

When Conceptual Modeling meets Databases

Enrico Franconi

KRDB Research Centre for Knowledge and Data
Free University of Bozen-Bolzano, Italy

<http://krdb.eu/franconi>

EROSS 2020

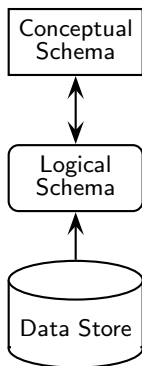
Summary of this talk

- ▶ What is a Conceptual Schema
- ▶ The role of a Conceptual Schema in Databases
- ▶ Querying via a Conceptual Schema
 - ▶ The classical Database case
 - ▶ Lossless transformations
 - ▶ Lossy transformations
 - ▶ DBoxes
 - ▶ ABoxes
 - ▶ Bad news and fixes

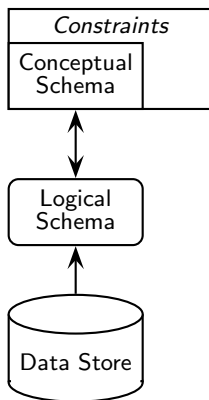
What is a Conceptual Schema

- ▶ A semantic layer gives an abstract perspective to a system, different from its original representation.
- ▶ Its purpose is to introduce a **vocabulary** understandable by the agents that need to interact with the system.
- ▶ A semantic layer is a formal conceptualisation of the agents' perspective: a **conceptual schema** (**ontology**, **logical theory**).
- ▶ A conceptual schema includes a set of **constraints** (the *semantics*) over the given vocabulary, which specifies what should hold in any possible configuration of the system.
- ▶ Any possible configuration of the system (an **instance**) conforms to the constraints expressed by the conceptual schema.
- ▶ Given a conceptual schema, a **legal instance** is a finite possible configuration of the system satisfying the constraints.
- ▶ We focus on database systems, their conceptual schemas, and their legal database instances.

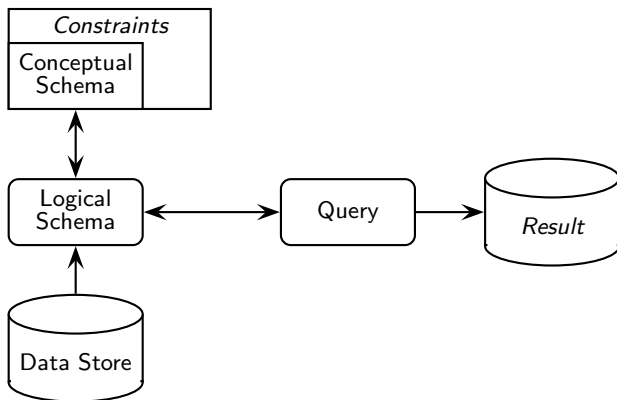
The role of a Conceptual Schema in Databases



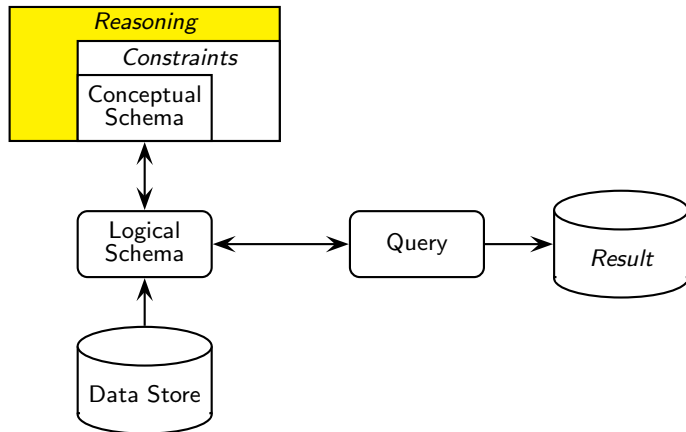
The role of a Conceptual Schema in Databases



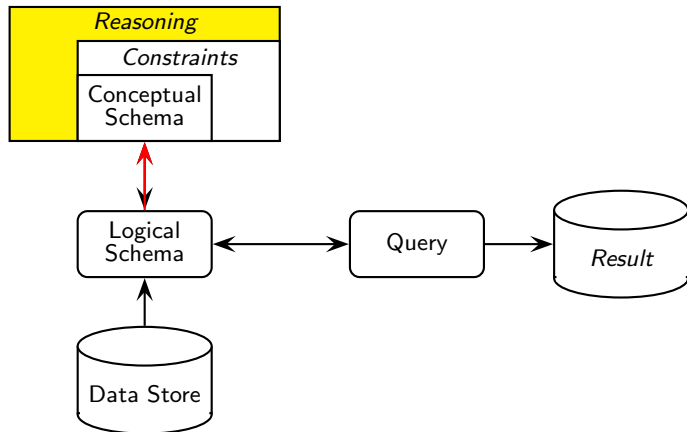
The role of a Conceptual Schema in Databases



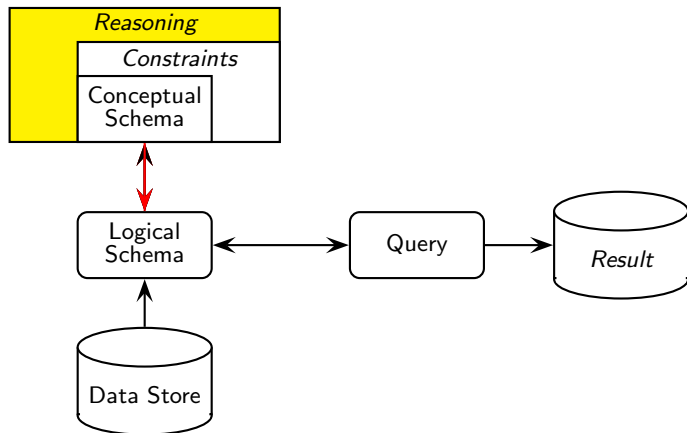
The role of a Conceptual Schema in Databases



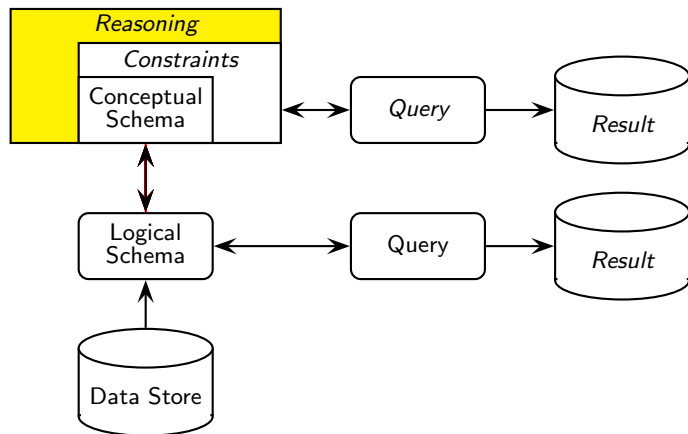
The role of a Conceptual Schema in Databases



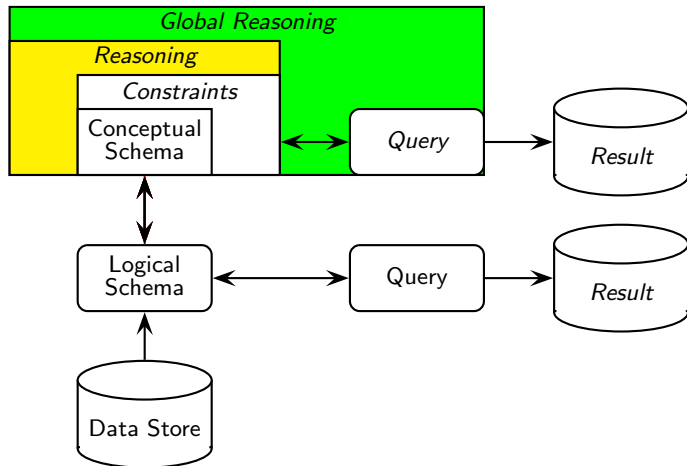
The role of a Conceptual Schema in Databases



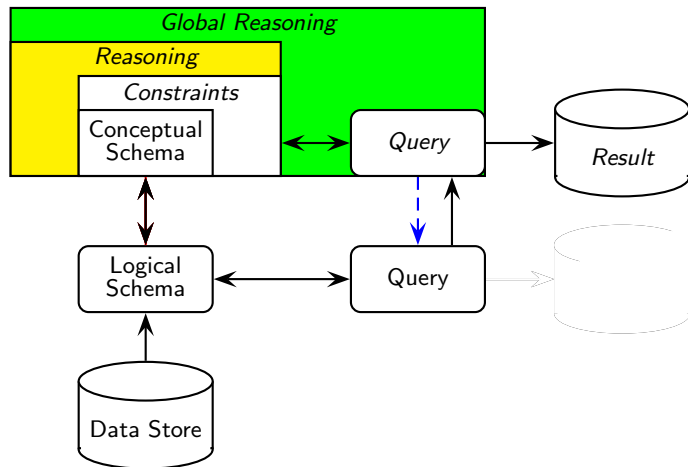
The role of a Conceptual Schema in Databases



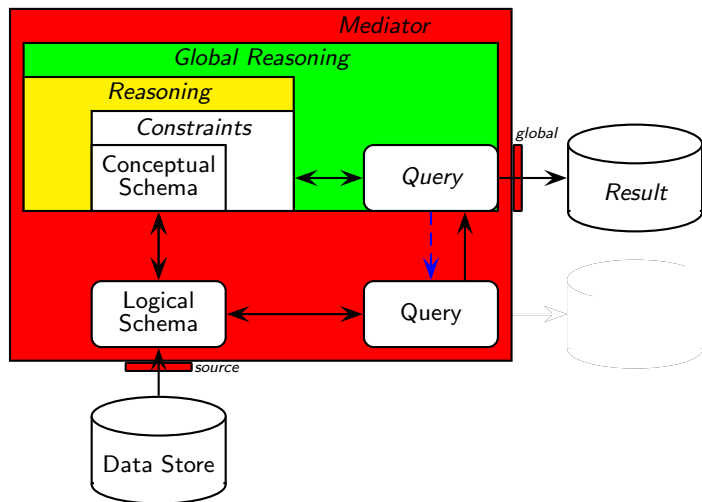
The role of a Conceptual Schema in Databases



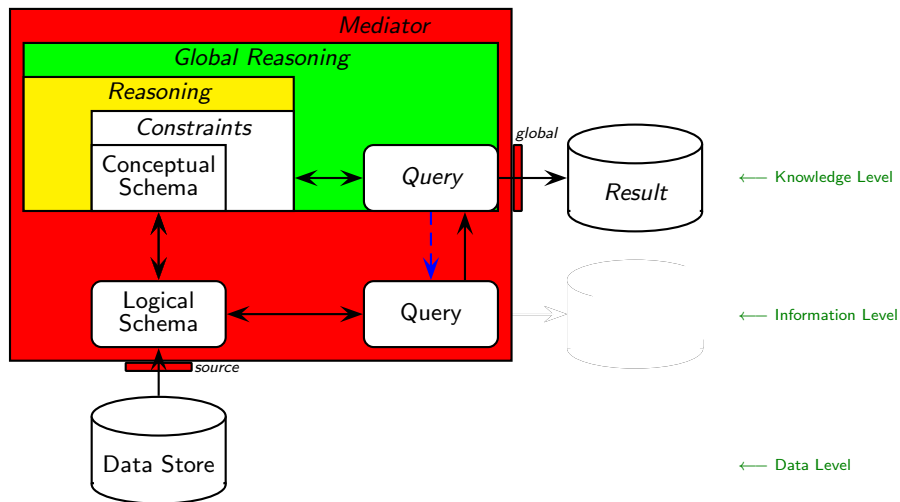
The role of a Conceptual Schema in Databases



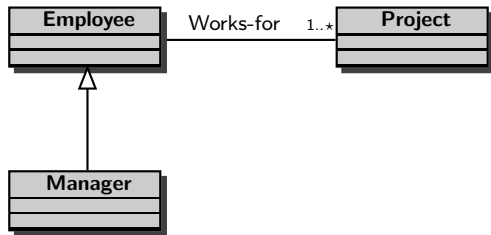
The role of a Conceptual Schema in Databases



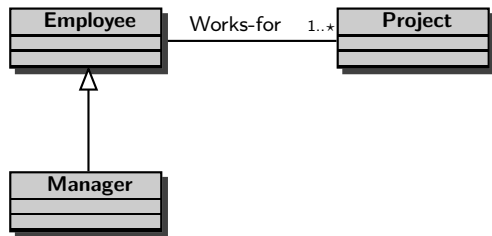
The role of a Conceptual Schema in Databases



Queries via Conceptual Schemas: the DB case



Queries via Conceptual Schemas: the DB case



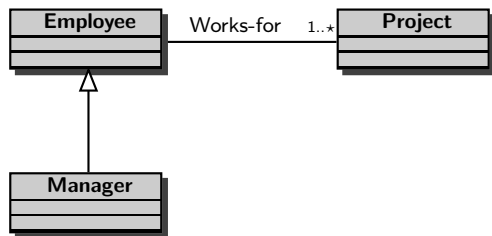
Employee = { John, Mary, Paul }

Manager = { John, Paul }

Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

Project = { Prj-A, Prj-B }

Queries via Conceptual Schemas: the DB case



Employee = { John, Mary, Paul }

Manager = { John, Paul }

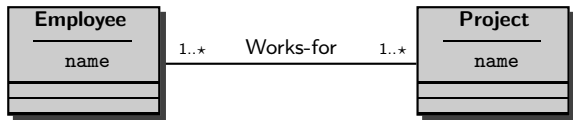
Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

Project = { Prj-A, Prj-B }

$Q(x) :- \exists y. \text{Manager}(x) \wedge \text{Works-for}(x,y) \wedge \text{Project}(y)$

$\implies \{ \text{John} \}$

(Conceptual) Schemas with different vocabularies



Employee = { John, Mary }

EmployeeName = { <John, "John">, <Mary, "Mary"> }

Project = { Prj-A, Prj-B }

ProjectName = { <Prj-A, "Prj-A">, <Prj-B, "Prj-B"> }

Works-for = { <John, Prj-A>, <John, Prj-B>, <Mary, Prj-A> }

Employee-table: (EmpOID, Ename, ProjOID, Pname)

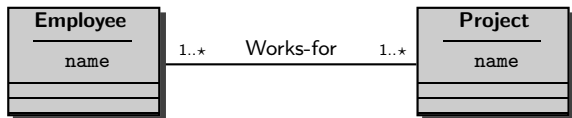
FD: ProjOID \rightarrow Pname

Employee-table = { <John, "John", Prj-A, "Prj-A">, <John, "John", Prj-B, "Prj-B">, <mary, "Mary", Prj-A, "Prj-A"> }

Lossless Transformations (aka Equivalence)

- ▶ Two (conceptual) schemas with different vocabularies are lossless transformations of one another (i.e., they are equivalent) if there exist **two total injective mappings** from legal database instances in one schema to legal database instances in the other schema such that **their composition is the identity**.
- ▶ Query answering across equivalent schemas involves expanding those mappings as views.
- ▶ This notion is the formal foundation of **Database Design**, of **Database Normalisation**, and of **Database Reverse Engineering**.

Lossless Transformations (aka Equivalence)



Employee = π_{EmpOID} Employee-table

EmployeeName = $\pi_{\text{EmpOID, Ename}}$ Employee-table

Project = π_{ProjOID} Employee-table

ProjectName = $\pi_{\text{ProjOID, Pname}}$ Employee-table

Works-for = $\pi_{\text{EmpOID, ProjOID}}$ Employee-table

Employee-table =

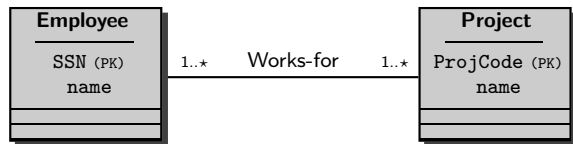
\bowtie (EmployeeName, ProjectName, Works-for)

Employee-table: (EmpOID, Ename, ProjOID, Pname)

FD: ProjID \rightarrow Pname

Object Identifiers

- ▶ A special lossless transformation: elimination of the **Object Identifiers** using **Primary Keys** or (better) **Reference Scheme Modelling** (see Terry Halpin's talk in this EROSS series).

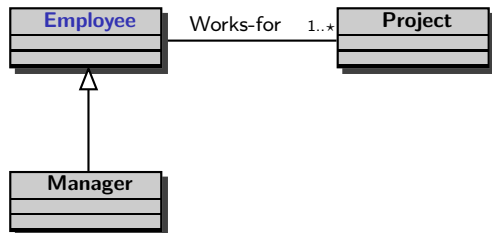


Employee-table: (SSN, Ename, ProjCode, Pname)

FD: ProjCode \rightarrow Pname

Remarkable Lossy Transformations

Queries via Conceptual Schemas: the DBox case

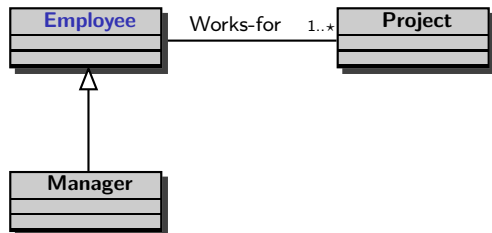


Manager = { John, Paul }

Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

Project = { Prj-A, Prj-B }

Queries via Conceptual Schemas: the DBox case



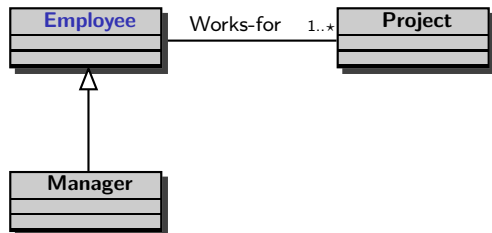
Manager = { John, Paul }

Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

Project = { Prj-A, Prj-B }

$Q(x) :- \text{Employee}(x)$

Queries via Conceptual Schemas: the DBox case



Manager = { John, Paul }

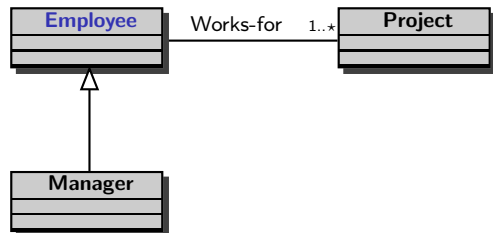
Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

Project = { Prj-A, Prj-B }

$Q(x) :- \text{Employee}(x)$

$\implies \{ \text{John, Paul, Mary} \}$ *certain answer*

Queries via Conceptual Schemas: the DBox case



Manager = { John, Paul }

Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-B⟩ }

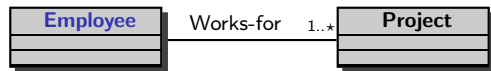
Project = { Prj-A, Prj-B }

$Q(x) :- \text{Employee}(x)$

\Rightarrow { John, Paul, Mary } *certain answer*

\Rightarrow $Q'(x) :- \text{Manager}(x) \vee \exists y. \text{Works-for}(x,y)$

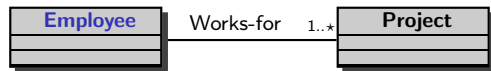
Queries via Conceptual Schemas: the ABox case



Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$

Project = { Prj-A, Prj-B, ... }

Queries via Conceptual Schemas: the ABox case

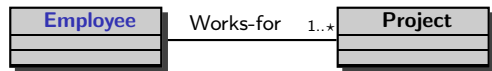


Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$

Project = { Prj-A, Prj-B, ... }

$Q(y) :- \exists x. \text{Works-for}(x,y)$ *certain answer*

Queries via Conceptual Schemas: the ABox case



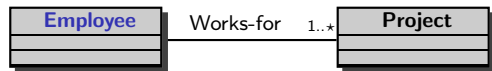
Works-for = { ⟨John, Prj-A⟩, ⟨Mary, Prj-A⟩, ... }

Project = { Prj-A, Prj-B, ... }

$Q(y) :- \exists x. \text{Works-for}(x,y)$ *certain answer*

$\implies \{ \text{Prj-A}, \text{Prj-B} \}$

Queries via Conceptual Schemas: the ABox case



Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$

Project = { Prj-A, Prj-B, ... }

$Q(y) :- \exists x. \text{Works-for}(x,y)$ *certain answer*

$\implies \{ \text{Prj-A}, \text{Prj-B} \}$

$\implies Q'(y) :- \text{Project}(y) \vee \exists x. \text{Works-for}(x,y)$

Bad news

Query answering with **certain answer semantics**
with DBoxes or ABoxes
with **expressive** conceptual modelling languages
is coNP-hard in data complexity.

Bad news

Query answering with **certain answer semantics**
with DBoxes or ABoxes
with **expressive** conceptual modelling languages
is coNP-hard in data complexity.

Fixes:

1. reduce expressivity of conceptual modelling languages and queries (OBDA approach);
2. use exact answer semantics allowing only determined queries (reduction to lossless transformations).

Bad news

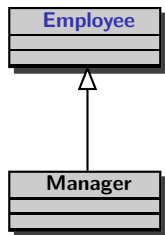
Query answering with **certain answer semantics**
with DBoxes or ABoxes
with **expressive** conceptual modelling languages
is coNP-hard in data complexity.

Fixes:

1. reduce expressivity of conceptual modelling languages and queries (OBDA approach);
2. use exact answer semantics allowing only determined queries (reduction to lossless transformations).

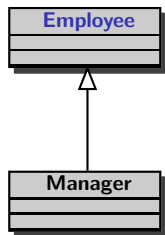
More bad news follow.

Queries with certain answer semantics



Manager = { John, Paul }

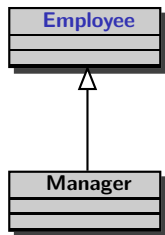
Queries with certain answer semantics



Manager = { John, Paul }

$Q(x) :- \text{Employee}(x)$

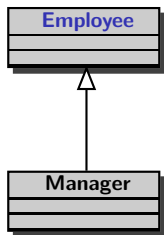
Queries with certain answer semantics



Manager = { John, Paul }

$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

Queries with certain answer semantics

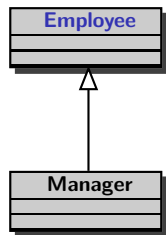


Manager = { John, Paul }

$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

So: are Managers and Employees the same?

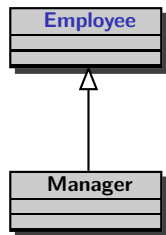
Queries with certain answer semantics



`Manager = { John, Paul }`

`Q(x) :- Employee(x) \implies { John, Paul }` *certain answer*

Queries with certain answer semantics



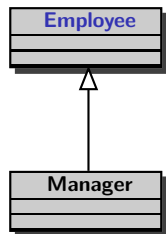
Manager = { John, Paul }

$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

$Q^1 :- \text{Manager}(\text{George})$

$Q^2 :- \text{Employee}(\text{George})$

Queries with certain answer semantics



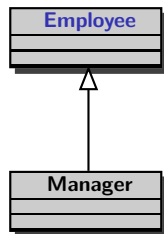
Manager = { John, Paul }

$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

$Q^1 :- \text{Manager}(\text{George}) \implies \text{FALSE}$

$Q^2 :- \text{Employee}(\text{George})$

Queries with certain answer semantics



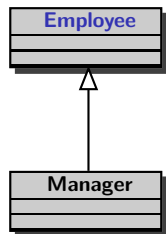
Manager = { John, Paul }

$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

$Q^1 :- \text{Manager}(\text{George}) \implies \text{FALSE}$

$Q^2 :- \text{Employee}(\text{George}) \implies \text{DON'T KNOW}$

Queries with certain answer semantics



Manager = { John, Paul }

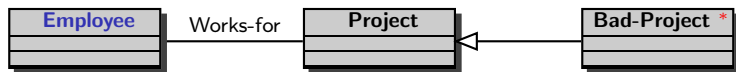
$Q(x) :- \text{Employee}(x) \implies \{ \text{John, Paul} \}$ *certain answer*

$Q^1 :- \text{Manager}(\text{George}) \implies \text{FALSE}$

$Q^2 :- \text{Employee}(\text{George}) \implies \text{DON'T KNOW}$

The result of the query can not be materialised (as a view)

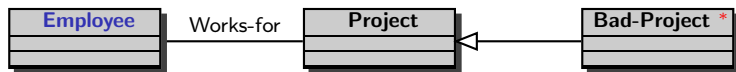
An ABox changes the DB semantics



Bad-Project defined with a rule/view:

* $\text{Bad-Project}(x) \leftrightarrow \text{Project}(x) \wedge \neg \exists y. \text{Works-for}(y, x)$

An ABox changes the DB semantics



Bad-Project defined with a rule/view:

* $\text{Bad-Project}(x) \leftrightarrow \text{Project}(x) \wedge \neg \exists y. \text{Works-for}(y, x)$

► DBox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B} \}$

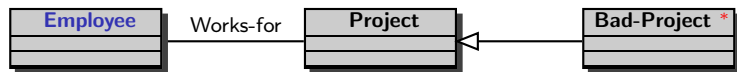
► $Q(X) :- \text{Bad-Project}(X)$

► ABox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B}, \dots \}$

► $Q(X) :- \text{Bad-Project}(X)$

An ABox changes the DB semantics



Bad-Project defined with a rule/view:

* $\text{Bad-Project}(x) \leftrightarrow \text{Project}(x) \wedge \neg \exists y. \text{Works-for}(y, x)$

► DBox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B} \}$

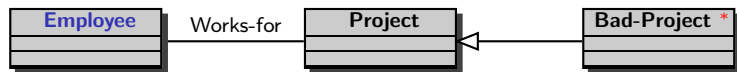
► $Q(X) :- \text{Bad-Project}(X)$
 $\implies \{ \text{Prj-B} \}$

► ABox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B}, \dots \}$

► $Q(X) :- \text{Bad-Project}(X)$

An ABox changes the DB semantics



Bad-Project defined with a rule/view:

* $\text{Bad-Project}(x) \leftrightarrow \text{Project}(x) \wedge \neg \exists y. \text{Works-for}(y, x)$

► DBox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B} \}$

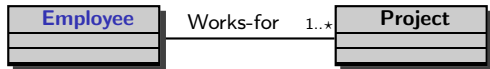
► $Q(X) :- \text{Bad-Project}(X)$
 $\implies \{ \text{Prj-B} \}$

► ABox:

$\text{Works-for} = \{ \langle \text{John}, \text{Prj-A} \rangle, \langle \text{Mary}, \text{Prj-A} \rangle, \dots \}$
 $\text{Project} = \{ \text{Prj-A}, \text{Prj-B}, \dots \}$

► $Q(X) :- \text{Bad-Project}(X)$
 $\implies \{ \}$

Compositionality of queries

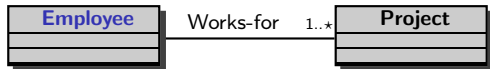


► ABox:

Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \dots$ }

Project = { Prj-A, Prj-B, ... }

Compositionality of queries



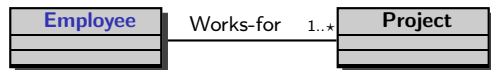
▶ ABox:

Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \dots$ }

Project = { Prj-A, Prj-B, ... }

▶ $Q(x) :- \exists y. \text{Works-for}(y,x)$ $Q = \pi_2 \text{ Works-for}$

Compositionality of queries



▶ ABox:

Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \dots$ }

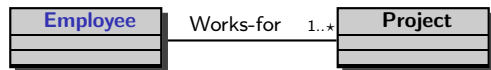
Project = { Prj-A, Prj-B, ... }

▶ $Q(x) :- \exists y. \text{Works-for}(y,x)$ $Q = \pi_2 \text{ Works-for}$

▶ $Q = \text{EVAL}(\pi_2 \text{ Works-for})$

▶ $Q = \pi_2 (\text{EVAL}(\text{Works-for}))$

Compositionality of queries



▶ ABox:

Works-for = { \langle John, Prj-A \rangle , ... }

Project = { Prj-A, Prj-B, ... }

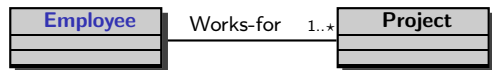
▶ $Q(x) :- \exists y. \text{Works-for}(y,x)$ $Q = \pi_2 \text{ Works-for}$

▶ $Q = \text{EVAL}(\pi_2 \text{ Works-for})$

$\implies \{ \text{Prj-A}, \text{Prj-B} \}$

▶ $Q = \pi_2 (\text{EVAL}(\text{Works-for}))$

Compositionality of queries



▶ ABox:

Works-for = { $\langle \text{John}, \text{Prj-A} \rangle, \dots$ }

Project = { Prj-A, Prj-B, ... }

▶ $Q(x) :- \exists y. \text{Works-for}(y,x)$ $Q = \pi_2 \text{ Works-for}$

▶ $Q = \text{EVAL}(\pi_2 \text{ Works-for})$
 $\implies \{ \text{Prj-A}, \text{Prj-B} \}$

▶ $Q = \pi_2 (\text{EVAL}(\text{Works-for}))$
 $\implies \{ \text{Prj-A} \}$

Conclusions

Conclusions

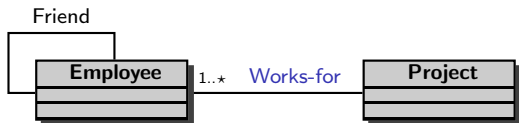
Do you want to exploit conceptual schema knowledge
(i.e., constraints or an ontology)
in your data intensive application?

Conclusions

Do you want to exploit conceptual schema knowledge
(i.e., constraints or an ontology)
in your data intensive application?

Pay attention!

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

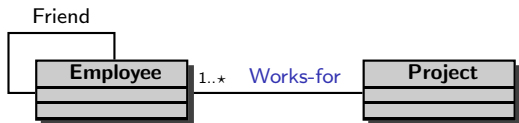
Employee = {John, Mary, etc}

Project = { Prj-A, Prj-B, Prj-C }

$Q :- \text{Works-for}(e_1, p) \wedge \text{Works-for}(e_2, p) \wedge \text{Friend}(e_1, e_2)$

Are there two friends working for the same project?

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

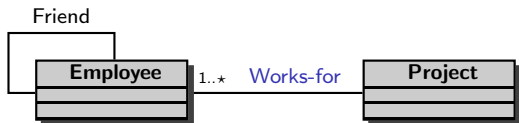
Employee = {John, Mary, etc}

Project = { Prj-A, Prj-B, Prj-C }

$Q :- \text{Works-for}(e_1, p) \wedge \text{Works-for}(e_2, p) \wedge \text{Friend}(e_1, e_2)$

Are there two friends working for the same project?

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

Employee = {John, Mary, etc}

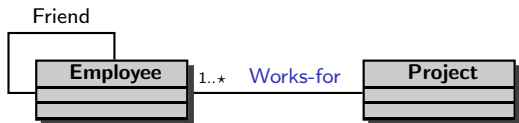
Project = { Prj-A, Prj-B, Prj-C }

Q :- Works-for(e_1, p) \wedge Works-for(e_2, p) \wedge Friend(e_1, e_2)

Are there two friends working for the same project?

- ▶ **YES:** in any legal database instance, there are at least two friends working for the same project.

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

Employee = {John, Mary, etc}

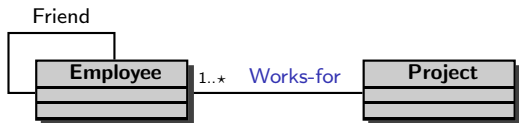
Project = { Prj-A, Prj-B, Prj-C }

Q :- Works-for(e_1, p) \wedge Works-for(e_2, p) \wedge Friend(e_1, e_2)

Are there two friends working for the same project?

- ▶ **YES**: in any legal database instance, there are at least two friends working for the same project.
- ▶ **NO**: there is at least a legal database instance in which no two friends work for the same project.

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

Employee = {John, Mary, etc}

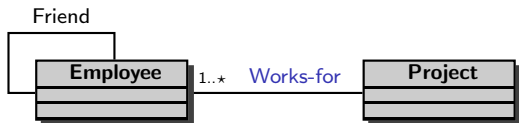
Project = { Prj-A, Prj-B, Prj-C }

Q :- Works-for(e_1, p) \wedge Works-for(e_2, p) \wedge Friend(e_1, e_2)

Are there two friends working for the same project?

- ▶ **YES**: in any legal database instance, there are at least two friends working for the same project.
- ▶ **NO**: there is at least a legal database instance in which no two friends work for the same project.
- ▶ **NOTE**: with *ABox semantics* the answer is always **NO**, since there is at least a legal database instance with *enough* distinct projects so that no two friends work for the same project.

Data Complexity of Query Answering with DBoxes



Friend = {⟨John, Mary⟩, etc}

Employee = {John, Mary, etc}

Project = { Prj-A, Prj-B, Prj-C }

Q :- Works-for(e_1, p) \wedge Works-for(e_2, p) \wedge Friend(e_1, e_2)

Are there two friends working for the same project?

- ▶ **YES**: in any legal database instance, there are at least two friends working for the same project.
- ▶ **NO**: there is at least a legal database instance in which no two friends work for the same project.
- ▶ **NOTE**: with *ABox semantics* the answer is always **NO**, since there is at least a legal database instance with *enough* distinct projects so that no two friends work for the same project.

This encodes 3-colorability of maps \rightarrow co-NP-hard