

Bridging a gap between static analysis and ontology-based reasoning

Filip Murlak
University of Warsaw

Workshop on Foundations of Databases and AI 2025 @ TU Wien

Databases

Data (tables, trees, graphs) – relational structures/models

Queries (SQL, XPath, Cypher) – formulas with free variables

Metadata (schemas, integrity constraints) – theories

Goal Evaluate queries over data, guided by metadata.

Databases

Data (tables, trees, graphs) – relational structures/models

Queries (SQL, XPath, Cypher) – formulas with free variables

Metadata (schemas, integrity constraints) – theories

Goal Evaluate queries over data, guided by metadata.

Problem (Query containment modulo schema)

Given queries P and Q , and a schema \mathcal{S} , decide if $P \subseteq_{\mathcal{S}} Q$; that is,

for every database D , if $D \models \mathcal{S}$ and $D \models P$, then $D \models Q$.

Knowledge representation

Facts (ABoxes) - ground atomic formulas (a CQ with no variables)

Ontologies (TBoxes, rules) - theories

Queries (SPARQL) - formulas with free variables

Goal Answer queries based on facts derived using ontologies.

Knowledge representation

Facts (ABoxes) - ground atomic formulas (a CQ with no variables)

Ontologies (TBoxes, rules) - theories

Queries (SPARQL) - formulas with free variables

Goal Answer queries based on facts derived using ontologies.

Problem (Query entailment)

Given facts \mathcal{A} , an ontology \mathcal{T} , and a query Q , decide if $\mathcal{T}, \mathcal{A} \models Q$; that is,
for every interpretation I , if $I \models \mathcal{T}$ and $I \models \mathcal{A}$, then $I \models Q$.

Two problems, or one?

Problem (Query containment modulo schema)

Given queries P and Q , and a schema \mathcal{S} , decide if $P \subseteq_{\mathcal{S}} Q$; that is,

for every database D , if $D \models \mathcal{S}$ and $D \models P$, then $D \models Q$.

Problem (Query entailment)

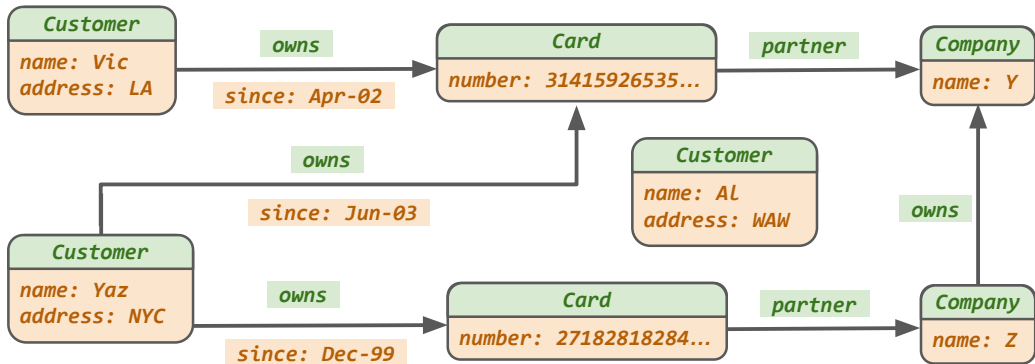
Given facts \mathcal{A} , an ontology \mathcal{T} , and a query Q , decide if $\mathcal{T}, \mathcal{A} \models Q$; that is,

for every interpretation I , if $I \models \mathcal{T}$ and $I \models \mathcal{A}$, then $I \models Q$.

$$D \approx I \quad P \approx \mathcal{A} \quad \mathcal{S} \approx \mathcal{T}$$

[Calvanese, De Giacomo, Lenzerini '98]

Graphs



Graphs

We work with labelled graphs, modelled as relational structures:

- ▶ unary predicates = node labels = concept names A, B, \dots
- ▶ binary predicates = edge labels = role names r, s, \dots

That is,

- ▶ nodes have multiple labels;
- ▶ edges have single labels;
- ▶ parallel edges with different labels are allowed;
- ▶ in a subgraph, nodes may have fewer labels.

Queries

- ▶ Conjunctive queries (CQs), unions of CQs (UCQs)

$$\exists x \exists y A(x) \wedge r(x, y) \wedge \bar{A}(y) \quad \vee \quad \exists x \exists y \exists z r(x, y) \wedge r(y, z) \wedge r(z, x)$$

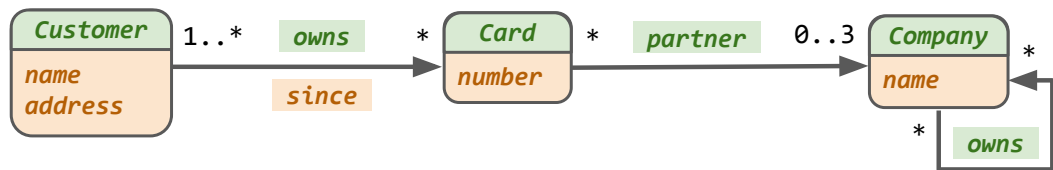
The core of relational query languages, such as SQL.

- ▶ Conjunctive regular path queries (CRPQs), unions of CRPQs (UCRPQs)

$$\exists x r^+(x, x) \quad \vee \quad \exists x \exists y A(x) \wedge (r^* \cup s)(x, y) \wedge (p \cdot B)^*(y, x)$$

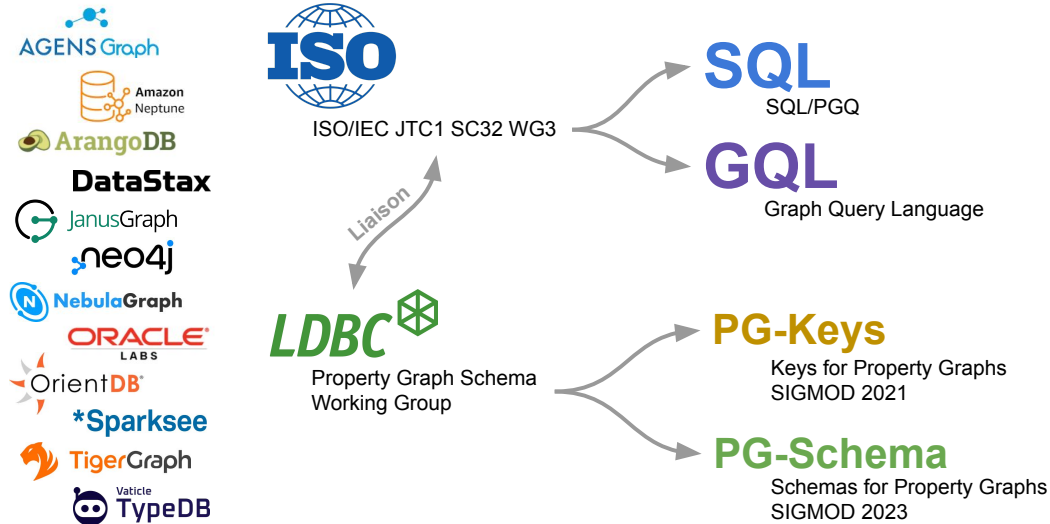
Graph query languages have reachability/RPQs at the core (SPARQL, Neo4j's Cypher, SQL/PGQ, GQL).

Schemas



What kind of schemas for graph data?

RDF has SHACL and ShEx, for property graphs the picture is less clear.



Common Foundations for SHACL, ShEx, and PG-Schema

Shqiponja Ahmetaj
shqiponjaahmetaj@tuwien.ac.at
TU Wien
Vienna, Austria

Katja Hose
katja.hose@tuwien.ac.at
TU Wien
Vienna, Austria

Wim Martens
wim.martens@uni-bayreuth.de
University of Bayreuth
Bayreuth, Germany

Cem Okulmus
okulmus@mail.uni-paderborn.de
Paderborn University
Paderborn, Germany

Mantas Šimkus
mantas.simkus@tuwien.ac.at
TU Wien
Vienna, Austria

Iovka Boneva
iovka.boneva@univ-lille.fr
Univ. Lille, CNRS, Inria, Centrale Lille,
UMR 9189 CRIStAL
Lille, France

Maxime Jakubowski
maxime.jakubowski@tuwien.ac.at
TU Wien
Vienna, Austria

Fabio Mogavero
fabio.mogavero@unina.it
Università di Napoli Federico II
Naples, Italy

Axel Polleres
axel.polleres@wu.ac.at
WU Wien
Vienna, Austria
CSH Vienna
Vienna, Austria

Dominik Tomaszuk
d.tomaszuk@uwb.edu.pl
TU Wien
Vienna, Austria
University of Białystok
Białystok, Poland

Jan Hidders
j.hidders@bbk.ac.uk
Birkbeck, University of London
London, UK

Jose Emilio Labra Gayo
labra@uniovi.es
University of Oviedo
Oviedo, Spain

Filip Murlak
f.murlak@uw.edu.pl
University of Warsaw
Warsaw, Poland

Ognjen Savković
ognjen.savkovic@unibz.it
Free University of Bolzano
Bolzano, Italy

ABSTRACT

Graphs have emerged as an important foundation for a variety of applications, including capturing and reasoning over factual knowledge, semantic data integration, social networks, and providing factual knowledge for machine learning algorithms. To formalise certain properties of the data and to ensure data quality, there is a need to describe the schema of such graphs. Because of the breadth of applications and availability of different data models, such as RDF and property graphs, both the Semantic Web and the database community have independently developed graph schema languages: SHACL, ShEx, and PG-Schema. Each language has its unique approach to defining constraints and validating graph data, leaving

potential users in the dark about their commonalities and differences. In this paper, we provide formal, concise definitions of the core components of each of these schema languages. We employ a uniform framework to facilitate a comprehensive comparison between the languages and identify a common set of functionalities, shedding light on both overlapping and distinctive features of the three languages.

ACM Reference Format:

Shqiponja Ahmetaj, Iovka Boneva, Jan Hidders, Katja Hose, Maxime Jakubowski, Jose Emilio Labra Gayo, Wim Martens, Fabio Mogavero, Filip Murlak, Cem Okulmus, Axel Polleres, Ognjen Savković, Mantas Šimkus, and Dominik Tomaszuk. 2025. Common Foundations for SHACL, ShEx, and PG-Schema. In *Proceedings of the ACM Web Conference 2025 (WWW '25)*, April 28–May 2, 2025, Sydney, NSW, Australia. ACM, New York, NY, USA, 23 pages. <https://doi.org/10.1145/XXXXXX.XXXXXX>

1 INTRODUCTION

Driven by the unprecedented growth of inter connected data, graph-based data representations have emerged as an expressive and versatile framework for modelling and analysing connections in data sets [46]. This rapid growth however, has led to a proliferation of diverse approaches, each with its own identity and perspective.

Ahmetaj, Boneva, Hidders, Hose, Jakubowski, Labra Gayo, Martens, Mogavero, Murlak, Okulmus, Polleres, Savković, Šimkus, Tomaszuk, *Common Foundations for SHACL, ShEx, and PG-Schema*, WWW 2025 (to appear).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions.acm.org.
WWW '25, April 28–May 2, 2025, Sydney, NSW, Australia.
© 2025 Copyright held by the owner/authors. Publication rights licensed to ACM.
ACM ISBN 978-0-407-42744-2/25/04
<https://doi.org/10.1145/XXXXXX.XXXXXX>

Description logics

Basic description logic \mathcal{ALC} has statements of the form

$$C \sqsubseteq D$$

where C, D are *complex concepts* build according to the following grammar by

$$C, D ::= \perp \mid \top \mid A \mid C \sqcup D \mid C \sqcap D \mid \neg C \mid \exists r.C \mid \forall r.C.$$

For example, $\text{Person} \sqsubseteq \exists \text{childOf. Person}$ and $\text{Male} \sqcap \exists \text{childOf. Person} \sqsubseteq \text{Son}$.

\mathcal{ALC} in normal form:

$$A_1 \sqcap A_2 \sqcap \dots \sqcap A_n \sqsubseteq B_1 \sqcup B_2 \sqcup \dots \sqcup B_m$$

$$A \sqsubseteq \exists r.B$$

$$A \sqsubseteq \forall r.B$$

$$\text{empty} \sqcap = \top$$

$$\text{empty} \sqcup = \perp$$

A zoo of description logics

Features can be added to *ALC*, giving logics like *ALCHOIQ* or *SOQ* (*S* is shorthand for *ALCS*).

O: use constants as singleton concepts $\{a\}$

$A \sqsubseteq \exists r. \{a\}$

I: use inverse roles r^- anywhere

$A \sqsubseteq \exists r^-. B$

F: declare role r to be a partial function

$\text{fun}(r)$

Q: use counting quantifiers $\exists^{\leq n}$, $\exists^{\geq n}$

$A \sqsubseteq \exists^{\leq 5} r. B$

S: declare role r to be transitive

$\text{tra}(r)$

H: declare role r to be contained in role s

$r \sqsubseteq s$

ALCQI can express EER, and a lot of SHACL and PG-Schema.

Finite vs unrestricted models

Problem (Query containment modulo schema)

Given queries P and Q , and a schema \mathcal{S} , decide if $P \subseteq_{\mathcal{S}} Q$; that is,

for every *finite* database D , if $D \models \mathcal{S}$ and $D \models P$, then $D \models Q$.

Problem (Query entailment)

Given facts \mathcal{A} , an ontology \mathcal{T} , and a query Q , decide if $\mathcal{T}, \mathcal{A} \models Q$; that is,

for every *possibly infinite* interpretation I , if $I \models \mathcal{T}$ and $I \models \mathcal{A}$, then $I \models Q$.

A paradox

Does $\mathcal{T} = \{\text{Person} \sqsubseteq \exists \text{childOf. Person}\}$ model reality well?

A paradox

Does $\mathcal{T} = \{\text{Person} \sqsubseteq \exists \text{childOf. Person}\}$ model reality well?

Suppose we know that at least one person exists: $\mathcal{A} = \{\text{Person}(\text{filip})\}$.

A paradox

Does $\mathcal{T} = \{\text{Person} \sqsubseteq \exists \text{childOf. Person}\}$ model reality well?

Suppose we know that at least one person exists: $\mathcal{A} = \{\text{Person}(\text{filip})\}$.

Then, over finite models we can conclude that somebody is their own ancestor...

$$\mathcal{T}, \mathcal{A} \models_{\text{fin}} \exists x \text{childOf}^+(x, x)$$

Over unrestricted models we are safe:

$$\mathcal{T}, \mathcal{A} \not\models \exists x \text{childOf}^+(x, x)$$

Query entailment

Landscape of query entailment

Entailment of CQs is **EXPTIME**-complete for

- ▶ $ALCH$ [Ortiz, Šimkus, Eiter 2008], $ALCHQ$ [Lutz 2008]

2EXPTIME-complete for

- ▶ $ALCI$ [Lutz 2008], $ALCO$ [Ngo, Ortiz, Šimkus 2016]
- ▶ SH [Eiter, Lutz, Ortiz, Šimkus '10], S [Ibáñez-García, Jung, Michielini, M. '25],
- ▶ $SHIQ^r$ [Calvanese, Eiter, Ortiz 2007] [Glimm, Lutz, Horrocks, Sattler 2008]
- ▶ $SHOQ^r$ [Glimm, Horrocks, Sattler 2008]
- ▶ SOQ^u [Gogacz, Gutiérrez-Basulto, Ibáñez-García, Jung, Murlak 2019]

decidable for

- ▶ $ALCHOIQb$ [Glimm, Rudolph 2010]

undecidable for

- ▶ SHQ^u [Horrocks, Sattler, Tobies 2000], SIQ^u [Kazakov, Sattler, Zolin 2007]
- ▶ $SHOIQ^r$ and $SHOIF$ [Rudolph 2016]

Landscape of query entailment: finite controllability

Finite controllability: $\models = \models_{\text{fin}}$

GF [Bárány, Gottlob, Otto 2014] (covers *ALCHOIb*)

ALCOF [Gogacz, Ibáñez-García, Murlak 2018]

Landscape of query entailment: finite controllability

Finite controllability: $\models = \models_{\text{fin}}$

\mathcal{GF} [Bárány, Gottlob, Otto 2014] (covers $\mathcal{ALCHOIb}$)

\mathcal{ALCOF} [Gogacz, Ibáñez-García, Murlak 2018]

No finite controllability: $\models \neq \models_{\text{fin}}$

\mathcal{ALCIF}

$A(a)$

$\top \sqsubseteq \exists r. \neg A$

$\top \sqsubseteq \exists^{\leq 1} r. \top$

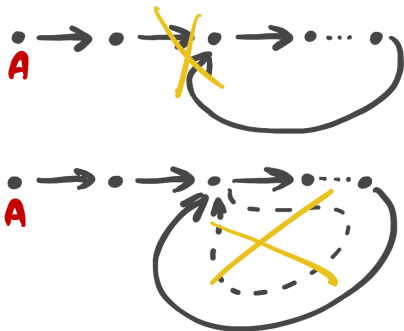
\mathcal{S}

$A(a)$

$\top \sqsubseteq \exists r. \neg A$

$\text{tra}(r)$

$\exists x r(x, x)$



What about \mathcal{ALCOQ} and up to $\mathcal{ALCHOQb}$?

Landscape of finite entailment: conjunctive queries

Finite entailment of CQs is $2EXPTIME$ -complete for

- ▶ GC^2 [Pratt-Hartmann 2009]
- ▶ GF [Bárány, Gottlob, Otto 2014]
- ▶ SOI and SIF [Gogacz, Ibáñez-García, Murlak 2018]
- ▶ $SHOI^\square$ [Danielski, Kieroński 2019]
- ▶ SOQ^u [Gogacz, Gutiérrez-Basulto, Ibáñez-García, Jung, Murlak 2019]

undecidable for

- ▶ SHQ^u and SIQ^u [Kazakov, Sattler, Zolin 2007]
- ▶ $SHOIF$ [Rudolph 2016]

a challenge for

- ▶ $ALCOIF$

Landscape of finite entailment: UCRPQs and UC2RPQs

Finite entailment of UCRPQ is $2EXPTIME$ -complete for

- ▶ $ALCI$ and $ALCQ$ [Gutiérrez-Basulto, Gutowski, Ibáñez-García, M. '22,'24]
 - ▶ should extend to $ALCOI$ and $ALCOQ$

a challenge for

- ▶ $ALCIQ$

Finite entailment of UC2RPQs (two-way UCRPQs) is **undecidable** for

- ▶ $ALCOIF$ [Rudolph 2016]

a challenge for

- ▶ ALC

Query containment

Query containment without schema is well understood

UCQs: NP-complete

[Chandra, Merlin '77]

UC2RPQs: EXPSPACE-complete

[Florescu, Levy, Suciu '98]

[Calvanese, De Giacomo, Lenzerini, Vardi '00]

Fragments complete for NP, co-NP^{NP}, PSPACE

[Deutsch, Tannen '02]

[Figueira, Godbole, Krishna, Martens, Niewerth, Trautner '20]

Dichotomy between EXPSPACE-hard and PSPACE-easy

[Figueira '20]

Query containment modulo schema is mostly open

Containment of UCQs modulo full dependencies [Beeri, Vardi '84]

Containment of UCQs w. reachability modulo full dependencies w. reachability
[Deutsch, Tannen '01]

Containment of UC2RPQs in acyclic UC2RPQs modulo Horn *ALC1Q*
[Boneva, Groz, Hidders, Murlak, Staworko '23]

Strong results on query containment modulo schema/constraints over infinite graphs.
[Calvanese, De Giacomo, Lenzerini '98,'08] [Calvanese, Ortiz, Šimkus '11]

Step 1: reduce containment to finite entailment

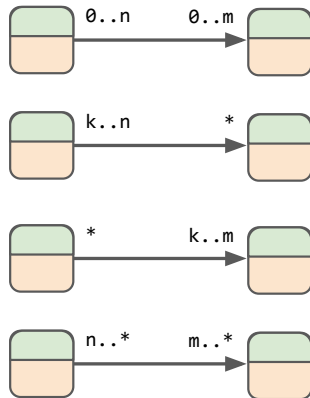
Theorem (Gutiérrez-Basulto, Gutowski, Ibáñez-García, Murlak 2024)

*Containment of UC2RPQs modulo **one-way** schemas reduces to finite entailment.*

- ▶ One-way schemas do not mix forward at-least and backward at-most constraints (and vice versa).
- ▶ Solving one instance of containment requires multiple instances of entailment.
- ▶ All instances of entailment involve single-node input graphs.

Corollary (from the proof)

*Containment of UC2RPQs modulo schemas **without at-least constraints** is decidable in 2EXPTIME .*

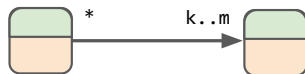


Step 2: use finite entailment

Theorem (Gutiérrez-Basulto, Gutowski, Ibáñez-García, Murlak 2024)

Query containment is **2EXPTIME**-complete for

- (a) *UCRPQs and one-way schemas,*
- (b) *simple UC2RPQs and forward schemas.*



- ▶ Forward schemas do not use backward at-least and at-most constraints.
- ▶ Simple UC2RPQs do not use concatenation in regular expressions.

Conclusion

Summary and take-away

- ▶ Query containment and finite entailment are closely related problems.
- ▶ Rich body of techniques and results in DLs that can be potentially reused, because DLs are pretty good at capturing relevant schema information.
- ▶ Capturing more refined schemas requires fancier logics, new results are needed. Looks challenging, but this is what we like, isn't it?
- ▶ Complexity is high in general, but there's space for tractable special cases. The more we know about user's needs, the better we can tailor the algorithms.
- ▶ Other ways to build a bridge: closed predicates and mixed models.

Sometimes infinity is an oversimplification

