



Tree-Shaped Aggregate Questions over Ontologies

Camilo Thorne, Diego Calvanese

KRDB Research Centre

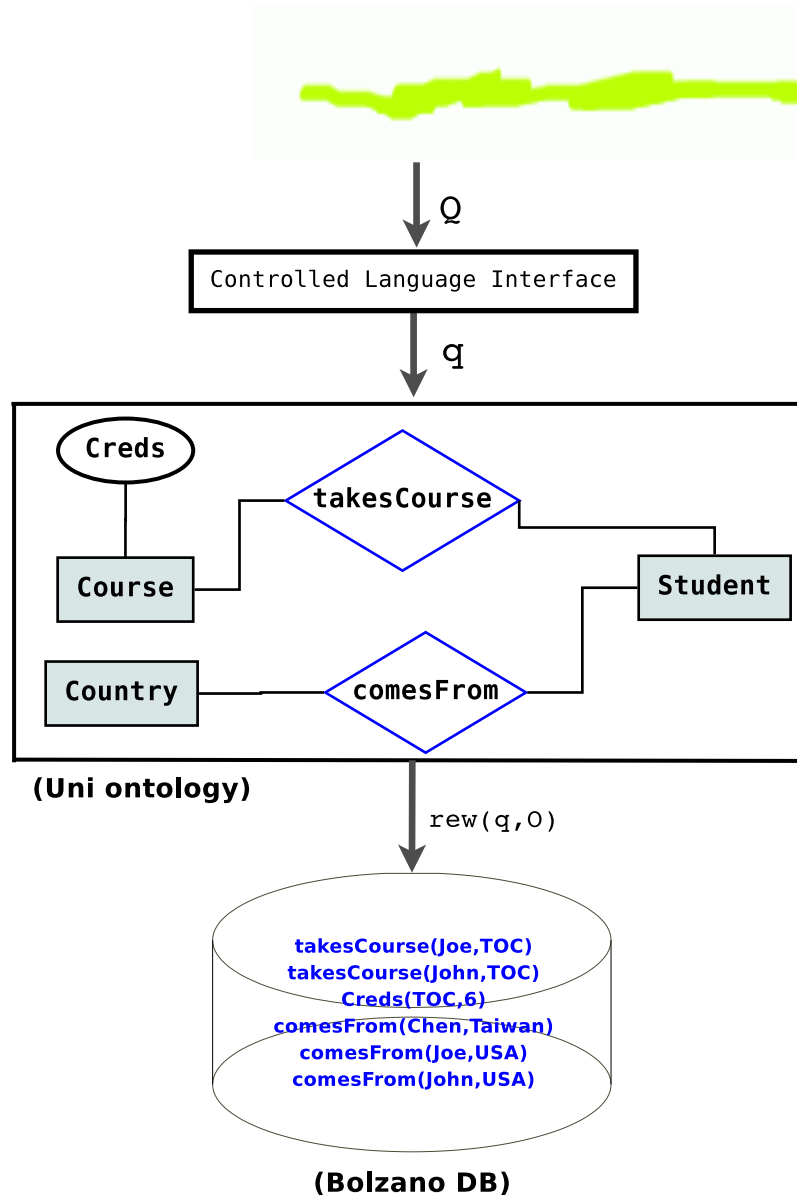
Free University of Bozen-Bolzano

Flexible Query Answering Systems - FQAS 2009

CL Interfaces

- ⑥ **Controlled language interfaces** (CLIs) to information systems have been proposed to address the usability problem [Sowa 2004, Fuchs et al. 2006]
- ⑥ Controlled languages (CLs) are ambiguity-free fragments of natural language which
 - symbolically translate into formal ontology and/or query languages
 - the translation is compositional and polynomial in the size of the input
- ⑥ More recently, CLIs have been proposed for ontology-based data access systems (OBDASs) [Bernstein et al. 2005, 2007]

OBDASs



In an ontology based data access system $(\mathcal{O}, \mathcal{D})$ data stored in databases \mathcal{D} is queried through an ontology \mathcal{O}

Questions are translated into queries, expanded w.r.t. \mathcal{O} and shipped to \mathcal{D}

Measuring CLI Performance

- ⑥ An important measure of performance of OBDASs is the **data complexity** of query evaluation [Vardi 1982]
- ⑥ CL constructs in formal ontology/query languages have an impact on data complexity:
 - which combinations of CL constructs give way to **tractable** (**P**TIME or less) data complexity?
 - which combinations of CL constructs give rise to **intractable** (**coNP**-hard) data complexity?

Measuring CLI Performance

- ⑥ An important measure of performance of OBDASs is the **data complexity** of query evaluation [Vardi 1982]
- ⑥ CL constructs in formal ontology/query languages have an impact on data complexity:
 - which combinations of CL constructs give way to **tractable** (**P**TIME or less) data complexity?
 - which combinations of CL constructs give rise to **intractable** (**coNP**-hard) data complexity?
- ⑥ We can study this impact by looking at what happens with the formal query/ontology languages they translate into!
- ⑥ Target: SQL aggregate queries

Tree-Shaped Conditions and Tree Shaped Queries

- ⑥ We consider logic-based declarative counterparts of SQL
- ⑥ A tree-shaped condition (TSC) $\varphi(x)$ is a FOL formula

$$\begin{aligned}\varphi(x) ::= & A(x) \mid \exists y R(x, y) \\ & \mid \exists y (R(x, y) \wedge \varphi(y)) \\ & \mid \varphi(x) \wedge \varphi'(x)\end{aligned}$$

- ⑥ A tree-shaped query (TSQ) q is a query

$$q(x) \leftarrow \varphi(x)$$

with $\varphi(x)$ a TSC

Tree-Shaped Conditions and Tree Shaped Queries

- ⑥ We consider logic-based declarative counterparts of SQL
- ⑥ A tree-shaped condition (TSC) $\varphi(x)$ is a FOL formula

$$\begin{aligned}\varphi(x) ::= & A(x) \mid \exists y R(x, y) \\ & \mid \exists y (R(x, y) \wedge \varphi(y)) \\ & \mid \varphi(x) \wedge \varphi'(x)\end{aligned}$$

- ⑥ A tree-shaped query (TSQ) q is a query

$$q(x) \leftarrow \varphi(x)$$

with $\varphi(x)$ a TSC

- ⑥ In domains with numerical data we want also to aggregate!

Aggregate Tree-Shaped Queries

⑥ We consider SQL aggregation functions: **max**, **min**, **count**, **cntd**, **sum** and **avg**, denoted α

⑥ An aggregate tree-shaped query (ATQ) q is a query

$$q(x, \alpha(y)) \leftarrow \varphi_1(x) \wedge R(x, y) \wedge \varphi_2(y)$$

- x is the grouping variable
- y is the aggregation variable
- $\alpha(y)$ is the aggregation term
- $\varphi_1(x)$ and $\varphi_2(y)$ are TSCs
- $R(x, y)$ is a FOL atom

⑥ Note: $\varphi_1(x)$ expresses GROUP-BY clauses

Ontology Languages

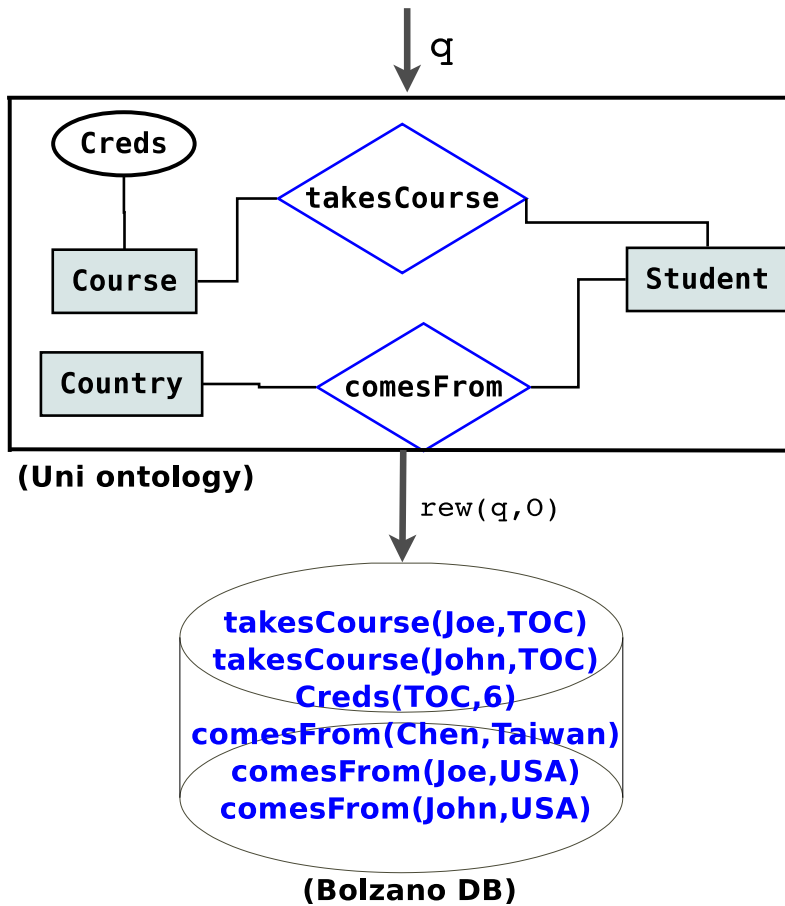
- ⑥ We are interested in the **OWL** ontology language and its fragments
- ⑥ OWL is formally underpinned by **description logics** (DLs)
- ⑥ **DL-Lite**: captures the fundamental features of conceptual modelling languages (e.g., ER-diagrams)
 - IS-A, disjointness, keys
 - role typing, mandatory participation
- ⑥ **ALCI**: the least subset of OWL containing *DL-Lite* and closed under boolean operations

Certain Answers

- ⑥ Aggregate functions are interpreted as in SQL
- ⑥ Given $q(x, \alpha(y)) \leftarrow \psi$ and $(\mathcal{O}, \mathcal{D})$:
 - a certain assignments γ is
 - a FOL **satisfying assignment** for ψ
over each $\mathcal{D}' \supseteq \mathcal{D}$ s.t. $\mathcal{D}' \models \mathcal{O}$
 - the certain group $H_{\gamma(x)}$ collects
 - the **values** $\gamma'(y)$ of y
for all γ s.t. $\gamma'(x) = \gamma(x)$
 - the set of certain answers $\text{cert}(q, \mathcal{O}, \mathcal{D})$ collects
 - the **tuples** $(\gamma(x), \alpha(H_{\gamma(x)}))$

Example

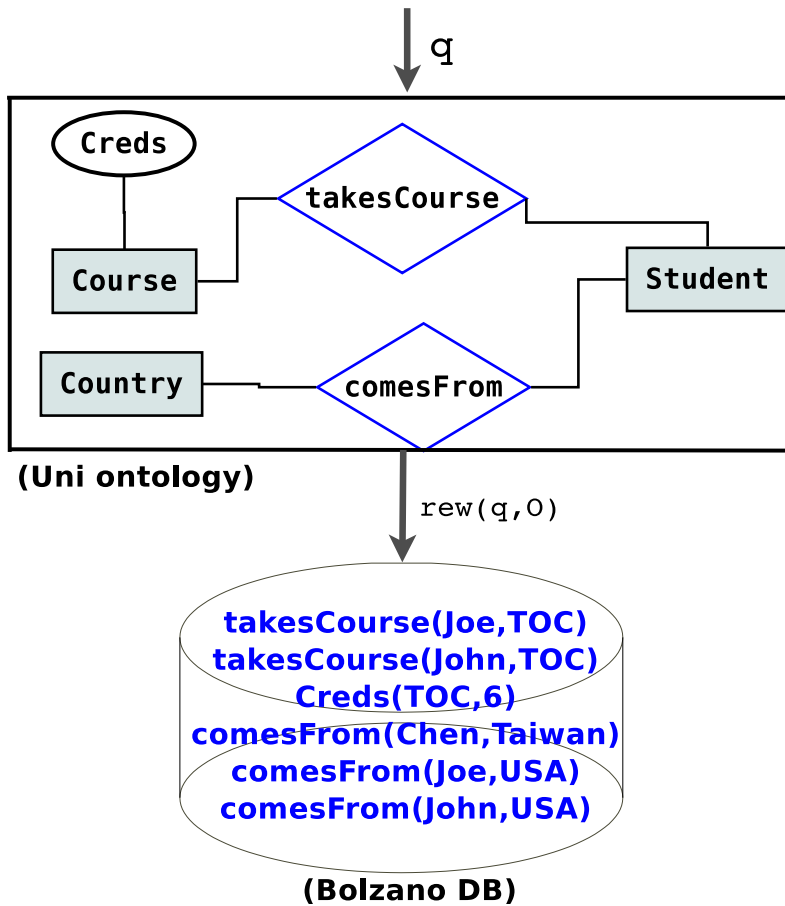
Which is the number of students, per country, who study something?



$$q(x, count(z)) \leftarrow Student(x) \wedge \exists y (takesCourse(x, y)) \wedge comesFrom(x, z) \wedge Country(z)$$

Example

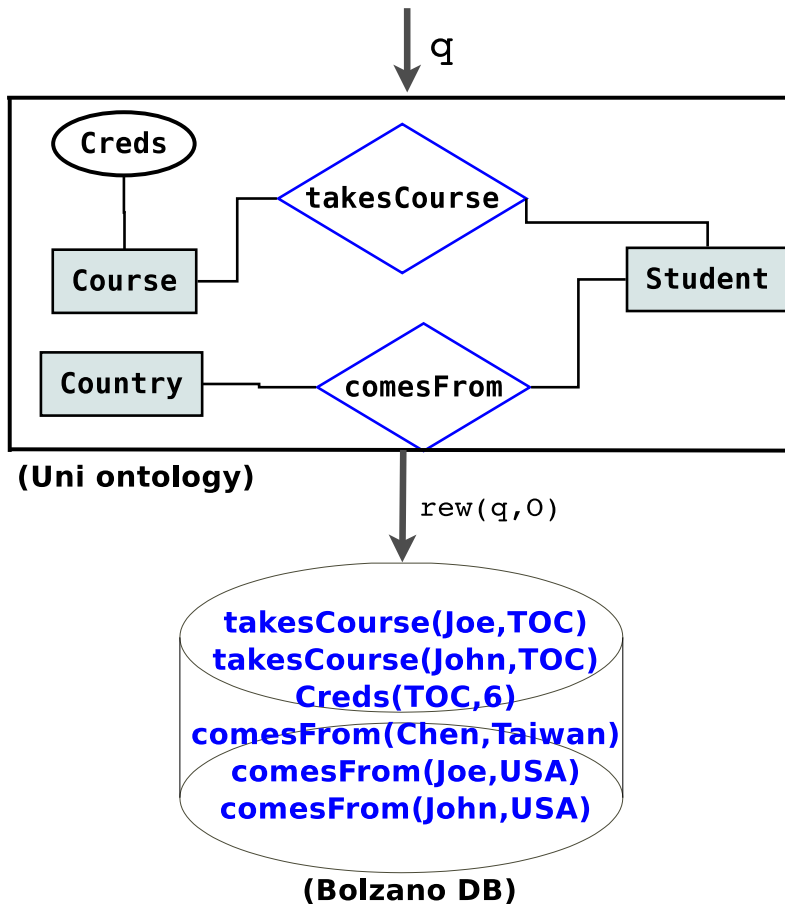
Which is the number of students, per country, who study something?



$$q(x, \mathit{count}(z)) \leftarrow \mathit{Student}(x) \wedge \exists y(\mathit{takesCourse}(x, y)) \wedge \mathit{comesFrom}(x, z) \wedge \mathit{Country}(z))$$

Example

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



$$q(x, count(z)) \leftarrow Student(x) \wedge \exists y(takesCourse(x, y)) \wedge comesFrom(x, z) \wedge Country(z)$$

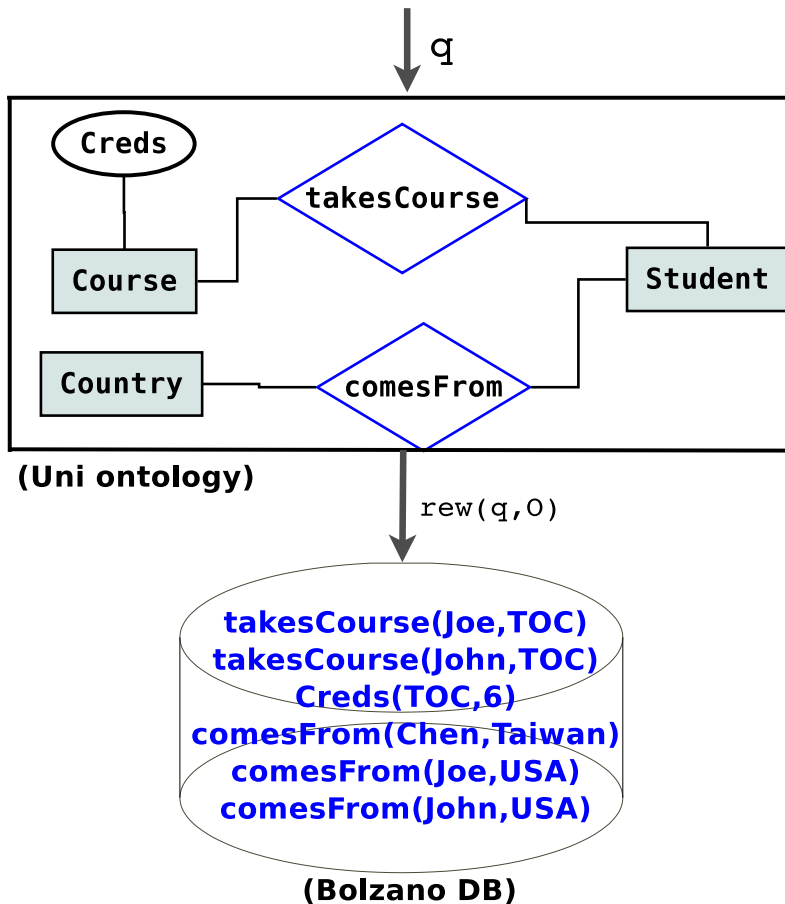
After evaluating and grouping

$$H_{USA} = \{John, Joe\}$$

$$H_{Taiwan} = \{\}$$

Example

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



$$q(x, count(z)) \leftarrow Student(x) \wedge \exists y (takesCourse(x, y) \wedge comesFrom(x, z) \wedge Country(z))$$

Therefore

$$cert(q, \mathcal{O}_{Uni}, \mathcal{D}_{Bz}) = \{(USA, 2), (Taiwan, 0)\}$$

ATQ-English and Formal Semantics

- ⑥ To express ATQs we use computational semantics compositional translations $\tau(\cdot)$

ATQ-English and Formal Semantics

- ⑥ To express ATQs we use computational semantics compositional translations $\tau(\cdot)$
- ⑥ Given $\{e, t, \mathbb{N}, \mathbb{Q}\}$ as basic types we define aggregate determiners like

$$\tau(\text{the average number of}) := \lambda P^{e \rightarrow \mathbb{N}}. \text{avg}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q}$$

$$\tau(\text{the number of}) := \lambda P^{e \rightarrow \mathbb{N}}. \text{count}(P) : (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q}$$

- ⑥ We make use of semantically annotated grammars, i.e., context-free grammars with rules like

$$Q_{wh} \rightarrow I_i S_{g_i} ? \quad \tau(Q_{wh}) := \lambda \bar{z}. \tau(I_i)(\lambda i. \tau(S_{g_i}))$$

$$S_{g_i} \rightarrow NP_{g_i} VP_i \quad \tau(S_{g_i}) := \tau(NP_{g_i})(\tau(VP_i))$$

ATQ-English: Sample Lexicon

$\tau(\text{the number of}) := \lambda P^{e \rightarrow \mathbb{N}}. \text{count}(P): (e \rightarrow \mathbb{N}) \rightarrow \mathbb{Q}$

$\tau(\text{who}) := \lambda P^{e \rightarrow \mathbb{N}} \lambda Q^{e \rightarrow \mathbb{N}}. \lambda x^e. (P(x) \wedge Q(x)): (e \rightarrow \mathbb{N}) \rightarrow ((e \rightarrow \mathbb{N}) \rightarrow (e \rightarrow \mathbb{N}))$

$\tau(\text{something}) := \lambda P^{e \rightarrow t}. \exists y^e P(y): (e \rightarrow \mathbb{N}) \rightarrow t$

$\tau(\text{per country}) := \lambda P^{e \rightarrow \mathbb{N}}. \lambda y^e. (P(y) \wedge \text{Country}(z) \wedge \text{comesFrom}(y, z)): (e \rightarrow \mathbb{N}) \rightarrow (e \rightarrow \mathbb{N})$

$\tau(\text{per country}) := \lambda Q^{e \rightarrow \mathbb{N}}. \lambda P^{e \rightarrow \mathbb{N}} \lambda y^e. (P(y) \wedge \text{Country}(z) \wedge \text{comesFrom}(y, z) \wedge Q(z)):$
 $(e \rightarrow \mathbb{N}) \rightarrow ((e \rightarrow \mathbb{N}) \rightarrow (e \rightarrow \mathbb{N}))$

$\tau(\text{which}) := \lambda R^{e \rightarrow (\mathbb{Q} \rightarrow t)}. \lambda n^{\mathbb{Q}}. R(n, m): (e \rightarrow (\mathbb{Q} \rightarrow t)) \rightarrow (\mathbb{Q} \rightarrow t)$

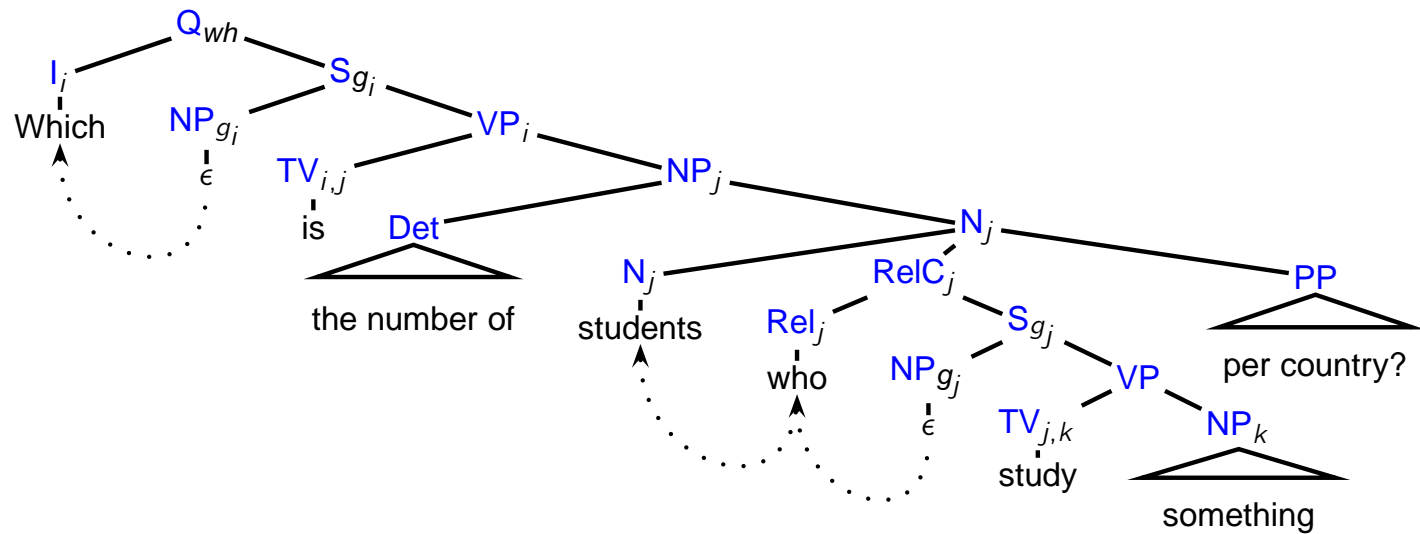
$\tau(\text{is}) := \lambda n^{\mathbb{Q}}. \lambda m^{\mathbb{Q}}. n \approx m: \mathbb{Q} \rightarrow (\mathbb{Q} \rightarrow t)$

$\tau(\epsilon) := \lambda P^{\mathbb{Q} \rightarrow t}. P(n): (\mathbb{Q} \rightarrow t) \rightarrow t$

$\tau(\text{students}) := \lambda x^e. \text{Student}(x): e \rightarrow \mathbb{N}$

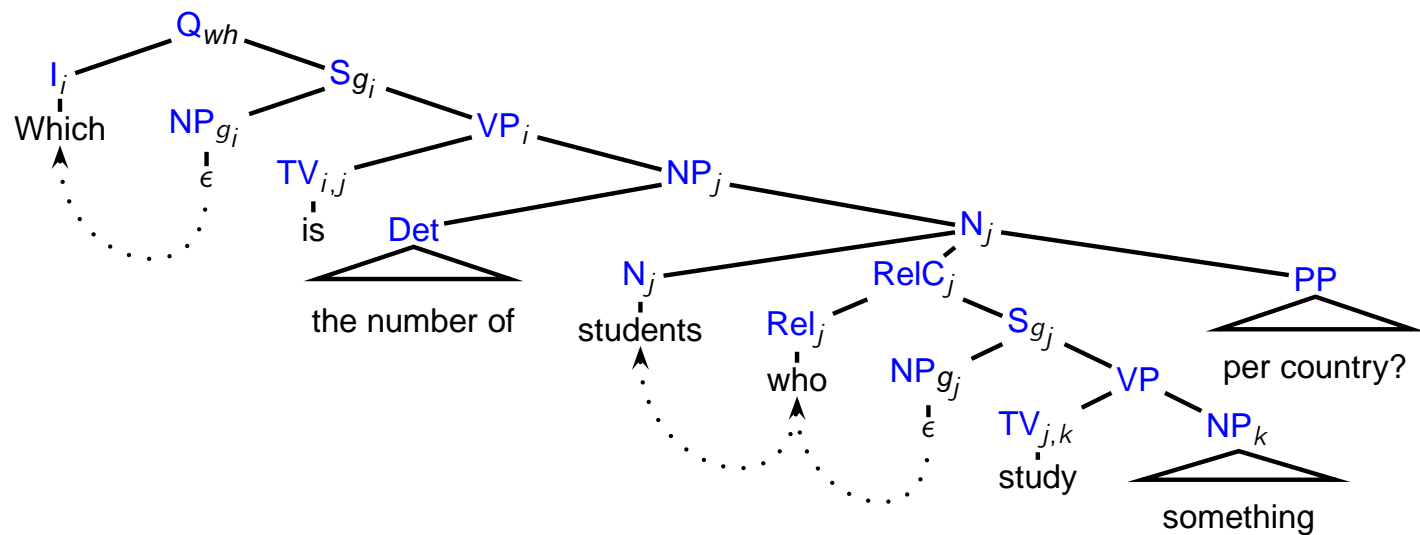
$\tau(\text{study}) := \lambda \alpha^{(e \rightarrow \mathbb{N}) \rightarrow \mathbb{N}}. \lambda x^e. \alpha(\lambda y^e. \text{takesCourse}(x, y)): ((e \rightarrow \mathbb{N}) \rightarrow \mathbb{N}) \rightarrow (e \rightarrow \mathbb{N})$

ATQ-English: Expressing Aggregations and Grouping



$$\lambda z^e . \lambda m^{\mathbb{Q}} . m \approx \mathit{count}(\lambda x^e . \mathit{Student}(x) \wedge \exists y^e (\mathit{takesCourse}(x, y)) \wedge \mathit{comesFrom}(x, z) \wedge \mathit{Country}(z)) : e \rightarrow (\mathbb{Q} \rightarrow t)$$

ATQ-English: Expressing Aggregations and Grouping



$$q(x, \textit{count}(z)) \leftarrow \textit{Student}(x) \wedge \exists y(\textit{takesCourse}(x, y)) \\ \wedge \textit{comesFrom}(x, z) \wedge \textit{Country}(z))$$

Data Complexity of ATQ-English

- ⑥ Which is the number of students, per country, who study something?

```
SELECT s1.CName, COUNT(s.SName)
FROM comesFrom s1
WHERE EXISTS (SELECT *
              FROM takesCourse s2
              WHERE s1.SName = s2.SName)
GROUP BY s1.CName
```

Data Complexity of ATQ-English

- ⑥ Which is the number of students, per country, who study something?

```
SELECT s1.CName, COUNT(s.SName)
FROM comesFrom s1
WHERE EXISTS (SELECT *
              FROM takesCourse s2
              WHERE s1.SName = s2.SName)
GROUP BY s1.CName
```

- ⑥ Answering an ATQ is in **LOGSPACE** in data complexity when combined with ontologies from the *DL-Lite* family

Beyond ATQs and TSQs

- ⑥ Disjunctions $\psi(x) \vee \psi'(x)$ of TSCs do not increase data complexity either

Beyond ATQs and TSQs

- ⑥ Disjunctions $\psi(x) \vee \psi'(x)$ of TSCs do not increase data complexity either
- ⑥ But negation does:

Which is the number of students who do **not** come from Italy?

```
SELECT COUNT(s1.SName)
FROM comesFrom s1
WHERE s1.CName != "Italy"
```

Beyond ATQs and TSQs

- ⑥ Disjunctions $\psi(x) \vee \psi'(x)$ of TSCs do not increase data complexity either
- ⑥ Complexity jumps to **coNP**-hard with
 - negations of TSCs, $\neg\psi(x)$
 - universally restricted TSCs, $\forall x(R(x, y) \rightarrow \psi(y))$
 - comparisons, $t \theta t'$, $\theta \in \{=, >, <, \geq, \leq\}$

Conclusions

- ⑥ We have proposed to measure the performance of CLIs through the data complexity of data access
- ⑥ We have expressed in CL ATQs
- ⑥ We have shown that
 - aggregations, disjunctions, conjunctions and existential quantifications are **cheap** to use in CLIs to OBDASs
 - extending coverage to negation, comparisons and universal restrictions is **costly**
- ⑥ We can capture in CL and answer efficiently a significant class of SQL aggregate queries over OBDASs

Data Complexity

	\exists -ATQs/TSQs	\leq -ATQs/TSQs	\forall -ATQs/TSQs	\neg -ATQs/TSQs
<i>DL-Lite</i>	in LOGSPACE	coNP -complete	coNP -complete	coNP -complete
<i>ALCI</i>	coNP -complete	coNP -hard	coNP -complete	coNP -complete
OWL	coNP -hard	coNP -hard	coNP -hard	coNP -hard

- ⑥ We derive a **coNP** upper bound when \mathcal{O} and q are in the 2-variable fragment of FOL [Pratt 2009]

Data Complexity

	\exists -ATQs/TSQs	\leq -ATQs/TSQs	\forall -ATQs/TSQs	\neg -ATQs/TSQs
<i>DL-Lite</i>	in LOGSPACE	coNP -complete	coNP -complete	coNP -complete
<i>ALCI</i>	coNP -complete	coNP -hard	coNP -complete	coNP -complete
OWL	coNP -hard	coNP -hard	coNP -hard	coNP -hard

- ⑥ We derive a **coNP** upper bound when \mathcal{O} and q are in the 2-variable fragment of FOL [Pratt 2009]
- ⑥ But evidence from question corpora suggests that such **coNP**-hard constructs are infrequent [Bernardi et al. 2007]

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

Encode Φ into DB \mathcal{D}_Φ

$$\begin{array}{c} \vdots \\ P(c_i, l_{i,1}), P(c_i, l_{i,2}), N(c_i, n_{i,1}), N(c_i, n_{i,2}) \\ \vdots \\ A_t(\text{true}) \end{array}$$

Consider \mathcal{O}_Φ

$$\forall x (A_f(x) \rightarrow \neg A_t(x))$$

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

Ask TCQ q_Φ of body ψ_Φ

$$\begin{aligned} & \exists c \exists l_1 \exists l_2 \exists l_3 \exists l_4 (\\ & P_1(c, l_1) \wedge \forall v_1 (\text{Val}(l_1, v_1) \rightarrow A_f(v_1)) \wedge \\ & P_2(c, l_2) \wedge \forall v_2 (\text{Val}(l_2, z_2) \rightarrow A_f(v_2)) \wedge \\ & N_1(c, l_3) \wedge \exists v_3 (\text{Val}(l_3, v_3) \wedge A_t(v_3)) \wedge \\ & N_2(c, l_4) \wedge \exists v_4 (\text{Val}(l_4, v_4) \wedge A_t(v_4))) \end{aligned}$$

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

We claim

$(\) \notin \text{cert}(q_\Phi, \mathcal{O}_\Phi, \mathcal{D}_\Phi)$ iff Φ has no model

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

" \Rightarrow " given a \mathcal{D} s.t.

- $\mathcal{D} \models \mathcal{O}_\Phi$
- $\mathcal{D}_\Phi \subseteq \mathcal{D}$
- $\mathcal{D} \not\models \psi_\Phi$

Then

$$\begin{aligned} \forall v (\mathbf{Val}(l, v) \wedge A_f(v)) &\Rightarrow_{\mathcal{D}} \forall v (\mathbf{Val}(l, v) \wedge \neg A_f(v)) \\ &\equiv_{\mathcal{D}} \neg \exists v_i (\mathbf{Val}(l, v) \wedge A_t(v)) \end{aligned}$$

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

Define $\delta(\cdot)$ by putting

$$\delta(I) = \text{true iff } (I, \text{true}) \in \text{Val}^{\mathcal{D}}$$

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

" \Leftarrow ": Symmetrical argument

coNP-hardness of \forall -TSQs: Reduction from 2+2-SAT

A conjunction of 2+2 clauses is a conjunction $\Phi := C_1 \wedge \cdots \wedge C_k$

$$C_i := p_{i1} \vee p_{i2} \vee \neg n_{i1} \vee \neg n_{i2}$$

SAT for 2+2 clauses in **NP**-complete [Scharf, 1994]

Therefore:

$() \notin \text{cert}(q_\Phi, \mathcal{O}_\Phi, \mathcal{D}_\Phi)$ iff Φ has no model

Moral: conjunction, universal and existential quantification in queries and disjointness in ontologies blows up data complexity!