Referring Expressions in Articial Intelligence and **Knowledge Representation Systems**

David Toman‡

(joint work with Alexander Borgida[†] and Grant Weddell[‡])



[†]Department of Computer Science Rutgers University, New Brunswick, USA borgida@cs.rutgers.edu



Waterloo [‡]Cheriton School of Computer Science University of Waterloo, Canada {david, qweddell}@uwaterloo.ca

IDENTIFYING AND COMMUNICATING REFERENCES

(TO OBJECTS/ENTITIES)





(Real world) Entities vs. (Computer) Representation(s)

Problem

- Information systems store information about entities
- Computers store (arrays of) ints and strings

How do we bridge the GAP?



(Real world) Entities vs. (Computer) Representation(s)

Problem

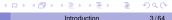
- Information systems store information about *entities*
- Computers store (arrays of) ints and strings

How do we bridge the GAP?

Typical solutions:

- OIDs (proxying entity *identity* by a number uniformly in the whole system)
 - ⇒ typically managed by *The System* (OO languages), or
- Keys (proxying entity identity by a unique combination of values (local))
 - ⇒ typically declared/managed by user (Relational DBMS).





a.k.a. proxying identities by values in a data type (say int)





4/64

a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

every object WILL have an OID (say 64 bits)

 \Rightarrow storage/performance overhead (need to be generated/managed)

can we proxy by (storage) address?

what about memory/storage reuse and/or garbage collection??

what about data replication??





a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

Information Integration: The CORBA Case

What happens to an *object* stored in *different ORBs*??

⇒ what does CORBA::Object::is_equivalent(in Object) do??



a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

Information Integration: The CORBA Case

What happens to an *object* stored in *different ORBs*??

⇒ what does CORBA::Object::is_equivalent(in Object) do??

... and before someone mentions URL/URI/IRIs:





a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

Information Integration: The CORBA Case

Unintuitive Answers: RDF/Freebase/... Cases

Freebase The (object id of the) "Synchronicity" album by "The Police" is /quid/9202a8c04000641f8000000002f9e349

(as of April, 2015.)

W3C URI/IRI/... do not improve the situation

⇒ and RDF *introduces* additional internal identifiers!





a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

Information Integration: The CORBA Case

Unintuitive Answers: RDF/Freebase/... Cases

Missing (implied) Answers: The OBDA Case

In the presence of *background knowledge* we may *know* that certain objects exist, but we cannot identify/report them due to lack of an *explicit identifier* (example later)





a.k.a. proxying identities by values in a data type (say int)

Performance: The PROTEL2 Case

Information Integration: The CORBA Case

Unintuitive Answers: RDF/Freebase/... Cases

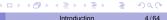
Missing (implied) Answers: The OBDA Case

Alternative Preferred Answers

Internal (computer) addresses vs. physical locations of equipment

- ⇒ programs need electronic address (to route the electric signals)
- ⇒ technicians need physical location (to find the equipment)





Relational Keys





Goal of the Tutorial

Goal

Introduce *referring expressions* as an uniform approach to identification of entities in information systems.



Goal of the Tutorial

Goal

Introduce *referring expressions* as an uniform approach to identification of entities in information systems.



Outline

- Referring Expressions in Philosophy/Linguistics
- Logical Foundations: Single Interpretations vs. Models of Theories
- Use of Referring Expressions in Information Systems
 - Referring Expressions in Answers to Queries over Knowledge Bases
 - Referring Expressions for recording Ground Knowledge
 - 3 Referring Expressions in Conceptual Design
- Summary and Open Problems

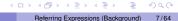




REFERRING EXPRESSIONS

(BACKGROUND)





What is an Referring Expression?

Referring Expression

A referring expression in linguistics is any noun phrase identifying an object in a way that will be useful to interlocutors.

What is an Referring Expression?

Referring Expression

A referring expression in linguistics is any noun phrase identifying an object in a way that will be useful to interlocutors.

Russell: "On Denoting," Mind, New Series, Vol.14, No.56, pp. 479–493, 1905.

A definite description "the F is a G" is understood to have the form

$$\exists x. F(x) \land \forall y (F(y) \rightarrow x = y) \land G(x)$$

A definite description is a denoting phrase in the form of "the F" where F is a noun-phrase or a singular common noun. The definite description is proper if F applies to a unique individual or object.



What is an Referring Expression?

Referring Expression

A referring expression in linguistics is any noun phrase identifying an object in a way that will be useful to interlocutors.

Russell: "On Denoting," Mind, New Series, Vol.14, No.56, pp. 479–493, 1905.

A definite description "the F is a G" is understood to have the form

$$\exists x. F(x) \land \forall y (F(y) \rightarrow x = y) \land G(x)$$

A definite description is a denoting phrase in the form of "the F" where F is a noun-phrase or a singular common noun. The definite description is proper if F applies to a unique individual or object.

The discussion of *definite* and *indefinite* descriptions (in English, phrases of the form 'the F' and 'an F') has been at the centre of analytic philosophy for over a century (so we won't go there today!).





Issues and Criticisms

Referring to Non-existing Object:

"The King of Kentucky (is...)" [Strawson] (object does NOT exist in this interpretation? or *in principle*?)

Referring to Object in Context:

"The table (is covered with books)" (non-unique reference without assuming additional context)

Multiple Reference:

"The Morning Star" vs. "The Evening Star" [Frege] (multiple distinct references to the same object)

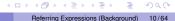
Rigidity:

Should referring expressions identify the same object in all possible worlds? [Kripke, S.: Identity and Necessity, In Identity and Individuation. NYU Press, pp. 135-164 (1971)]





REFERRING EXPRESSIONS AND (LOGICAL) THEORIES



How do we *communicate* Results of Queries?

Typical solution: tuples of *constant symbols* that, when substituted for free variables, make a query logically implied by the Knowledge Base.

How do we *communicate* Results of Queries?

Typical solution: tuples of *constant symbols* that, when substituted for free variables, make a query logically implied by the Knowledge Base.

- only explicitly named objects are returned as certain answers
- 2 often system-generated ids (that aren't too user-friendly)

How do we *communicate* Results of Queries?

Typical solution: tuples of *constant symbols* that, when substituted for free variables, make a query logically implied by the Knowledge Base.

- only explicitly named objects are returned as certain answers
- 2 often system-generated ids (that aren't too user-friendly)

Example (Freebase)

The (object id of the) "Synchronicity" album by "The Police" is /quid/9202a8c04000641f8000000002f9e349 (as of April, 2015.)



How do we *communicate* Results of Queries?

Typical solution: tuples of *constant symbols* that, when substituted for free variables, make a query *logically implied* by the Knowledge Base.

- only explicitly named objects are returned as certain answers
- often system-generated ids (that aren't too user-friendly)

Example (Freebase)

The (object id of the) "Synchronicity" album by "The Police" is /guid/9202a8c04000641f8000000002f9e349 (as of April, 2015.)

Referring Expressions

More answers (e.g., objects *without* explicit name), and/or more informative/*preferred* answers, e.g.:

 $ALBUM(x) \land (title(x) = "Synchronicity") \land (band(x) = "The Police")$



How do we communicate Results of Queries?

Typical solution: tuples of *constant symbols* that, when substituted for free variables, make a query *logically implied* by the Knowledge Base.

- only explicitly named objects are returned as certain answers
- often system-generated ids (that aren't too user-friendly)

Example (Freebase)

The (object id of the) "Synchronicity" album by "The Police" is /guid/9202a8c04000641f8000000002f9e349 (as of April, 2015.)

Referring Expressions

More answers (e.g., objects *without* explicit name), and/or more informative/*preferred* answers, e.g.:

 $ALBUM \sqcap (title = "Synchronicity") \sqcap (band = "The Police")$





Bottom Line

Referring Expressions

Formulæ $\phi\{x\}$ (in the language of the Knowledge Base)

- with *exactly one free variable* (x) that are
- *singular* with respect to a Knowledge Base \mathcal{K} , i.e.,

$$|\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| = 1$$

for all models \mathcal{I} of \mathcal{K} .

⇒ this intuition may be refined w.r.t. queries (e.g., singular among answers)



Bottom Line

Referring Expressions

Formulæ $\phi\{x\}$ (in the language of the Knowledge Base)

- with exactly one free variable (x) that are
- ${f z}$ singular with respect to a Knowledge Base ${\cal K}$, i.e.,

$$|\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| = 1$$

for all models \mathcal{I} of \mathcal{K} .

⇒ this intuition may be refined w.r.t. queries (e.g., singular *among answers*)

Why not terms?

Terms (with the standard FO semantics) suffer from *totality*

⇒ must denote *something* in *every* interpretation





Denoting/Non-denoting Referring Expressions

Singularity Revisited

$$\begin{split} |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 1 \quad \text{a denoting refering expression} \\ |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 0 \quad \text{a non-denoting refering expression} \\ \text{for all } \mathcal{I} \models \mathcal{K}; \text{ these ought to be the only two possibilities!} \end{split}$$

Denoting/Non-denoting Referring Expressions

Singularity Revisited

$$\begin{split} |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 1 \quad \text{a denoting refering expression} \\ |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 0 \quad \text{a non-denoting refering expression} \\ & \quad \text{for all } \mathcal{I} \models \mathcal{K}; \text{ these ought to be the only two possibilities!} \end{split}$$

Terms vs. Formulas Revisited [Artale et al., 2021]

Free (description) logics allow terms to be partial functions

 $\Rightarrow \iota \phi$ coerces a unary formula ϕ to a (partial) term that is defined iff ϕ is *denoting* (and *singular*)



Denoting/Non-denoting Referring Expressions

Singularity Revisited

$$\begin{split} |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 1 \quad \text{a denoting refering expression} \\ |\{o \mid \mathcal{I}, [x \mapsto o] \models \phi\}| &= 0 \quad \text{a non-denoting refering expression} \\ & \quad \text{for all } \mathcal{I} \models \mathcal{K}; \text{ these ought to be the only two possibilities!} \end{split}$$

Terms vs. Formulas Revisited [Artale et al., 2021]

Free (description) logics allow terms to be partial functions

 $\Rightarrow \iota \phi$ coerces a unary formula ϕ to a (partial) term

that is defined iff ϕ is *denoting* (and *singular*)

How do we Guarantee Singularity?

- 1 (fresh) nominals [Artale et al., 2021]
- 2 functionality [Borgida et al., 2016]



Single Interpretations/Models

Generating Referring Expressions (GRE)

Task: given an *interpretation*, find formulæ (referring expressions) that denote (selected) single objects.

Carlos Areces, Santiago Figueira, Daniel Gorín: Using Logic in the Generation of Referring Expressions. Logical Aspects of Computational Linguistics 2011.

Carlos Areces, Alexander Koller, Kristina Striegnitz: Referring Expressions as Formulas of Description Logic. International Natural Language Generation Conference 2008.



Logical Theories and Knowledge Bases

Russell's *Definite Descriptions* ... denote exactly *one* object

What happens if we consider *logical theories* rather than a *particular model*?

- constant symbols (similar for function/predicate symbols)
- ... can be interpreted by different individuals in different models



Logical Theories and Knowledge Bases

Russell's *Definite Descriptions* . . . denote exactly *one* object

What happens if we consider *logical theories* rather than a *particular model*?

- constant symbols (similar for function/predicate symbols)
 ... can be interpreted by different individuals in different models
- \Rightarrow (standard) constants don't quite satisfy Russell's/Kripke's requirements \Rightarrow *rigid designators* (symbols interpreted identically in all models)?





Rigidity and Genericity: DB Theory Way

Database (theory) Approach

- Database Instances (aka models) expect constants to be rigid ⇒ but constraints/queries do not *know*
- Database Queries are required to be *generic* ⇒ invariant under permutations of the underlying domain

Rigidity and Genericity: DB Theory Way

Database (theory) Approach

- Database Instances (aka models) expect constants to be rigid ⇒ but constraints/queries do not know
- Database Queries are required to be *generic*⇒ invariant under permutations of the underlying domain

Certain Answers (to $\varphi\{x\}$ in \mathcal{K})

- **11** Logical Definition: $\{a \mid \mathcal{K} \models \varphi[a/x]\}$
- 2 DB Definition: $\bigcap_{l\models K}\{a\mid \mathcal{I}, [x\mapsto a]\models \varphi\}$ (conflates constants with domain elements)





Rigidity and Genericity: DB Theory Way

Database (theory) Approach

- Database Instances (aka models) expect constants to be rigid ⇒ but constraints/queries do not know
- Database Queries are required to be *generic* ⇒ invariant under permutations of the underlying domain

Certain Answers (to $\varphi\{x\}$ in \mathcal{K})

- Logical Definition: $\{a \mid \mathcal{K} \models \varphi[a/x]\}$
- **2** DB Definition: $\bigcap_{l \models K} \{ a \mid \mathcal{I}, [x \mapsto a] \models \varphi \}$

(conflates constants with domain elements)

... for generic (and domain-independent) queries the result is *the same!*



Referring to Objects (fine print)

The rest of the presentation is based on

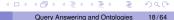
- KR16 Alexander Borgida, David Toman, and Grant E. Weddell: On Referring Expressions in Query Answering over First Order Knowledge Bases. Proc. International Conference on Principles of Knowledge Representation and Reasoning KR 2016, 319-328, 2016.
- ER16 Alexander Borgida, David Toman, and Grant Weddell: On Referring Expressions in Information Systems Derived from Conceptual Modelling. Proc, International Conference on Conceptual Modeling ER 2016, 183-197, 2016.
- Al16 David Toman, and Grant Weddell: Ontology Based Data Access with Referring Expressions for Logics with the Tree Model Property. Proc. Australasian Joint Conference on Artificial Intelligence, 2016.
- EKAW18 Weicong Ma, C. Maria Keet, Wayne Oldford, David Toman, and Grant Weddell: The Utility of the Abstract Relational Model and Attribute Paths in SQL. Proc. *International Conference on Knowledge Engineering and Knowledge Management*, 195-211, EKAW 2018.
- DL18 David Toman and Grant E. Weddell: Identity Resolution in Conjunctive Querying over DL-based Knowledge Bases. Proc. Description Logics DL 2018, 2018 (to appear in PRICAI 2019).
- DL19 David Toman, Grant E. Weddell: Exhaustive Query Answering via Referring Expressions. Proc. Description Logics DL 2019, 2019.
- DL22 Alexander Borgida, Enrico Franconi, David Toman, and Grant E. Weddell: Accessing Document Data Sources using Referring Expression Types. Proc. Description Logics DL 2022, 2022 (next week).



ONTOLOGY BASED DATA ACCESS

(BETTER QUERY ANSWERS WHEN QUERYING KNOWLEDGE BASES)





Queries and Ontologies

Ontology-based Data Access

Enriches (query answers over) explicitly represented data using background knowledge (captured using an ontology.)



Queries and Ontologies

Ontology-based Data Access

Enriches (query answers over) explicitly represented data using background knowledge (captured using an ontology.)

Example

Bob is a BOSS

■ Every BOSS is an EMPloyee

List all EMPloyees ⇒ {Bob}

(explicit data)

(ontology)

(query)

19/64

Goal: compute all certain answers

 \Rightarrow answers *common* in all models of KB (aka. answers *logically implied* by KB)





Approaches to Ontology-based Data Access

Main Task

INPUT: Ontology (\mathcal{T}), Data (\mathcal{A}), and a Query (\mathcal{Q})

Knowledge Base(\mathcal{K})

OUTPUT: $\{a \mid \mathcal{K} \models Q[a]\}$

- Reduction to *standard reasoning* (e.g., satisfiability)
- Reduction to *querying a relational database*
 - \Rightarrow very good at $\{a \mid A \models Q[a]\}$ for range restricted Q
 - \Rightarrow what to do with \mathcal{T} ??
 - incorporate into Q (perfect rewriting for DL-Lite et al. (AC 0 logics)); or
 - incorporate into \mathcal{A} (combined approach for \mathcal{EL} (PTIME-complete logics)); or sometimes both (\mathcal{CFDI} or Horn- $\mathcal{ALC}*$ logics).



"David is a UWaterloo Employee" and "every Employee has a Phone"

Question: Does David have a Phone?

Answer: YES

"David is a UWaterloo Employee" and "every Employee has a Phone"

Question: Does David have a Phone?

Answer: YES

Question: OK, tell me about David's Phone!

Answer: {}

"David is a UWaterloo Employee" and "every Employee has a Phone"

Question: Does David have a Phone?

Answer: YES

Question: OK, tell me about David's Phone!

Answer: {}



"David is a UWaterloo Employee" and "every Employee has a Phone"

Question: Does David have a Phone?

Answer: YES

Question: OK, tell me about David's Phone!

Answer: {}



Better Answers (possibly)

- it is a phone with phone # +1(519) 888-4567x34447;
- it is a *UWaterloo* phone with an extension x34447;
- it is a phone in the Davis Centre, Office 3344;
- 4 it is a Waterloo phone attached to port 0x0123abcd;
- it is a Waterloo CS phone with inventory # 100034447;
- it is *David's* phone (??)



Definition (Singular Referring Expression)

is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

Definition (Singular Referring Expression)

... is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

- it is a phone with phone # "+1(519) 888-4567x34447";
- it is a *UWaterloo* phone with extension x34447;
- it is a phone in the Davis Centre, Office 3344;
- it is a Waterloo phone attached to port 0x0123abcd;
- it is a Waterloo CS phone with inventory # 100034447;
- 6 it is *David's* phone;
- it is the red phone;



Definition (Singular Referring Expression)

... is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

- it is a phone with phone # "+1(519) 888-4567x34447";
- it is a *UWaterloo* phone with extension x34447;
- it is a phone in the Davis Centre, Office 3344;
- it is a Waterloo phone attached to port 0x0123abcd;
- it is a Waterloo CS phone with inventory # 100034447;
- 6 it is *David's* phone;
- it is the red phone;



Definition (Singular Referring Expression)

... is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

- it is a phone with phone # "+1(519) 888-4567x34447";
- it is a *UWaterloo* phone with extension x34447;
- it is a phone in the Davis Centre, Office 3344;
- it is a Waterloo phone attached to port 0x0123abcd;
- it is a *Waterloo CS* phone *with inventory # 100034447*;
- 6 it is *David's* phone;
- it is the *red phone*;



Definition (Singular Referring Expression)

... is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

```
    it is a phone with phone # "+1(519) 888-4567x34447";
    it is a UWaterloo phone with extension x34447;
```

- it is a phone in the Davis Centre, Office 3344;
- 4 it is a Waterloo phone attached to port 0x0123abcd;
- it is a Waterloo priorie attacheu to port 0x0123abcd,
- it is a Waterloo CS phone with inventory # 100034447;
- 6 it is *David's* phone;
- it is the *red phone*;



Definition (Singular Referring Expression)

... is a noun phrase that, when used as a query answer, identifies a particular object in this query answer.

```
it is a phone with phone # "+1(519) 888-4567x34447";
it is a UWaterloo phone with extension x34447;
  it is a phone in the Davis Centre, Office 3344;
  it is a Waterloo phone attached to port 0x0123abcd;
it is a Waterloo CS phone with inventory # 100034447;
6 it is David's phone;
it is the red phone;
```



Definition (Singular Referring Expression)

... is a unary formula that, when used as a query answer, identifies a particular object in this query answer.

```
it is a phone x s.t. PhoneNo(x, "+1(519) 888-4567x34447") holds;
```

- it is a phone x s.t. UWPhone(x) \land PhoneExt(x, "x34447") holds;
- it is a phone x s.t. UWRoom(x, "DC3344") holds;
 - it is a phone x s.t. UWPhone(x) \land PhonePort(x, 0x0123abcd) holds;
- it is a phone x s.t. UWCSPhone(x) \wedge InvNo(x, "100034447") holds;
- it is a phone x s.t. IsOwner("David", x) holds;
- it is the phone x s.t. Colour(x, "red") holds;



From Query Answers to Referring Expressions [KR16]

(Certain) Query Answers

Given a query $\psi\{x_1,\ldots,x_k\}$ and a KB \mathcal{K} ;

Classical answers: substitutions

$$\theta = \{x_1 \mapsto a_1, \dots, x_k \mapsto a_k\}$$

that map free variables of ψ to constants that appear in \mathcal{K} and $\mathcal{K} \models \psi \theta$.



From Query Answers to Referring Expressions [KR16]

(Certain) Query Answers

Given a query $\psi\{x_1,\ldots,x_k\}$ and a KB \mathcal{K} ;

Classical answers: substitutions

$$\theta = \{x_1 \mapsto a_1, \dots, x_k \mapsto a_k\}$$

that map free variables of ψ to constants that appear in \mathcal{K} and $\mathcal{K} \models \psi \theta$.

■ Referring Expression-based answers: *R-substitutions*

$$\theta = \{\mathbf{x}_1 \mapsto \phi_1\{\mathbf{x}_1\}, \dots, \mathbf{x}_k \mapsto \phi_k\{\mathbf{x}_k\}\}\$$

where $\phi_i\{x_i\}$ are unary formulæ in the language of K such that

(soundness)

$$\exists x_1,\ldots,x_k.(\phi_1\wedge\ldots\wedge\phi_k)\wedge\psi$$

(existence)

$$\exists \forall x_1, \dots, x_k, y_i. \phi_1 \wedge \dots \wedge \phi_k \wedge \psi \wedge \phi_i[x_i/y_i] \wedge \psi[x_i/y_i] \rightarrow x_i = y_i$$
 (singularity)

 \ldots are logically implied by \mathcal{K} .



Example KB

```
\mathcal{T} = \{ \text{ fatherof}(x, y) \rightarrow (\text{Father}(x) \land \text{Person}(y)), \}
           Father(x) \rightarrow Person(x),
           Father(x) \rightarrow \exists y.fatherof(x, y),
           Person(x) \rightarrow \exists y.fatherof(y, x)
A = \{ Father(fred), Person(mary) \}
```

Example KB

```
\mathcal{T} = \{ \text{ fatherof}(x, y) \rightarrow (\text{Father}(x) \land \text{Person}(y)), \}
            Father(x) \rightarrow Person(x),
            Father(x) \rightarrow \exists y.fatherof(x, y),
            Person(x) \rightarrow \exists y.fatherof(y, x)
\mathcal{A} = \{ \text{ Father(fred)}, \text{Person(mary)} \}
```

Query: Father(x)?

Answers: x = fred

Example KB

```
\mathcal{T} = \{ \text{ fatherof}(x, y) \to (\text{Father}(x) \land \text{Person}(y)), \\ \text{Father}(x) \to \text{Person}(x), \\ \text{Father}(x) \to \exists y. \text{fatherof}(x, y), \\ \text{Person}(x) \to \exists y. \text{fatherof}(y, x) \\ \}
\mathcal{A} = \{ \text{ Father}(\text{fred}), \text{Person}(\text{mary}) \}
```

Query: Father(x)?

Answers: x = fred, fatherof(x, mary), $\exists y.\text{fatherof}(x, y) \land \text{fatherof}(y, \text{mary})$, ...

Example KB

```
\mathcal{T} = \{ \text{ fatherof}(x, y) \to (\text{Father}(x) \land \text{Person}(y)), \\ \text{Father}(x) \to \text{Person}(x), \\ \text{Father}(x) \to \exists y. \text{fatherof}(x, y), \\ \text{Person}(x) \to \exists y. \text{fatherof}(y, x) \\ \text{fatherof}(x, z) \land \text{fatherof}(y, z) \to x = y \}
\mathcal{A} = \{ \text{ Father}(\text{fred}), \text{Person}(\text{mary}) \}
```

Query: Father(x)?

```
Answers: x = \text{fred}, \text{fatherof}(x, \text{mary}), \exists y.\text{fatherof}(x, y) \land \text{fatherof}(y, \text{mary}), ... \text{fatherof}(x, \text{fred}), \exists y.\text{fatherof}(x, y) \land \text{fatherof}(y, \text{fred}), ...
```



Example KB

```
\mathcal{T} = \{ \text{ fatherof}(x, y) \to (\text{Father}(x) \land \text{Person}(y)), \\ \text{Father}(x) \to \text{Person}(x), \\ \text{Father}(x) \to \exists y. \text{fatherof}(x, y), \\ \text{Person}(x) \to \exists y. \text{fatherof}(y, x) \\ \text{fatherof}(x, z) \land \text{fatherof}(y, z) \to x = y \}
\mathcal{A} = \{ \text{ Father}(\text{fred}), \text{Person}(\text{mary}) \}
```

Query: Father(x)?

```
Answers: x = \text{fred}, fatherof(x, mary), \exists y. fatherof(x, y) \land fatherof(y, mary), ... fatherof(x, fred), \exists y. fatherof(x, y) \land fatherof(y, fred), ...
```

Query: Person(x)?

Answers: x = mary, x = fred, $\frac{\text{fatherof(fred}, x)}{\text{fatherof(}x, \text{mary)}}$, $\frac{\text{fatherof(}x, \text{fred)}}{\text{fatherof(}x, \text{mary)}}$



Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \\ \text{ spouse}(x, z) \land \text{ spouse}(y, z) \rightarrow x = y \\ \text{ spouse}(x, y) \rightarrow x \neq y \\ \mathcal{A} = \{ \text{ spouse}(\text{mary}, \text{fred}) \}
```

Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \}
           spouse(x, z) \land spouse(y, z) \rightarrow x = y
           spouse(x, y) \rightarrow x \neq y
\mathcal{A} = \{ \text{ spouse(mary, fred)} \}
```

Query: spouse(x, mary)?

Answers: x = fred

Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \\ \text{ spouse}(x, z) \land \text{ spouse}(y, z) \rightarrow x = y \\ \text{ spouse}(x, y) \rightarrow x \neq y \\ \mathcal{A} = \{ \text{ spouse}(\text{mary}, \text{fred}) \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred})$, ...



Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \}
           spouse(x, z) \land spouse(y, z) \rightarrow x = y
           spouse(x, y) \rightarrow x \neq y
\mathcal{A} = \{ \text{ spouse(mary, fred)} \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred}), \dots$

How many distinct answers to $\exists y.\text{spouse}(x, y)$?



Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \}
           spouse(x, z) \land spouse(y, z) \rightarrow x = y
           spouse(x, y) \rightarrow x \neq y
\mathcal{A} = \{ \text{ spouse(mary, fred)} \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred}), \dots$

How many distinct answers to $\exists y.\text{spouse}(x, y)$?

$$fred = spouse(x, mary) = \exists y.spouse(x, y) \land spouse(y, fred) = \dots$$



Example KB

```
\mathcal{T} = \{ \text{ spouse}(x, y) \rightarrow \text{ spouse}(y, x), \}
           spouse(x, z) \land spouse(y, z) \rightarrow x = y
           spouse(x, y) \rightarrow x \neq y
\mathcal{A} = \{ \text{ spouse(mary, fred)} \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred}), \dots$

How many distinct answers to $\exists y.\text{spouse}(x, y)$?

```
fred = spouse(x, mary) = \exists y.spouse(x, y) \land spouse(y, fred) = \dots
mary = spouse(x, fred) = \exists y.spouse(x, y) \land spouse(y, mary) = \dots
```



Example KB

```
\mathcal{T} = \{ \text{spouse}(x, y) \rightarrow \text{spouse}(y, x), \\ \text{spouse}(x, z) \land \text{spouse}(y, z) \rightarrow x = y \\ \text{spouse}(x, y) \rightarrow x \neq y \}
\mathcal{A} = \{ \text{spouse}(\text{mary}, \text{fred}) \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred})$, ...

How many *distinct* answers to $\exists y.spouse(x, y)$?

```
fred = spouse(x, mary) = \exists y.spouse(x, y) \land spouse(y, fred) = ... mary = spouse(x, fred) = \exists y.spouse(x, y) \land spouse(y, mary) = ... mary \neq fred (last constraint!)
```





Example KB

```
\mathcal{T} = \{ \text{spouse}(x, y) \rightarrow \text{spouse}(y, x), \\ \text{spouse}(x, z) \land \text{spouse}(y, z) \rightarrow x = y \\ \text{spouse}(x, y) \rightarrow x \neq y \}
\mathcal{A} = \{ \text{spouse}(\text{mary}, \text{fred}) \}
```

Query: spouse(x, mary)?

Answers: x = fred, spouse(x, mary), $\exists y.\text{spouse}(x, y) \land \text{spouse}(y, \text{fred})$, ...

How many *distinct* answers to $\exists y.spouse(x, y)$?

```
\begin{array}{l} \text{fred} = \text{spouse}(x, \text{mary}) = \exists y. \text{spouse}(x, y) \land \text{spouse}(y, \text{fred}) = \dots \\ \text{mary} = \text{spouse}(x, \text{fred}) = \exists y. \text{spouse}(x, y) \land \text{spouse}(y, \text{mary}) = \dots \\ \text{mary} \neq \text{fred (last constraint!)} \Rightarrow \text{exactly 2 distinct objects} \end{array}
```



Controlling the number of Answers: Finite Representation

How do we deal with multiple referring expression answers/preferences/...?

- potentially too many implied answers (infinitely many!)
- potentially too many ways to refer to the same object





Controlling the number of Answers: Finite Representation

How do we deal with multiple referring expression answers/preferences/...?

- potentially too many implied answers (infinitely many!)
- potentially too many ways to refer to the same object

Can we (somehow) get ALL answers to Q over K?

Yes (for logics with *recursively enumerable* logical consequence):

- for all (tuples of) unary formulas $\varphi(x)$
- do test if $\varphi(x)$ is a singular certain answer to Q in K.





Controlling the number of Answers: Finite Representation

How do we deal with multiple referring expression answers/preferences/...?

- potentially too many implied answers (infinitely many!)
- potentially too many ways to refer to the same object

Can we (somehow) get ALL answers to Q over K?

Yes (for logics with *recursively enumerable* logical consequence):

- for all (tuples of) unary formulas $\varphi(x)$
- do test if $\varphi(x)$ is a singular certain answer to Q in \mathcal{K} .
 - ⇒ this does NOT guarantee decidability [Artale et al., 2021]





Controlling the number of Answers: Finite Representation

How do we deal with multiple referring expression answers/preferences/...?

- potentially too many implied answers (infinitely many!)
- potentially too many ways to refer to the same object

Can we (somehow) get ALL answers to Q over K?

Yes (for logics with *recursively enumerable* logical consequence):

for all (tuples of) unary formulas $\varphi(x)$

Borgida, Toman, and Weddel

- do test if $\varphi(x)$ is a singular certain answer to Q in \mathcal{K} .
 - ⇒ this does NOT guarantee decidability [Artale et al., 2021]
 - \Rightarrow is there a *finite representation* of all answers (and what is "all")?





Example: Horn Logics with Tree Models [DL19]

What to do \mathcal{EL}^{\perp} (and Horn- \mathcal{ALC})?

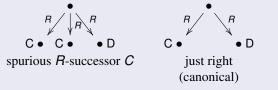
■ singularity requires role functionality (not expressible in \mathcal{EL}^{\perp} /Horn- \mathcal{ALC})

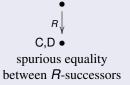


Example: Horn Logics with Tree Models [DL19]

What to do \mathcal{EL}^{\perp} (and Horn- \mathcal{ALC})?

- singularity requires role functionality (not expressible in \mathcal{EL}^{\perp} /Horn- \mathcal{ALC})
- (Tree) Models of $a : \exists R.C \sqcap \exists R.D$:





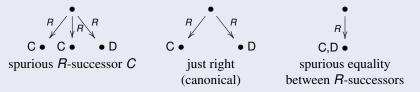
⇒ singular certain answers: singular in a canonical model



Example: Horn Logics with Tree Models [DL19]

What to do \mathcal{EL}^{\perp} (and Horn- \mathcal{ALC})?

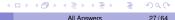
- **singularity** requires *role functionality* (not expressible in \mathcal{EL}^{\perp} /Horn- \mathcal{ALC})
- (Tree) Models of $a : \exists R.C \sqcap \exists R.D$:



⇒ singular certain answers: singular in a canonical model

⇒ coincide with singular answers in DLs with functional roles (FunDL).





How Does it Work?

Base Case: Instance Retrieval B(x) over \mathcal{T} and $A = \{a : A\}$

Looping automaton-like construction

- \Rightarrow only non-redundant successors in matching tuples
- \Rightarrow preserves complexity bounds for both logics





How Does it Work?

Base Case: Instance Retrieval B(x) over \mathcal{T} and $A = \{a : A\}$

Looping automaton-like construction

- \Rightarrow only non-redundant successors in matching tuples
- ⇒ preserves complexity bounds for both logics

Generalizations&Limitations

- General ABoxes and Conjunctive Queries
 - ⇒ lots of case analysis followed by existing approaches
- Finite representation of answers (succinctness??)
- More Expressive Logics
 - ⇒ this will NOT work with at-least restrictions (functionality is fine)
- Mon-Horn Logics
 - ⇒ non-unique canonical models
 - ⇒ disjunctions in referring expressions (questionable)





Controlling the number of Answers: Typing Restrictions

How do we deal with multiple referring expression answers/preferences/...?

- potentially too many implied answers (infinitely many!)
- potentially too many ways to refer to the same object

Referring Expression Types and Typed Queries

Types:
$$Rt ::= Pd = \{?\} \mid Rt_1 \land Rt_2 \mid T \rightarrow Rt \mid Rt_1; Rt_2 \Rightarrow \text{ each type induces a set of unary formulæ;}$$

Queries: select $x_1 : Rt_1, \dots, x_k : Rt_k$ where ψ

 $\Rightarrow x_1 : Rt_1, \dots, x_k : Rt_k$ is called the head, ψ is the body.





Referring Expression Types

Desiderata: only Referring Expressions that are Singular

Given \blacksquare a KB \mathcal{K} (the "background knowledge"),

2 a query $\psi\{x_1,\ldots,x_k\}$, and

3 types Rt_1, \ldots, Rt_k for sets of unary formulæ S_1, \ldots, S_k

We ask whether, for every \mathcal{K}' (the "data") consistent with \mathcal{K} and an answer

$$\theta = \{x_1 \mapsto \phi_1\{x_1\}, \dots, x_k \mapsto \phi_k\{x_k\}\}\$$

to ψ with respect to $\mathcal{K} \cup \mathcal{K}'$ such that $\phi_i \in S_i$, it is the case that θ is singular.

Referring Expression Types

Desiderata: only Referring Expressions that are Singular

Given \blacksquare a KB \mathcal{K} (the "background knowledge"),

 $extbf{2}$ a query $\psi\{x_1,\ldots,x_k\}$, and

3 types Rt_1, \ldots, Rt_k for sets of unary formulæ S_1, \ldots, S_k

We ask whether, for every \mathcal{K}' (the "data") consistent with \mathcal{K} and an answer

$$\theta = \{x_1 \mapsto \phi_1\{x_1\}, \dots, x_k \mapsto \phi_k\{x_k\}\}\$$

to ψ with respect to $\mathcal{K} \cup \mathcal{K}'$ such that $\phi_i \in S_i$, it is the case that θ is singular.

Theorem (Weak Identification; paraphrased)

Given a query ψ with a head H and a KB K, the question "are all answers to ψ conforming to H over any $K \cup K'$ singular?" reduces to logical implication in the underlying logic of K.





Reference via a Single-Attribute Key

"The ssn# of any person with phone 1234567"

 $select x : ssn\# = \{?\}$

where $Person(x) \land phone\#(x, 1234567)$



Reference via a Single-Attribute Key

Reference by a Multi-Attribute Key

"The title and publisher of any journals"

```
select x : title = \{?\} \land publishedBy = \{?\} where Journal(x)
```

Reference via a Single-Attribute Key

Reference by a Multi-Attribute Key

Choice of Identification in a Heterogeneous Set

```
"Any legal entity"
```

```
select x : Person \rightarrow ssn\# = \{?\};
            Company \rightarrow tickerSymbol = \{?\}
where LegalEntity(x)
```

```
answers: \{x \mapsto Person(x) \land ssn\#(x,7654)\}
            \{x \mapsto Company(x) \land tickerSymbol(x, "IBM")\}.
```



Reference via a Single-Attribute Key

Reference by a Multi-Attribute Key

Choice of Identification in a Heterogeneous Set

Preferred Identification

"Any publication, identified by its most specific identifier, when available."

```
select x : Journal \rightarrow (title = \{?\} \land publisher = \{?\});
             EditedCollection \rightarrow isbn# = \{?\}; \{?\}
where Publication(x)
```

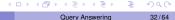
```
answers: \{x \mapsto Journal(x) \land title(x, "AIJ") \land publisher(x, "Elsevier")\}
            \{x \mapsto EditedCollection(x) \land isbn\#(x, 123456789)\}
            \{x \mapsto x = /\text{guid}/9202a8c04000641f8000000...\}
```



REQA (Referring Expression-based QA)

GOAL: reduce REQA to standard OBDA (used as an oracle)





REQA (outline, unary queries only)

GOAL: reduce REQA to standard OBDA (used as an oracle)

Input: K (background knowledge), K' (data), $\psi\{x\}$ (query), H (query head)

1 Normalize H to $H_1; \ldots; H_\ell$, each of the form

$$T_i \to Pd_{i,1} = \{?\} \land \ldots \land Pd_{i,k_i} = \{?\};$$

2 Create queries $\psi_i\{x, y_1, \dots, y_{k_i}\}$ as

$$\psi \wedge T_i(x) \wedge Pd_{i,1}(x, y_1) \wedge \ldots \wedge Pd_{i,k_i}(x, y_{k_i});$$

- 3 Create K_i with a witnesses for x when no such witness exists;
- 4 Evaluate $\mathcal{K} \cup \mathcal{K}' \cup \mathcal{K}_i \models \psi_i$ (OBDA oracle);
- **5** Resolve preferences (based on value of x); and
- **6** Reconstruct a referring expression from the values of y_1, \ldots, y_{k_i}

... extends naturally to higher arity queries: (more) messy



The Tractable (practical) Cases

Lite Description Logics

```
DL-Lite_{core}^{\mathcal{F}}(idc):
```

- Weak identification sequence of KB consistency tests
- Query answering → REQA
 - + Witnesses for x w.r.t. H + Perfect Reformulation

$\mathcal{CFDI}_{no}^{\forall}$:

- Weak identification sequence of logical implications
- Query answering → REQA
 - + Combined Combined Approach

Logics with Tree Models (outside of an ABox) [Al16]

The witnesses for anonymous objects (step (3))

→ *last* named individual on a path *towards* the anonymous object





RECORDING/REPRESENTING FACTUAL DATA





Referring Expressions for Ground Knowledge

Standard approach: constant symbols \sim objects (and values!)

⇒ needs a constant symbol for *every individual* (Skolems?)





Referring Expressions for Ground Knowledge

Standard approach: constant symbols \sim objects (and values!)

⇒ needs a constant symbol for every individual (Skolems?)

How are *external* objects identified in a KB?

■ Two PERSON objects, o_1 and o_2 , identified by their ssn value:

PERSON
$$\sqcap \exists ssn. \{123\}$$
 and PERSON $\sqcap \exists ssn. \{456\}$.

Role (feature) assertions of the form $mother(o_1) = o_2$ can then be captured as:

```
PERSON \sqcap \exists ssn. \{123\} \sqcap \exists mother. (PERSON \sqcap \exists ssn. \{345\}).
```



Referring Expressions for Ground Knowledge

Standard approach: constant symbols \sim objects (and values!)

⇒ needs a constant symbol for every individual (Skolems?)

How are external objects identified in a KB?

■ Two PERSON objects, o_1 and o_2 , identified by their ssn value:

PERSON
$$\sqcap \exists ssn. \{123\}$$
 and PERSON $\sqcap \exists ssn. \{456\}$.

■ Role (feature) assertions of the form $mother(o_1) = o_2$ can then be captured as:

PERSON
$$\sqcap \exists ssn. \{123\} \sqcap \exists mother. (PERSON $\sqcap \exists ssn. \{345\}).$$$

Issues:

- admissibility: what descriptions qualify here? ⇒ singularity!
- minimality: is the description succinct? (similar to keys/superkeys issues)



Heterogeneous Data Integration (example)

Example

```
■ TBox \mathcal{T} = \{ FRIEND \sqsubseteq PERSON,
                      FRIEND \square PERSON : fname \rightarrow id,
                      MATRIARCH \square PERSON,
                      MATRIARCH \square PERSON : Iname \rightarrow id,
                      PERSON \square PERSON : fname, lname \rightarrow id, ... }
■ CBox C = \{ FRIEND \sqcap \exists fname. \{ \text{"Mary"} \}, \}
                      PERSON \sqcap (\exists fname.{"Mary"}) \sqcap (\exists lname.{"Smith"}),
                      MATRIARCH \sqcap \exists lname. \{\text{"Smith"}\}, \dots \}
```

Heterogeneous Data Integration (example)

Example

```
■ TBox \mathcal{T} = \{ FRIEND \sqsubseteq PERSON,
                      FRIEND \square PERSON : fname \rightarrow id,
                      MATRIARCH \square PERSON,
                      MATRIARCH \square PERSON : Iname \rightarrow id,
                      PERSON \square PERSON : fname, lname \rightarrow id, ... }
■ CBox C = \{ FRIEND \sqcap \exists fname. \{ "Mary"\},
                      PERSON \sqcap (\exists fname. {"Mary"}) \sqcap (\exists lname. {"Smith"}),
                      MATRIARCH \sqcap \exists lname. \{\text{``Smith''}\}, \dots \}
```

Heterogeneous Identification

```
"FRIEND \sqcap \exists fname. \{\text{"Mary"}\}" identifies the same object as
"PERSON \sqcap (\exists fname.{"Mary"}) \sqcap (\exists lname.{"Smith"})" and in turn as
"MATRIARCH □ ∃Iname.{"Smith"}"
```



Heterogeneous Data Integration (example)

Example

```
■ TBox \mathcal{T} = \{ FRIEND \sqsubseteq PERSON, FRIEND \sqsubseteq PERSON : fname \rightarrow id, MATRIARCH \sqsubseteq PERSON, MATRIARCH \sqsubseteq PERSON : fname \rightarrow id, PERSON \sqsubseteq PERSON : fname, fname \rightarrow id, ... \}

■ CBox \mathcal{C} = \{ FRIEND \sqcap \exists fname. {"Mary"}, PERSON \sqcap (\exists fname. {"Mary"}) \sqcap (\exists fname. {"Smith"}), MATRIARCH \sqcap \exists fname. {"Smith"}, ... \}
```

Heterogeneous Identification

```
"FRIEND \sqcap \exists fname. \{\text{"Mary"}\}" identifies the same object as "PERSON \sqcap (\exists fname. \{\text{"Mary"}\}) \sqcap (\exists fname. \{\text{"Smith"}\})" and in turn as "MATRIARCH \sqcap \exists fname. \{\text{"Smith"}\}"
```

... and thus is an answer to $\{x \mid MATRIARCH(x)\}$.



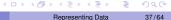
Minimality

IDEA: minimal referring expressions (ala Candidate Keys)

C is a referring expression singular w.r.t. a TBox \mathcal{T} (e.g., a *superkey*)

- C's subconcepts A, $\{a\}$, $\exists f. \top$, $\exists f^{-1}. \top$, and $\top \sqcap \top$ are *leaves* of C.
- $C[L \mapsto \top]$ is a description C in which a leaf L was replaced by \top .
- "first-leaf" and "next-leaf" successively enumerate all leaves of C.
 - 1. L := first-leaf(C);
 - 2. **while** $C[L \mapsto \top]$ is singular w.r.t. \mathcal{T} **do**
 - 3. $C := C[L \mapsto \top]; L := \text{next-leaf}(C);$
 - 4. done
 - 5. return C;





Minimality

IDEA: minimal referring expressions (ala Candidate Keys)

C is a referring expression singular w.r.t. a TBox \mathcal{T} (e.g., a *superkey*)

- C's subconcepts A, $\{a\}$, $\exists f. \top$, $\exists f^{-1}. \top$, and $\top \sqcap \top$ are *leaves* of C.
- $C[L \mapsto \top]$ is a description C in which a leaf L was replaced by \top .
- "first-leaf" and "next-leaf" successively enumerate all leaves of C.
 - 1. L := first-leaf(C);
 - 2. **while** $C[L \mapsto \top]$ is singular w.r.t. \mathcal{T} **do**
 - 3. $C := C[L \mapsto \top]; L := \text{next-leaf}(C);$
 - 4. done
 - 5. return C;
 - \Rightarrow computes a syntactically-minimal co-referring expression for C.
 - \Rightarrow order of enumeration \rightarrow variant minimal co-referring expressions.



Reasoning and QA with CBoxes [DL18]

Theorem (CBox Admissibility)

Let \mathcal{T} be a $\mathcal{CFDI}_{nc}^{\forall}$ TBox and C a concept description. Then C is a singular referring expression w.r.t. \mathcal{T} if and only if the knowledge base

$$(\mathcal{T} \cup \{A \sqsubseteq \neg B\}, \mathsf{Simp}(a : C) \cup \mathsf{Simp}(b : C) \cup \{a : A, b : B\})$$

is inconsistent, where a and b are distinct constant symbols, and A and B are primitive concepts not occurring in $\mathcal T$ and $\mathcal C$.

Theorem (Satisfiability of KBs with CBoxes)

Let $\mathcal{K}=(\mathcal{T},\mathcal{C})$ be a knowledge base with an admissible CBox \mathcal{C} . Then \mathcal{K} is consistent iff $(\mathcal{T},\mathsf{Simp}(\mathcal{C}))$ is consistent.

Theorem (Query Answering)

Let $\mathcal{K} = (\mathcal{T}, \mathcal{C})$ be a consistent knowledge base and $Q = \{(x_1, \dots, x_k) : \varphi\}$ a conjunctive query over \mathcal{K} . Then (C_1, \dots, C_k) is a certain answer to Q in \mathcal{K} if and only if $(a_{C_1}, \dots, a_{C_k})$ is a certain answer to Q over $(\mathcal{T}, \text{Simp}(\mathcal{C}))$.

Documents and Ontologies

Ontologies for Documents: Goals

- 1 to capture class mambership of entities captured in a document, and
- 2 to establish how entities are identified in a document.



Documents and Ontologies

Ontologies for Documents: Goals

- 1 to capture class mambership of entities captured in a document, and
- 2 to establish how entities are identified in a document.

IDEA: Documents as Concepts, Semantics as Ontology

- Syntactical document structure captured as a concept in FunDL
 - ⇒ similar to the IBM IMS hierarchical data model
- Ontology adds meaning to this concept and its subconcepts
 - identifies class membership of entities described by subdocuments,
 - discovers subdocuments pertaining to the same entity, and
 - drives document normalization.





Example: JSON Document

```
"collection": "person",
"data" : [
  { "fname": "John", "lname": "Smith", "age": 25,
    "wife": { "fname" : "Mary" },
    "phone": [
        {"colour": "red", "dnum": "212 555-1234"}
    "fname": "Mary", "lname": "Jones",
    "salary": "$150,000 (CAD)",
    "spouse": { "fname": "John" },
    "phone": [
        {"loc": "home", "dnum": "212 555-1234"},
        {"loc": "work", "dnum": "212 666-4567"}
```

Example: JSON as a FunDL Concept

```
\exists collection.\{"person"\} \sqcap
∃data.∃dom<sup>-</sup>.∃ran(
              \exists fname.{"John"} \sqcap \exists lname.{"Smith"} \sqcap
              \existsage.{"25"}\sqcap\existswife.\existsfname.{"Mary"}\sqcap
              \exists phone(\exists dom^{-}.\exists ran(
                           \exists colour. \{ "red" \} \sqcap \exists dnum. \{ "212 555-1234" \} ) ) ) \sqcap
          ∃dom-.∃ran(
              ∃fname.{"Mary"}∏∃lname.{"Jones"}∏
              \exists salary.{"$150000CAD"} \sqcap \exists spouse.\exists fname.{"John"} \sqcap
              ∃phone(∃dom-.∃ran(
                             \exists loc. \{"home"\} \sqcap \exists dnum. \{"212 555-1234"\}) \sqcap
                          ∃dom-.∃ran(
                             \exists loc. \{ "work" \} \sqcap \exists dnum. \{ "212 666-4567" \} ) ) )
```



Example: Ontology

1 The TBox:

```
(\exists \texttt{collection.T}) \sqcap (\exists \texttt{data.T}) \sqsubseteq DOCUMENT \\ (\exists \texttt{fname.T}) \sqcap (\exists \texttt{lname.T}) \sqsubseteq PERSON \\ \exists \texttt{dnum.T} \sqsubseteq PHONE \\ DOCUMENT \sqsubseteq DOCUMENT : \texttt{collection} \rightarrow \textit{id} \\ PERSON \sqsubseteq PERSON : \texttt{fname,lname} \rightarrow \textit{id} \\ PHONE \sqsubseteq PHONE : \texttt{dnum} \rightarrow \textit{id} \\ PERSON \sqsubseteq \exists \texttt{wife.PERSON} \\ \exists \text{wife.PERSON}
```

The Referring Expression Type Assignment:

```
\begin{array}{lll} \mathsf{RTA}(\mathit{DOCUMENT}) &=& \mathit{DOCUMENT} \sqcap \exists \mathtt{collection}. \{?\} \\ \mathsf{RTA}(\mathit{PERSON}) &=& \mathit{PERSON} \sqcap \exists \mathtt{lname}. \{?\} \sqcap \exists \mathtt{fname}. \{?\} \\ \mathsf{RTA}(\mathit{PHONE}) &=& \mathit{PHONE} \sqcap \exists \mathtt{dnum}. \{?\} \end{array}
```



Example: Normalized CBox/Document

```
DOCUMENT | Gollection. { "person" } | Gollection |
       \exists dom^-.\exists ran(PERSON \sqcap \exists fname.{"John"} \sqcap \exists lname.{"Smith"}) \sqcap
       ∃dom<sup>-</sup>.∃ran(PERSON □ ∃fname.{"Mary"} □ ∃lname.{"Jones"}))
PERSON \sqcap \exists fname. \{ "John" \} \sqcap \exists lname. \{ "Smith" \} \sqcap
            ∃age.{"25"}∏∃wife.∃fname{"Mary"}∏
            \exists phone(\exists dom^-.\exists ran(PHONE \sqcap \exists dnum{"212 555-1234"}))
PERSON □ ∃ fname.{"Mary"} □ ∃lname.{"Jones"} □
            ∃salary.{"$150000CAD"}∏∃spouse.∃fname{"John"}∏
            \exists phone(\exists dom^-.\exists ran(PHONE \sqcap \exists dnum{"212 555-1234"}) \sqcap
                      \exists dom^{-}.\exists ran(PHONE \sqcap \exists dnum{"212 666-4567"}))
```

```
PHONE □ ∃dnum{"212 555-1234"} □ ∃loc.{"home"} □ ∃colour.{"red"}
PHONE □ ∃dnum{"212 555-4567"} □ ∃loc.{"work"}
```



Representing Data

Example: Normalized CBox/Document

```
DOCUMENT | Gollection. { "person" } | Gollection.
       ∃dom<sup>-</sup>.∃ran(PERSON □∃fname.{"John"} □∃lname.{"Smith"}) □
       ∃dom<sup>-</sup>.∃ran(PERSON □ ∃fname.{"Mary"} □ ∃lname.{"Jones"}))
PERSON \sqcap \exists fname. \{ "John" \} \sqcap \exists lname. \{ "Smith" \} \sqcap
            ∃age.{"25"}∏∃wife.∃fname{"Mary"}∏
            \exists phone(\exists dom^-.\exists ran(PHONE \sqcap \exists dnum{"212 555-1234"}))
PERSON □ ∃ fname.{"Mary"} □ ∃lname.{"Jones"} □
            \existssalary.{"$150000CAD"}\sqcap \existsspouse.\existsfname{"John"}\sqcap
            \exists phone(\exists dom^-.\exists ran(PHONE \sqcap \exists dnum{"212 555-1234"}) \sqcap
                     \exists dom^{-}.\exists ran(PHONE \sqcap \exists dnum{"212 666-4567"}))
```

```
PHONE □ ∃dnum{"212 555-1234"} □ ∃loc.{"home"} □ ∃colour.{"red"}
```

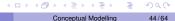
 $\underline{PHONE} \sqcap \exists \operatorname{dnum} \{"212 555-4567"\} \sqcap \exists \operatorname{loc.} \{"\operatorname{work"}\}$



CONCEPTUAL MODELLING

(Decoupling *modelling* from *identification* issues)





Thesis:

Modeling of *Entities* and their *Relationships* should be decoupled from issues of managing the identity of such entities.

Thesis:

Modeling of *Entities* and their *Relationships* **should be decoupled** from issues of *managing the identity* of such entities.

Weak Entities and dominant entity identification

Example (ROOM within BUILDING)

For the entity set ROOM with attributes room-number and capacity

- ⇒ natural attributes are insufficient to identify ROOMs
- \Rightarrow need for a *key* of dominant set, such as <code>BUILDING</code>





Thesis:

Modeling of *Entities* and their *Relationships* should be decoupled from issues of managing the identity of such entities.

Weak Entities and dominant entity identification

Preferred Identification in sub/super-classes

Example (PERSON and FAMOUS-PERSON)

For the entity set FAMOUS-PERSON a sub-entity of PERSON

- ⇒ choice of key (ssn) for PERSON forces the same key for FAMOUS-PERSON
- ⇒ we may prefer to use name in this case (e.g., Eric Clapton or The Edge)





Thesis:

Modeling of *Entities* and their *Relationships* should be decoupled from issues of managing the identity of such entities.

Weak Entities and dominant entity identification

Preferred Identification in sub/super-classes

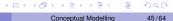
Generalizations and heterogeneity

Example (LEGAL-ENTITY: PERSON or COMPANY)

For the entity set LEGAL-ENTITY a generalization of PERSON and COMPANY

- ⇒ commonly required to create an artificial attribute le-num
- ⇒ despite the fact that all entities are already identified by the (more) natural ssn and (name, city) identifiers.





Conceptual Modeling and Identification [ER16]

Thesis:

Modeling of *Entities* and their *Relationships* **should be decoupled** from issues of *managing the identity* of such entities.

Weak Entities and dominant entity identification

Preferred Identification in sub/super-classes

Generalizations and heterogeneity

Contributions

- Methodology that allows decoupling identification from modeling;
- Referring Expressions that subsequently resolve identity issues; and
- Compilation-based technology that makes further translation to a pure relational model seamless.





Abstract (Relational) Model ARM

A simple conceptual model $\mathcal C$

Common features of so-called "attribute-based" semantic models

 \Rightarrow class hierarchies, disjointness, coverage, attributes and typing, functional dependencies, \dots

Example (DMV)

```
class PERSON (ssn: INT, name: STRING,
  isa LEGAL-ENTITY, disjoint with VEHICLE)
class COMPANY (name: STRING, city: STRING,
  isa LEGAL-ENTITY)
class LEGAL-ENTITY (covered by PERSON, COMPANY)
class VEHICLE (vin: INT, make: STRING,
  owned-by: LEGAL-ENTITY)
class CAN-DRIVE (driver: PERSON, driven: VEHICLE)
```



Abstract (Relational) Model ARM

A simple conceptual model ARM

Common features of so-called "attribute-based" semantic models

⇒ class hierarchies, disjointness, coverage, attributes and typing, functional dependencies, ...

Example (DMV and Relational Understanding)

```
table PERSON (self: OID, ssn: INT, name: STRING,
  isa LEGAL-ENTITY, disjoint with VEHICLE)
table COMPANY (self: OID, name: STRING, city: STRING,
  isa LEGAL-ENTITY)
table LEGAL-ENTITY (covered by PERSON, COMPANY)
table VEHICLE (self: OID, vin: INT, make: STRING,
  owned-by: LEGAL-ENTITY)
table CAN-DRIVE (self: OID, driver: PERSON, driven: VEHICLE)
```



Abstract Relational Queries

SQLP

(pretty) standard select-from-where-union-except SQL syntax ... with extensions to ARM: abstract attributes (OID) and attribute paths



Abstract Relational Queries

SQLP

(pretty) standard select-from-where-union-except SQL syntax ... with extensions to ARM: abstract attributes (OID) and attribute paths

■ The name of anyone who can drive a vehicle made by Honda:

```
select d.driver.name from CAN-DRIVE d
where d.driven.make = 'Honda'
```

attribute paths in the select and where clauses

■ The owners of Mitsubishi vehicles:

Borgida, Toman, and Weddel

```
select v.owned-by from VEHICLE v
where v.make = 'Mitsubishi'
```

retrieving abstract attributes may yield

heterogeneous results (PERSONs and COMPANIES)





Abstract Relational Queries

SQLP

(pretty) standard select-from-where-union-except SQL syntax ... with extensions to ARM: abstract attributes (OID) and attribute paths

■ The name of anyone who can drive a vehicle made by Honda:

```
select d.driver.name from CAN-DRIVE d
where d.driven.make = 'Honda'
```

attribute paths in the select and where clauses

■ The owners of Mitsubishi vehicles:

```
select v.owned-by from VEHICLE v
where v.make = 'Mitsubishi'
```

retrieving abstract attributes may yield

heterogeneous results (PERSONs and COMPANIES)

Note that queries **do NOT** rely on *(external) identification* of entities/objects.



How to Make this Technology Succeed?

- ARM/SQLP Helps Users (User Study) [EKAW18]
- ARM/SQLP Can be Efficiently Implemented [ER16]
 - Mapping to standard relational model with the help of referring expressions
 ⇒ and WITHOUT introducing explicit, material OIDs
 - Reverse-Engineering ARM from Legacy Relational Schemata





Experimental Design (HCI experiments)

Hypotheses

 H_t : no difference between RM/SQL and ARM/SQLP in the mean time taken

 H_c : no difference between RM/SQL and ARM/SQLP in the mean correctness

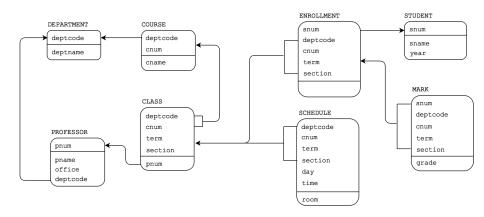
Methods

- Undergraduate (9) and Graduate (15) UW students
- Protocol
 - 1 Instructions (5") and Examples of SQL/SQLP (10")
 - 2 Six Questions (Q1–Q6), no time limit
 - 3 Subjects recorded start/end times for each Question
- Performance Assessment
 - 3 assessors
 - 2 agreed upon grading scale



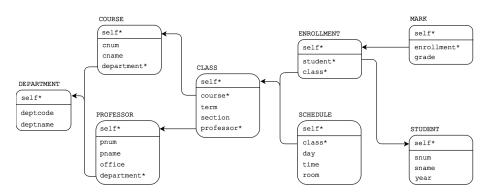


Course Enrollment as an RM Schema





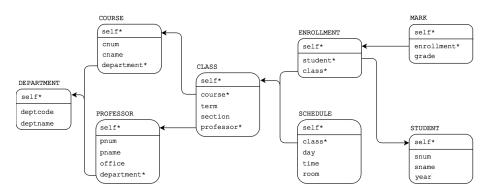
Course Enrolment as an ARM Schema







Course Enrolment as an ARM Schema



ARM *completely frees* domain experts/users from the need to understand how entities are *identified* in an information system.



Example Queries

Query: Names of students who have been taught by Prof. 'Alan John'

RM/SQL:

```
select distinct s.sname as name
from STUDENT s, ENROLLMENT e, CLASS c, PROFESSOR p
where e.snum = s.snum
and e.deptcode = c.deptcode and e.cnum = c.cnum
and e.term = c.term and e.section = c.section
and c.pnum = p.pnum and p.pname = 'Alan John'
```

Example Queries

Query: Names of students who have been taught by Prof. 'Alan John'

RM/SQL:

```
select distinct s.sname as name from STUDENT s, ENROLLMENT e, CLASS c, PROFESSOR p where e.snum = s.snum \,
```



and e.deptcode = c.deptcode and e.cnum = c.cnum
and e.term = c.term and e.section = c.section
and c.pnum = p.pnum and p.pname = 'Alan John'

Domain expert needs to understand structure of PK/FKs: BAD!!



Example Queries

Query: Names of students who have been taught by Prof. 'Alan John'

RM/SQL:

```
select distinct s.sname as name from STUDENT s, ENROLLMENT e, CLASS c, PROFESSOR p where e.snum = s.snum \,
```



and e.deptcode = c.deptcode and e.cnum = c.cnum
and e.term = c.term and e.section = c.section
and c.pnum = p.pnum and p.pname = 'Alan John'

Domain expert needs to understand structure of PK/FKs: BAD!!

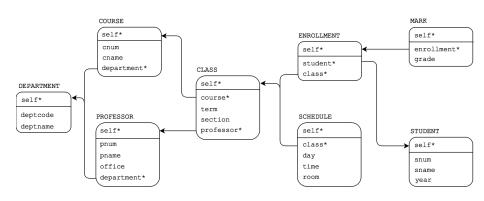
ARM/SQLP:

```
select distinct e.student.sname as name
from ENROLLMENT e
where e.class.professor.pname = 'Alan John'
```



ARM Schema and Path Navigation

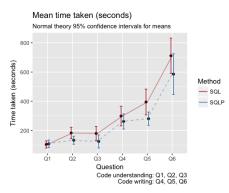
```
select distinct e.student.sname as name
from ENROLLMENT e
where e.class.professor.pname = 'Alan John'
```

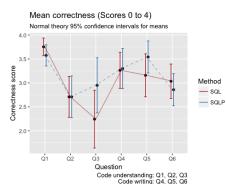




Experiments: Results

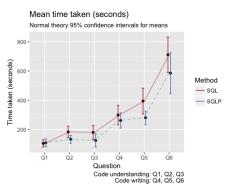
Mean performance for all subjects: SQL solid; SQLP dashed.

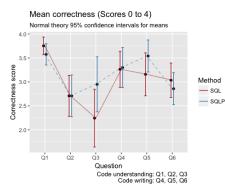




Experiments: Results

Mean performance for all subjects: SQL solid; SQLP dashed.





- SQLP outperforms SQL in time taken
- No significant difference in correctness (Q3, Q5 almost significant)



How to make the Technology Succeed?

- 1 ARM/SQLP Helps Users (User Study)
- 2 ARM/SQLP Can be Efficiently Implemented [ER16]
 - Mapping to standard relational model with the help of referring expressions
 - Reverse-Engineering ARM from Legacy Relational Schemata





Example (How to refer to LEGAL-ENTITY)

■ invent a *new attribute for this purpose* (will be *inherited* by subclasses)

Example (How to refer to LEGAL-ENTITY)

- invent a new attribute for this purpose (will be inherited by subclasses)
- use (a combination of) the identities of *generalized* entities, e.g., ssn for PERSON and (name, city) for COMPANY.



Example (How to refer to LEGAL-ENTITY)

- invent a new attribute for this purpose (will be inherited by subclasses)
- use (a combination of) the identities of generalized entities, e.g., ssn for PERSON and (name, city) for COMPANY.
 - ⇒ but what happens to objects that are both a PERSON and a COMPANY??

Example (How to refer to LEGAL-ENTITY)

- invent a new attribute for this purpose (will be inherited by subclasses)
- use (a combination of) the identities of generalized entities, e.g., ssn for PERSON and (name, city) for COMPANY.

```
but what happens to chicate that are both a DEDGON and a COMPANY 9
```

- \Rightarrow but what happens to objects that are both a PERSON and a COMPANY??
- ⇒ we need to resolve the *preferred* identification:

```
PERSON \rightarrow ssn=?; COMPANY \rightarrow (name=?, city=?).
```



Example (How to refer to LEGAL-ENTITY)

- invent a new attribute for this purpose (will be inherited by subclasses)
- use (a combination of) the identities of generalized entities, e.g., ssn for PERSON and (name, city) for COMPANY.
 - ⇒ but what happens to objects that are both a PERSON and a COMPANY??
 - ⇒ we need to resolve the preferred identification:

```
PERSON \rightarrow ssn=?; COMPANY \rightarrow (name=?, city=?).
```

Goal(s)

- Flexible assignment of *Referring Expression Types* to classes,
- Automatic check(s) for sanity of such an assignment, and
- 3 Compilation of gueries (updates) over ARM to ones over concrete tables.

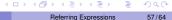




IDEA

Assign a referring expression type RTA(T) to each table T in Σ .





IDEA

Assign a referring expression type RTA(T) to each table T in Σ .

Example

Is every RTA(.) assignment "good"? Consider the SQLP query

select X.self from PERSON X, COMPANY Y where X.self = Y.self

- assignment: RTA(PERSON) = (ssn = ?), RTA(COMPANY) = (name = ?, city = ?)
 - \Rightarrow the ability to compare the OID values is lost \Rightarrow BAD RTA!;



IDEA

Assign a referring expression type RTA(T) to each table T in Σ .

Example

Is every RTA(.) assignment "good"? Consider the SQLP query

select x.self from PERSON x, COMPANY y where x.self = y.self

- **1 assignment:** RTA(PERSON) = (ssn = ?), RTA(COMPANY) = (name = ?, city = ?)
 - \Rightarrow the ability to compare the OID values is lost \Rightarrow BAD RTA!;
- (modified) assignment:

```
\mathsf{RTA}(\mathsf{COMPANY}) = (\mathsf{PERSON} \to \mathsf{ssn} = ?); (\mathsf{name} = ?, \mathsf{city} = ?)
```

⇒ the ability to compare the OID values is preserved as COMPANY objects are identified by ssn values when also residing in PERSON.



IDEA

Assign a referring expression type RTA(T) to each table T in Σ .

Definition (Identity-resolving RTA(.))

Let Σ be a ARM schema and RTA a referring type assignment for Σ . Given a linear order $\mathcal{O} = (T_{i_1}, \dots, T_{i_n})$ on the set Tables(Σ), define $\mathcal{O}(\mathsf{RTA})$ as the referring expression type RTA(T_{i_1}); . . . ; RTA(T_{i_2}).

We say that RTA is *identity resolving* if there is some linear order \mathcal{O} such that the following conditions hold for each $T \in \text{Tables}(\Sigma)$:

 \blacksquare RTA(T) = Prune($\mathcal{O}(RTA), T$),

Borgida, Toman, and Weddel

- $\Sigma \models (\text{covered by } \{T_1, ..., T_n\}) \in T$, and
- for each component $T_i \to (\mathsf{Pf}_{i,1} = ?, \dots, \mathsf{Pf}_{i,k_i} = ?)$ of RTA(T), the following also holds:
 - (i) Pf_{i,i} is well defined for T_i , for $1 < i < k_i$, and
 - (ii) $\Sigma \models (\text{pathfd } \mathsf{Pf}_{i,1}, \dots, \mathsf{Pf}_{i,k_i} \to id) \in T_i$.



IDEA

Assign a referring expression type RTA(T) to each table T in Σ .

Definition (Identity-resolving RTA(.))

The definition achieves the following:

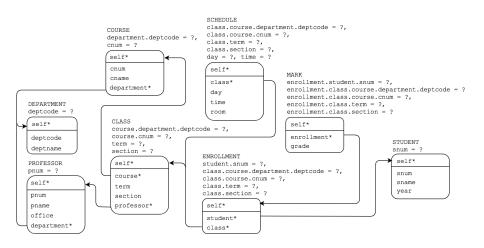
- Referring expression types assigned to classes (tables) that can share objects must guarantee that a particular object is *uniquely* identified;
- Referring expression types for disjoint classes/tables can be assigned independently;

Consequences:

- Referring expressions serve as a sound&complete proxy for entity/object (OID) equality;
- Referring expression can be *coerced* to a least common supertype.



Course Enrollment as an ARM Schema







Concrete Relational Back-end

- 1 Every abstract attribute and its referring expression type
 - ⇒ a concrete relational representation (denoted by Rep(.)): essentially a discriminated variant record;
- Representations can be coerced to a common supertype
 - the ability to compare the representations a sound and complete proxy for comparing object ids;
- A SQLP query is compiled to a standard SQL query over the concrete representation of an abstract instance in such a way that:

Theorem

Let Σ be a ARM schema and let RTA an identity resolving type assignment for $\Sigma.$ For any SQLP $\ query \ Q \ over \ \Sigma$

$$\mathsf{Rep}(Q(I), \Sigma) = (\mathsf{C}^{\Sigma,\mathsf{RTA}}(Q))(\mathsf{Rep}(I, \Sigma))$$

for every database instance I of Σ .





Obtaining an Initial ARM Schema (legacy setting)

RM2ARM Algorithm (highlights; see [EKAW18])

For every table in RM:

- add "self OID" (as a new primary key)
- replace foreign keys with unary ones and discard original FK attributes
 - ⇒ what if original FK overlaps with primary key attributes?
 - ⇒ how about *cycles* between (overlapping) PKs and FKs?
- add ISA constraints (and remove corresponding FKs)
 - ⇒ from PK to PK foreign keys in RM
- 4 add *disjointness* constraints
 - ⇒ for tables with different PKs





Obtaining an Initial ARM Schema (legacy setting)

RM2ARM Algorithm (highlights; see [EKAW18])

For every table in RM:

- 1 add "self OID" (as a new primary key)
- replace foreign keys with unary ones and discard original FK attributes
 - ⇒ what if original FK overlaps with primary key attributes?
 - ⇒ how about cycles between (overlapping) PKs and FKs?
- add ISA constraints (and remove corresponding FKs)
 - ⇒ from PK to PK foreign keys in RM
- 4 add *disjointness* constraints
 - ⇒ for tables with different PKs
- 5 generate *referring expressions* (so the ARM2RM mapping works)



SUMMARY





Summary

Contributions

Referring expressions allow one to get more/better (certain) answers . . .

- General approach to OBDA-style query answering
 - ⇒ Ability to refer to implicit individuals/entities
- 2 General approach to representing data
 - ⇒ without need for *object id invention*;
- Methodology that allows decoupling identification from modeling
 - ⇒ Referring Expressions resolve identity issues and
 - ⇒ Compilation to *pure relational model*.





Future work&Extensions

- Strong Identification (distinct referring expr's refer to distinct objects);
- 2 More complex referring expression types;
- Replacing types by other *preferred way* to chose among referring expressions (e.g., *length/formula complexity/...* measure);
- 4 Alternatives to concrete representations;
- More general/axiomatic definition of identity resolving RTA(.)s;



Message from our Sponsors

Data Systems Group at the University of Waterloo



- Data Systems Group
- ~15 professors, affiliated faculty, postdocs, ~50 grads,
- Wide range of research interests
 - Advanced guery processing/Knowledge representation
 - System aspects of database systems and Distributed data management
 - Data quality/Managing uncertain data/Data mining
 - Information Retrieval and "big data"
 - New(-ish) domains (text, streaming, graph data/RDF, OLAP)
- Research sponsored by governments, and local/global companies NSERC/CFI/OIT and Google, IBM, SAP, OpenText, ...
- Part of a School of CS with 90+ professors, 450+ grad students, etc. Al&ML, Algorithms&Data Structures, PL, Theory, Systems, ...

... and we are always looking for good graduate students (MMath/PhD)



