

Logical Approach to Physical Data Independence and Query Compilation

Classical OBDA and Data Exchange

David Toman

D.R. Cheriton School of Computer Science

University of

Waterloo



OBDA AND LITE LOGICS

Setup

Setting

- Input:
- (1) Schema Σ (set of integrity constraints);
 - (2) Data $D = \{R_1, \dots, R_k\}$ (instance of **access paths**); and
 - (3) Query φ (a formula)

Setup

Setting

- Input:
- (1) Schema Σ (set of integrity constraints);
 - (2) Data $D = \{R_1, \dots, R_k\}$ (instance of **access paths**); and
 - (3) Query φ (a formula)

Definition (Certain Answers)

$$\text{cert}_{\Sigma, D}(\varphi) = \{\vec{a} \mid \Sigma \cup D \models \varphi(\vec{a})\} = \bigcap_{I \models \Sigma \cup D} \{\vec{a} \mid I \models \varphi(\vec{a})\}$$

Setup

Setting

- Input: (1) Schema Σ (set of integrity constraints);
(2) Data $D = \{R_1, \dots, R_k\}$ (instance of **access paths**); and
(3) Query φ (a formula)

Definition (Certain Answers)

$$\text{cert}_{\Sigma, D}(\varphi) = \{\vec{a} \mid \Sigma \cup D \models \varphi(\vec{a})\} = \bigcap_{I \models \Sigma \cup D} \{\vec{a} \mid I \models \varphi(\vec{a})\}$$

Convention: ABox \mathcal{A} vs. database $D_{\mathcal{A}}$

We assume that for every access path $R_{\text{AP}}(\vec{x})$ in $D_{\mathcal{A}}$ there is

- a logical predicate $R(\vec{x})$ (with the same arity), and
- a constraint $\forall \vec{x}. R_{\text{AP}}(\vec{x}) \rightarrow R(\vec{x})$.

Can this be Done Efficiently at all?

Question

Can there be a *non-trivial* schema language for which *query answering* (under certain answer semantics) is *tractable* (in data complexity)?

Can this be Done Efficiently at all?

Question

Can there be a *non-trivial* schema language for which *query answering* (under certain answer semantics) is *tractable* (in data complexity)?

YES: *Conjunctive queries* (or positive) and “*lite*” *Description Logics*:

- 1 The DL-Lite family
 - ⇒ conjunction, \perp , domain/range, unqualified \exists , role inverse, UNA
 - ⇒ certain answers in AC_0 for data complexity (maps to SQL)
- 2 The \mathcal{EL} family
 - ⇒ conjunction, qualified \exists
 - ⇒ certain answers *PTIME-complete* for data complexity
- 3 The \mathcal{CFD} family
 - ⇒ qualified \forall (over total functions), functional dependencies
 - ⇒ certain answers *PTIME-complete* for data complexity

DL-Lite Family of DLs

Definition (DL-Lite family: Schemata/TBoxes)

- 1 Roles R and concepts C as follows:

$$R ::= P \mid P^- \quad C ::= \perp \mid A \mid \exists R$$

- 2 Schemata are represented as TBoxes: a finite set \mathcal{T} of *constraints*

$$C_1 \sqcap \dots \sqcap C_n \sqsubseteq C \quad R_1 \sqsubseteq R_2$$

Access paths (data) \Rightarrow ABox \mathcal{A} (recall the “convention” about access paths!)

DL-Lite Family of DLs

Definition (DL-Lite family: Schemata/TBoxes)

- 1 Roles R and concepts C as follows:

$$R ::= P \mid P^- \quad C ::= \perp \mid A \mid \exists R$$

- 2 Schemata are represented as TBoxes: a finite set \mathcal{T} of constraints

$$C_1 \sqcap \dots \sqcap C_n \sqsubseteq C \quad R_1 \sqsubseteq R_2$$

Access paths (data) \Rightarrow ABox \mathcal{A} (recall the “convention” about access paths!)

How to compute answers to CQs?

IDEA: incorporate *schematic knowledge* into the query.

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^{-} \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^{-} \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^{\dagger} = (\exists y. Works(x, y) \wedge Project(y)) \vee$$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^- \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^\dagger = (\exists y. Works(x, y) \wedge Project(y)) \vee (\exists y, z. Works(x, y) \wedge Works(z, y)) \vee$$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^- \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^\dagger = (\exists y. Works(x, y) \wedge Project(y)) \vee$$
$$(\exists y, z. Works(x, y) \wedge Works(z, y)) \vee$$
$$(\exists y. Works(x, y)) \vee$$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^{-} \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^{\dagger} = (\exists y. Works(x, y) \wedge Project(y)) \vee$$
$$(\exists y, z. Works(x, y) \wedge Works(z, y)) \vee$$
$$(\exists y. Works(x, y)) \vee$$
$$(Employee(x))$$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^- \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^\dagger = (\exists y. Works_{AP}(x, y) \wedge Project_{AP}(y)) \vee \\ (\exists y, z. Works_{AP}(x, y) \wedge Works_{AP}(z, y)) \vee \\ (\exists y. Works_{AP}(x, y)) \vee \\ (Employee_{AP}(x))$$

Query Execution:

$$Q^\dagger \left(\left\{ \begin{array}{l} Employee(bob), \\ Works(sue, slides) \end{array} \right\} \right)$$

Example

TBox (Schema): $Employee \sqsubseteq \exists Works$
 $\exists Works^- \sqsubseteq Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Rewriting:

$$Q^\dagger = (\exists y. Works_{AP}(x, y) \wedge Project_{AP}(y)) \vee \\ (\exists y, z. Works_{AP}(x, y) \wedge Works_{AP}(z, y)) \vee \\ (\exists y. Works_{AP}(x, y)) \vee \\ (Employee_{AP}(x))$$

Query Execution:

$$Q^\dagger \left(\left\{ \begin{array}{l} Employee(bob), \\ Works(sue, slides) \end{array} \right\} \right) = \{bob, sue\}$$

QuOnto: Rewriting Approach [Calvanese et al.]

Input: Conjunctive query Q , DL-Lite TBox Σ

$R = \{Q\};$

repeat

foreach query $Q' \in R$ **do**

foreach axiom $\alpha \in \Sigma$ **do**

if α is applicable to Q' **then**

$R = R \cup \{Q'[\text{lhs}(\alpha)/\text{rhs}(\alpha)]\}$

foreach two atoms D_1, D_2 in Q' **do**

if D_1 and D_2 unify **then**

$\sigma = \text{MGU}(D_1, D_2); R = R \cup \{\lambda(Q', \sigma)\};$

until no query unique up to variable renaming can be added to R ;

return $Q^\dagger := (\bigvee R)$

QuOnto: Rewriting Approach [Calvanese et al.]

Input: Conjunctive query Q , DL-Lite TBox Σ

$R = \{Q\};$

repeat

foreach query $Q' \in R$ **do**

foreach axiom $\alpha \in \Sigma$ **do**

if α is applicable to Q' **then**

$R = R \cup \{Q'[\text{lhs}(\alpha)/\text{rhs}(\alpha)]\}$

foreach two atoms D_1, D_2 in Q' **do**

if D_1 and D_2 unify **then**

$\sigma = \text{MGU}(D_1, D_2); R = R \cup \{\lambda(Q', \sigma)\};$

until no query unique up to variable renaming can be added to R ;

return $Q^\dagger := (\bigvee R)$

Theorem

$\Sigma \cup \mathcal{A} \models Q(\vec{a})$ if and only if $D_{\mathcal{A}} \models Q^\dagger(\vec{a})$

QuOnto: Rewriting Approach [Calvanese et al.]

Input: Conjunctive query Q , DL-Lite TBox Σ

$R = \{Q\};$

repeat

foreach query $Q' \in R$ **do**

foreach axiom $\alpha \in \Sigma$ **do**

if α is applicable to Q' **then**

$R = R \cup \{Q'[\text{lhs}(\alpha)/\text{rhs}(\alpha)]\}$

foreach two atoms D_1, D_2 in Q' **do**

if D_1 and D_2 unify **then**

$\sigma = \text{MGU}(D_1, D_2); R = R \cup \{\lambda(Q', \sigma)\};$

until no query unique up to variable renaming can be added to R ;

return $Q^\dagger := (\bigvee R)$

Theorem

$\Sigma \cup \mathcal{A} \models Q(\vec{a})$ if and only if $D_{\mathcal{A}} \models Q^\dagger(\vec{a})$ \leftarrow can be VERY large

Definition (\mathcal{EL} -Lite family: Schemata and TBoxes)

- 1 Concepts C as follows:

$$C ::= A \mid \top \mid \perp \mid C \sqcap C \mid \exists R.C$$

- 2 Schemata are represented as TBoxes: a finite set \mathcal{T} of *constraints*

$$C_1 \sqsubseteq C_2 \qquad R_1 \sqsubseteq R_2$$

Access paths (data) \Rightarrow ABox \mathcal{A} (recall the “convention” about access paths!)

Definition (\mathcal{EL} -Lite family: Schemata and TBoxes)

- 1 Concepts C as follows:

$$C ::= A \mid \top \mid \perp \mid C \sqcap C \mid \exists R.C$$

- 2 Schemata are represented as TBoxes: a finite set \mathcal{T} of *constraints*

$$C_1 \sqsubseteq C_2 \qquad R_1 \sqsubseteq R_2$$

Access paths (data) \Rightarrow ABox \mathcal{A} (recall the “convention” about access paths!)

How to compute answers to CQs?

IDEA: incorporate *schematic knowledge* into the data.

Combined Approach

Can an approach based on *rewriting* be used for \mathcal{EL} ?

Combined Approach

Can an approach based on *rewriting* be used for \mathcal{EL} ?

NO: \mathcal{EL} is PTIME-complete (data complexity).

Combined Approach

Can an approach based on *rewriting* be used for \mathcal{EL} ?

NO: \mathcal{EL} is PTIME-complete (data complexity).

Combined Approach

We effectively transform

- 1 the ABox (access paths) \mathcal{A} to a *canonical structure* $D_{\mathcal{A}}^*$ utilizing Σ ,
- 2 the conjunctive query Q to a relational query Q^\ddagger .

... both *polynomial* in the input(s).

Combined Approach

Can an approach based on *rewriting* be used for \mathcal{EL} ?

NO: \mathcal{EL} is PTIME-complete (data complexity).

Combined Approach

We effectively transform

- 1 the ABox (access paths) \mathcal{A} to a *canonical structure* $D_{\mathcal{A}}^*$ utilizing Σ ,
- 2 the conjunctive query Q to a relational query Q^\ddagger .

... both *polynomial* in the input(s).

Theorem (Lutz, __, Wolter: IJCAI'09)

$\Sigma \cup \mathcal{A} \models Q(\vec{a})$ if and only if $D_{\mathcal{A}}^* \models Q^\ddagger(\vec{a})$

Example (with almost DL-Lite schema)

TBox (Schema): $Employee \sqsubseteq \exists Works.Project$
 $\exists Works.\top \sqsubseteq \exists Works.Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Data: $\{Employee(bob), Works(sue, slides)\}$

Example (with almost DL-Lite schema)

TBox (Schema): $Employee \sqsubseteq \exists Works.Project$
 $\exists Works.\top \sqsubseteq \exists Works.Project$

Conjunctive Query: $\exists y. Works(x, y) \wedge Project(y)$

Data: $\{Employee(bob), Works(sue, slides)\}$

Rewriting:

- 1 $D_{\mathcal{A}}^* = \{ Employee(bob), Works(bob, cWorks), Works(sue, slides), Works(sue, cWorks), Project(cWorks), \}$
- 2 $Q^\ddagger = Q \wedge (x \neq cWorks)$

Example (with almost DL-Lite schema)

TBox (Schema): $Employee \sqsubseteq \exists Works.Project$
 $\exists Works.\top \sqsubseteq \exists Works.Project$

Conjunctive Query: $\exists y.Works(x, y) \wedge Project(y)$

Data: $\{Employee(bob), Works(sue, slides)\}$

Rewriting:

- 1 $D_{\mathcal{A}}^* = \{ Employee(bob), Works(bob, cWorks), Works(sue, slides), Works(sue, cWorks), Project(cWorks), \}$
- 2 $Q^\ddagger = Q \wedge (x \neq cWorks)$

Query Execution:

$$Q^\ddagger(D_{\mathcal{A}}^*) = \{bob, sue\}$$

A Combined Approach and DL-Lite

Can the *exponential size* of rewriting be avoided for DL-Lite?

A Combined Approach and DL-Lite

Can the *exponential size* of rewriting be avoided for DL-Lite?

Yes: using the Combined Approach

...but query rewriting is much more involved due to *inverse roles*;

A Combined Approach and DL-Lite

Can the *exponential size* of rewriting be avoided for DL-Lite?

Yes: using the Combined Approach

...but query rewriting is much more involved due to *inverse roles*;

Theorem (Konchatov, Lutz, __, Wolter, KR10)

$$\Sigma \cup A \models Q(\vec{a}) \text{ if and only if } D_{\mathcal{A}}^* \models Q^{\ddagger}(\vec{a})$$

(... still exponential for *role hierarchies*.)

A Combined Approach and DL-Lite

Can the *exponential size* of rewriting be avoided for DL-Lite?

Yes: using the Combined Approach

... but query rewriting is much more involved due to *inverse roles*;

Theorem (Konchatov, Lutz, __, Wolter, KR10)

$$\Sigma \cup A \models Q(\vec{a}) \text{ if and only if } D_{\mathcal{A}}^* \models Q^{\ddagger}(\vec{a})$$

(... still exponential for *role hierarchies*.)

Theorem (Lutz, Seylan, __, Wolter, ISWC13)

$$\Sigma \cup A \models Q(\vec{a}) \text{ if and only if } D_{\mathcal{A}}^* \models Q^{\text{filter}}(\vec{a})$$

(... polynomial in $|\mathcal{H}|$, but uses UDF feature of DB2.)

Definition (CFD_{nc} : Schemata and TBoxes)

- 1 Syntax formed from *path functions* Pf and *concepts* C, D as follows:

$$C ::= A \mid \forall Pf.C$$

$$D ::= A \mid \neg C \mid \forall Pf.C \mid C : Pf_1, \dots, Pf_k \rightarrow Pf$$

- 2 Schemata are represented as a TBox:

finite set \mathcal{T} of *constraints* $C \sqsubseteq D$.

- 3 Data is represented as an ABox (recall again the AP “convention”):

finite set \mathcal{A} of *concept* ($A(a)$) and *equational* ($Pf(a) = Pf'(b)$) assertions.

Definition (CFD_{nc} : Schemata and TBoxes)

- 1 Syntax formed from *path functions* Pf and *concepts* C, D as follows:

$$C ::= A \mid \forall Pf.C$$

$$D ::= A \mid \neg C \mid \forall Pf.C \mid C : Pf_1, \dots, Pf_k \rightarrow Pf$$

- 2 Schemata are represented as a TBox:

finite set \mathcal{T} of *constraints* $C \sqsubseteq D$.

- 3 Data is represented as an ABox (recall again the AP “convention”):

finite set \mathcal{A} of *concept* ($A(a)$) and *equational* ($Pf(a) = Pf'(b)$) assertions.

Definition (CFD_{nc} : Schemata and TBoxes)

- 1 Syntax formed from *path functions* Pf and *concepts* C, D as follows:

$$C ::= A \mid \forall \text{Pf}. C$$

$$D ::= A \mid \neg C \mid \forall \text{Pf}. C \mid C : \text{Pf}_1, \dots, \text{Pf}_k \rightarrow \text{Pf}$$

- 2 Schemata are represented as a TBox:

finite set \mathcal{T} of *constraints* $C \sqsubseteq D$.

- 3 Data is represented as an ABox (recall again the AP “convention”):

finite set \mathcal{A} of *concept* ($A(a)$) and *equational* ($\text{Pf}(a) = \text{Pf}'(b)$) assertions.

Rewriting Approach: can't work—reachability in ABox (PTIME-c)

Combined Approach: can't work—too many *types* (anon. completion too big)

Definition (CFD_{nc} : Schemata and TBoxes)

- 1 Syntax formed from *path functions* Pf and *concepts* C, D as follows:

$$C ::= A \mid \forall Pf.C$$

$$D ::= A \mid \neg C \mid \forall Pf.C \mid C : Pf_1, \dots, Pf_k \rightarrow Pf$$

- 2 Schemata are represented as a TBox:

finite set \mathcal{T} of *constraints* $C \sqsubseteq D$.

- 3 Data is represented as an ABox (recall again the AP “convention”):

finite set \mathcal{A} of *concept* ($A(a)$) and *equational* ($Pf(a) = Pf'(b)$) assertions.

Query Answering: The Perfect Combined Approach

IDEA: incorporate

- reachability induced by *schematic knowledge* into the data, and
- types induced by *schematic knowledge* into the query.

DATA EXCHANGE

Setup

Schema Mapping

- source schema (signature) S_P and (closed) data;
- target schema (signature) S_L ;
- mapping constraints: *s-t TGDs*—formulas of the form

$\forall \vec{x}.\varphi(\vec{x}) \rightarrow \exists \vec{y}.\psi(\vec{x}, \vec{y})$ where φ is a CQ over S_P and ψ a CQ over S_L .

The general setting of data exchange is this:



[Arenas et al: Foundations of Data Exchange]

Setup

Schema Mapping

- source schema (signature) S_P and (closed) data;
- target schema (signature) S_L ;
- mapping constraints: *s-t TGDs*—formulas of the form

$\forall \vec{x}. \varphi(\vec{x}) \rightarrow \exists \vec{y}. \psi(\vec{x}, \vec{y})$ where φ is a CQ over S_P and ψ a CQ over S_L .

The general setting of data exchange is this:



[Arenas et al: Foundations of Data Exchange]

Definition

J (over S_L) is a *solution* for I (over S_P) w.r.t. Σ if $(I, J) \models \Sigma$.

... too many solutions (TGDs imply open world @ S_L !)

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions
⇒ variables (marked nulls): *representation system* [Imielinski&Lipski'84]

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions
 - ⇒ variables (marked nulls): *representation system* [Imielinski&Lipski'84]
 - ⇒ can be used to answer CQ/UCQ (how and why?)

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions
 - ⇒ variables (marked nulls): *representation system* [Imielinski&Lipski'84]
 - ⇒ can be used to answer CQ/UCQ (how and why?)
- a smallest universal solution—the *core*.

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions
 - ⇒ variables (marked nulls): *representation system* [Imielinski&Lipski'84]
 - ⇒ can be used to answer CQ/UCQ (how and why?)
- a smallest universal solution—the *core*.
- core can be constructed using the *chase* (in PTIME);

Universal Solutions and Cores

Problem(s):

Multiple *solutions* (target instances) for single *closed world* source

⇒ how to answer queries over target? *certain answers* w.r.t. all solutions.

IDEA:

Find the *best* solution: one that can be used instead of every other solution.

- an *universal solution*: homomorphism to all other solutions
 - ⇒ variables (marked nulls): *representation system* [Imielinski&Lipski'84]
 - ⇒ can be used to answer CQ/UCQ (how and why?)
- a smallest universal solution—the *core*.

- core can be constructed using the *chase* (in PTIME);
- what happens if we have additional constraints on the target (S_L)?

LIMITS AND ISSUES WITH POSSIBLE WORLDS

Certain Answers: What is the Price?

High Computational Cost even for mild deviation from *Lite* Logics (and CQ)

coNP-hard for *DATA COMPLEXITY*

Example

- Schema&Data:

$$\Sigma = \left\{ \begin{array}{l} \forall x, y. ColNode(x, y) \leftrightarrow Node(x), \\ \forall x, y. ColNode(x, y) \leftrightarrow Colour(y) \end{array} \right\}$$

$$D = \left\{ \begin{array}{l} Edge = \{(n_i, n_j)\}, Node = \{n_1, \dots, n_m\}, \\ Colour = \{r, g, b\} \end{array} \right\}$$

Certain Answers: What is the Price?

High Computational Cost even for mild deviation from *Lite* Logics (and CQ)

coNP-hard for **DATA COMPLEXITY**

Example

- Schema&Data:

$$\Sigma = \left\{ \begin{array}{l} \forall x, y. \text{ColNode}(x, y) \leftrightarrow \text{Node}(x), \\ \forall x, y. \text{ColNode}(x, y) \leftrightarrow \text{Colour}(y) \end{array} \right\}$$

$$D = \left\{ \begin{array}{l} \text{Edge} = \{(n_i, n_j)\}, \text{Node} = \{n_1, \dots, n_m\}, \\ \text{Colour} = \{r, g, b\} \end{array} \right\}$$

- Query: $\exists x, y, c. \text{Edge}(x, y) \wedge \text{ColNode}(x, c) \wedge \text{ColNode}(y, c)$

Certain Answers: What is the Price?

High Computational Cost even for mild deviation from *Lite* Logics (and CQ)

coNP-hard for *DATA COMPLEXITY*

Example

- Schema&Data:

$$\Sigma = \left\{ \begin{array}{l} \forall x, y. ColNode(x, y) \leftrightarrow Node(x), \\ \forall x, y. ColNode(x, y) \leftrightarrow Colour(y) \end{array} \right\}$$

$$D = \left\{ \begin{array}{l} Edge = \{(n_i, n_j)\}, Node = \{n_1, \dots, n_m\}, \\ Colour = \{r, g, b\} \end{array} \right\}$$

- Query: $\exists x, y, c. Edge(x, y) \wedge ColNode(x, c) \wedge ColNode(y, c)$
 \Rightarrow the graph $(Node, Edge)$ is NOT 3-colourable.

Certain Answers: What is the Price?

High Computational Cost even for mild deviation from *Lite* Logics (and CQ)

coNP-hard for *DATA COMPLEXITY*

Example

- Schema&Data:

$$\Sigma = \left\{ \begin{array}{l} \forall x, y. ColNode(x, y) \leftrightarrow Node(x), \\ \forall x, y. ColNode(x, y) \leftrightarrow Colour(y) \end{array} \right\}$$

$$D = \left\{ \begin{array}{l} Edge = \{(n_i, n_j)\}, Node = \{n_1, \dots, n_m\}, \\ Colour = \{r, g, b\} \end{array} \right\}$$

- Query: $\exists x, y, c. Edge(x, y) \wedge ColNode(x, c) \wedge ColNode(y, c)$
 \Rightarrow the graph $(Node, Edge)$ is NOT 3-colourable.

... coNP-complete for all DLs between \mathcal{AL} and \mathcal{SHIQ} .

Certain Answers: What is the Price?

High Computational Cost even for mild deviation from *Lite* Logics (and CQ)

coNP-hard for *DATA COMPLEXITY*

Example

- Schema&Data:

$$\Sigma = \left\{ \begin{array}{l} \forall x, y. ColNode(x, y) \leftrightarrow Node(x), \\ \forall x, y. ColNode(x, y) \leftrightarrow Colour(y) \end{array} \right\}$$

$$D = \left\{ \begin{array}{l} Edge = \{(n_i, n_j)\}, Node = \{n_1, \dots, n_m\}, \\ Colour = \{r, g, b\} \end{array} \right\}$$

- Query: $\exists x, y, c. Edge(x, y) \wedge ColNode(x, c) \wedge ColNode(y, c)$
 \Rightarrow the graph $(Node, Edge)$ is NOT 3-colourable.

... coNP-complete for all DLs between \mathcal{AL} and \mathcal{SHIQ} .

OBDA-Lite can only say $Colour \supseteq \{r, g, b\}$ (due to OWA)

Data Exchange cannot say $\forall x, y. ColNode(x, y) \rightarrow Colour(y)$ (not an s-t TGD)

Certain Answers: What about more complex Queries?

(safe) Negation, Inequality

Theorem (Gutiérrez-Basulto et al., RR13)

OBDA for CQ with single inequality or with safe negated atoms over DL-Lite^ℋ is undecidable.

Aggregation

- ⇒ *count/sum* aggregate functions do not play nicely with *certain answers*
- epistemic operators (count the number of *known* answers)
[Calvanese et al., ONISW08]
 - range/lower bounds semantics (at least so many)
[Kostylev and Reutter, AAAI13]
- ... and it is (data complexity-wise) hard in all cases.

Certain Answers??

Example (Unintuitive Behaviour of Queries:)

1 $\exists x. \textit{Phone}(\text{"John"}, x)?$

2 $\textit{Phone}(\text{"John"}, x)?$

under $\Sigma = \{\forall x. \textit{Person}(x) \rightarrow \exists y. \textit{Phone}(x, y)\}$
and $D = \{\textit{Person}(\text{"John"})\}$.

Certain Answers??

Example (Unintuitive Behaviour of Queries:)

1 $\exists x. \text{Phone}(\text{"John"}, x)?$

2 $\text{Phone}(\text{"John"}, x)?$

under $\Sigma = \{\forall x. \text{Person}(x) \rightarrow \exists y. \text{Phone}(x, y)\}$
and $D = \{\text{Person}(\text{"John"})\}$.

Embedded SQL-like Example

```
if " $\exists x. \text{Phone}(\text{"John"}, x)$ " then  
begin  
  x := " $\text{Phone}(\text{"John"}, x)$ ";  
  print "John's phone number is:" x  
end
```


Certain Answers??

Example (Unintuitive Behaviour of Queries:)

1 $\exists x. \text{Phone}(\text{"John"}, x)? \Rightarrow \text{YES}$

2 $\text{Phone}(\text{"John"}, x)? \Rightarrow \{\}$

under $\Sigma = \{\forall x. \text{Person}(x) \rightarrow \exists y. \text{Phone}(x, y)\}$
and $D = \{\text{Person}(\text{"John"})\}$.

Embedded SQL-like Example

```
if " $\exists x. \text{Phone}(\text{"John"}, x)$ " then
begin
  x := " $\text{Phone}(\text{"John"}, x)$ ";
  print "John's phone number is:" x
end
```

Certain Answers??

Example (Unintuitive Behaviour of Queries:)

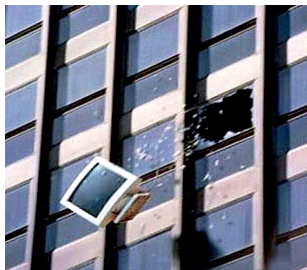
1 $\exists x. \text{Phone}(\text{"John"}, x)? \Rightarrow \text{YES}$

2 $\text{Phone}(\text{"John"}, x)? \Rightarrow \{\}$

under $\Sigma = \{\forall x. \text{Person}(x) \rightarrow \exists y. \text{Phone}(x, y)\}$
and $D = \{\text{Person}(\text{"John"})\}$.

Embedded SQL-like Example

```
if " $\exists x. \text{Phone}(\text{"John"}, x)$ " then
begin
  x := " $\text{Phone}(\text{"John"}, x)$ ";
  print "John's phone number is:" x
end
```



Summary

- certain answers are tractable only for *Lite schemata* and *Conjunctive/UC Queries*
- pretty much any extension leads to complexity (decidability) issues

Summary

- certain answers are tractable only for *Lite schemata* and *Conjunctive/UC Queries*
- pretty much any extension leads to complexity (decidability) issues

Next time: THE DATABASE EMPIRE STRIKES BACK