

$sc(?X, ?Y), sc(?Y, ?Z) \rightarrow sc(?X, ?Z)$

- Reason about **disjointness constraints**:

$disj(?X, ?Y), sc(?Z, ?X), sc(?W, ?Y) \rightarrow disj(?Z, ?W)$

The Program τ_{OWL2QL}

- Reason about **membership assertions**:

$triple_1(?X, ?U, ?Y), sp(?U, ?V) \rightarrow triple_1(?X, ?V, ?Y)$

$triple_1(?X, ?U, ?Y), inv(?U, ?V) \rightarrow triple_1(?Y, ?V, ?X)$

$type(?X, ?Y), rest(?Y, ?U) \rightarrow \exists ?Z triple_1(?X, ?U, ?Z)$

$type(?X, ?Y), sc(?Y, ?Z) \rightarrow type(?X, ?Z)$

$type(?X, ?Y) \rightarrow triple_1(?X, rdf:type, ?Y)$

$triple_1(?X, ?U, ?Y), rest(?Z, ?U) \rightarrow type(?X, ?Z)$

$type(?X, ?Y), type(?X, ?Z), disj(?Y, ?Z) \rightarrow \perp$

Active Domain Semantics Revisited

$G = \{(dog, rdf:type, animal), (animal, rdfs:subClassOf, \exists eats)\}$

Dropping the active domain semantics in SPARQL
is non-trivial:

Consider the query $(?X, rdfs:subClassOf, ?Y)$

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

$(c, \text{rdfs:subClassOf}, \exists p)$

$(\exists p^-, \text{rdfs:subClassOf}, c)$

$(a, \text{rdf:type}, c)$

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

$(c, \text{rdfs:subClassOf}, \exists p)$

$(\exists p^-, \text{rdfs:subClassOf}, c)$

$(a, \text{rdf:type}, c)$

$c(?X) \rightarrow \exists ?Y p(?X, ?Y)$

$p(?X, ?Y) \rightarrow c(?Y)$

$c(a)$

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

(c, rdfs:subClassOf, $\exists p$)

($\exists p^-$, rdfs:subClassOf, c)

(a, rdf:type, c)

SPARQL query Q:

SELECT ?X ((?X, p, ?U) AND (?U, p, ?V)) UNION

SELECT ?X ?Y ((?X, p, ?W) AND (?W, p, ?Y))

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

(c, rdfs:subClassOf, $\exists p$)

($\exists p^-$, rdfs:subClassOf, c)

(a, rdf:type, c)

SPARQL query Q:

SELECT ?X ((?X, p, ?U) AND (?U, p, ?V)) UNION

SELECT ?X ?Y ((?X, p, ?W) AND (?W, p, ?Y))

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

(c, rdfs:subClassOf, $\exists p$)

Answer: a

($\exists p^+$, rdfs:subClassOf, c)

(a, rdf:type, c)

SPARQL query Q:

SELECT ?X ((?X, p, ?U) AND (?U, p, ?V)) UNION

SELECT ?X ?Y ((?X, p, ?W) AND (?W, p, ?Y))

Is TriQ-Lite really necessary?

Proof sketch:

RDF graph G:

(c, rdfs:subClassOf, $\exists p$)

($\exists p^{-}$, rdfs:subClassOf, c)

(a, rdf:type, c)

Answer: a

In a Datalog[\neg s, \perp] program: a is an answer if and only if there exists a URI b such that (a,b) is an answer

SPARQL query Q:

SELECT ?X ((?X, p, ?U) AND (?U, p, ?V)) UNION

SELECT ?X ?Y ((?X, p, ?W) AND (?W, p, ?Y))

Program Expressive Power

Pep: captures the expressive power of a program

Program Expressive Power

Pep: captures the expressive power of a program

$\text{Pep}_\Omega[\Pi] = \{ (D, \Lambda, p(a_1, \dots, a_k)) \mid (\Pi, \Lambda) \text{ is a program in } \Omega \text{ and } p(a_1, \dots, a_k) \text{ is in } \text{ans}((\Pi, \Lambda), D) \}$

$\text{Pep}[\Omega] = \{ \text{Pep}_\Omega[\Pi] \mid \Pi \text{ is a program in } \Omega \}$

Ω_1 is more expressive than Ω_2 if $\text{Pep}[\Omega_2]$ is a proper subset of $\text{Pep}[\Omega_1]$

Program Expressive Power

Pep: captures the expressive power of a program

$\text{Pep}_\Omega[\Pi] = \{ (D, \Lambda, p(a_1, \dots, a_k)) \mid (\Pi, \Lambda) \text{ is a program in } \Omega \text{ and } p(a_1, \dots, a_k) \text{ is in } \text{ans}((\Pi, \Lambda), D) \}$

$\text{Pep}[\Omega] = \{ \text{Pep}_\Omega[\Pi] \mid \Pi \text{ is a program in } \Omega \}$

Ω_1 is more expressive than Ω_2 if $\text{Pep}[\Omega_2]$ is a proper subset of $\text{Pep}[\Omega_1]$

Theorem: weakly-guarded constant-join Datalog $[\exists, \neg s, \perp]$ is more expressive than Datalog $[\neg s, \perp]$