

Data Complexity of Natural and Controlled Languages

Camilo Thorne¹

¹ KRDB Research Centre for Knowledge and Data
Free University of Bozen-Bolzano
cthorne@inf.unibz.it
<http://www.inf.unibz.it/~cathorne>

DCC UdC - Concepción - 5/8/2011

Outline

- 1 Motivation
- 2 Controlled Languages (CLs)
 - Controlled Natural Language Fragments
 - Natural Language Fragments and Engendered Logics
 - Semantic Complexity
- 3 Semantic Complexity of the Fragments of English (FOEs)
 - The Fragments of English (FOEs)
 - Data and Combined Complexity
- 4 The IS-A_i Fragments and Aggregate Questions
 - The IS-A_i Fragments
 - Aggregate Questions
- 5 Conclusions and Further Work

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires **deep semantic analysis** of English

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope (PP-attachment ambiguity)

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope ⇒ one possibility!

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope ⇒ another possibility!

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope ⇒ two interpretations!

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope ⇒ two interpretations!

- 1 ambiguity blows up processing (exponentially many semantic representations)
- 2 worst case ⇒ computationally unsolvable

Motivation

- Suppose a system that on input

Every Italian loves pasta and football
Silvio is Italian

outputs

Silvio loves pasta

- Requires deep semantic analysis of English
- But, natural language is ambiguous

EX: Consider: John saw a man with a telescope ⇒ two interpretations!

- 1 ambiguity blows up processing (exponentially many semantic representations)
- 2 worst case ⇒ computationally unsolvable

QS: What if we remove ambiguity?

Controlled Languages (CL)

- **Controlled language** (CLs) are ambiguity-free subset of natural languages

- 1 polynomially translate to logic **meaning representations** φ
- 2 translation modelled by formal semantics **compositional translations** $\tau(\cdot)$

Controlled Languages (CL)

- **Controlled language** (CLs) are ambiguity-free subset of natural languages

- 1 polynomially translate to logic **meaning representations** φ
- 2 translation modelled by formal semantics **compositional translations** $\tau(\cdot)$

- Defined by constraining syntax, semantics and vocabulary
- Standard machinery of formal semantics [Montague 1970]
- Controlled languages interfaces proposed [Sowa 2004, Fuchs et al. 2006]

Controlled Languages (CL)

- **Controlled language** (CLs) are ambiguity-free subset of natural languages

- 1 polynomially translate to logic **meaning representations** φ
- 2 translation modelled by formal semantics **compositional translations** $\tau(\cdot)$

- Defined by constraining syntax, semantics and vocabulary
- Standard machinery of formal semantics [Montague 1970]
- Controlled languages interfaces proposed [Sowa 2004, Fuchs et al. 2006]

QS: Efficient translation, but how costly is **logical reasoning**?

Controlled Languages - COP [Pratt & Third 2006]

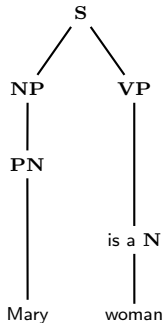
Syntax Rules	Semantics (= $\tau(\cdot)$)
$S \rightarrow NP VP$	$\tau(NP)(\tau(VP)) \triangleright \tau(S)$
$VP \rightarrow \text{is a } N$	$\tau(VP) = \tau(N)$
$VP \rightarrow \text{is not a } N$	$\tau(VP) = \neg\tau(N)$
$NP \rightarrow PN$	$\tau(NP) = \tau(PN)$
$NP \rightarrow \text{Det } N$	$\tau(\text{Det})(\tau(N)) \triangleright \tau(NP)$

Controlled Languages - COP [Pratt & Third 2006]

Syntax Rules	Semantics (= $\tau(\cdot)$)
$S \rightarrow NP VP$	$\tau(NP)(\tau(VP)) \triangleright \tau(S)$
$VP \rightarrow \text{is a } N$	$\tau(VP) = \tau(N)$
$VP \rightarrow \text{is not a } N$	$\tau(VP) = \neg\tau(N)$
$NP \rightarrow PN$	$\tau(NP) = \tau(PN)$
$NP \rightarrow \text{Det } N$	$\tau(\text{Det})(\tau(N)) \triangleright \tau(NP)$

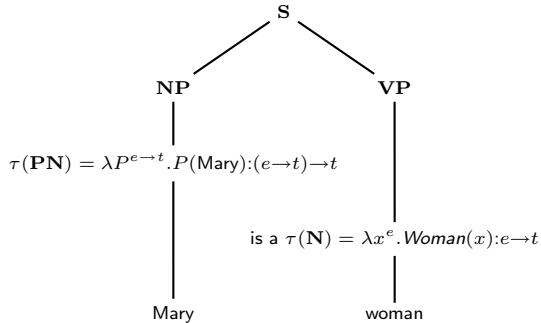
Lexicon	Semantics (= $\tau(\cdot)$)
$N \rightarrow \text{woman}$	$\tau(N) = \lambda x^e. \text{Woman}(x):e \rightarrow t$
$N \rightarrow \text{man}$	$\tau(N) = \lambda x^e. \text{Man}(x):e \rightarrow t$
$PN \rightarrow \text{Mary}$	$\tau(PN) = \lambda P^{e \rightarrow t}. P(\text{Mary}):(e \rightarrow t) \rightarrow t$
$\text{Det} \rightarrow \text{every}$	$\tau(\text{Det}) = \lambda P^{e \rightarrow t}. \lambda Q^{e \rightarrow t}. \forall x^e (P(x) \Rightarrow Q(x)):(e \rightarrow t) \rightarrow ((e \rightarrow t) \rightarrow t)$
$\text{Det} \rightarrow \text{some}$	$\tau(\text{Det}) = \lambda P^{e \rightarrow t}. \lambda Q^{e \rightarrow t}. \exists x^e (P(x) \wedge Q(x)):(e \rightarrow t) \rightarrow ((e \rightarrow t) \rightarrow t)$
$\text{Det} \rightarrow \text{no}$	$\tau(\text{Det}) = \lambda P^{e \rightarrow t}. \lambda Q^{e \rightarrow t}. \forall x^e (P(x) \Rightarrow \neg Q(x)):(e \rightarrow t) \rightarrow ((e \rightarrow t) \rightarrow t)$

Controlled Languages - COP (cntd.)



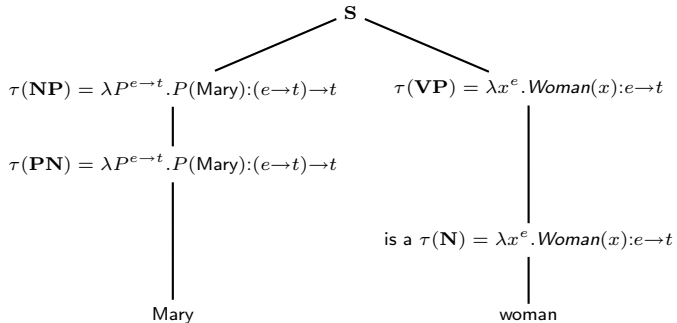
Parsing and interpreting "Mary is a woman"

Controlled Languages - COP (cntd.)



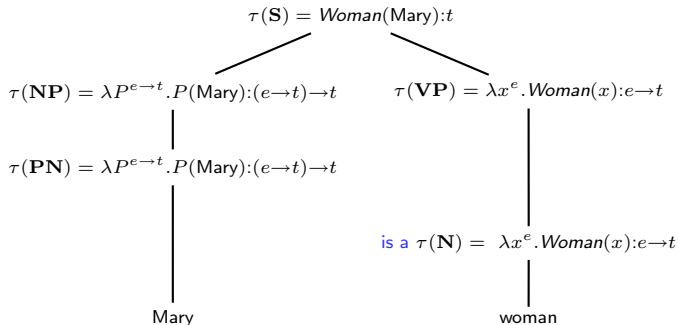
Parsing and interpreting "Mary is a woman"

Controlled Languages – COP (cntd.)



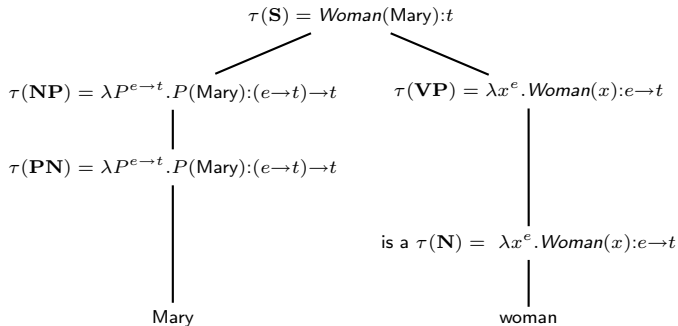
Parsing and interpreting “Mary is a woman”

Controlled Languages - COP (cntd.)



Parsing and interpreting "Mary is a woman"

Controlled Languages - COP (cntd.)



Parsing and interpreting “Mary is a woman”

NB: Computing $\tau(\cdot)$ is “easy” for context-free fragments

- 1 parsing: polynomial in the input utterance
- 2 interpreting: linear in the the parse tree

Engendered Logics and Complexity

- Modulo $\tau(\cdot)$ controlled fragments **express** a fragment of FO

EX: COP expresses

$Woman(Mary)$	Mary is a woman
$\neg Man(Mary)$	Mary is not a man.
$\forall x(Man(x) \Rightarrow Person(x))$	Every man is a person
$\forall x(Woman(x) \Rightarrow \neg Man(x))$	No woman is a man
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\exists x(Person(x) \wedge Woman(x))$	Some person is a woman
$\exists x(Person(x) \wedge \neg Woman(x))$	Some person is not a woman

Engendered Logics and Complexity

- Modulo $\tau(\cdot)$ controlled fragments **express** a fragment of FO

EX: COP expresses

$Woman(Mary)$	Mary is a woman
$\neg Man(Mary)$	Mary is not a man.
$\forall x(Man(x) \Rightarrow Person(x))$	Every man is a person
$\forall x(Woman(x) \Rightarrow \neg Man(x))$	No woman is a man
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\exists x(Person(x) \wedge Woman(x))$	Some person is a woman
$\exists x(Person(x) \wedge \neg Woman(x))$	Some person is not a woman

- By studying such FO fragments

1 exploit computational logic \Rightarrow **computational semantics**

Engendered Logics and Complexity

- Modulo $\tau(\cdot)$ controlled fragments **express** a fragment of FO

EX: COP expresses

$Woman(Mary)$	Mary is a woman
$\neg Man(Mary)$	Mary is not a man.
$\forall x(Man(x) \Rightarrow Person(x))$	Every man is a person
$\forall x(Woman(x) \Rightarrow \neg Man(x))$	No woman is a man
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\forall x(Person(x) \Rightarrow Human(x))$	Every person is a human
$\exists x(Person(x) \wedge Woman(x))$	Some person is a woman
$\exists x(Person(x) \wedge \neg Woman(x))$	Some person is not a woman

- By studying such FO fragments

- exploit computational logic \Rightarrow **computational semantics**
- understand semantic processing complexity \Rightarrow **semantic complexity**

A Survey of Controlled Languages

CL (English)	Maps to	Goal
FOEs [Pratt & Third 2005]	FO fragments	Knowledge representation
ACE [Fuchs 2005]	FO	Knowledge representation
ACE-OWL [Kaaljurand 2007]	OWL-DL	Ontology authoring, querying
PENG [Schwitter 2003]	OWL-DL	Ontology authoring, querying
SOS [Schwitter2008]	OWL-DL	Ontology authoring, querying
CLCE [Sowa2004]	FOL	Knowledge representation
AECMA [Unwalla 2005]	?	User specifications
English Query (EQ) [Blum 1999]	SQL	DB querying/management
OWL-CNL [Schwitter 2006]	OWL-DL	Ontology authoring
Easy English [Bernth 1998]	?	User specifications
λ -SQL [Winter 2006]	SQL	Database querying
nRQL [Schwitter 2008]	FO queries	Ontology querying
Rabbit [Schwitter2008]	OWL	Ontology authoring
ACE-PQL [Bernstein 2005]	PQL	Ontology querying
QE-III [Clifford 1987]	IL	Database querying

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

- 1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?
- 2 **QA**: Given quantified sentences S , facts D and question Q , does $\tau(S) \cup D$ **logically entail** $\tau(Q)$?

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

- 1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?
- 2 **QA**: Given quantified sentences S , facts D and question Q , does $\tau(S) \cup F$ **logically entail** $\tau(Q)$?

- Fine-grained analysis of semantic complexity possible

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

- 1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?
- 2 **QA**: Given quantified sentences S , facts D and question Q , does $\tau(S) \cup D$ **logically entail** $\tau(Q)$?

- Fine-grained analysis of semantic complexity possible

- 1 if measured only in D : **data complexity**

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

- 1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?
- 2 **QA**: Given quantified sentences S , facts D and question Q , does $\tau(S) \cup F$ **logically entail** $\tau(Q)$?

- Fine-grained analysis of semantic complexity possible

- 1 if measured only in D : **data complexity**
- 2 if measured in D , S and Q : **combined complexity**

Semantic Complexity

- We are interested in the **semantic complexity** of CLs
- Semantic complexity \Rightarrow **computational complexity** of decision problems

- 1 **SAT**: Given quantified sentences S and facts D , is $\tau(S) \cup D$ **satisfiable**?
- 2 **QA**: Given quantified sentences S , facts D and question Q , does $\tau(S) \cup F$ **logically entail** $\tau(Q)$?

- Fine-grained analysis of semantic complexity possible

- 1 if measured only in D : **data complexity**
- 2 if measured in D , S and Q : **combined complexity**

NB: Possible when semantic complexity is decidable

The Fragments of English (FOEs) [Pratt & Third 2006]

FOE	Constructs
COP	Copula, common and proper nouns, negation, universal, existential quantifiers
COP+Rel	COP plus relative pronouns
COP+TV	COP plus transitive verbs
COP+TV+DTV	COP+TV plus ditransitive verbs
COP+Rel+TV	COP+Rel plus transitive verbs
COP+Rel+TV+DTV	COP+Rel+TV plus ditransitive verbs
COP+Rel+TV+RA	COP+Rel+TV plus anaphoric pronouns (e.g., he, him, it, herself) of bounded scope
COP+Rel+TV+GA	COP+Rel+TV plus unbounded anaphoric pronouns
COP+Rel+TV+DTV+RA	COP+Rel+TV+DTV plus bounded anaphoric pronouns

Combined Complexity of the FOEs [Pratt & Third 2006]

FOE	SAT (combined)
COP	in NL
COP+TV	NL-complete
COP+TV+DTV	in P
COP+Rel	NP-complete
COP+TV+Rel	EXP-complete
COP+TV+Rel+RA	EXP-complete
COP+TV+DTV+Rel	NEXP-complete
COP+TV+Rel+GA	undecidable

Different function words yield different complexity!

Combined Complexity of the FOEs [Pratt & Third 2006]

FOE	SAT (combined)
COP	in NL
COP+TV	NL-complete
COP+TV+DTV	in P
COP+Rel	NP-complete
COP+TV+Rel	EXP-complete
COP+TV+Rel+RA	EXP-complete
COP+TV+DTV+Rel	NEXP-complete
COP+TV+Rel+GA	undecidable

Not "Boolean-closed" ⇒ tractable!

Combined Complexity of the FOEs [Pratt & Third 2006]

FOE	SAT (combined)
COP	in NL
COP+TV	NL-complete
COP+TV+DTV	in P
COP+Rel	NP-complete
COP+TV+Rel	EXP-complete
COP+TV+Rel+RA	EXP-complete
COP+TV+DTV+Rel	NEXP-complete
COP+TV+Rel+GA	undecidable

“Boolean-closed” ⇒ **intractable!**

Combined Complexity of the FOEs [Pratt & Third 2006]

FOE	SAT (combined)
COP	in NL
COP+TV	NL-complete
COP+TV+DTV	in P
COP+Rel	NP-complete
COP+TV+Rel	EXP-complete
COP+TV+Rel+RA	EXP-complete
COP+TV+DTV+Rel	NEXP-complete
COP+TV+Rel+GA	undecidable

“Boolean-closed” + restricted anaphora ⇒ [decidable!](#)

Combined Complexity of the FOEs [Pratt & Third 2006]

FOE	SAT (combined)
COP	in NL
COP+TV	NL-complete
COP+TV+DTV	in P
COP+Rel	NP-complete
COP+TV+Rel	EXP-complete
COP+TV+Rel+RA	EXP-complete
COP+TV+DTV+Rel	NEXP-complete
COP+TV+Rel+GA	undecidable

“Boolean-closed” + full anaphora ⇒ undecidable!

Data Complexity of the FOEs [Thorne 2011]

- Data complexity [Vardi 1982] measures whether reasoning **scales to large data** repositories
- Typical scenarios: interfaces to knowledge bases, ontologies
- It allows a substantial increase in expressivity:

	QA	SAT
COP	in L	in L
COP+TV	in P	in L
COP+TV+DTV	in coNP	in L
COP+Rel	coNP-complete	in L
COP+Rel+TV	coNP-complete	NP-complete
COP+Rel+DTV	coNP-complete	NP-complete
COP+Rel+DTV+TV	coNP-complete	NP-complete

NB: Questions express conjunction and existential quantification

Data Complexity of the FOEs [Thorne 2011]

- Data complexity [Vardi 1982] measures whether reasoning **scales to large data** repositories
- Typical scenarios: interfaces to knowledge bases, ontologies
- It allows a substantial increase in expressivity:

	QA	SAT
COP	in L	in L
COP+TV	in P	in L
COP+TV+DTV	in coNP	in L
COP+Rel	coNP-complete	in L
COP+Rel+TV	coNP-complete	NP-complete
COP+Rel+DTV	coNP-complete	NP-complete
COP+Rel+DTV+TV	coNP-complete	NP-complete

NB: Questions express conjunction and existential quantification

Data Complexity of the FOEs [Thorne 2011]

- Data complexity [Vardi 1982] measures whether reasoning **scales to large data** repositories
- Typical scenarios: interfaces to knowledge bases, ontologies
- But “Boolean-closedness” hits in quickly!

	QA	SAT
COP	in L	in L
COP+TV	in P	in L
COP+TV+DTV	in coNP	in L
COP+Rel	coNP-complete	in L
COP+Rel+TV	coNP-complete	NP-complete
COP+Rel+DTV	coNP-complete	NP-complete
COP+Rel+DTV+TV	coNP-complete	NP-complete

NB: Questions express conjunction and existential quantification

Ontology Languages

Semantic Web Language (OWL)

```
<owl:Class rdf:about="#Employee">  
  <rdfs:subClassOf>  
    <owl:Restriction>  
      <owl:onProperty  
        rdf:resource="#develops"/>  
      <owl:someValuesFrom  
        rdf:resource="#Project"/>  
    </owl:Restriction>  
  </rdfs:SubclassOf>  
</owl:Class>
```

Description Logics + CL

$Employee \sqsubseteq \exists \text{develops:Project}$

Every employee develops
some project

- OWL is a machine-readable language developed for the semantic web
- CLs are human-readable, yet as unambiguous as OWL

ALCI Ontologies and Ontology Languages

- Given concept names A and role names r , **roles** R and **left** C_l or **right** C_r **concepts** are defined by

$$R \rightarrow r \mid r^- \quad C_f \rightarrow A \mid \exists R:C'_f \mid \neg C'_f \mid C'_f \sqcap C''_f, f \in \{l, r\}$$

- An **assertion** is a formula $C_l \sqsubseteq C_r$
- An **ontology** \mathcal{O} is a set of assertions

ALCI Ontologies and Ontology Languages

- Given concept names A and role names r , **roles** R and **left** C_l or **right** C_r **concepts** are defined by

$$R \rightarrow r \mid r^- \quad C_f \rightarrow A \mid \exists R:C'_f \mid \neg C'_f \mid C'_f \sqcap C''_f, f \in \{l, r\}$$

- An **assertion** is a formula $C_l \sqsubseteq C_r$
- An **ontology** \mathcal{O} is a set of assertions

EX: $Manager \sqsubseteq Employee \Rightarrow$ “every manager is an employee”

- Open world semantics
 - 1 concepts denote sets
 - 2 roles denote relations
 - 3 assertions denote set inclusion

ALCI Ontologies and Ontology Languages

- Given concept names A and role names r , roles R and left C_l or right C_r concepts are defined by

$$R \rightarrow r \mid r^- \quad C_f \rightarrow A \mid \exists R:C'_f \mid \neg C'_f \mid C'_f \sqcap C''_f, f \in \{l, r\}$$

- An **assertion** is a formula $C_l \sqsubseteq C_r$
- An **ontology** \mathcal{O} is a set of assertions

EX: $Manager \sqsubseteq Employee \Rightarrow$ “every manager is an employee”

- Open world semantics
 - 1 concepts denote sets
 - 2 roles denote relations
 - 3 assertions denote set inclusion

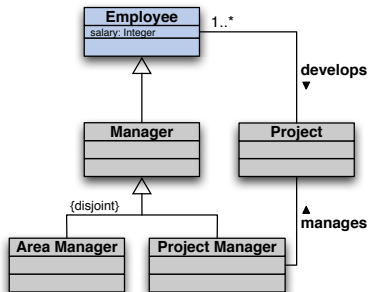
NB: If right or left concepts restricted \Rightarrow **ontology languages**

Conceptual models and OBDA

- Description logic ontologies capture conceptual data models (UML class diagrams, ER-diagrams, etc.)
- In [ontology-based data access](#) (OBDA), relational data is accessed via such ontologies

Conceptual models and OBDA

- Description logic ontologies capture conceptual data models (UML class diagrams, ER-diagrams, etc.)
- In **ontology-based data access** (OBDA), relational data is accessed via such ontologies



$Manager \sqsubseteq Employee$
 $AreaManager \sqsubseteq Manager$
 $ProjectManager \sqsubseteq Manager$
 $AreaManager \sqsubseteq \neg ProjectManager$

$\exists develops \sqsubseteq Employee$
 $\exists develops^- \sqsubseteq Project$
 $Project \sqsubseteq \exists develops^-$

$\exists manages \sqsubseteq TopManager$
 $\exists manages^- \sqsubseteq Project$
 $TopManager \sqsubseteq \exists manages$
 $Project \sqsubseteq \exists manages^-$

$Employee \sqsubseteq \exists has:Salary$

Tree Shaped Queries

- A **tree shaped query** (TSQ) is defined by

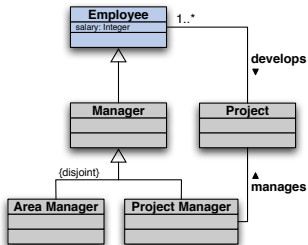
$$\phi(x) \rightarrow A(x) \mid \exists y R(x, y) \mid \phi(x) \wedge \phi'(x) \mid \exists y (R(x, y) \wedge \phi(x))$$

Tree Shaped Queries

- A **tree shaped query** (TSQ) is defined by

$$\phi(x) \rightarrow A(x) \mid \exists y R(x, y) \mid \phi(x) \wedge \phi'(x) \mid \exists y (R(x, y) \wedge \phi(y))$$

- Subset of SQL SELECT-PROJECT-JOIN queries



$$Manager(x) \wedge ProjectManager(x) \wedge \exists y (manages(x, y) \wedge Project(y))$$

```

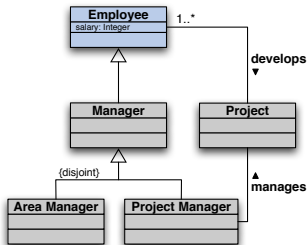
SELECT Manager.MName
FROM Manager, ProjectManager, manages, Project
WHERE Manager.MName=ProjectManager.MName
AND Manager.MName=manages.MName
AND Project.PName=manages.PName
    
```

Tree Shaped Queries

- A **tree shaped query** (TSQ) is defined by

$$\phi(x) \rightarrow A(x) \mid \exists y R(x, y) \mid \phi(x) \wedge \phi'(x) \mid \exists y (R(x, y) \wedge \phi(y))$$

- Subset of SQL SELECT-PROJECT-JOIN queries



$$Manager(x) \wedge ProjectManager(x) \wedge \exists y (manages(x, y) \wedge Project(y))$$

```
SELECT Manager.MName
FROM Manager, ProjectManager, manages, Project
WHERE Manager.MName=ProjectManager.MName
AND Manager.MName=manages.MName
AND Project.PName=manages.PName
```

NB: Semantics of query answering \Rightarrow form of logical entailment!

The IS-A_is [Thorne & Calvanese 2009]

QS: What is the data complexity of CLs that map to OWL?

The IS-A_is [Thorne & Calvanese 2009]

QS: What is the data complexity of CLs that map to OWL?

- We want to single out

- 1 the **maximal** CLs that are **tractable** w.r.t. data complexity
- 2 the **minimal** CLs that are **intractable** w.r.t. data complexity

- Utterances should translate into assertions $C_l \sqsubseteq C_r$
- We partition [Bernardi et al. 2007] **Ns** and **VPs** into

- 1 **left** components: **N_l** and **VP_l**
- 2 **right** components: **N_r** and **VP_r**

The IS-A_is [Thorne & Calvanese 2009]

QS: What is the data complexity of CLs that map to OWL?

- We want to single out

- 1 the **maximal** CLs that are **tractable** w.r.t. data complexity
- 2 the **minimal** CLs that are **intractable** w.r.t. data complexity

- Utterances should translate into assertions $C_l \sqsubseteq C_r$
- We partition [Bernardi et al. 2007] **Ns** and **VPs** into

- 1 **left** components: **N_l** and **VP_l**
- 2 **right** components: **N_r** and **VP_r**

- Engenders a family of CLs

The IS-A_is [Thorne & Calvanese 2009]

QS: What is the data complexity of CLs that map to OWL?

- We want to single out

- 1 the **maximal** CLs that are **tractable** w.r.t. data complexity
- 2 the **minimal** CLs that are **intractable** w.r.t. data complexity

- Utterances should translate into assertions $C_l \sqsubseteq C_r$
- We partition [Bernardi et al. 2007] **Ns** and **VPs** into

- 1 **left** components: **N_l** and **VP_l**
- 2 **right** components: **N_r** and **VP_r**

- Engenders a family of CLs

NB: A similar technique captures TSQs

The IS-A_is [Thorne & Calvanese 2009]

$$\overbrace{\lambda C_l. \lambda C_r. C_l \sqsubseteq C_r}^{\text{every}} \quad \widehat{C_l}^{N_r} \quad \widehat{C_r}^{VP_l}$$

$$\overbrace{\lambda C_l. \lambda C_r. C_l \sqsubseteq C_r}^{\text{everybody who}} \quad \widehat{C_l}^{VP_l} \quad \widehat{C_r}^{VP_r}$$

The IS-A_is [Thorne & Calvanese 2009]

$$\underbrace{\text{every}}_{\lambda C_l . \lambda C_r . C_l \sqsubseteq C_r} \quad \underbrace{N_r}_{C_l} \quad \underbrace{VP_l}_{C_r}$$

$$\underbrace{\text{everybody who}}_{\lambda C_l . \lambda C_r . C_l \sqsubseteq C_r} \quad \underbrace{VP_l}_{C_l} \quad \underbrace{VP_r}_{C_r}$$

Concept C_f	Constituent α_f	Grammar Rules
A	N_f, VP_f	$VP_f \rightarrow \text{is a } N_f \mid IV \mid \text{is Adj}$ $N_f \rightarrow N$
$\exists r : C_f$	TV only N_f	$VP_f \rightarrow TV NP_f$ $NP_f \rightarrow Det N_f$
$C_f \sqcap C_f$	Adj N_f	$VP_f \rightarrow Adj N_f$ $N_f \rightarrow Adj N_f$
		⋮

The IS-A_is [Thorne & Calvanese 2009]

$$\underbrace{\text{every}}_{\lambda C_l. \lambda C_r. C_l \sqsubseteq C_r} \underbrace{N_r}_{C_l} \underbrace{VP_l}_{C_r} \qquad \underbrace{\text{everybody who}}_{\lambda C_l. \lambda C_r. C_l \sqsubseteq C_r} \underbrace{VP_l}_{C_l} \underbrace{VP_r}_{C_r}$$

Concept C_f	Constituent α_f	Grammar Rules
A	N_f, VP_f	$VP_f \rightarrow \text{is a } N_f \mid IV \mid \text{is Adj}$ $N_f \rightarrow N$
$\exists r: C_f$	TV only N_f	$VP_f \rightarrow \text{TV } NP_f$ $NP_f \rightarrow \text{Det } N_f$
$C_f \sqcap C_f$	Adj N_f	$VP_f \rightarrow \text{Adj } N_f$ $N_f \rightarrow \text{Adj } N_f$

⋮

CL	Assertions	Sample Sentence(s)
IS-A ₁	$A \sqsubseteq \forall r: A'$	Every project manager manages only projects.
IS-A ₂	$A \sqcap A' \sqsubseteq \forall r: (A'' \sqcap A''')$	Every good manager manages only good projects.

⋮

Data Complexity [Thorne & Calvanese 2011]

	SAT (data)	QA (data)	QA (combined)
IS-A ₁	in P	NL-complete	in PSPACE
IS-A ₂	in P	P-complete	in PSPACE
IS-A ₃	in P	P-complete	in NEXP
IS-A ₄	in P	P-complete	in PSPACE
IS-A ₅	in P	CONP-complete	in NEXP
IS-A ₆	in P	CONP-complete	CONP-complete
IS-A ₇	in P	CONP-complete	CONP-complete

Data Complexity [Thorne & Calvanese 2011]

	SAT (data)	QA (data)	QA (combined)
IS-A ₁	in P	NL-complete	in PSPACE
IS-A ₂	in P	P-complete	in PSPACE
IS-A ₃	in P	P-complete	in NEXP
IS-A ₄	in P	P-complete	in PSPACE
IS-A ₅	in P	coNP-complete	in NEXP
IS-A ₆	in P	coNP-complete	coNP-complete
IS-A ₇	in P	coNP-complete	coNP-complete

Data Complexity [Thorne & Calvanese 2011]

	SAT (data)	QA (data)	QA (combined)
IS-A ₁	in P	NL-complete	in PSPACE
IS-A ₂	in P	P-complete	in PSPACE
IS-A ₃	in P	P-complete	in NEXP
IS-A ₄	in P	P-complete	in PSPACE
IS-A ₅	in P	coNP-complete	in NEXP
IS-A ₆	in P	coNP-complete	coNP-complete
IS-A ₇	in P	coNP-complete	coNP-complete

⇒ Data intractability is caused by “Boolean-closedness”

Data Complexity [Thorne & Calvanese 2011]

	SAT (data)	QA (data)	QA (combined)
IS-A ₁	in P	NL-complete	in PSPACE
IS-A ₂	in P	P-complete	in PSPACE
IS-A ₃	in P	P-complete	in NEXP
IS-A ₄	in P	P-complete	in PSPACE
IS-A ₅	in P	coNP-complete	in NEXP
IS-A ₆	in P	coNP-complete	coNP-complete
IS-A ₇	in P	coNP-complete	coNP-complete

⇒ Data intractability is caused by “Boolean-closedness”

NB: Maximal tractable CL obtained by eliminating negation!

Semantic Parsers and Aggregations

- Boxer [Bos 2005]: state-of-the art wide-coverage English semantic parser generating logic meaning representations

Which is the **number of** students who study something, **per country**?

$$\begin{aligned} \forall x_1 (& \mathbf{thing}(x_1) \Rightarrow \exists x_2 \exists x_3 \exists x_4 \exists x_5 \exists x_6 \\ & (\mathbf{number}(x_2) \wedge \mathbf{of}(x_2, x_3)) \wedge \\ & \mathbf{student}(x_3) \wedge \mathbf{study}(x_5) \wedge \mathbf{event}(x_5) \\ & \wedge \mathbf{agent}(x_3, x_5) \wedge \mathbf{patient}(x_6, x_4) \wedge \mathbf{thing}(x_6) \\ & \mathbf{per}(x_2, x_4) \wedge \mathbf{country}(x_4)) \wedge \\ & \exists x_7 (\mathbf{event}(x_7) \wedge x_1 = x_2)) \end{aligned}$$

Semantic Parsers and Aggregations

- Boxer [Bos 2005]: state-of-the art wide-coverage English semantic parser generating logic meaning representations

Which is the **number of** students who study something, **per country**?

$$\begin{aligned} \forall x_1 (& \mathbf{thing}(x_1) \Rightarrow \exists x_2 \exists x_3 \exists x_4 \exists x_5 \exists x_6 \\ & (\mathbf{number}(x_2) \wedge \mathbf{of}(x_2, x_3) \wedge \\ & \mathbf{student}(x_3) \wedge \mathbf{study}(x_5) \wedge \mathbf{event}(x_5) \\ & \wedge \mathbf{agent}(x_3, x_5) \wedge \mathbf{patient}(x_6, x_4) \wedge \mathbf{thing}(x_6) \\ & \mathbf{per}(x_2, x_4) \wedge \mathbf{country}(x_4)) \wedge \\ & \exists x_7 (\mathbf{event}(x_7) \wedge x_1 = x_2)) \end{aligned}$$

- 1 Aggregations **modify** nominals compound expressing TSQs
- 2 Aggregations are very common in English

Semantic Parsers and Aggregations

- Boxer [Bos 2005]: state-of-the art wide-coverage English semantic parser generating logic meaning representations

Which is the **number of** students who study something, **per country**?

$$\begin{aligned} \forall x_1 (& \mathbf{thing}(x_1) \Rightarrow \exists x_2 \exists x_3 \exists x_4 \exists x_5 \exists x_6 \\ & (\mathbf{number}(x_2) \wedge \mathbf{of}(x_2, x_3) \wedge \\ & \mathbf{student}(x_3) \wedge \mathbf{study}(x_5) \wedge \mathbf{event}(x_5) \\ & \wedge \mathbf{agent}(x_3, x_5) \wedge \mathbf{patient}(x_6, x_4) \wedge \mathbf{thing}(x_6) \\ & \mathbf{per}(x_2, x_4) \wedge \mathbf{country}(x_4)) \wedge \\ & \exists x_7 (\mathbf{event}(x_7) \wedge x_1 = x_2)) \end{aligned}$$

- 1 Aggregations **modify** nominals compound expressing TSQs
- 2 Aggregations are very common in English

QS: Which is the **data complexity** of CLs with aggregations?

Aggregate Determiners [Thorne & Calvanese 2009]

- **Aggregate determiners** combine with nominals to produce **NPs** denoting numbers

Det → the average
Det → the number of

$$\begin{aligned}\tau(\mathbf{Det}) &:= \lambda P^{\mathbb{Q} \rightarrow \mathbb{N}}.\mathbf{avg}(P): (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q} \\ \tau(\mathbf{Det}) &:= \lambda P^{e \rightarrow \mathbb{N}}.\mathbf{count}(P): (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q}\end{aligned}$$

Aggregate Determiners [Thorne & Calvanese 2009]

- **Aggregate determiners** combine with nominals to produce **NPs** denoting numbers

Det → the average
Det → the number of

$$\begin{aligned} \tau(\mathbf{Det}) &:= \lambda P^{\mathbb{Q} \rightarrow \mathbb{N}}. \mathbf{avg}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q} \\ \tau(\mathbf{Det}) &:= \lambda P^{e \rightarrow \mathbb{N}}. \mathbf{count}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q} \end{aligned}$$

- We use again **semantically annotated grammars**

$$\begin{array}{l} \mathbf{Q}_{wh} \rightarrow \mathbf{I}_i \mathbf{S}_{g_i} ? \\ \mathbf{S}_{g_i} \rightarrow \mathbf{NP}_{g_i} \mathbf{VP}_i \end{array} \quad \lambda \bar{z}. \tau(\mathbf{Intpro}_i)(\lambda i. \tau(\mathbf{S}_{g_i})) \triangleright \tau(\mathbf{Q}_{wh})$$

$$\tau(\mathbf{NP}_{g_i})(\tau(\mathbf{VP}_i)) \triangleright \tau(\mathbf{S}_{g_i})$$

Aggregate Determiners [Thorne & Calvanese 2009]

- Aggregate determiners combine with nominals to produce NPs denoting numbers

Det → the average
 Det → the number of

$$\begin{aligned} \tau(\text{Det}) &:= \lambda P^{\mathbb{Q} \rightarrow \mathbb{N}}. \text{avg}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q} \\ \tau(\text{Det}) &:= \lambda P^{e \rightarrow \mathbb{N}}. \text{count}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q} \end{aligned}$$

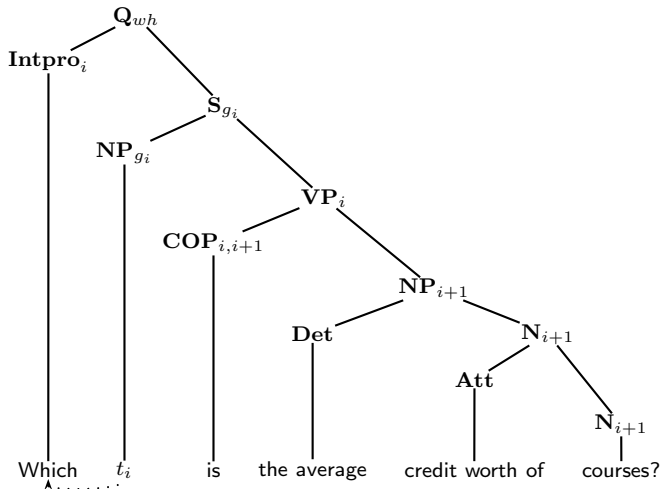
- We use again **semantically annotated grammars**

$$\begin{array}{l} \mathbf{Q}_{wh} \rightarrow \mathbf{I}_i \mathbf{S}_{g_i} ? \\ \mathbf{S}_{g_i} \rightarrow \mathbf{NP}_{g_i} \mathbf{VP}_i \end{array} \quad \lambda \bar{z}. \tau(\mathbf{Intpro}_i)(\lambda i. \tau(\mathbf{S}_{g_i})) \triangleright \tau(\mathbf{Q}_{wh})$$

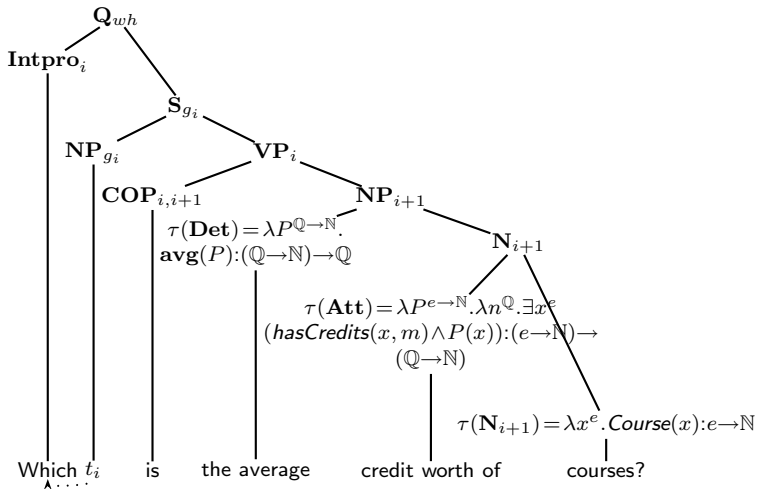
$$\tau(\mathbf{NP}_{g_i})(\tau(\mathbf{VP}_i)) \triangleright \tau(\mathbf{S}_{g_i})$$

- We consider **bag-typed** (i.e., $e \rightarrow \mathbb{N}$, $\mathbb{Q} \rightarrow \mathbb{N}$, etc.) constituents

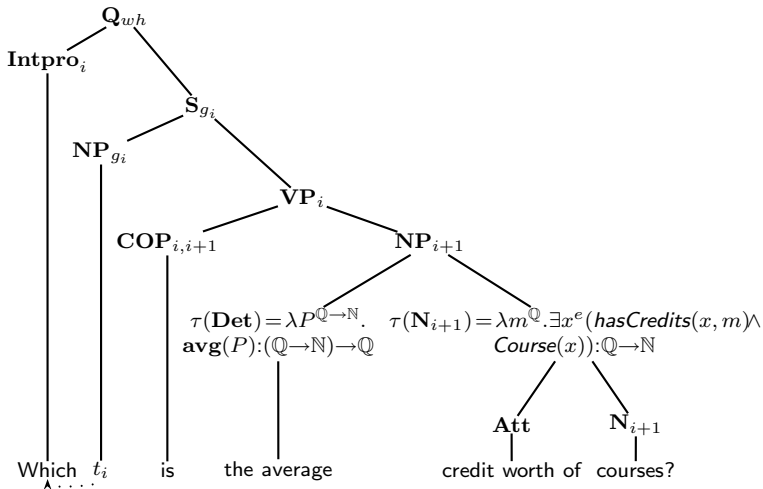
Aggregate Questions: Parsing Example



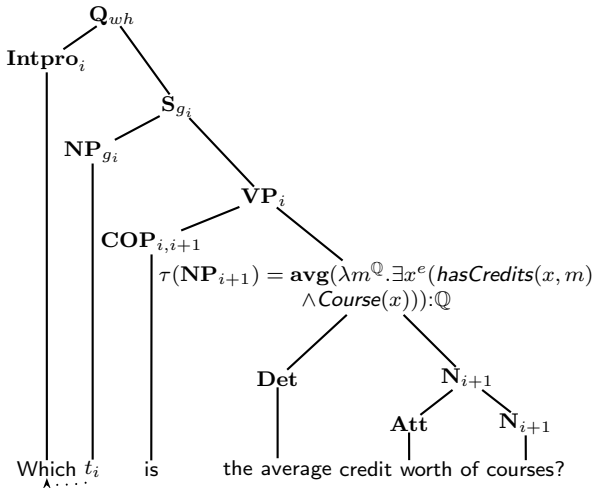
Aggregate Questions: Parsing Example



Aggregate Questions: Parsing Example

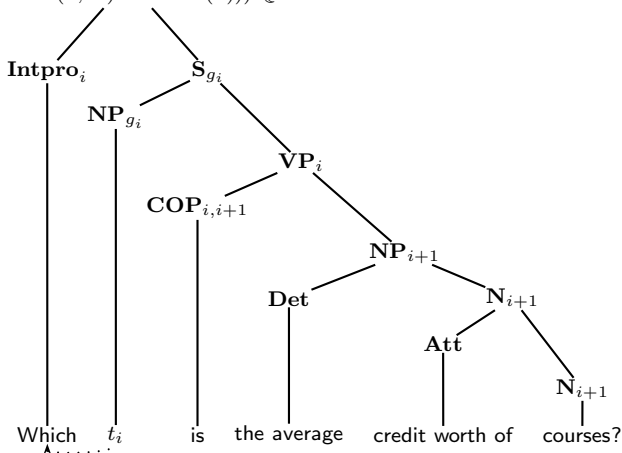


Aggregate Questions: Parsing Example

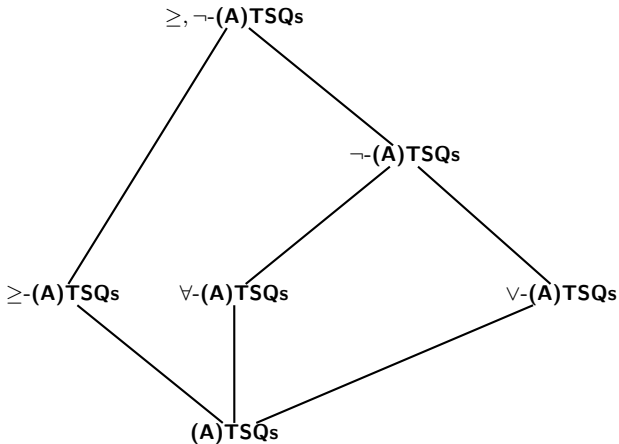


Aggregate Questions: Parsing Example

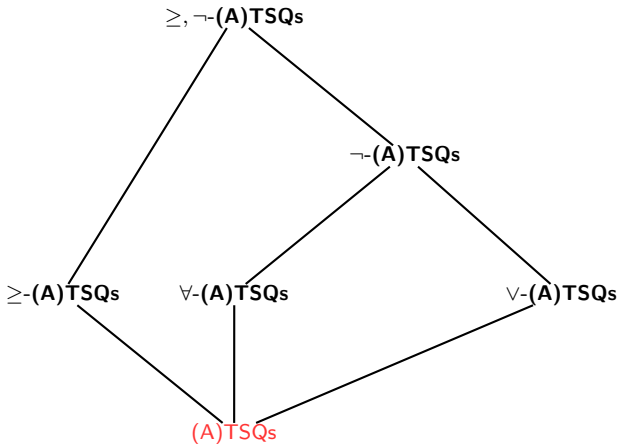
$$\tau(Q_{wh}) = \lambda n^Q. n \approx \text{avg}(\lambda m^Q. \exists x^e. (\text{hasCredits}(x, m) \wedge \text{Course}(x))): Q \rightarrow t$$



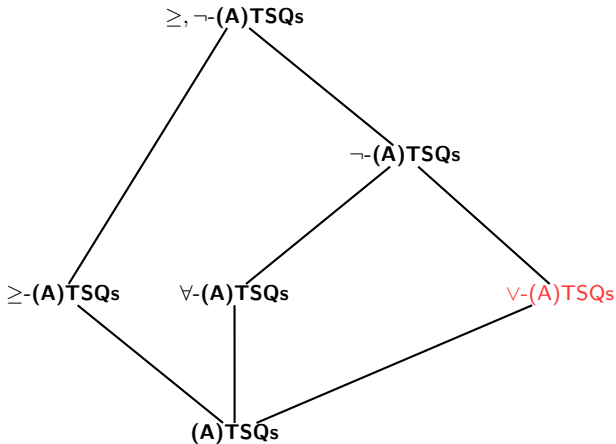
Aggregate Questions and Beyond [Thorne & Calvanese 2009]



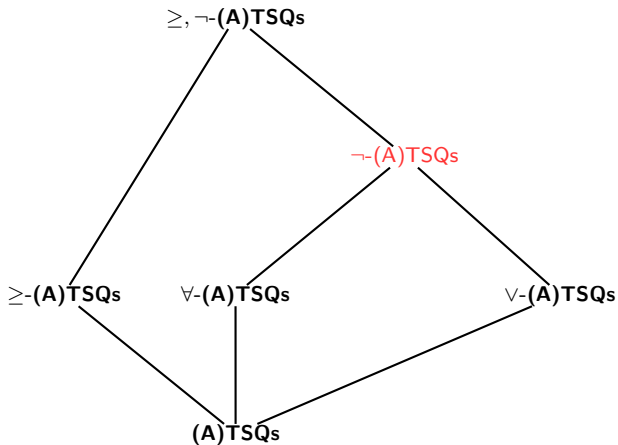
Aggregate Questions and Beyond [Thorne & Calvanese 2009]



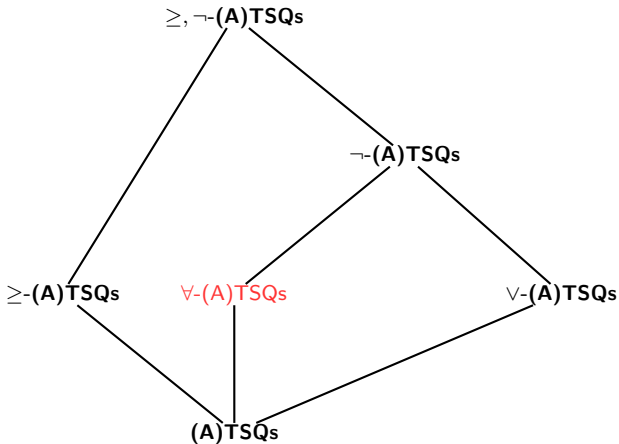
Aggregate Questions and Beyond [Thorne & Calvanese 2009]



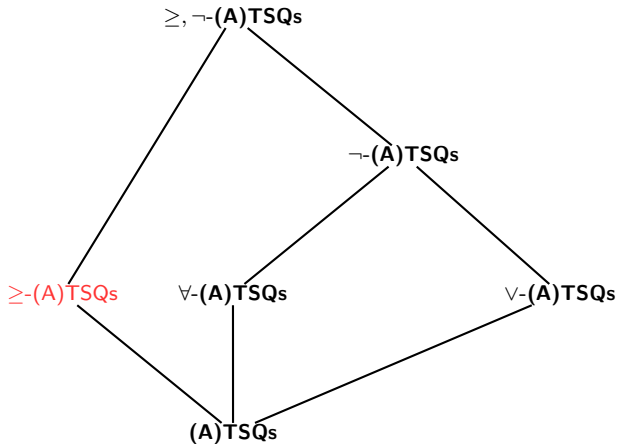
Aggregate Questions and Beyond [Thorne & Calvanese 2009]



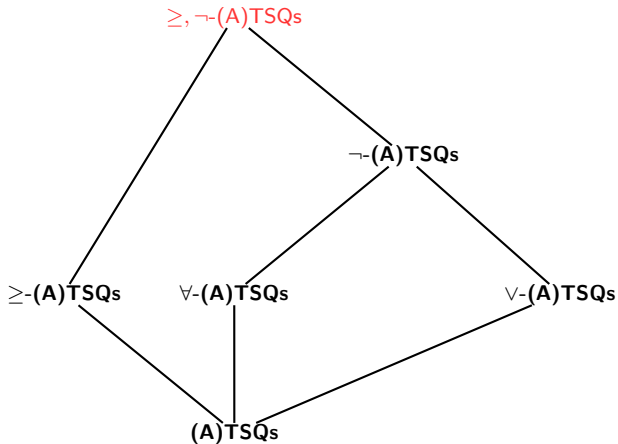
Aggregate Questions and Beyond [Thorne & Calvanese 2009]



Aggregate Questions and Beyond [Thorne & Calvanese 2009]



Aggregate Questions and Beyond [Thorne & Calvanese 2009]



Data Complexity [Thorne & Calvanese 2009]

Data Complexity [Thorne & Calvanese 2009]

- Complexity is tractable (in P) with

- 1 aggregates over TSQs (space logarithmic overhead)
- 2 disjunctions of TSQs, $\phi(x) \vee \phi'(x)$

Data Complexity [Thorne & Calvanese 2009]

- Complexity is **tractable** (in P) with

- 1 **aggregates** over TSQs, (space logarithmic overhead)
- 2 **disjunctions** of TSQs, $\phi(x) \vee \phi'(x)$

- Complexity is **intractable** (coNP-hard) with

- 1 **negations** of TSQs, $\neg\phi(x)$
- 2 **universally restricted** TSQs, $\forall x(R(x, y) \rightarrow \phi(y))$
- 3 **comparisons** in TSQs, $t \theta t'$, $\theta \in \{=, >, <, \geq, \leq\}$

- \geq, \neg (A)TSQs are in coNP

Data Complexity [Thorne & Calvanese 2009]

- Complexity is **tractable** (in P) with

- 1 **aggregates** over TSQs (space logarithmic overhead)
- 2 **disjunctions** of TSQs, $\phi(x) \vee \phi'(x)$

- Complexity is **intractable** (coNP-hard) with

- 1 **negations** of TSQs, $\neg\phi(x)$
- 2 **universally restricted** TSQs, $\forall x(R(x, y) \rightarrow \phi(y))$
- 3 **comparisons** in TSQs, $t \theta t'$, $\theta \in \{=, >, <, \geq, \leq\}$

- \geq, \neg (A)TSQs are in coNP

NB: We consider queries over **ontologies** (i.e., OBDA), not relational databases

Conclusions and Further Work

- Studied controlled languages and semantic complexity
- Argued that semantic complexity is an important issue
- Eliminating ambiguity makes translation efficient, but reasoning becomes intractable with simple fragments
- Pinpointed (data) tractable and intractable fragments

Conclusions and Further Work

- Studied controlled languages and semantic complexity
- Argued that semantic complexity is an important issue
- Eliminating ambiguity makes translation efficient, but reasoning becomes intractable with simple fragments
- Pinpointed (data) tractable and intractable fragments
- Further work: how can theoretical semantic complexity be
 - 1 empirically validated by behavior (processing time) of systems?
 - 2 empirically validated by behavior of automated reasoners?

Any Questions?

THANK YOU!



FREIE UNIVERSITÄT BOZEN
LIBERA UNIVERSITÀ DI BOLZANO
FREE UNIVERSITY OF BOZEN · BOLZANO



Some References

■ DLs and OBDA:

- Diego Calvanese et al., "Tractable reasoning and efficient query answering in description logics: the *DL-Lite* family", 2007, JAIR
- Diego Calvanese et al., "The *DL-Lite* Family and Relations", 2009, JAIR

■ CLs and semantic complexity:

- Norbert Fuchs et al., "Controlled English for reasoning on the semantic web", 2009, Semantic Techniques for the Web: the REVERSE Perspective
- Ian Pratt and Allan Third, "More Fragments of Language", 2006, Notre Dame Journal of Formal Logic

■ Data complexity of CLs:

- Camilo Thorne, "Querying the Fragments of English", 2011, Proceedings of WoLLIC 2011
- Camilo Thorne, Diego Calvanese, "Tractability and Intractability of Controlled Languages for Data Access", 2011, Studia Logica
- Camilo Thorne, Diego Calvanese, "The Data Complexity of the Syllogistic Fragments of English", Proceedings of AC 2009
- Camilo Thorne, Diego Calvanese, "Tree-shaped aggregate questions over ontologies", 2009, Proceedings of FQAS 2009
- Raffaella Bernardi et al., "Lite natural language", 2007, Proceedings of IWCS-7

(Data) Scalable and Non-Scalable Constructs [Thorne 2010]

	Declarations	Questions
Constructs that scale (P or less)	<ul style="list-style-type: none"> - Negation in predicate VPs, relatives in predicate VPs, conjunction in predicate VPs - Relatives and conjunction in subject NPs and predicate VPs, but no negation 	<ul style="list-style-type: none"> - Existential quantifiers, conjunction, relatives, disjunctions, aggregations
Constructs that do not (CONP-hard)	<ul style="list-style-type: none"> - Negation in subject NPs - Relatives and negation in subject NPs and predicate VPs 	<ul style="list-style-type: none"> - Full negation - Comparisons - Universal restrictions
Undecidable Constructs	<ul style="list-style-type: none"> - TVs, relatives, negation, existential and universal quantifiers, restricted anaphoric pronouns and indeterminate pronouns in subject NPs and predicate VPs, copula 	<ul style="list-style-type: none"> - TVs, existential indeterminate pronouns, relatives and restricted anaphoric pronouns

The IS- A_i s [Thorne & Calvanese 2009]

Concept C_f	Constituent α_f	Grammar Rules
A	N_f, VP_f	
$\exists P: A$	TV some N_f , TV somebody who VP_f	$VP_f \rightarrow \text{is a } N_f \mid \text{IV} \mid \text{is Adj}$
$\exists P^-: A$	TV by some N_f , TV by somebody who VP_f	$N_f \rightarrow N$
$\forall P: A$	TV only VP_f , TV only who VP_f	
$\exists P$	TV something, TV somebody	\emptyset
$A_1 \sqcap \dots \sqcap A_n$	Adj N_f, N_f who VP_f	$VP_f \rightarrow \text{is a } N_f \mid \text{IV} \mid \text{is Adj}$ $\mid VP_f \text{ and } VP_f$
	N_f and N_f, VP_f and VP_f	$N_f \rightarrow N \mid \text{Adj } N_f$ $\mid N_f \text{ and } N_f$
$A_1 \sqcup \dots \sqcup A_n$	VP_f or VP_f	$VP_f \rightarrow \text{is a } N_f \mid \text{IV} \mid \text{is Adj}$ $\mid VP_f \text{ and } VP_f$ $N_f \rightarrow N \mid N_f \text{ and } N_f$
$\neg A$	is not Adj , does not IV , is not a N_f	$N_f \rightarrow N$

The IS-A_is [Thorne & Calvanese 2009]

CL	Assertions	Sample Sentence(s)
IS-A ₁	$A \sqsubseteq \forall r:A'$	Every project manager manages only projects.
IS-A ₂	$A \sqcap A' \sqsubseteq \forall r:(A'' \sqcap A''')$	Every good manager manages only good projects.
IS-A ₃	$\exists r:A \sqsubseteq A' \sqcap A''$	Anybody who manages some project is an employee and is a manager.
	$\exists r^-:A \sqsubseteq A' \sqcap A''$	Anything that is managed by some important manager is a big project.
	$A \sqsubseteq \exists P$	Every manager manages something.
IS-A ₄	$A \sqcap A' \sqsubseteq A'' \sqcap A'''$	Every cruel manager is a bad manager.
	$\exists r:(A \sqcap A') \sqsubseteq A'' \sqcap A'''$	Anybody who manages some bankrupt project is a bad manager.
IS-A ₅	$\forall r:A \sqsubseteq A' \sqcap A''$	Anybody who manages only projects is a manager and a project manager.
IS-A ₆	$A \sqsubseteq A' \sqcup A''$	Every manager is a project manager or is an area manager.
IS-A ₇	$\neg A \sqsubseteq A' \sqcap A''$	Anybody who is not an area manager is an employee who is a project manager.

FO Semantics of \mathcal{ALCI} [Baader et al. 2004]

Consider FO interpretations $\mathcal{I} := (\Delta, \cdot^{\mathcal{I}})$:

$$\begin{aligned}
 A^{\mathcal{I}} &\subseteq \Delta \\
 (\exists R:C_f)^{\mathcal{I}} &:= \{d \mid \text{exists } d' \text{ s.t.} \\
 &\quad (d, d') \in R^{\mathcal{I}} \text{ and } d' \in C_f^{\mathcal{I}}\} \\
 (\neg C_f)^{\mathcal{I}} &:= \Delta \setminus C_f^{\mathcal{I}} \\
 (C_f \sqcap C'_f)^{\mathcal{I}} &:= C_f^{\mathcal{I}} \cap C'^{\mathcal{I}}_f \\
 r^{\mathcal{I}} &\subseteq \Delta \times \Delta \\
 (r^-)^{\mathcal{I}} &:= \{(d, d') \mid (d', d) \in r^{\mathcal{I}}\}
 \end{aligned}
 \quad
 \begin{aligned}
 \mathcal{I} \models C_l \sqsubseteq C_r &\text{ iff } C_l^{\mathcal{I}} \subseteq C_r^{\mathcal{I}} \\
 \mathcal{I} \models \mathcal{O} &\text{ iff for all } C_l \sqsubseteq C_r \in \mathcal{O}, \\
 &\quad \mathcal{I} \models C_l^{\mathcal{I}} \subseteq C_r^{\mathcal{I}}
 \end{aligned}$$

Open world semantics: \mathcal{O} may have many models!

Research Methodology [Thorne 2010]

We proposed the following strategy:

- 1 Define a **declarative** CL L_D and a grammar for it
- 2 Define a meaning representation $\tau(L_D)$
- 3 Define an **interrogative** CL L_Q and a grammar for it
- 4 Define a meaning representation $\tau(L_Q)$
- 5 Study the properties of
 - 1 the **ontology language** $\tau(L_D)$
 - 2 the **query language** $\tau(L_Q)$

⇒ infer data complexity by reasoning on $\tau(L_D)$ and $\tau(L_Q)$