

# Data Integration

Maurizio Lenzerini

Università di Roma “La Sapienza”

*DASI '06: Phd School on Data and Service Integration*  
Bertinoro, December 11–15, 2006

# Structure of the course

- 1 Introduction to data integration
  - Basic issues in data integration
  - Logical formalization
- 2 Query answering in the absence of constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
- 3 Query answering in the presence of constraints
  - The role of integrity constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
  - Inconsistency tolerance
- 4 Concluding remarks

# Structure of the course

- 1 Introduction to data integration
  - Basic issues in data integration
  - Logical formalization
- 2 Query answering in the absence of constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
- 3 Query answering in the presence of constraints
  - The role of integrity constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
  - Inconsistency tolerance
- 4 Concluding remarks

# Structure of the course

- 1 Introduction to data integration
  - Basic issues in data integration
  - Logical formalization
- 2 Query answering in the absence of constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
- 3 Query answering in the presence of constraints
  - The role of integrity constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
  - Inconsistency tolerance
- 4 Concluding remarks

# Structure of the course

- 1 Introduction to data integration
  - Basic issues in data integration
  - Logical formalization
- 2 Query answering in the absence of constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
- 3 Query answering in the presence of constraints
  - The role of integrity constraints
  - Global-as-view (GAV) setting
  - Local-as-view (LAV) and GLAV setting
  - Inconsistency tolerance
- 4 Concluding remarks

# Part I

## Introduction to data integration

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
  
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems













# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Architectures for integrated access to distributed data

- **Distributed databases**  
data sources are homogeneous databases under the control of the distributed database management system
- **Multidatabase or federated databases**  
data sources are autonomous, heterogeneous databases; procedural specification
- **(Mediator-based) data integration**  
access through a global schema mapped to autonomous and heterogeneous data sources; declarative specification
- **Peer-to-peer data integration**  
network of autonomous systems mapped one to each other, without a global schema; declarative specification





# Database federation tools: Characteristics

- **Physical transparency**, i.e., masking from the user the physical characteristics of the sources
- **Heterogeneity**, i.e., federating highly diverse types of sources
- **Extensibility**
- **Autonomy** of data sources
- **Performance**, through distributed query optimization

However, current tools do not (directly) support **logical (or conceptual) transparency**























# The modeling problem

## Basic questions:

- How to model the global schema
  - data model
  - constraints
- How to model the sources
  - data model (conceptual and logical level)
  - access limitations
  - data values (common vs. different domains)
- How to model the mapping between global schemas and sources
- How to verify the quality of the modeling process

**A word of caution:** Data modeling (in data integration) is an art.  
Theoretical frameworks can help humans, not replace them

# The modeling problem

## Basic questions:

- How to model the global schema
  - data model
  - constraints
- How to model the sources
  - data model (conceptual and logical level)
  - access limitations
  - data values (common vs. different domains)
- How to model the mapping between global schemas and sources
- How to verify the quality of the modeling process

**A word of caution:** Data modeling (in data integration) is an art. Theoretical frameworks can help humans, not replace them

# The querying problem

- A query expressed in terms of the global schema must be **reformulated** in terms of (a set of) queries over the sources and/or materialized views
- The computed sub-queries are shipped to the sources, and the results are collected and **assembled** into the final answer
- The computed query plan should guarantee
  - completeness of the obtained answers wrt the semantics
  - efficiency of the whole query answering process
  - efficiency in accessing sources
- This process heavily depends on the approach adopted for modeling the data integration system

This is the problem that we want to address in this course

# The querying problem

- A query expressed in terms of the global schema must be **reformulated** in terms of (a set of) queries over the sources and/or materialized views
- The computed sub-queries are shipped to the sources, and the results are collected and **assembled** into the final answer
- The computed query plan should guarantee
  - completeness of the obtained answers wrt the semantics
  - efficiency of the whole query answering process
  - efficiency in accessing sources
- This process heavily depends on the approach adopted for modeling the data integration system

**This is the problem that we want to address in this course**

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems



# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
  
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Formal framework for data integration

## Definition

A **data integration system**  $\mathcal{I}$  is a triple  $\langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ , where

- $\mathcal{G}$  is the global schema  
*i.e., a logical theory over a relational alphabet  $\mathcal{A}_{\mathcal{G}}$*
- $\mathcal{S}$  is the source schema  
*i.e., simply a relational alphabet  $\mathcal{A}_{\mathcal{S}}$  disjoint from  $\mathcal{A}_{\mathcal{G}}$*
- $\mathcal{M}$  is the mapping between  $\mathcal{S}$  and  $\mathcal{G}$   
*We consider different approaches to the specification of mappings*

# Semantics of a data integration system

Which are the dbs that satisfy  $\mathcal{I}$ , i.e., the logical models of  $\mathcal{I}$ ?

- We refer only to dbs over a **fixed infinite domain**  $\Delta$  of elements
- We start from the data present in the sources: these are modeled through a **source database**  $\mathcal{C}$  over  $\Delta$  (also called source model), fixing the extension of the predicates of  $\mathcal{A}_S$
- The dbs for  $\mathcal{I}$  are logical interpretations for  $\mathcal{A}_G$ , called **global dbs**

## Definition

The **set of databases for  $\mathcal{A}_G$  that satisfy  $\mathcal{I}$  relative to  $\mathcal{C}$**  is:

$$sem^{\mathcal{C}}(\mathcal{I}) = \{ \mathcal{B} \mid \mathcal{B} \text{ is a global database that is legal wrt } \mathcal{G} \\ \text{and that satisfies } \mathcal{M} \text{ wrt } \mathcal{C} \}$$

What it means to satisfy  $\mathcal{M}$  wrt  $\mathcal{C}$  depends on the nature of  $\mathcal{M}$



# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - **Relational calculus**
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Relational calculus: the basics

**Basic idea:** we use the language of first-order logic to express which tuples should be in the result to a query

- We assume to have a domain  $\Delta$  and a set  $\Sigma$  of constants, one for each element of  $\Delta$
- Let  $\mathcal{A}$  be a **relational alphabet**, i.e., a set of predicates, each with an associated arity (we assume a positional notation)
- A **database  $\mathcal{D}$  over  $\mathcal{A}$  and  $\Delta$**  is a set of relations, one for each predicate in  $\mathcal{A}$ , over the constants in  $\Sigma$  (in turn interpreted as elements of  $\Delta$ )
- Let  $\mathcal{L}_{\mathcal{A}}$  be the first-order language over
  - the constants in  $\Sigma$
  - the predicates of  $\mathcal{A}$  plus the built-in predicates of relational algebra (e.g.,  $<$ ,  $>$ ,  $\dots$ )
  - no function symbols

# Relational calculus: Syntax

## Definition

An **(domain) relational calculus query** over alphabet  $\mathcal{A}$  has the form

$$\{ (x_1, \dots, x_n) \mid \varphi \},$$

where

- $n \geq 0$  is the **arity** of the query
- $x_1, \dots, x_n$  are (not necessarily distinct) variables
- $\varphi$  is the **body** of the query, i.e., a formula of  $\mathcal{L}_{\mathcal{A}}$  whose free variables are exactly  $x_1, \dots, x_n$
- $(x_1, \dots, x_n)$  is called the **target list** of the query

If  $r$  is a predicate of arity  $k$ , an **atom** with predicate  $r$  has the form  $r(y_1, \dots, y_k)$ , where  $y_1, \dots, y_k$  are variables or constants

# Relational calculus: Semantics

Relational calculus queries are evaluated on particular interpretations

## Definition

A **correct interpretation** for relational calculus queries over  $\mathcal{A}$  is a pair  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$ , where  $\Delta$  is a domain, and  $\mathcal{D}$  is a database over  $\mathcal{A}$  and  $\Delta$

## Definition

The **value** of a relational calculus query  $q = \{(x_1, \dots, x_n) \mid \varphi\}$  in an interpretation  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$  is the set of tuples  $(c_1, \dots, c_n)$  of constants in  $\Sigma$  such that  $\langle \mathcal{I}, \mathcal{V} \rangle \models \varphi$ , where  $\mathcal{V}$  is the variable assignment that assigns  $c_i$  to  $x_i$

When the domain  $\Delta$  is clear, we can omit it, and write directly  $\langle \mathcal{D}, \mathcal{V} \rangle \models \varphi$ , instead of  $\langle \langle \Delta, \mathcal{D} \rangle, \mathcal{V} \rangle \models \varphi$



# Relational calculus: Semantics

Relational calculus queries are evaluated on particular interpretations

## Definition

A **correct interpretation** for relational calculus queries over  $\mathcal{A}$  is a pair  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$ , where  $\Delta$  is a domain, and  $\mathcal{D}$  is a database over  $\mathcal{A}$  and  $\Delta$

## Definition

The **value** of a relational calculus query  $q = \{(x_1, \dots, x_n) \mid \varphi\}$  in an interpretation  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$  is the set of tuples  $(c_1, \dots, c_n)$  of constants in  $\Sigma$  such that  $\langle \mathcal{I}, \mathcal{V} \rangle \models \varphi$ , where  $\mathcal{V}$  is the variable assignment that assigns  $c_i$  to  $x_i$

When the domain  $\Delta$  is clear, we can omit it, and write directly  $\langle \mathcal{D}, \mathcal{V} \rangle \models \varphi$ , instead of  $\langle \langle \Delta, \mathcal{D} \rangle, \mathcal{V} \rangle \models \varphi$

# Relational calculus: Semantics

Relational calculus queries are evaluated on particular interpretations

## Definition

A **correct interpretation** for relational calculus queries over  $\mathcal{A}$  is a pair  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$ , where  $\Delta$  is a domain, and  $\mathcal{D}$  is a database over  $\mathcal{A}$  and  $\Delta$

## Definition

The **value** of a relational calculus query  $q = \{(x_1, \dots, x_n) \mid \varphi\}$  in an interpretation  $\mathcal{I} = \langle \Delta, \mathcal{D} \rangle$  is the set of tuples  $(c_1, \dots, c_n)$  of constants in  $\Sigma$  such that  $\langle \mathcal{I}, \mathcal{V} \rangle \models \varphi$ , where  $\mathcal{V}$  is the variable assignment that assigns  $c_i$  to  $x_i$

When the domain  $\Delta$  is clear, we can omit it, and write directly  $\langle \mathcal{D}, \mathcal{V} \rangle \models \varphi$ , instead of  $\langle \langle \Delta, \mathcal{D} \rangle, \mathcal{V} \rangle \models \varphi$

# Result of relational calculus queries

## Definition

The **result of the evaluation** of a relational calculus query  $q = \{(x_1, \dots, x_n) \mid \varphi\}$  on a database  $\mathcal{D}$  over  $\mathcal{A}$  and  $\Delta$  is the relation  $q^{\mathcal{D}}$  such that

- the arity of  $q^{\mathcal{D}}$  is  $n$
- the extension of  $q^{\mathcal{D}}$  is the set of constants that constitute the value of the query  $q$  in the interpretation  $\langle \Delta, \mathcal{D} \rangle$

# Conjunctive queries

- are the most common kind of relational calculus queries
- also known as **select-project-join** SQL queries
- allow for easy optimization in relational DBMSs

## Definition

A **conjunctive query** (CQ) is a relational calculus query of the form

$$\{ (\vec{x}) \mid \exists \vec{y}. r_1(\vec{x}_1, \vec{y}_1) \wedge \cdots \wedge r_m(\vec{x}_m, \vec{y}_m) \}$$

where

- $\vec{x}$  is the union of the  $\vec{x}_i$ 's, and  $\vec{y}$  is the union of the  $\vec{y}_i$ 's
- $r_1, \dots, r_m$  are relation symbols (not built-in predicates)

We use the following abbreviation:  $\{ (\vec{x}) \mid r_1(\vec{x}_1, \vec{y}_1), \dots, r_m(\vec{x}_m, \vec{y}_m) \}$

# Conjunctive queries

- are the most common kind of relational calculus queries
- also known as **select-project-join** SQL queries
- allow for easy optimization in relational DBMSs

## Definition

A **conjunctive query** (CQ) is a relational calculus query of the form

$$\{ (\vec{x}) \mid \exists \vec{y}. r_1(\vec{x}_1, \vec{y}_1) \wedge \cdots \wedge r_m(\vec{x}_m, \vec{y}_m) \}$$

where

- $\vec{x}$  is the union of the  $\vec{x}_i$ 's, and  $\vec{y}$  is the union of the  $\vec{y}_i$ 's
- $r_1, \dots, r_m$  are relation symbols (not built-in predicates)

We use the following abbreviation:  $\{ (\vec{x}) \mid r_1(\vec{x}_1, \vec{y}_1), \dots, r_m(\vec{x}_m, \vec{y}_m) \}$

# Complexity of relational calculus

We consider the complexity of the **recognition problem**, i.e., checking whether a tuple of constants is in the answer to a query:

- measured wrt the size of the database  $\rightsquigarrow$  **data complexity**
- measured wrt the size of the query and the database  $\rightsquigarrow$  **combined complexity**

## Complexity of relational calculus

- data complexity: polynomial, actually in LOGSPACE
- combined complexity: PSPACE-complete

## Complexity of conjunctive queries

- data complexity: in LOGSPACE
- combined complexity: NP-complete

# Complexity of relational calculus

We consider the complexity of the **recognition problem**, i.e., checking whether a tuple of constants is in the answer to a query:

- measured wrt the size of the database  $\rightsquigarrow$  **data complexity**
- measured wrt the size of the query and the database  $\rightsquigarrow$  **combined complexity**

## Complexity of relational calculus

- data complexity: polynomial, actually in LOGSPACE
- combined complexity: PSPACE-complete

## Complexity of conjunctive queries

- data complexity: in LOGSPACE
- combined complexity: NP-complete

# Complexity of relational calculus

We consider the complexity of the **recognition problem**, i.e., checking whether a tuple of constants is in the answer to a query:

- measured wrt the size of the database  $\rightsquigarrow$  **data complexity**
- measured wrt the size of the query and the database  $\rightsquigarrow$  **combined complexity**

## Complexity of relational calculus

- data complexity: polynomial, actually in LOGSPACE
- combined complexity: PSPACE-complete

## Complexity of conjunctive queries

- data complexity: in LOGSPACE
- combined complexity: NP-complete



# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Queries to a data integration system $\mathcal{I}$

- The domain  $\Delta$  is fixed, and we do not distinguish an element of  $\Delta$  from the constant denoting it  $\rightsquigarrow$  **standard names**
- Queries to  $\mathcal{I}$  are relational calculus queries over the alphabet  $\mathcal{A}_g$  of the global schema
- When “evaluating”  $q$  over  $\mathcal{I}$ , we have to consider that for a **given source database  $\mathcal{C}$** , there may be **many global databases  $\mathcal{B}$**  in  $sem^{\mathcal{C}}(\mathcal{I})$
- We consider those answers to  $q$  that hold for **all** global databases in  $sem^{\mathcal{C}}(\mathcal{I})$   
 $\rightsquigarrow$  **certain answers**

# Semantics of queries to $\mathcal{I}$

## Definition

Given  $q$ ,  $\mathcal{I}$ , and  $\mathcal{C}$ , the set of **certain answers to  $q$  wrt  $\mathcal{I}$  and  $\mathcal{C}$**  is

$$\mathit{cert}(q, \mathcal{I}, \mathcal{C}) = \{ (c_1, \dots, c_n) \in q^{\mathcal{B}} \mid \text{for all } \mathcal{B} \in \mathit{sem}^{\mathcal{C}}(\mathcal{I}) \}$$

- Query answering is **logical implication**
- Complexity is measured mainly *wrt the size of the source db  $\mathcal{C}$* , i.e., we consider **data complexity**
- We consider the problem of deciding whether  $\vec{c} \in \mathit{cert}(q, \mathcal{I}, \mathcal{C})$ , for a given  $\vec{c}$

# Databases with incomplete information, or knowledge bases

- **Traditional database:** one model of a first-order theory  
Query answering means **evaluating** a formula in the model
- **Database with incomplete information, or knowledge base:** set of models (specified, for example, as a restricted first-order theory)  
Query answering means computing the tuples that satisfy the query in **all** the models in the set

There is a **strong connection** between query answering in data integration and query answering in databases with incomplete information under constraints (or, query answering in knowledge bases)

# Databases with incomplete information, or knowledge bases

- **Traditional database:** one model of a first-order theory  
Query answering means **evaluating** a formula in the model
- **Database with incomplete information, or knowledge base:** set of models (specified, for example, as a restricted first-order theory)  
Query answering means computing the tuples that satisfy the query in **all** the models in the set

There is a **strong connection** between query answering in data integration and query answering in databases with incomplete information under constraints (or, query answering in knowledge bases)

# Query answering with incomplete information

- [Reiter '84]: relational setting, databases with incomplete information modeled as a first order theory
- [Vardi '86]: relational setting, complexity of reasoning in closed world databases with unknown values
- Several approaches both from the DB and the KR community
- [van der Meyden '98]: survey on logical approaches to incomplete information in databases

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - **Formalizing the mapping**
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# The mapping

How is the mapping  $\mathcal{M}$  between  $\mathcal{S}$  and  $\mathcal{G}$  specified?

- Are the sources defined in terms of the global schema?  
Approach called **source-centric**, or **local-as-view**, or **LAV**
- Is the global schema defined in terms of the sources?  
Approach called **global-schema-centric**, or **global-as-view**, or **GAV**
- A mixed approach?  
Approach called **GLAV**



# GAV vs. LAV – Example

## Global schema:

*movie(Title, Year, Director)*

*european(Director)*

*review(Title, Critique)*

## Source 1:

*r<sub>1</sub>(Title, Year, Director)* since 1960, european directors

## Source 2:

*r<sub>2</sub>(Title, Critique)* since 1990

## Query: Title and critique of movies in 1998

$\{ (t, r) \mid \exists d. \text{movie}(t, 1998, d) \wedge \text{review}(t, r) \}$ , abbreviated

$\{ (t, r) \mid \text{movie}(t, 1998, d), \text{review}(t, r) \}$

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - **Formalizing GAV data integration systems**
  - Formalizing LAV data integration systems

# Formalization of GAV

In GAV (with **sound sources**), the mapping  $\mathcal{M}$  is a set of assertions:

$$\phi_S \rightsquigarrow g$$

one for each element  $g$  in  $\mathcal{A}_G$ , with  $\phi_S$  a **query** over  $S$  of the arity of  $g$

Given a source db  $C$ , a db  $B$  for  $G$  satisfies  $\mathcal{M}$  wrt  $C$  if for each  $g \in G$ :

$$\phi_S^C \subseteq g^B$$

In other words, the assertion means  $\forall \vec{x}. \phi_S(\vec{x}) \rightarrow g(\vec{x})$

Given a source database,  $\mathcal{M}$  **provides direct information** about which data satisfy the elements of the global schema

*Relations in  $G$  are views, and queries are expressed over the views.*

Thus, it **seems** that we can simply evaluate the query over the data satisfying the global relations (as if we had a single database at hand)

# GAV – Example

**Global schema:**  $\text{movie}(\textit{Title}, \textit{Year}, \textit{Director})$   
 $\text{european}(\textit{Director})$   
 $\text{review}(\textit{Title}, \textit{Critique})$

**GAV:** to each relation in the global schema,  $\mathcal{M}$  associates a view over the sources:

$$\begin{aligned} \{ (t, y, d) \mid r_1(t, y, d) \} &\rightsquigarrow \text{movie}(t, y, d) \\ \{ (d) \mid r_1(t, y, d) \} &\rightsquigarrow \text{european}(d) \\ \{ (t, r) \mid r_2(t, r) \} &\rightsquigarrow \text{review}(t, r) \end{aligned}$$

Logical formalization:

$$\begin{aligned} \forall t, y, d. r_1(t, y, d) &\rightarrow \text{movie}(t, y, d) \\ \forall d. (\exists t, y. r_1(t, y, d)) &\rightarrow \text{european}(d) \\ \forall t, r. r_2(t, r) &\rightarrow \text{review}(t, r) \end{aligned}$$

# GAV – Example of query processing

The query

$$\{ (t, r) \mid \text{movie}(t, 1998, d), \text{review}(t, r) \}$$

is processed by means of **unfolding**, i.e., by expanding each atom according to its associated definition in  $\mathcal{M}$ , so as to come up with source relations

In this case:

$$\begin{array}{ccc} \{ (t, r) \mid \text{movie}(t, 1998, d), \text{review}(t, r) \} & & \\ \text{unfolding} & \downarrow & \downarrow \\ \{ (t, r) \mid r_1(t, 1998, d), r_2(t, r) \} & & \end{array}$$

# GAV – Example of constraints

## Global schema containing constraints:

movie(*Title*, *Year*, *Director*)

european(*Director*)

review(*Title*, *Critique*)

european\_movie\_60s(*Title*, *Year*, *Director*)

$\forall t, y, d. \text{european\_movie\_60s}(t, y, d) \rightarrow \text{movie}(t, y, d)$

$\forall d. \exists t, y. \text{european\_movie\_60s}(t, y, d) \rightarrow \text{european}(d)$

## GAV mappings:

$\{ (t, y, d) \mid r_1(t, y, d) \} \rightsquigarrow \text{european\_movie\_60s}(t, y, d)$

$\{ (d) \mid r_1(t, y, d) \} \rightsquigarrow \text{european}(d)$

$\{ (t, r) \mid r_2(t, r) \} \rightsquigarrow \text{review}(t, r)$

# Outline

- 1 Basic issues in data integration
  - The problem of data integration
  - Variants of data integration
  - Problems in data integration
- 2 Data integration: Logical formalization
  - Semantics of a data integration system
  - Relational calculus
  - Queries to a data integration system
  - Formalizing the mapping
  - Formalizing GAV data integration systems
  - Formalizing LAV data integration systems

# Formalization of LAV

In LAV (with **sound sources**), the mapping  $\mathcal{M}$  is a set of assertions:

$$s \rightsquigarrow \phi_{\mathcal{G}}$$

one for each source element  $s$  in  $\mathcal{A}_{\mathcal{S}}$ , with  $\phi_{\mathcal{G}}$  a **query** over  $\mathcal{G}$

Given source db  $\mathcal{C}$ , a db  $\mathcal{B}$  for  $\mathcal{G}$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  if for each  $s \in \mathcal{S}$ :

$$s^{\mathcal{C}} \subseteq \phi_{\mathcal{G}}^{\mathcal{B}}$$

In other words, the assertion means  $\forall \vec{x}. s(\vec{x}) \rightarrow \phi_{\mathcal{G}}(\vec{x})$

The mapping  $\mathcal{M}$  and the source database  $\mathcal{C}$  do **not** provide direct information about which data satisfy the global schema

*Sources are views, and we have to answer queries on the basis of the available data in the views*



# LAV – Example

**Global schema:**    *movie*(*Title*, *Year*, *Director*)  
                           *european*(*Director*)  
                           *review*(*Title*, *Critique*)

**LAV:** to each **source relation**,  $\mathcal{M}$  associates a **view** over the global schema:

$$\begin{aligned} r_1(t, y, d) &\rightsquigarrow \{ (t, y, d) \mid \text{movie}(t, y, d), \text{european}(d), y \geq 1960 \} \\ r_2(t, r) &\rightsquigarrow \{ (t, r) \mid \text{movie}(t, y, d), \text{review}(t, r), y \geq 1990 \} \end{aligned}$$

The query  $\{ (t, r) \mid \text{movie}(t, 1998, d), \text{review}(t, r) \}$  is processed by means of an inference mechanism that aims at re-expressing the atoms of the global schema in terms of atoms at the sources.

In this case:

$$\{ (t, r) \mid r_2(t, r), r_1(t, 1998, d) \}$$

# GAV and LAV – Comparison

**GAV:** (e.g., Carnot, SIMS, Tsimmis, IBIS, Momis, DisAtDis, ...)

- Quality depends on how well we have compiled the sources into the global schema through the mapping
- Whenever a source changes or a new one is added, the global schema needs to be reconsidered
- Query processing can be based on some sort of unfolding (query answering looks easier – without constraints)

**LAV:** (e.g., Information Manifold, DWQ, Pícel)

- Quality depends on how well we have characterized the sources
- High modularity and extensibility (if the global schema is well designed, when a source changes, only its definition is affected)
- Query processing needs reasoning (query answering complex)

# Beyond GAV and LAV: GLAV

In GLAV (with **sound sources**), the mapping  $\mathcal{M}$  is a set of assertions:

$$\phi_S \rightsquigarrow \phi_G$$

with  $\phi_S$  a **query** over  $\mathcal{S}$ , and  $\phi_G$  a **query** over  $\mathcal{G}$  of the same arity as  $\phi_S$

Given source db  $\mathcal{C}$ , a db  $\mathcal{B}$  for  $\mathcal{G}$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  if for each  $\phi_S \rightsquigarrow \phi_G$  in  $\mathcal{M}$ :

$$\phi_S^{\mathcal{C}} \subseteq \phi_G^{\mathcal{B}}$$

In other words, the assertion means  $\forall \vec{x}. \phi_S(\vec{x}) \rightarrow \phi_G(\vec{x})$

As for LAV, the mapping  $\mathcal{M}$  does **not** provide direct information about which data satisfy the global schema

To answer a query  $q$  over  $\mathcal{G}$ , we have to **infer** how to use  $\mathcal{M}$  in order to access the source database  $\mathcal{C}$

# GLAV – Example

**Global schema:**  $\text{work}(\text{Person}, \text{Project}), \quad \text{area}(\text{Project}, \text{Field})$

**Source 1:**  $\text{hasjob}(\text{Person}, \text{Field})$

**Source 2:**  $\text{teaches}(\text{Professor}, \text{Course}), \quad \text{in}(\text{Course}, \text{Field})$

**Source 3:**  $\text{get}(\text{Researcher}, \text{Grant}), \quad \text{for}(\text{Grant}, \text{Project})$

**GLAV mapping:**

$$\begin{aligned} \{(r, f) \mid \text{hasjob}(r, f)\} &\rightsquigarrow \{(r, f) \mid \text{work}(r, p), \text{area}(p, f)\} \\ \{(r, f) \mid \text{teaches}(r, c), \text{in}(c, f)\} &\rightsquigarrow \{(r, f) \mid \text{work}(r, p), \text{area}(p, f)\} \\ \{(r, p) \mid \text{get}(r, g), \text{for}(g, p)\} &\rightsquigarrow \{(r, f) \mid \text{work}(r, p)\} \end{aligned}$$

# GLAV – A technical observation

In GLAV (with **sound sources**), the mapping  $\mathcal{M}$  is constituted by a set of assertions:

$$\phi_S \rightsquigarrow \phi_G$$

Each such assertion can be rewritten wlog by introducing a **new predicate**  $r$  (not to be used in the queries) of the same arity as the two queries and replace the assertion with the following two:

$$\phi_S \rightsquigarrow r \qquad r \rightsquigarrow \phi_G$$

In other words, we replace  $\forall \vec{x}. \phi_S(\vec{x}) \rightarrow \phi_G(\vec{x})$   
with  $\forall \vec{x}. \phi_S(\vec{x}) \rightarrow r(\vec{x})$  and  $\forall \vec{x}. r(\vec{x}) \rightarrow \phi_G(\vec{x})$







































































































# Outline

- 3 Query answering in GAV without constraints
  - Retrieved global database
  - Query answering via unfolding
  - Universal solutions
  
- 4 Query answering in (G)LAV without constraints
  - (G)LAV and incompleteness
  - Approaches to query answering in (G)LAV
  - (G)LAV: Direct methods (aka view-based query answering)
  - (G)LAV: Query answering by (view-based) query rewriting

# (G)LAV – Basic technique

From [Duschka & Genesereth PODS'97]:

$$\begin{aligned}
 r_1(t) &\rightsquigarrow \{ (t) \mid \text{movie}(t, y, d) \wedge \text{european}(d) \} \\
 r_2(t, v) &\rightsquigarrow \{ (t, v) \mid \text{movie}(t, y, d) \wedge \text{review}(t, v) \}
 \end{aligned}$$

$$\begin{aligned}
 \forall t. r_1(t) &\rightarrow \exists y, d. \text{movie}(t, y, d) \wedge \text{european}(d) \\
 \forall t, v. r_2(t, v) &\rightarrow \exists y, d. \text{movie}(t, y, d) \wedge \text{review}(t, v)
 \end{aligned}$$

$$\begin{aligned}
 \text{movie}(t, f_1(t), f_2(t)) &\leftarrow r_1(t) \\
 \text{european}(f_2(t)) &\leftarrow r_1(t) \\
 \text{movie}(t, f_4(t, v), f_5(t, v)) &\leftarrow r_2(t, v) \\
 \text{review}(t, v) &\leftarrow r_2(t, v)
 \end{aligned}$$

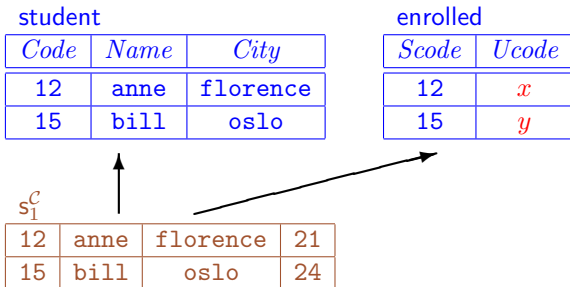
*Answering a query means evaluating a goal wrt to this nonrecursive logic program* (which can be transformed into a union of CQs)

$\rightsquigarrow$  **Data complexity is polynomial** (actually LOGSPACE)



## (G)LAV – Example of canonical model

$$\{ (c, n, ci) \mid s_1(c, n, ci, a) \} \rightsquigarrow \{ (c, n, ci) \mid \text{student}(c, n, ci) \wedge \text{enrolled}(c, u) \}$$



Example of source db  $\mathcal{C}$  and corresponding canonical model  $\mathcal{M}(\mathcal{C}) \downarrow$



# (G)LAV – Universal solution

Let  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$  be a (G)LAV data integration system without constraints in the global schema, and  $\mathcal{C}$  a source database

## Theorem

Then  $\mathcal{M}(\mathcal{C}) \downarrow$  is a **universal solution** for  $\mathcal{I}$  relative to  $\mathcal{C}$  (follows from [Fagin&al. ICDT'03])

It follows that:

- If  $q$  is a conjunctive query, then  $\vec{t} \in \text{cert}(q, \mathcal{I}, \mathcal{C})$  iff  $\vec{t} \in q^{\mathcal{M}(\mathcal{C}) \downarrow}$
- Complexity of query answering is **polynomial**, actually LOGSPACE





# (G)LAV – More expressive queries?

- More expressive **source queries** in the mapping?
  - Same results hold if we use **any computable query** as source query in the mapping assertions
- More expressive **queries over the global schema** in the mapping?
  - Already positive queries lead to intractability
- More expressive **user queries**?
  - Same results hold if we use **Datalog queries** as user queries
  - Even the simplest form of negation (inequalities) leads to intractability

# (G)LAV – More expressive queries?

- More expressive **source queries** in the mapping?
  - Same results hold if we use **any computable query** as source query in the mapping assertions
  
- More expressive **queries over the global schema in the mapping?**
  - Already positive queries lead to intractability
  
- More expressive **user queries?**
  - Same results hold if we use **Datalog queries** as user queries
  - Even the simplest form of negation (inequalities) leads to intractability

# (G)LAV – More expressive queries?

- More expressive **source queries** in the mapping?
  - Same results hold if we use **any computable query** as source query in the mapping assertions
  
- More expressive **queries over the global schema** in the mapping?
  - Already positive queries lead to intractability
  
- More expressive **user queries**?
  - Same results hold if we use **Datalog queries** as user queries
  - Even the simplest form of negation (inequalities) leads to intractability

# (G)LAV – Intractability for positive views

From [van der Meyden TCS'93], by reduction from 3-colorability

We define the following LAV data integration system  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ :

$$\begin{aligned} \mathcal{G} : & \text{edge}(x, y), \quad \text{color}(x, c) & \mathcal{S} : & s_E(x, y), \quad s_N(x) \\ \mathcal{M} : & s_E(x, y) \rightsquigarrow \text{edge}(x, y) \\ & s_N(x) \rightsquigarrow \text{color}(x, \text{RED}) \vee \text{color}(x, \text{BLUE}) \vee \text{color}(x, \text{GREEN}) \end{aligned}$$

Given a graph  $G = (N, E)$ , we define the following source database  $\mathcal{C}$ :

$$s_E^{\mathcal{C}} = \{ (a, b), (b, a) \mid (a, b) \in E \} \quad s_N^{\mathcal{C}} = \{ (a) \mid a \in N \}$$

Consider the boolean query  $q: \exists x, y, c. \text{edge}(x, y) \wedge \text{color}(x, c) \wedge \text{color}(y, c)$  describing mismatched edge pairs:

- If  $G$  is 3-colorable, then  $\exists \mathcal{B}$  s.t.  $q^{\mathcal{B}} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- If  $G$  is not 3-colorable, then  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{true}$

## Theorem

*Data complexity is coNP-hard for positive views and conjunctive queries*

## (G)LAV – Intractability for positive views

From [van der Meyden TCS'93], by reduction from 3-colorability

We define the following LAV data integration system  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ :

$$\begin{aligned} \mathcal{G} : & \text{edge}(x, y), \quad \text{color}(x, c) & \mathcal{S} : & \text{s}_E(x, y), \quad \text{s}_N(x) \\ \mathcal{M} : & \text{s}_E(x, y) \rightsquigarrow \text{edge}(x, y) \\ & \text{s}_N(x) \rightsquigarrow \text{color}(x, \text{RED}) \vee \text{color}(x, \text{BLUE}) \vee \text{color}(x, \text{GREEN}) \end{aligned}$$

Given a graph  $G = (N, E)$ , we define the following source database  $\mathcal{C}$ :

$$\text{s}_E^{\mathcal{C}} = \{ (a, b), (b, a) \mid (a, b) \in E \} \quad \text{s}_N^{\mathcal{C}} = \{ (a) \mid a \in N \}$$

Consider the boolean query  $q: \exists x, y, c. \text{edge}(x, y) \wedge \text{color}(x, c) \wedge \text{color}(y, c)$  describing mismatched edge pairs:

- If  $G$  is 3-colorable, then  $\exists \mathcal{B}$  s.t.  $q^{\mathcal{B}} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- If  $G$  is not 3-colorable, then  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{true}$

### Theorem

*Data complexity is coNP-hard for positive views and conjunctive queries*

# (G)LAV – Intractability for positive views

From [van der Meyden TCS'93], by reduction from 3-colorability

We define the following LAV data integration system  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ :

$$\begin{aligned} \mathcal{G} : & \text{edge}(x, y), \quad \text{color}(x, c) & \mathcal{S} : & \text{s}_E(x, y), \quad \text{s}_N(x) \\ \mathcal{M} : & \text{s}_E(x, y) \rightsquigarrow \text{edge}(x, y) \\ & \text{s}_N(x) \rightsquigarrow \text{color}(x, \text{RED}) \vee \text{color}(x, \text{BLUE}) \vee \text{color}(x, \text{GREEN}) \end{aligned}$$

Given a graph  $G = (N, E)$ , we define the following source database  $\mathcal{C}$ :

$$\text{s}_E^{\mathcal{C}} = \{ (a, b), (b, a) \mid (a, b) \in E \} \quad \text{s}_N^{\mathcal{C}} = \{ (a) \mid a \in N \}$$

Consider the boolean query  $q: \exists x, y, c. \text{edge}(x, y) \wedge \text{color}(x, c) \wedge \text{color}(y, c)$  describing mismatched edge pairs:

- If  $G$  is 3-colorable, then  $\exists \mathcal{B}$  s.t.  $q^{\mathcal{B}} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- If  $G$  is not 3-colorable, then  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{true}$

## Theorem

*Data complexity is coNP-hard for positive views and conjunctive queries*

# (G)LAV – Intractability for positive views

From [van der Meyden TCS'93], by reduction from 3-colorability

We define the following LAV data integration system  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ :

$$\begin{aligned} \mathcal{G} : & \text{edge}(x, y), \quad \text{color}(x, c) & \mathcal{S} : & \text{s}_E(x, y), \quad \text{s}_N(x) \\ \mathcal{M} : & \text{s}_E(x, y) \rightsquigarrow \text{edge}(x, y) \\ & \text{s}_N(x) \rightsquigarrow \text{color}(x, \text{RED}) \vee \text{color}(x, \text{BLUE}) \vee \text{color}(x, \text{GREEN}) \end{aligned}$$

Given a graph  $G = (N, E)$ , we define the following source database  $\mathcal{C}$ :

$$\text{s}_E^{\mathcal{C}} = \{ (a, b), (b, a) \mid (a, b) \in E \} \quad \text{s}_N^{\mathcal{C}} = \{ (a) \mid a \in N \}$$

Consider the boolean query  $q: \exists x, y, c. \text{edge}(x, y) \wedge \text{color}(x, c) \wedge \text{color}(y, c)$  describing mismatched edge pairs:

- If  $G$  is 3-colorable, then  $\exists \mathcal{B}$  s.t.  $q^{\mathcal{B}} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- If  $G$  is not 3-colorable, then  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{true}$

## Theorem

*Data complexity is coNP-hard for positive views and conjunctive queries*



# (G)LAV – Intractability for positive views

From [van der Meyden TCS'93], by reduction from 3-colorability

We define the following LAV data integration system  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ :

$$\begin{aligned} \mathcal{G} : & \text{edge}(x, y), \quad \text{color}(x, c) & \mathcal{S} : & \text{s}_E(x, y), \quad \text{s}_N(x) \\ \mathcal{M} : & \text{s}_E(x, y) \rightsquigarrow \text{edge}(x, y) \\ & \text{s}_N(x) \rightsquigarrow \text{color}(x, \text{RED}) \vee \text{color}(x, \text{BLUE}) \vee \text{color}(x, \text{GREEN}) \end{aligned}$$

Given a graph  $G = (N, E)$ , we define the following source database  $\mathcal{C}$ :

$$\text{s}_E^{\mathcal{C}} = \{ (a, b), (b, a) \mid (a, b) \in E \} \quad \text{s}_N^{\mathcal{C}} = \{ (a) \mid a \in N \}$$

Consider the boolean query  $q: \exists x, y, c. \text{edge}(x, y) \wedge \text{color}(x, c) \wedge \text{color}(y, c)$  describing mismatched edge pairs:

- If  $G$  is 3-colorable, then  $\exists \mathcal{B}$  s.t.  $q^{\mathcal{B}} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- If  $G$  is not 3-colorable, then  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{true}$

## Theorem

*Data complexity is coNP-hard for positive views and conjunctive queries*

# (G)LAV – In coNP for positive views and queries

- $\vec{t} \notin \text{cert}(q, \mathcal{I}, \mathcal{C})$  if and only if there is a database  $\mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and such that  $\vec{t} \notin q^{\mathcal{B}}$
- The mapping  $\mathcal{M}$  has the form:

$$\forall \vec{x}. \phi_{\mathcal{S}}(\vec{x}) \rightarrow \exists \vec{y}_1. \alpha_1(\vec{x}, \vec{y}_1) \vee \dots \vee \exists \vec{y}_h. \alpha_h(\vec{x}, \vec{y}_h)$$

Hence, each tuple in  $\mathcal{C}$  forces the existence of  $k$  tuples in any database that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , where  $k$  is the maximal length of conjunctions  $\alpha_i(\vec{x}, \vec{y}_i)$  in  $\mathcal{M}$

- If  $\mathcal{C}$  has  $n$  tuples, then there is a db  $\mathcal{B}' \subseteq \mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  with at most  $n \cdot k$  tuples. Since  $q$  is monotone,  $\vec{t} \notin q^{\mathcal{B}'}$
- Checking whether  $\mathcal{B}'$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and checking whether  $\vec{t} \notin q^{\mathcal{B}'}$  can be done in PTIME wrt the size of  $\mathcal{B}'$

## Theorem

For positive views and queries, query answ. is *coNP in data complexity*

## (G)LAV – In coNP for positive views and queries

- $\vec{t} \notin \text{cert}(q, \mathcal{I}, \mathcal{C})$  if and only if there is a database  $\mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and such that  $\vec{t} \notin q^{\mathcal{B}}$
- The mapping  $\mathcal{M}$  has the form:

$$\forall \vec{x}. \phi_{\mathcal{S}}(\vec{x}) \rightarrow \exists \vec{y}_1. \alpha_1(\vec{x}, \vec{y}_1) \vee \dots \vee \exists \vec{y}_h. \alpha_h(\vec{x}, \vec{y}_h))$$

Hence, each tuple in  $\mathcal{C}$  forces the existence of  $k$  tuples in any database that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , where  $k$  is the maximal length of conjunctions  $\alpha_i(\vec{x}, \vec{y}_i)$  in  $\mathcal{M}$

- If  $\mathcal{C}$  has  $n$  tuples, then there is a db  $\mathcal{B}' \subseteq \mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  with at most  $n \cdot k$  tuples. Since  $q$  is monotone,  $\vec{t} \notin q^{\mathcal{B}'}$
- Checking whether  $\mathcal{B}'$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and checking whether  $\vec{t} \notin q^{\mathcal{B}'}$  can be done in PTIME wrt the size of  $\mathcal{B}'$

### Theorem

For positive views and queries, query answ. is *coNP in data complexity*

# (G)LAV – In coNP for positive views and queries

- $\vec{t} \notin \text{cert}(q, \mathcal{I}, \mathcal{C})$  if and only if there is a database  $\mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and such that  $\vec{t} \notin q^{\mathcal{B}}$
- The mapping  $\mathcal{M}$  has the form:

$$\forall \vec{x}. \phi_{\mathcal{S}}(\vec{x}) \rightarrow \exists \vec{y}_1. \alpha_1(\vec{x}, \vec{y}_1) \vee \dots \vee \exists \vec{y}_h. \alpha_h(\vec{x}, \vec{y}_h)$$

Hence, each tuple in  $\mathcal{C}$  forces the existence of  $k$  tuples in any database that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , where  $k$  is the maximal length of conjunctions  $\alpha_i(\vec{x}, \vec{y}_i)$  in  $\mathcal{M}$

- If  $\mathcal{C}$  has  $n$  tuples, then there is a db  $\mathcal{B}' \subseteq \mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  with at most  $n \cdot k$  tuples. Since  $q$  is monotone,  $\vec{t} \notin q^{\mathcal{B}'}$
- Checking whether  $\mathcal{B}'$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and checking whether  $\vec{t} \notin q^{\mathcal{B}'}$  can be done in PTIME wrt the size of  $\mathcal{B}'$

## Theorem

*For positive views and queries, query answ. is coNP in data complexity*

# (G)LAV – In coNP for positive views and queries

- $\vec{t} \notin \text{cert}(q, \mathcal{I}, \mathcal{C})$  if and only if there is a database  $\mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and such that  $\vec{t} \notin q^{\mathcal{B}}$
- The mapping  $\mathcal{M}$  has the form:

$$\forall \vec{x}. \phi_{\mathcal{S}}(\vec{x}) \rightarrow \exists \vec{y}_1. \alpha_1(\vec{x}, \vec{y}_1) \vee \cdots \vee \exists \vec{y}_h. \alpha_h(\vec{x}, \vec{y}_h)$$

Hence, each tuple in  $\mathcal{C}$  forces the existence of  $k$  tuples in any database that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , where  $k$  is the maximal length of conjunctions  $\alpha_i(\vec{x}, \vec{y}_i)$  in  $\mathcal{M}$

- If  $\mathcal{C}$  has  $n$  tuples, then there is a db  $\mathcal{B}' \subseteq \mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  with at most  $n \cdot k$  tuples. Since  $q$  is monotone,  $\vec{t} \notin q^{\mathcal{B}'}$
- Checking whether  $\mathcal{B}'$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and checking whether  $\vec{t} \notin q^{\mathcal{B}'}$  can be done in PTIME wrt the size of  $\mathcal{B}'$

## Theorem

*For positive views and queries, query answ. is coNP in data complexity*

# (G)LAV – In coNP for positive views and queries

- $\vec{t} \notin \text{cert}(q, \mathcal{I}, \mathcal{C})$  if and only if there is a database  $\mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and such that  $\vec{t} \notin q^{\mathcal{B}}$
- The mapping  $\mathcal{M}$  has the form:

$$\forall \vec{x}. \phi_{\mathcal{S}}(\vec{x}) \rightarrow \exists \vec{y}_1. \alpha_1(\vec{x}, \vec{y}_1) \vee \dots \vee \exists \vec{y}_h. \alpha_h(\vec{x}, \vec{y}_h)$$

Hence, each tuple in  $\mathcal{C}$  forces the existence of  $k$  tuples in any database that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , where  $k$  is the maximal length of conjunctions  $\alpha_i(\vec{x}, \vec{y}_i)$  in  $\mathcal{M}$

- If  $\mathcal{C}$  has  $n$  tuples, then there is a db  $\mathcal{B}' \subseteq \mathcal{B}$  for  $\mathcal{I}$  that satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  with at most  $n \cdot k$  tuples. Since  $q$  is monotone,  $\vec{t} \notin q^{\mathcal{B}'}$
- Checking whether  $\mathcal{B}'$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$ , and checking whether  $\vec{t} \notin q^{\mathcal{B}'}$  can be done in PTIME wrt the size of  $\mathcal{B}'$

## Theorem

For positive views and queries, query answ. is **coNP in data complexity**

## (G)LAV – Conjunctive user queries with inequalities

Consider  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ , and source db  $\mathcal{C}$  (see [Fagin & al. ICDT'03]):

$$\begin{aligned} \mathcal{G} &: g(x, y) & \mathcal{S} &: s(x, y) \\ \mathcal{M} &: s(x, y) \rightsquigarrow \{ (x, y) \mid g(x, z) \wedge g(z, y) \} \\ \mathcal{C} &: \{ s(\mathbf{a}, \mathbf{a}) \} \end{aligned}$$

- $\mathcal{B}_1 = \{ g(\mathbf{a}, \mathbf{a}) \}$  is a solution
- If  $\mathcal{B}$  is a universal solution, then both  $g(\mathbf{a}, x)$  and  $g(x, \mathbf{a})$  are in  $\mathcal{B}$ , with  $x \neq \mathbf{a}$  (otherwise  $g(\mathbf{a}, \mathbf{a})$  would be *true* in every solution)

Let  $q = \{ () \mid g(x, y) \wedge x \neq y \}$

- $q^{\mathcal{B}_1} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- But  $q^{\mathcal{B}} = \text{true}$  for every universal solution  $\mathcal{B}$  for  $\mathcal{I}$  relative to  $\mathcal{C}$

Hence, the notion of universal solution is not the right tool

## (G)LAV – Conjunctive user queries with inequalities

Consider  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$ , and source db  $\mathcal{C}$  (see [Fagin & al. ICDT'03]):

$$\begin{aligned} \mathcal{G} : & \quad g(x, y) & \quad \mathcal{S} : & \quad s(x, y) \\ \mathcal{M} : & \quad s(x, y) \rightsquigarrow \{ (x, y) \mid g(x, z) \wedge g(z, y) \} \\ \mathcal{C} : & \quad \{ s(\mathbf{a}, \mathbf{a}) \} \end{aligned}$$

- $\mathcal{B}_1 = \{ g(\mathbf{a}, \mathbf{a}) \}$  is a solution
- If  $\mathcal{B}$  is a universal solution, then both  $g(\mathbf{a}, x)$  and  $g(x, \mathbf{a})$  are in  $\mathcal{B}$ , with  $x \neq \mathbf{a}$  (otherwise  $g(\mathbf{a}, \mathbf{a})$  would be *true* in every solution)

Let  $q = \{ () \mid g(x, y) \wedge x \neq y \}$

- $q^{\mathcal{B}_1} = \text{false}$ , hence  $\text{cert}(q, \mathcal{I}, \mathcal{C}) = \text{false}$
- But  $q^{\mathcal{B}} = \text{true}$  for every universal solution  $\mathcal{B}$  for  $\mathcal{I}$  relative to  $\mathcal{C}$

Hence, the notion of universal solution is not the right tool



## (G)LAV – Conjunctive user queries with inequalities

- coNP algorithm: guess equalities on variables in the canonical retrieved global database
- coNP-hard already for a conjunctive user query with two inequalities (and conjunctive view definitions) [Fagin & al. ICDT'03], see Phokion's course

### Theorem

*For conjunctive user queries with inequalities, (G)LAV query answering is **coNP-complete in data complexity***

*Note: inequalities in the view definitions do not affect expressive power and complexity (in fact, they can be removed)*

## (G)LAV – Conjunctive user queries with inequalities

- coNP algorithm: guess equalities on variables in the canonical retrieved global database
- coNP-hard already for a conjunctive user query with two inequalities (and conjunctive view definitions) [Fagin & al. ICDT'03], see Phokion's course

### Theorem

*For conjunctive user queries with inequalities, (G)LAV query answering is **coNP-complete in data complexity***

*Note:* inequalities in the view definitions do not affect expressive power and complexity (in fact, they can be removed)

# Outline

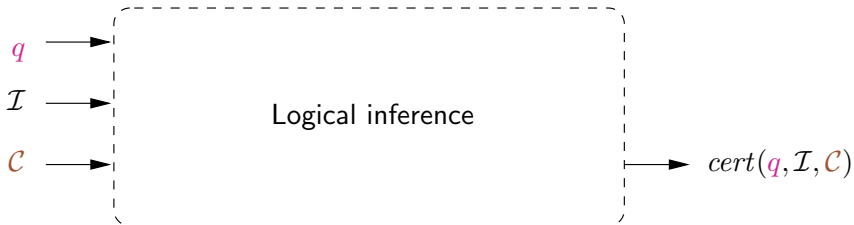
- 3 Query answering in GAV without constraints
  - Retrieved global database
  - Query answering via unfolding
  - Universal solutions
- 4 Query answering in (G)LAV without constraints
  - (G)LAV and incompleteness
  - Approaches to query answering in (G)LAV
  - (G)LAV: Direct methods (aka view-based query answering)
  - (G)LAV: Query answering by (view-based) query rewriting

# (G)LAV – View-based query rewriting

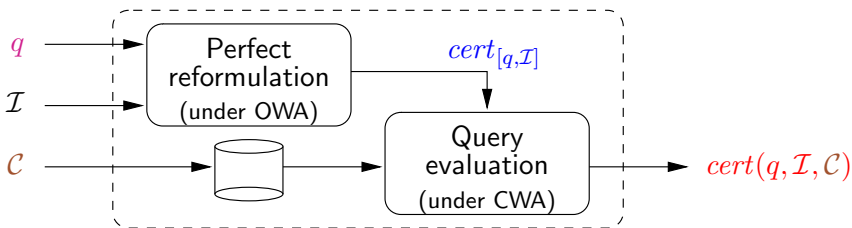
Query answering is divided in two steps:

- 1 Re-express the query in terms of a **given query language** over the alphabet of  $\mathcal{A}_S$
- 2 Evaluate the rewriting over the source database  $\mathcal{C}$

# Query answering

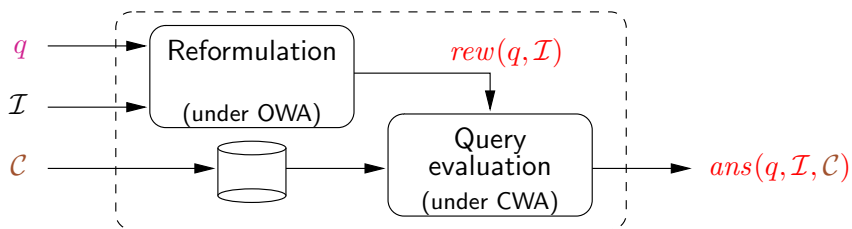


# Query answering: reformulation + evaluation



The query  $cert_{[q, \mathcal{I}]}$  could be expressed in an **arbitrary query language**

# Query rewriting



The language of  $rew(q, \mathcal{I})$  is chosen a priori!

# (G)LAV – Connection to rewriting

## Query answering by rewriting:

- ① Given  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$  and a query  $q$  over  $\mathcal{G}$ , rewrite  $q$  into a query, called  $rew(q, \mathcal{I})$ , over the alphabet  $\mathcal{A}_{\mathcal{S}}$  of the sources
- ② Evaluate the rewriting  $rew(q, \mathcal{I})$  over the source database  $\mathcal{C}$

We are interested in rewritings that are:

- **sound**, i.e., compute only tuples in  $cert(q, \mathcal{I}, \mathcal{C})$  for every  $\mathcal{C}$
- expressed in a given **query language**  $\mathcal{L}$
- **maximal** for the class of queries expressible in  $\mathcal{L}$

We may be interested in **exact** rewritings, i.e., rewritings that are logically equivalent to the query, modulo  $\mathcal{M}$

Exact rewritings may not exist



# (G)LAV – Connection to rewriting

## Query answering by rewriting:

- ① Given  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$  and a query  $q$  over  $\mathcal{G}$ , rewrite  $q$  into a query, called  $rew(q, \mathcal{I})$ , over the alphabet  $\mathcal{A}_{\mathcal{S}}$  of the sources
- ② Evaluate the rewriting  $rew(q, \mathcal{I})$  over the source database  $\mathcal{C}$

We are interested in rewritings that are:

- **sound**, i.e., compute only tuples in  $cert(q, \mathcal{I}, \mathcal{C})$  for every  $\mathcal{C}$
- expressed in a given **query language**  $\mathcal{L}$
- **maximal** for the class of queries expressible in  $\mathcal{L}$

We may be interested in **exact** rewritings, i.e., rewritings that are logically equivalent to the query, modulo  $\mathcal{M}$

Exact rewritings may not exist

# (G)LAV – Connection to rewriting

## Query answering by rewriting:

- ① Given  $\mathcal{I} = \langle \mathcal{G}, \mathcal{S}, \mathcal{M} \rangle$  and a query  $q$  over  $\mathcal{G}$ , rewrite  $q$  into a query, called  $rew(q, \mathcal{I})$ , over the alphabet  $\mathcal{A}_{\mathcal{S}}$  of the sources
- ② Evaluate the rewriting  $rew(q, \mathcal{I})$  over the source database  $\mathcal{C}$

We are interested in rewritings that are:

- **sound**, i.e., compute only tuples in  $cert(q, \mathcal{I}, \mathcal{C})$  for every  $\mathcal{C}$
- expressed in a given **query language**  $\mathcal{L}$
- **maximal** for the class of queries expressible in  $\mathcal{L}$

We may be interested in **exact** rewritings, i.e., rewritings that are logically equivalent to the query, modulo  $\mathcal{M}$

Exact rewritings may not exist

# (G)LAV – Example of maximal rewriting

$G$ :  $\text{nonstop}(Airline, Num, From, To)$

$S$ :  $\text{flightsByUnited}(From, To)$   
 $\text{flightsFromSFO}(Airline, Num, To)$

$M$ :  $\text{flightsByUnited}(From, To) \rightsquigarrow$   
 $\text{nonstop}(UA, Num, From, To)$   
 $\text{flightsFromSFO}(Airline, Num, To) \rightsquigarrow$   
 $\text{nonstop}(Airline, Num, SFO, To)$

$q$ :  $\{ (airline, num) \mid \text{nonstop}(airline, num, LAX, PHX) \}$

A maximal (wrt positive queries) rewriting of  $q$  is:

$\{ (UA, num) \mid \text{flightsByUnited}(num, LAX, PHX) \}$

# (G)LAV – Example of maximal rewriting

$\mathcal{G}$ :  $\text{nonstop}(Airline, Num, From, To)$

$\mathcal{S}$ :  $\text{flightsByUnited}(From, To)$   
 $\text{flightsFromSFO}(Airline, Num, To)$

$\mathcal{M}$ :  $\text{flightsByUnited}(From, To) \rightsquigarrow$   
 $\text{nonstop}(UA, Num, From, To)$   
 $\text{flightsFromSFO}(Airline, Num, To) \rightsquigarrow$   
 $\text{nonstop}(Airline, Num, SFO, To)$

$q$ :  $\{ (airline, num) \mid \text{nonstop}(airline, num, LAX, PHX) \}$

A maximal (wrt positive queries) rewriting of  $q$  is:

$\{ (UA, num) \mid \text{flightsByUnited}(num, LAX, PHX) \}$

# Perfect rewriting

What is the relationship between answering by rewriting and certain answers? [Calvanese & al. ICDT'05]:

- When does the (maximal) rewriting compute **all** certain answers?
- What do we gain or loose by focusing on a given class of queries?

Let's try to consider the “**best possible**” rewriting

Define  $cert_{[q, \mathcal{I}]}(\cdot)$  to be the function that, with  $q$  and  $\mathcal{I}$  fixed, given source database  $\mathcal{C}$ , computes the certain answers  $cert(q, \mathcal{I}, \mathcal{C})$ .

- $cert_{[q, \mathcal{I}]}$  can be seen as a query on the alphabet  $\mathcal{A}_{\mathcal{S}}$
- $cert_{[q, \mathcal{I}]}$  is a (sound) rewriting of  $q$  wrt  $\mathcal{I}$
- No sound rewriting exists that is better than  $cert_{[q, \mathcal{I}]}$

Hence,  $cert_{[q, \mathcal{I}]}$  is called the **perfect rewriting** of  $q$  wrt  $\mathcal{I}$

# Perfect rewriting

What is the relationship between answering by rewriting and certain answers? [Calvanese & al. ICDT'05]:

- When does the (maximal) rewriting compute **all** certain answers?
- What do we gain or loose by focusing on a given class of queries?

Let's try to consider the “**best possible**” rewriting

Define  $cert_{[q, \mathcal{I}]}(\cdot)$  to be the function that, with  $q$  and  $\mathcal{I}$  fixed, given source database  $\mathcal{C}$ , computes the certain answers  $cert(q, \mathcal{I}, \mathcal{C})$ .

- $cert_{[q, \mathcal{I}]}$  can be seen as a query on the alphabet  $\mathcal{A}_{\mathcal{S}}$
- $cert_{[q, \mathcal{I}]}$  is a (sound) rewriting of  $q$  wrt  $\mathcal{I}$
- No sound rewriting exists that is better than  $cert_{[q, \mathcal{I}]}$

Hence,  $cert_{[q, \mathcal{I}]}$  is called the **perfect rewriting** of  $q$  wrt  $\mathcal{I}$

# Perfect rewriting

What is the relationship between answering by rewriting and certain answers? [Calvanese & al. ICDT'05]:

- When does the (maximal) rewriting compute **all** certain answers?
- What do we gain or loose by focusing on a given class of queries?

Let's try to consider the “**best possible**” rewriting

Define  $cert_{[q, \mathcal{I}]}(\cdot)$  to be the function that, with  $q$  and  $\mathcal{I}$  fixed, given source database  $\mathcal{C}$ , computes the certain answers  $cert(q, \mathcal{I}, \mathcal{C})$ .

- $cert_{[q, \mathcal{I}]}$  can be seen as a query on the alphabet  $\mathcal{A}_{\mathcal{S}}$
- $cert_{[q, \mathcal{I}]}$  is a (sound) rewriting of  $q$  wrt  $\mathcal{I}$
- No sound rewriting exists that is better than  $cert_{[q, \mathcal{I}]}$

Hence,  $cert_{[q, \mathcal{I}]}$  is called the **perfect rewriting** of  $q$  wrt  $\mathcal{I}$

# Properties of the perfect rewriting

- Can the perfect rewriting be expressed in a certain query language?
- For a given class of queries, what is the relationship between a maximal rewriting and the perfect rewriting?
  - From a semantical point of view
  - From a computational point of view
- Which is the computational complexity of finding the perfect rewriting, and how big is it?
- Which is the computational complexity of evaluating the perfect rewriting?



## (G)LAV – The case of conjunctive queries

Let  $q$  and the queries in  $\mathcal{M}$  be conjunctive queries (CQs)

Let  $q'$  be the **union of all maximal rewritings of  $q$  for the class of CQs**

Theorem (Levy & al. PODS'95, Abiteboul & Duschka PODS'98)

- $q'$  is the maximal rewriting for the class of unions of conjunctive queries (UCQs)
- $q'$  is the perfect rewriting of  $q$  wrt  $\mathcal{I}$
- $q'$  is a PTIME query
- $q'$  is an exact rewriting (equivalent to  $q$  for each database  $\mathcal{B}$  of  $\mathcal{I}$ ), if an exact rewriting exists

Does this “ideal situation” carry on to cases where  $q$  and  $\mathcal{M}$  allow for union?

## (G)LAV – The case of conjunctive queries

Let  $q$  and the queries in  $\mathcal{M}$  be conjunctive queries (CQs)

Let  $q'$  be the **union of all maximal rewritings of  $q$  for the class of CQs**

Theorem (Levy & al. PODS'95, Abiteboul & Duschka PODS'98)

- $q'$  is the maximal rewriting for the class of unions of conjunctive queries (UCQs)
- $q'$  is the perfect rewriting of  $q$  wrt  $\mathcal{I}$
- $q'$  is a PTIME query
- $q'$  is an exact rewriting (equivalent to  $q$  for each database  $\mathcal{B}$  of  $\mathcal{I}$ ), if an exact rewriting exists

Does this “ideal situation” carry on to cases where  $q$  and  $\mathcal{M}$  allow for union?

















































































































































































































































# Outline

- 5 The role of global integrity constraints
- 6 Query answering in GAV with constraints
  - Incompleteness and inconsistency in GAV systems
  - Query answering in GAV under inclusion dependencies
  - Rewriting CQs under inclusion dependencies in GAV
  - Query answering in GAV under IDs and KDs
  - Query answering in GAV under IDs, KDs, and EDs
- 7 Query answering in LAV with constraints
  - LAV systems and integrity constraints
  - Query answering in (G)LAV under inclusion dependencies
  - Query answering in (G)LAV under IDs and EDs
  - LAV systems and key dependencies
- 8 Inconsistency tolerance







































































































## Part IV

### Concluding remarks

# Outline

## 9 Concluding remarks

# Outline

## 9 Concluding remarks

## Further issues and open problems

- Further forms of constraints, e.g.,
  - KDs with restricted forms of key-conflicting IDs
  - ontology languages, description logics, RDF  
[Calvanese & al. PODS'98, Calvanese & al., KR'06]
- Semistructured data and XML
  - constraints (DTDs, XML Schema, ...)
  - query languages (transitive closure)
- Finite models vs. unrestricted models [Rosati, PODS'06]
- Data exchange and materialization

# Acknowledgements





- Andrea Cali
- Diego Calvanese
- Giuseppe De Giacomo
- Domenico Lembo
- Riccardo Rosati
- Moshe Y. Vardi







# References I

-  S. Abiteboul and O. Duschka.  
Complexity of answering queries using materialized views.  
*In Proc. of PODS'98*, pages 254–265, 1998.
-  A. Calì, D. Calvanese, G. De Giacomo, and M. Lenzerini.  
On the expressive power of data integration systems.  
*In Proc. of ER 2002*, 2002.
-  A. Calì, D. Calvanese, G. De Giacomo, and M. Lenzerini.  
On the role of integrity constraints in data integration.  
*IEEE Bull. on Data Engineering*, 25(3):39–45, 2002.
-  A. Calì, D. Lembo, and R. Rosati.  
Query rewriting and answering under constraints in data integration systems.  
*In Proc. of IJCAI 2003*, pages 16–21, 2003.

# References II

-  D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.  
Data complexity of query answering in description logics.  
*In Proc. of KR 2006*, 2006.
-  D. Calvanese, G. De Giacomo, and M. Lenzerini.  
On the decidability of query containment under constraints.  
*In Proc. of PODS'98*, pages 149–158, 1998.
-  D. Calvanese, G. De Giacomo, M. Lenzerini, and M. Y. Vardi.  
View-based query processing: On the relationship between rewriting, answering and losslessness.  
*In Proc. of ICDT 2005*, volume 3363 of *LNCS*, pages 321–336. Springer, 2005.
-  D. Calvanese and R. Rosati.  
Answering recursive queries under keys and foreign keys is undecidable.  
*In Proc. of KRDB 2003*. CEUR Electronic Workshop Proceedings,  
<http://ceur-ws.org/Vol-79/>, 2003.

# References III

-  O. M. Duschka, M. R. Genesereth, and A. Y. Levy.  
Recursive query plans for data integration.  
*J. of Logic Programming*, 43(1):49–73, 2000.
-  A. Fuxman and R. J. Miller.  
First-order query rewriting for inconsistent databases.  
In *Proc. of ICDT 2005*, volume 3363 of *LNCS*, pages 337–351. Springer, 2005.
-  G. Grahne and A. O. Mendelzon.  
Tableau techniques for querying information sources through global schemas.  
In *Proc. of ICDT'99*, volume 1540 of *LNCS*, pages 332–347. Springer, 1999.
-  A. Y. Halevy.  
Answering queries using views: A survey.  
*VLDB Journal*, 10(4):270–294, 2001.

# References IV



M. Lenzerini.

Data integration: A theoretical perspective.

In *Proc. of PODS 2002*, pages 233–246, 2002.



N. Leone, T. Eiter, W. Faber, M. Fink, G. Gottlob, G. Greco, E. Kalka, G. Ianni, D. Lembo, V. Lio, B. Nowicki, R. Rosati, M. Ruzzi, W. Staniszki, and G. Terracina.

Boosting information integration: The INFOMIX system.

In *Proc. of SEBD 2005*, pages 55–66, 2005.



A. Y. Levy, A. Rajaraman, and J. J. Ordille.

Query answering algorithms for information agents.

In *Proc. of AAAI'96*, pages 40–47, 1996.

# References V



R. Rosati.

On the decidability and finite controllability of query processing in databases with incomplete information.

In *Proc. of PODS 2006*, 2006.