

# Logical foundations for data integration

*Maurizio Lenzerini*

**Dipartimento di Informatica e Sistemistica “Antonio Ruberti”  
Università di Roma “La Sapienza”**

*Joint work with D. Calvanese, G. De Giacomo, D. Lembo, R. Rosati*

**Invited talk at SOFSEM 2005**

*Liptovsky Jan, Slovak Republic – January 27, 2005*

## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- Conclusions

## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- Conclusions

## Three data integration architectures

- **Centralized data integration**

The traditional architecture for centralized, virtual data integration

- **Data exchange**

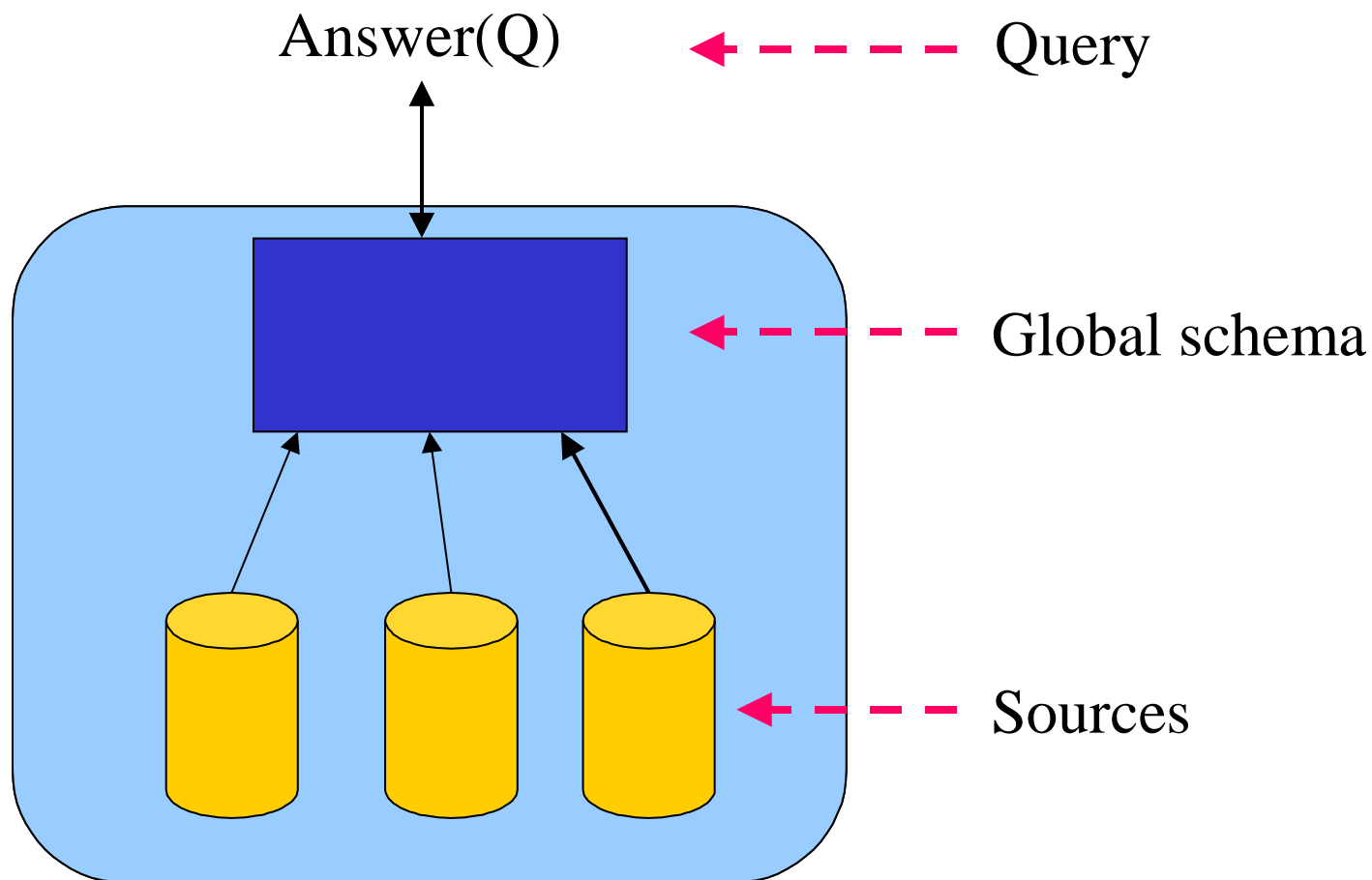
Materialization of data from a source database to a target database

- **Peer-to-peer data integration**

Decentralized, dynamic, data-centric coordination between autonomous organizations

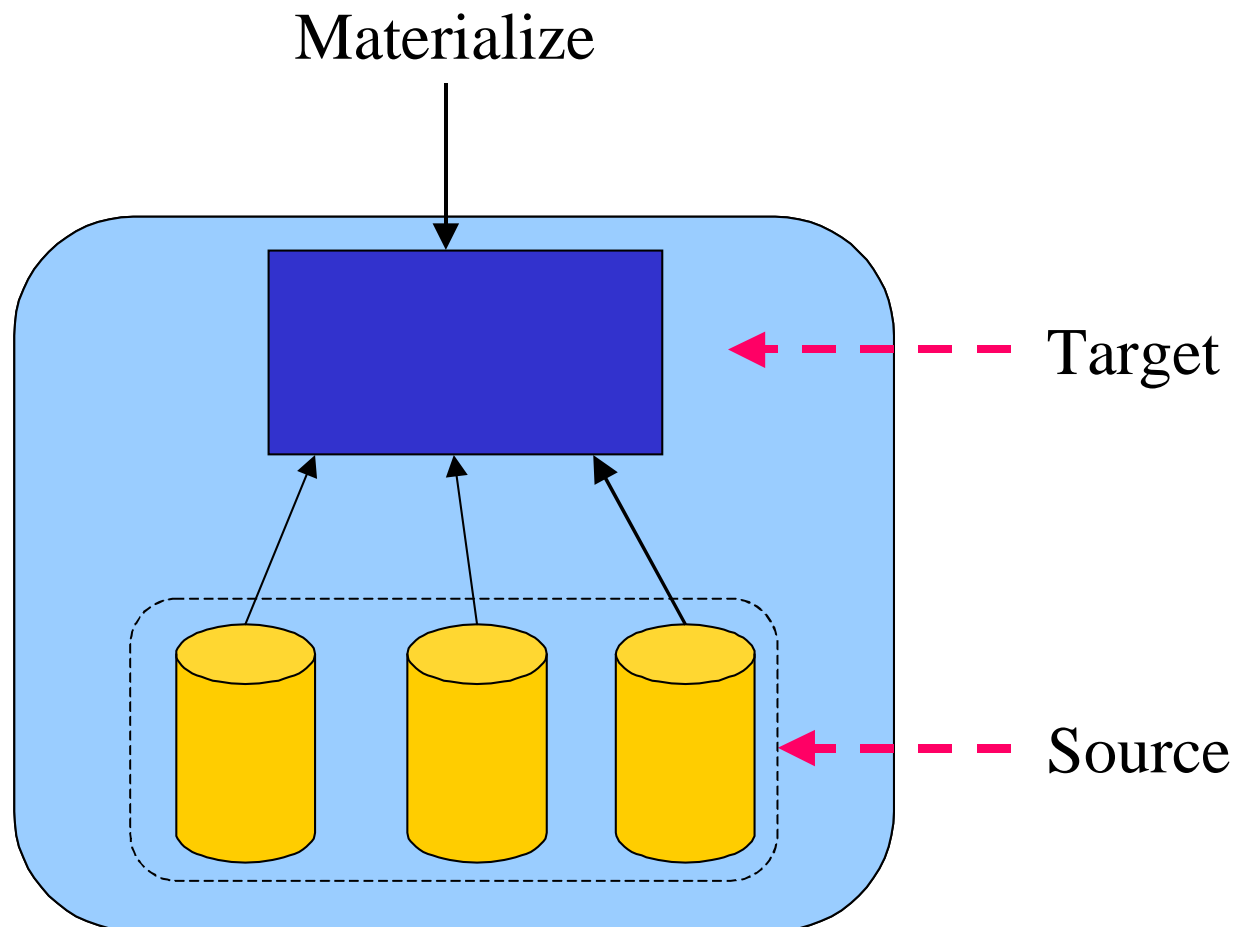
## Centralized data integration

- Mapping between sources and global schema
- Queries over the global schema



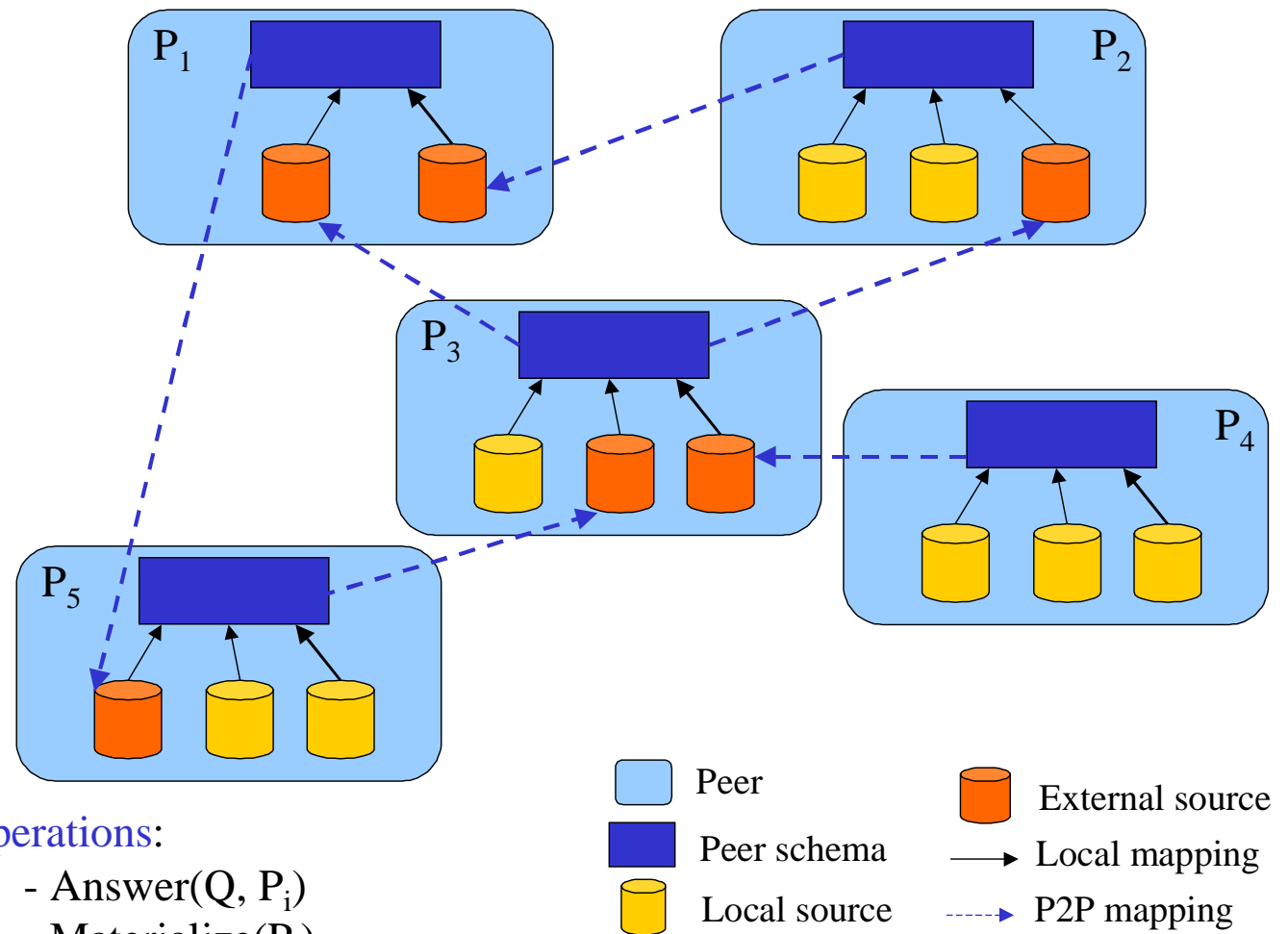
## Data exchange

- Mapping between sources and target schema
- Materialization according to the target schema



# Peer-to-peer data integration

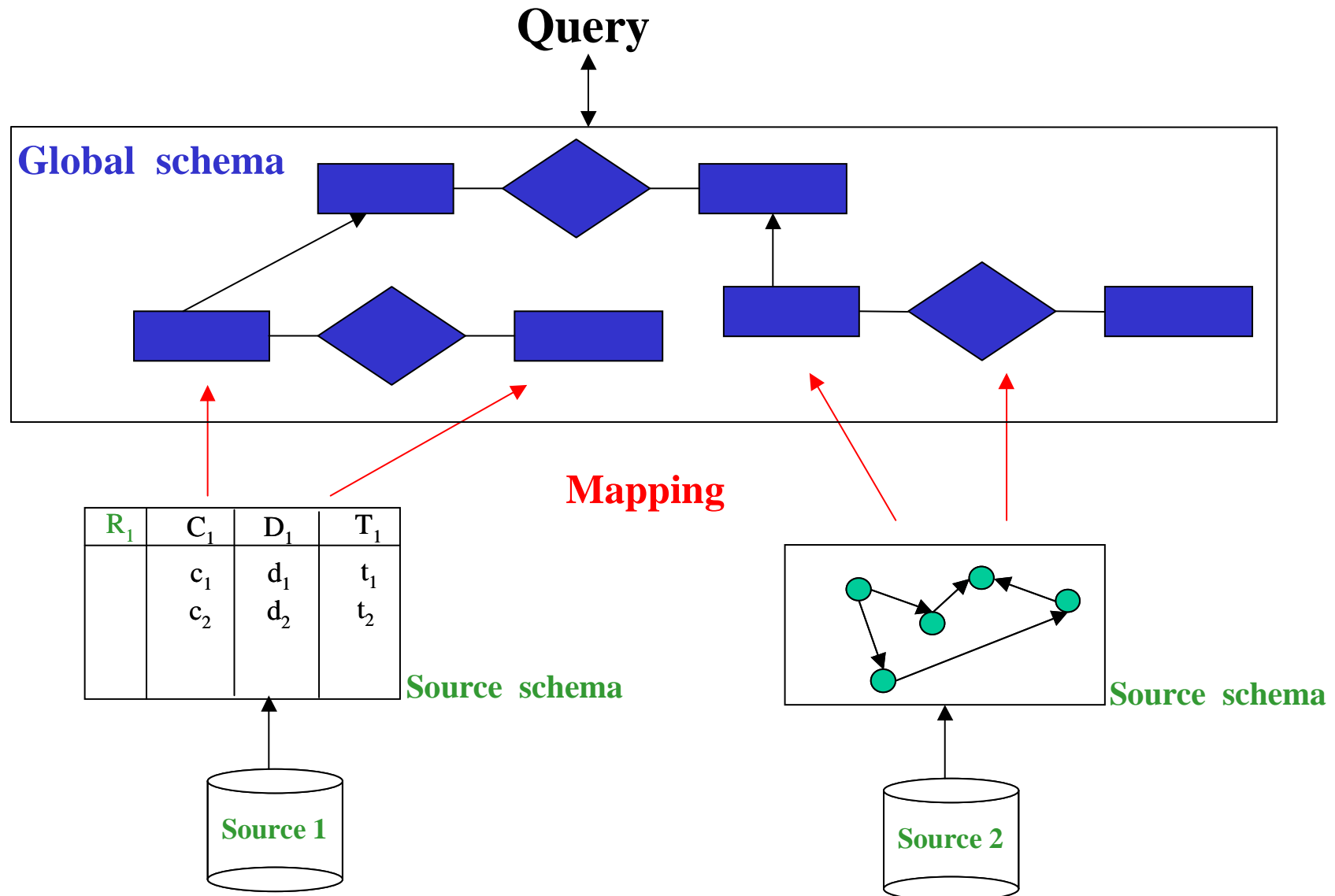
- Several peers
- Local mappings and P2P mappings
- Each query over one peer
- Dynamic mappings



## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- Conclusions

# Centralized data integration



## Formal framework for data integration

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

- $\mathcal{G}$  is the global schema

The global schema is a logical theory over an alphabet  $\mathcal{A}_{\mathcal{G}}$

- $\mathcal{S}$  is the source schema

The source schema is constituted simply by an alphabet  $\mathcal{A}_{\mathcal{S}}$  disjoint from  $\mathcal{A}_{\mathcal{G}}$

- $\mathcal{M}$  is the mapping between  $\mathcal{S}$  and  $\mathcal{G}$

Different approaches to the specification of mapping

## Semantics of a data integration system

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

Let  $\mathcal{D}$  be a source database (also called source model) over a fixed infinite domain  $\Gamma$  of constants, fixing the extension of the predicates of  $\mathcal{A}_S$ .

The set of models of (i.e., databases for  $\mathcal{A}_G$  that satisfy)  $\mathcal{I}$  relative to  $\mathcal{D}$  is:

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

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

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

A **query**  $q$  of arity  $n$  is a formula with  $n$  free variables.

If  $q$  is a query of arity  $n$  posed to a data integration system  $\mathcal{I}$  (i.e., a formula over  $\mathcal{A}_{\mathcal{G}}$  with  $n$  free variables), then the set of **certain answers to  $q$  wrt  $\mathcal{I}$  and  $\mathcal{D}$**  is

$$ans(q, \mathcal{I}, \mathcal{D}) = \{ \langle c_1, \dots, c_n \rangle \in q^{\mathcal{B}} \mid \forall \mathcal{B} \in sem^{\mathcal{D}}(\mathcal{I}) \}$$

Note: query answering is **logical implication**.

Note: complexity will be mainly measured **wrt the size of the source database  $\mathcal{D}$** , and will refer to the problem of deciding whether  $\vec{c} \in ans(q, \mathcal{I}, \mathcal{D})$ , 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).*

## 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:**  $\text{movie}(Title, Year, Director)$

$\text{european}(Director)$

$\text{review}(Title, Critique)$

**Source 1:**  $r_1(Title, Year, Director)$  since 1960, european directors

**Source 2:**  $r_2(Title, Critique)$  since 1990

**Query:** Title and critique of movies in 1998

$\exists D. \text{movie}(T, 1998, D) \wedge \text{review}(T, R)$ , written

$\{ (T, R) \mid \text{movie}(T, 1998, D) \wedge \text{review}(T, R) \}$

## Formalization of LAV

In LAV, the mapping  $\mathcal{M}$  is constituted by a set of assertions:

$$s \rightsquigarrow \phi_{\mathcal{G}} \text{ (sound source)} \quad \forall \vec{x} (s(\vec{x}) \rightarrow \phi_{\mathcal{G}}(\vec{x}))$$

one for each source element  $s$  in  $\mathcal{A}_{\mathcal{S}}$ , where  $\phi_{\mathcal{G}}$  is a query over  $\mathcal{G}$ . Given source database  $\mathcal{C}$ , a database  $\mathcal{B}$  for  $\mathcal{G}$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  if for each  $s \in \mathcal{S}$ :

$$s^{\mathcal{C}} \subseteq \phi_{\mathcal{G}}^{\mathcal{B}} \text{ (sound source)}$$

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:** associated to source relations we have **views** over the global schema

$r_1(T, Y, D) \rightsquigarrow \{ (T, Y, D) \mid \text{movie}(T, Y, D) \wedge \text{european}(D) \wedge Y \geq 1960 \}$   
 $r_2(T, R) \rightsquigarrow \{ (T, R) \mid \text{movie}(T, Y, D) \wedge \text{review}(T, R) \wedge Y \geq 1990 \}$

The query  $\{ (T, R) \mid \text{movie}(T, 1998, D) \wedge \text{review}(T, R) \}$  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:

$$\{ (T, R) \mid r_2(T, R) \wedge r_1(T, 1998, D) \}$$

## Formalization of GAV

In GAV, the mapping  $\mathcal{M}$  is constituted by a set of assertions:

$$g \rightsquigarrow \phi_{\mathcal{S}} \text{ (sound source)} \quad \forall \vec{x} (\phi_{\mathcal{S}}(\vec{x}) \rightarrow g(\vec{x}))$$

one for each element  $g$  in  $\mathcal{A}_{\mathcal{G}}$ , where  $\phi_{\mathcal{S}}$  is a query over  $\mathcal{S}$ . Given source database  $\mathcal{C}$ , a database  $\mathcal{B}$  for  $\mathcal{G}$  satisfies  $\mathcal{M}$  wrt  $\mathcal{C}$  if for each  $g \in \mathcal{G}$ :

$$g^{\mathcal{B}} \supseteq \phi_{\mathcal{S}}^{\mathcal{C}} \text{ (sound source)}$$

Given a source database,  $\mathcal{M}$  provides direct information about which data satisfy the elements of the global schema. Relations in  $\mathcal{G}$  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:**  $\text{movie}(Title, Year, Director)$   
 $\text{european}(Director)$   
 $\text{review}(Title, Critique)$

**GAV:** associated to relations in the global schema we have **views** over the sources

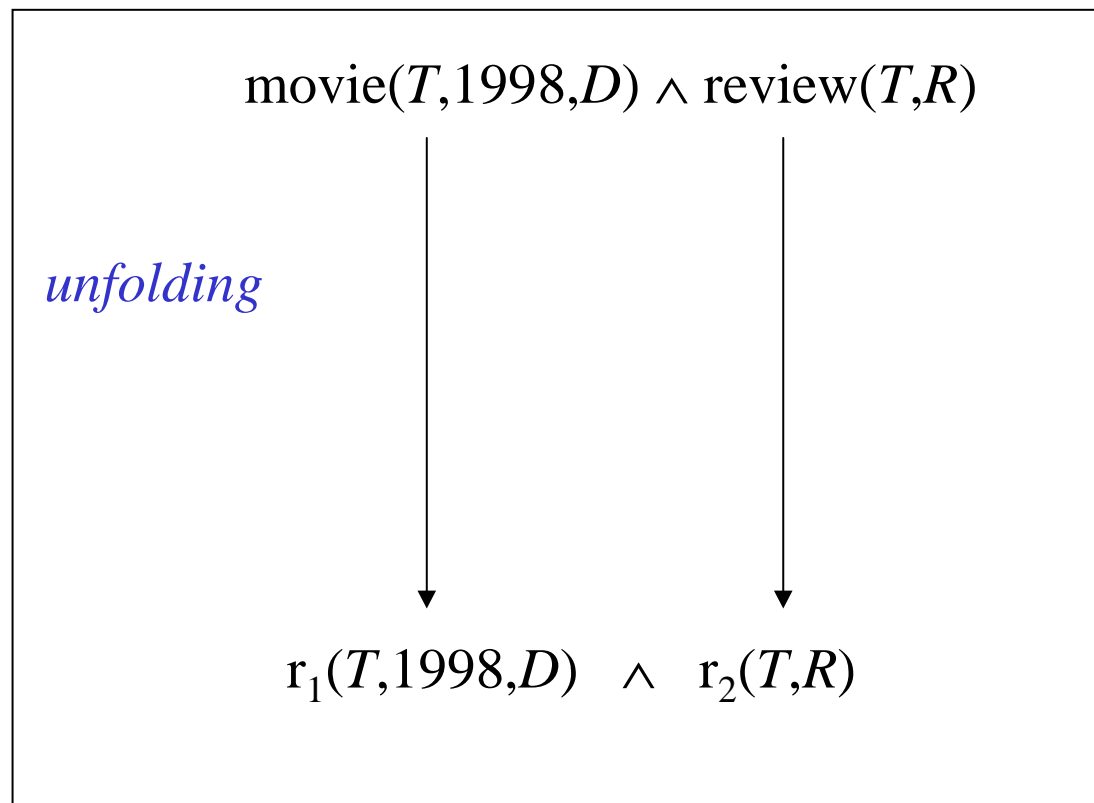
$\text{movie}(T, Y, D) \rightsquigarrow \{ (T, Y, D) \mid r_1(T, Y, D) \}$

$\text{european}(D) \rightsquigarrow \{ (D) \mid r_1(T, Y, D) \}$

$\text{review}(T, R) \rightsquigarrow \{ (T, R) \mid r_2(T, R) \}$

## GAV – example of query processing

The query  $\{ (T, R) \mid \text{movie}(T, 1998, D) \wedge \text{review}(T, R) \}$  is processed by means of unfolding, i.e., by expanding the atoms according to their definitions, so as to come up with source relations. In this case:



## GAV and LAV – comparison

**LAV:** (Information Manifold, DWQ, Picstel)

- 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 reformulation complex)

**GAV:** (Carnot, SIMS, Tsimmis, IBIS, Picstel, 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 reformulation looks easier)

For more details, see [Ullman TCS'00], [Halevy VLDBJ'01], [Lenzerini PODS'02].

## Beyond GAV and LAV: GLAV

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

$$\phi_S \rightsquigarrow \phi_G$$

where

- $\phi_S$  is a **query** over  $\mathcal{S}$ , and
- $\phi_G$  is a **query** over  $\mathcal{G}$  of the arity  $\phi_S$ .

Given source database  $\mathcal{D}$ , a database  $\mathcal{B}$  satisfies  $\mathcal{M}$  wrt  $\mathcal{D}$  if for each assertion in  $\mathcal{M}$  we have  $\phi_S^{\mathcal{D}} \subseteq \phi_G^{\mathcal{B}}$ , i.e., the assertion means  $\forall \vec{x} (\phi_S(\vec{x}) \rightarrow \phi_G(\vec{x}))$ .

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

## Example of GLAV

Global schema:  $Work(Person, Project)$ ,  $Area(Project, Field)$

Source 1:  $HasJob(Person, Field)$

Source 2:  $Teach(Professor, Course)$ ,  $In(Course, Field)$

Source 3:  $Get(Researcher, Grant)$ ,  $For(Grant, Project)$

GLAV mapping:

$$\{ (r, f) \mid HasJob(r, f) \} \quad \rightsquigarrow \quad \{ (r, f) \mid Work(r, p) \wedge Area(p, f) \}$$
$$\{ (r, f) \mid Teach(r, c) \wedge In(c, f) \} \quad \rightsquigarrow \quad \{ (r, f) \mid Work(r, p) \wedge Area(p, f) \}$$
$$\{ (r, p) \mid Get(r, g) \wedge For(g, p) \} \quad \rightsquigarrow \quad \{ (r, p) \mid Work(r, p) \}$$

# Query answering in different approaches

The problem of query answering comes in different forms, depending on several parameters:

- **Global schema**
  - **without** constraints (i.e., empty theory)
  - **with** constraints
- **Mapping**
  - **GAV**
  - **LAV**
  - **GLAV**
- **Class of queries allowed**
  - **client** queries
  - queries in the **mapping**

## Incompleteness and inconsistency

Constraints in $\mathcal{G}$	Type of mapping	Incompleteness	Inconsistency
no	GAV	yes/no	no
no	(G)LAV	yes	no
yes	GAV	yes	yes
yes	(G)LAV	yes	yes

## LAV without constraints: basic technique

Consider conjunctive queries and conjunctive views.

$$r_1(T) \quad \rightsquigarrow \quad \{ (T) \mid \text{movie}(T, Y, D) \wedge \text{european}(D) \}$$

$$r_2(T, V) \quad \rightsquigarrow \quad \{ (T, V) \mid \text{movie}(T, Y, D) \wedge \text{review}(T, V) \}$$

$$Q(X, Y) \quad \leftarrow \quad \text{movie}(X, 1990, D) \wedge \text{review}(X, Y) \wedge \text{european}(D)$$

$$\text{movie}(T, f_1(T), f_2(T)) \quad \leftarrow \quad r_1(T)$$

$$\text{european}(f_2(T)) \quad \leftarrow \quad r_1(T)$$

$$\text{movie}(T, f_4(T, V), f_5(T, V)) \quad \leftarrow \quad r_2(T, V)$$

$$\text{review}(T, V) \quad \leftarrow \quad r_2(T, V)$$

Answering query  $Q$  means evaluating the goal  $Q$  wrt to this nonrecursive logic program, i.e., this logic program is a **perfect rewriting**.

## GAV with constraints: incompleteness and inconsistency

Let us consider a system with a global schema with constraints, and with a GAV mapping  $\mathcal{M}$  with sound sources, whose assertions  $g \rightsquigarrow \phi_S$  have the logical form

$$\forall \vec{x} \phi_S(\vec{x}) \rightarrow g(\vec{x})$$

where  $\phi_S$  is a conjunctive query, and  $g$  is an element of  $\mathcal{G}$ .

Basic observation: since  $\mathcal{G}$  does have constraints, the global database obtained by “applying” the mapping may not be legal for  $\mathcal{G}$ .

## GAV with constraints: example

Global schema  $\mathcal{G}$ :

$\text{student}(Scode, Sname, Scity), \quad \text{key}\{Scode\}$

$\text{university}(Ucode, Uname), \quad \text{key}\{Ucode\}$

$\text{enrolled}(Scode, Ucode), \quad \text{key}\{Scode, Ucode\}$

$\text{enrolled}[Scode] \subseteq \text{student}[Scode]$

$\text{enrolled}[Ucode] \subseteq \text{university}[Ucode]$

**Sources  $\mathcal{S}$ :** database relations  $\mathbf{s}_1(X, Y, Z), \mathbf{s}_2(X, Y), \mathbf{s}_3(X, Y)$

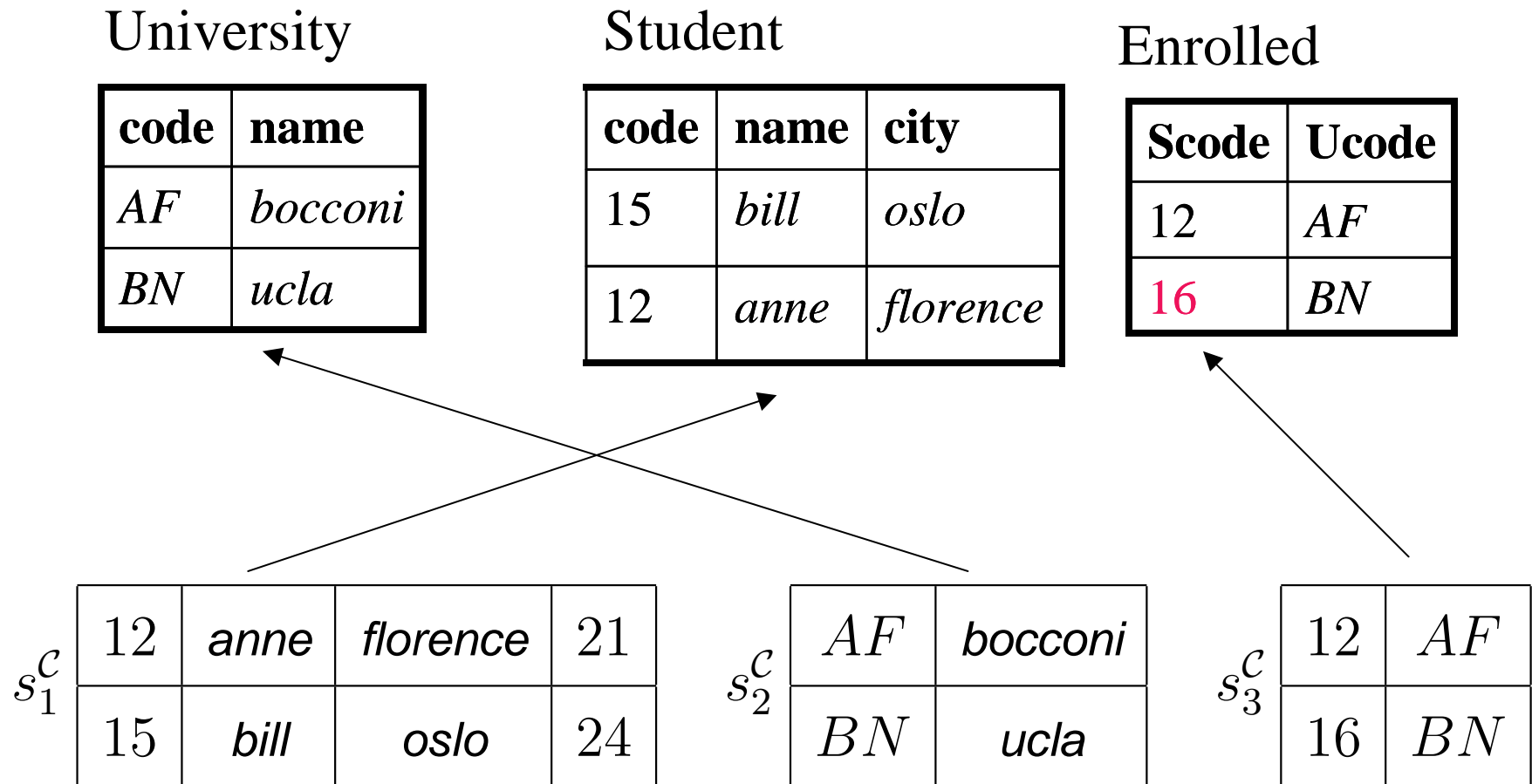
**Mapping  $\mathcal{M}$ :**

$\text{student}(X, Y, Z) \rightsquigarrow \{ (X, Y, Z) \mid \mathbf{s}_1(X, Y, Z, W) \}$

$\text{university}(X, Y) \rightsquigarrow \{ (X, Y) \mid \mathbf{s}_2(X, Y) \}$

$\text{enrolled}(X, Y) \rightsquigarrow \{ (X, Y) \mid \mathbf{s}_3(X, Y) \}$

## GAV with constraints: example



*Example of source database and corresponding retrieved global database*

## GAV with constraints: example of incompleteness

Source database  $\mathcal{C}$ :

 $s_1^{\mathcal{C}}$ 

12	<i>anne</i>	<i>florence</i>	21
15	<i>bill</i>	<i>oslo</i>	24

 $s_2^{\mathcal{C}}$ 

<i>AF</i>	<i>bocconi</i>
<i>BN</i>	<i>ucla</i>

 $s_3^{\mathcal{C}}$ 

12	<i>AF</i>
16	<i>BN</i>

$s_3^{\mathcal{C}}(16, BN)$  and the mapping imply  $\text{enrolled}^{\mathcal{B}}(16, BN)$ , for all  $\mathcal{B} \in \text{sem}^{\mathcal{C}}(\mathcal{I})$ .

Due to the integrity constraints in the global schema, **16 is the code of some student** in all  $\mathcal{B} \in \text{sem}^{\mathcal{C}}(\mathcal{I})$ .

Since  $\mathcal{C}$  says nothing about the name and the city of the student with code 16, we must accept as legal for  $\mathcal{I}$  wrt  $\mathcal{C}$  all virtual global databases that differ in such attributes.

# GAV with constraints: unfolding is not sufficient

**Mapping  $\mathcal{M}$ :**

$\text{student}(X, Y, Z) \rightsquigarrow \{ (X, Y, Z) \mid s_1(X, Y, Z, W) \}$

$\text{university}(X, Y) \rightsquigarrow \{ (X, Y) \mid s_2(X, Y) \}$

$\text{enrolled}(X, Y) \rightsquigarrow \{ (X, Y) \mid s_3(X, Y) \}$

 $s_1^{\mathcal{C}}$ 

12	<i>anne</i>	<i>florence</i>	21
15	<i>bill</i>	<i>oslo</i>	24

 $s_2^{\mathcal{C}}$ 

<i>AF</i>	<i>bocconi</i>
<i>BN</i>	<i>ucla</i>

 $s_3^{\mathcal{C}}$ 

12	<i>AF</i>
16	<i>BN</i>

Query:  $\{ (X) \mid \text{student}(X, Y, Z), \text{enrolled}(X, W) \}$

Unfolding wrt  $\mathcal{M}$ :  $\{ (X) \mid s_1(X, Y, Z, V), s_3(X, W) \}$

retrieves only the answer  $\{12\}$  from  $\mathcal{C}$ , although  $\{12, 16\}$  is the correct answer. The simple unfolding strategy is **not sufficient** in our context.

**Most GAV systems use the simple unfolding strategy!**

## GAV with constraints: example of inconsistency

Source database  $\mathcal{C}$ :

 $s_1^{\mathcal{C}}$ 

12	<i>anne</i>	<i>florence</i>	21
12	<i>bill</i>	<i>oslo</i>	24

 $s_2^{\mathcal{C}}$ 

<i>AF</i>	<i>bocconi</i>
<i>BN</i>	<i>ucla</i>

 $s_3^{\mathcal{C}}$ 

12	<i>AF</i>
16	<i>BN</i>

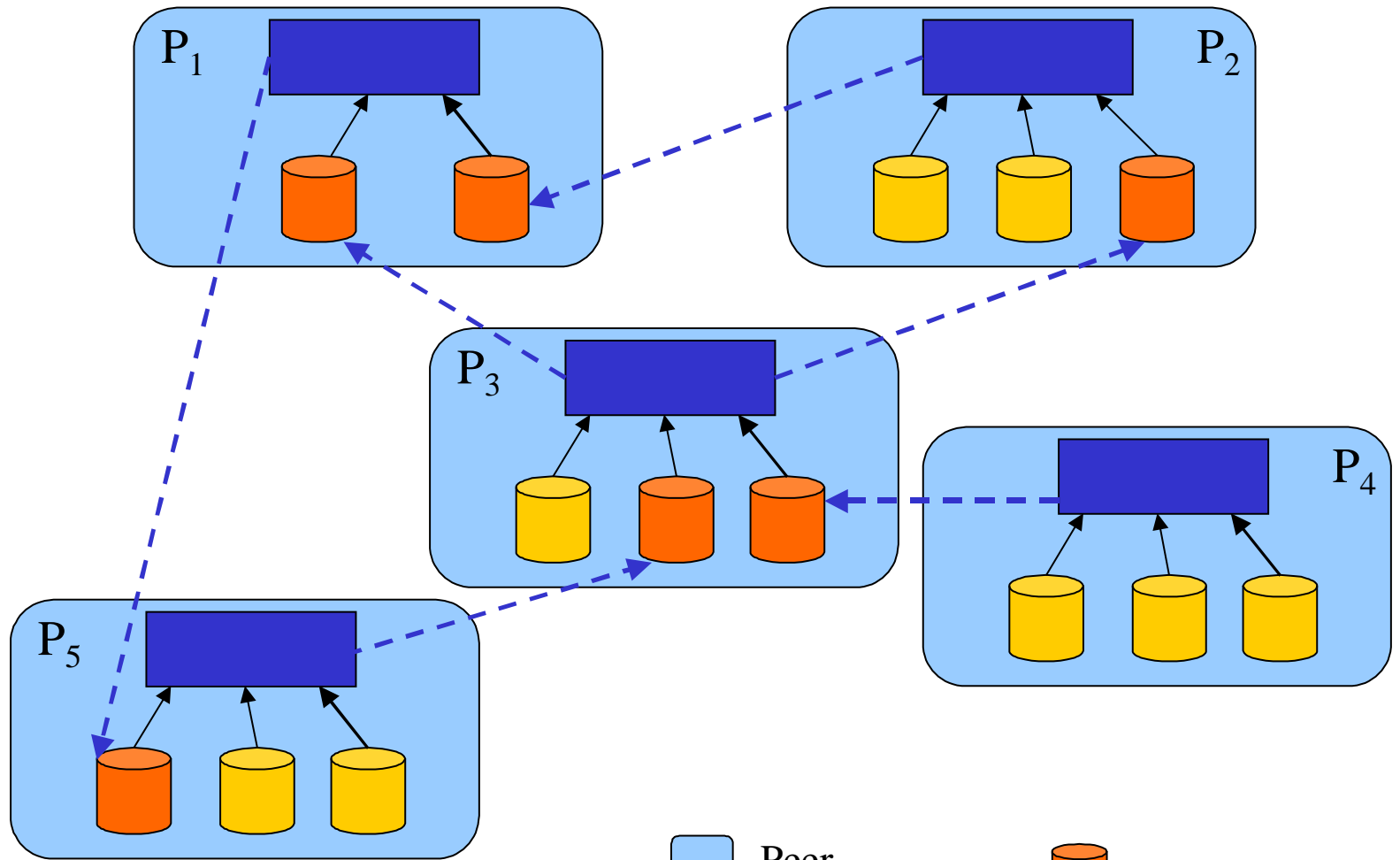
$s_1^{\mathcal{C}}$  imply  $\text{student}^{\mathcal{B}}(12, \textit{anne}, \textit{florence}, 21)$ , and  $\text{student}^{\mathcal{B}}(12, \textit{bill}, \textit{oslo}, 24)$ , for all  $\mathcal{B}$  that satisfies the mapping.

Due to the integrity constraints in the global schema, it follows that there is **no** database that satisfies both the mapping and the global schema, i.e.,  $\text{sem}^{\mathcal{C}}(\mathcal{I}) = \emptyset$ .

## Outline

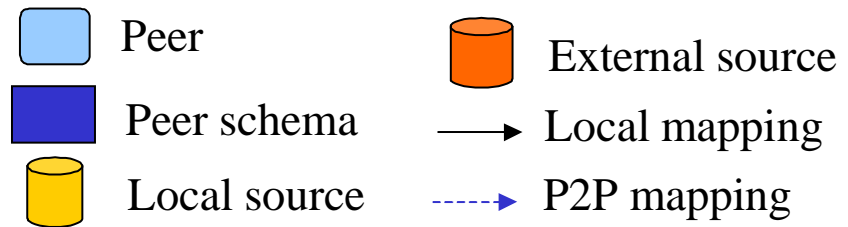
- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- Conclusions

# P2P data integration: the general picture



## Operations:

- Answer(Q, P<sub>i</sub>)
- Materialize(P<sub>i</sub>)



## P2P data integration: what's in a peer

A P2P system  $\Pi$  is a set  $\{P_1, \dots, P_n\}$  of peers, where each peer  $P_i = (G, S, L, M)$  models an autonomous information site, that

- exports its information content in terms of a peer schema  $G$ , whose alphabet is considered disjoint from that of the other peers
- represents its data as a set of sources  $S$  (local sources model its own data, and external sources model data coming from other peers)
- relates sources to global schema by means of local mappings  $L$
- is related to other peers in  $\Pi$  by means of a set of P2P mappings  $M$ , where each P2P mapping is a schema level assertion relating data of another peer  $P_j$  to one external source in  $P_i$

Inspired by [Catarci&Lenzerini COOPIS '92], [Halevy&al. ICDE'03]. Related work: [Ghidini&Serafini FCS '98], [Bernstein&al. WebDB '02], [Franconi&al. P2PDBIS '03].

## P2P data integration: local and P2P mappings

In a peer  $P_i = (G, S, L, M)$

- each local mapping in  $L$  has the form

$$ep_S \rightsquigarrow cq_G$$

where  $ep_S$  is an extraction program on the sources  $S$  and  $cq_G$  is a conjunctive queries over  $G$ , respectively

- each P2P mapping assertion in  $M$  has the form

$$cq \rightsquigarrow s$$

where:

- $cq$  is a conjunctive query over one of the other peers in  $\Pi$
- $s$  is an external source of the peer  $P$
- $cq$  and  $s$  are of the same arity

## Extraction programs

- The notion of **extraction program** aims at modeling computations done in order to
  - **extract**
  - **clean**
  - **transform**
  - **reconcile**

data coming from (local and external) data sources

- We assume that, given the extensions of the sources, an extraction program extracts a set of tuples (of the same arity as the arity of the program)
- We do **not** deal with extraction programs, but we point out that they are accomodated in the framework

## Importance of the semantics

The client sees the whole collection of peers through the eyes of one peer, and (s)he conceives the distributed information system as a unique database

- What does this database provide to the client?
- Can the client trust the answers to queries computed by system?
- Can we prove that it is sound and/or complete in some sense?

No answers to these questions without semantics!

## Semantics of one peer

For each peer  $P = (G, S, L, M)$  we define a FOL theory  $T_P$  as follows:

- The **alphabet** of  $T_P$  is obtained as union of the alphabets of the schema  $G$  and of the sources  $S$
- The **formulas** of  $T_P$  are as follows:
  - all FOL formulas in the schema  $G$
  - for each local mapping assertion  $\{\vec{x} \mid ep_S(\vec{x})\} \rightsquigarrow \{\vec{x} \mid \exists \vec{z} \varphi_G(\vec{x}, \vec{z})\}$  in  $L$ , one formula of the form

$$\forall \vec{x} (ep_S(\vec{x}) \supset \exists \vec{z} \varphi_G(\vec{x}, \vec{z}))$$

Notice that  $T_P$  does not consider the P2P mappings in  $M$

It follows that we are modeling each peer  $P$  as a **GLAV data integration system**, in turn modeled as a FOL theory  $T_P$  (ignoring the P2P mappings  $M$ )

## Semantics of a P2P system

- A **source database**  $\mathcal{D}$  for  $\Pi$  is the disjoint union of a set of local source databases, one local source database for each peer  $P_i$  in  $\Pi$
- Given a source database  $\mathcal{D}$  for  $\Pi$ , the **set of models of  $\Pi$  relative to  $\mathcal{D}$**  is:

$$sem^{\mathcal{D}}(\Pi) = \left\{ \mathcal{I} \mid \begin{array}{l} \mathcal{I} \text{ is a model of all peer theories } T_{P_i} \text{ relative to } \mathcal{D}, \text{ and} \\ \mathcal{I} \text{ satisfies all P2P mapping assertions} \end{array} \right\}$$

The meaning of  $\mathcal{I}$  satisfying a P2P mapping assertion may vary in the various approaches

- Given a **query**  $Q$  of arity  $n$  posed to a peer  $P_i$  of  $\Pi$ , and a source database  $\mathcal{D}$ , the **certain answers to  $Q$  relative to  $\mathcal{D}$**  are

$$ans(Q, \Pi, \mathcal{D}) = \left\{ \vec{t} \in \Gamma^n \mid \vec{t} \in Q^{\mathcal{I}}, \text{ for every } \mathcal{I} \in sem^{\mathcal{D}}(\Pi) \right\}$$

## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- **Hyper: epistemic semantics for P2P data integration**
- Dealing with inconsistencies in Hyper
- Conclusions

## Possible formalizations of P2P mappings

We consider two alternatives for specifying the semantics of P2P mappings:

- **Based on First-Order Logic**

P2P mappings are considered as material implications in logic

- **Based on Epistemic Logic**

P2P mappings are considered as specifications of exchange of certain answers

## First-Order Logic semantics of P2P mappings

The semantics of P2P mapping assertions is usually given in terms of **First-Order Logic**, e.g., [Halevy&al. ICDE'03]

An interpretation  $\mathcal{I}$  satisfies a P2P mapping assertion

$$\{\vec{x} \mid \exists \vec{y} \varphi(\vec{x}, \vec{y})\} \rightsquigarrow s(\vec{x})$$

if it satisfies the FOL formula

$$\forall \vec{x} (\exists \vec{y} \varphi(\vec{x}, \vec{y}) \equiv s(\vec{x}))$$

which is equivalent to the condition

$$\{\vec{x} \mid \exists \vec{y} \varphi_1(\vec{x}, \vec{y})\}^{\mathcal{I}} = (s(\vec{x}))^{\mathcal{I}}$$

## Inadequacy of FOL semantics of P2P mappings

The FOL semantics is not adequate for P2P data integration:

- **Lack of modularity**

- the system is modeled by a flat FOL theory, with no formal separation between the various peers
- the modular structure of the system is not reflected in the semantics

- **Bad computational properties**

Computing the set of certain answers to a conjunctive query  $Q$  posed to a peer is **undecidable**, even when all peer schemas are empty [Halevy&al. ICDE'03], [Koch FOIKS'02]

- **Lack of generality**

To recover decidability, one has to limit the expressive power of P2P mappings (e.g., assume acyclicity) [Halevy&al. ICDE'03]

## Epistemic semantics for P2P mappings: objectives

A new semantics for P2P mappings, with the following aims:

- Peers in our context are to be considered **autonomous sites** that exchange information
- We do not want to limit a-priori the **topology** of the mapping assertions among the peers in the system
- Defining a setting where query answering is decidable, and possibly, **polynomially tractable**

## Epistemic semantics for P2P mappings: basic idea

The new semantics is based on multi-agent epistemic logic [Reiter TARK'88], [Fagin&al 1995]

- A P2P mapping  $cq_i \rightsquigarrow s_j$  (with  $cq_i$  over  $P_i$  and  $s_j$  external source of  $P_j$ ) is interpreted as an epistemic formula which imposes that **only the certain answers** to  $cq_i$  in  $P_i$  (i.e., the facts that are **known** by  $P_i$ ) are transferred to  $P_j$  as facts **known** to satisfy  $s_j$ .

In other words, peer  $P_i$  communicates to peer  $P_j$  only facts that are known, i.e., true in every model of the P2P system

- The modular structure of the system is now reflected in the semantics (by virtue of the modal semantics of epistemic logics)
- Good computational properties: computing the certain answers to a conjunctive query  $Q$  relative to a source database  $\mathcal{D}$  is **polynomial time** in the size of  $\mathcal{D}$ , even for cyclic mappings

## Epistemic logic: basic notions

In multi-agent epistemic logic we want to model the **knowledge of  $n$  agents**

From a language point of views, we have **a new form of atom**, namely ( $\varphi$  is again a formula, and  $i \in \{1, \dots, n\}$ ):

$$\mathbf{K}_i \varphi$$

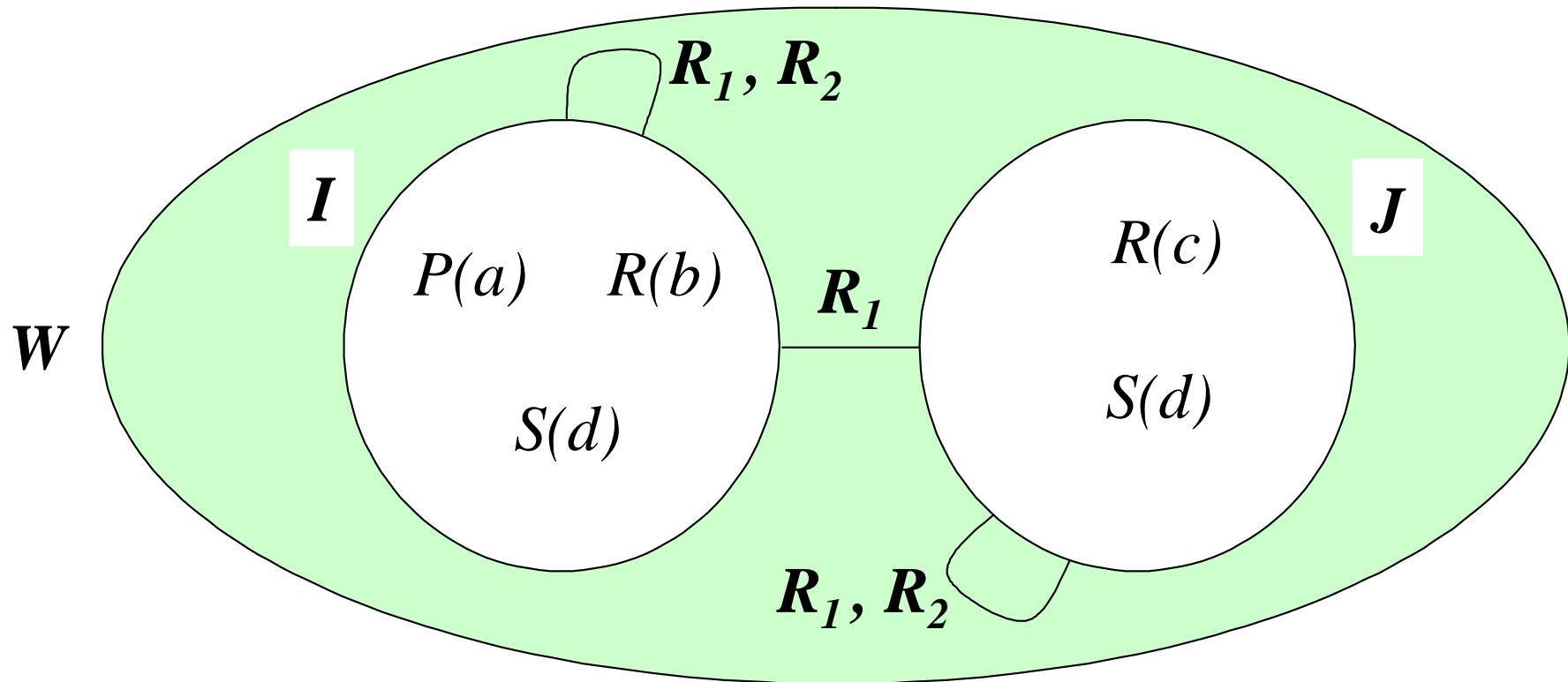
An **epistemic interpretation** is a tuple  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$ , where

- $\mathcal{W}$  is a set of FOL interpretations and  $\mathcal{I} \in \mathcal{W}$
- each  $\mathcal{R}_i$  is a binary relation over  $\mathcal{W}$

We use the  **$S5_n$  modal system**, and therefore, we also impose that

- **each  $\mathcal{R}_i$  is an equivalence relation**, i.e., it is reflexive, symmetric, and transitive

## Epistemic interpretation: example



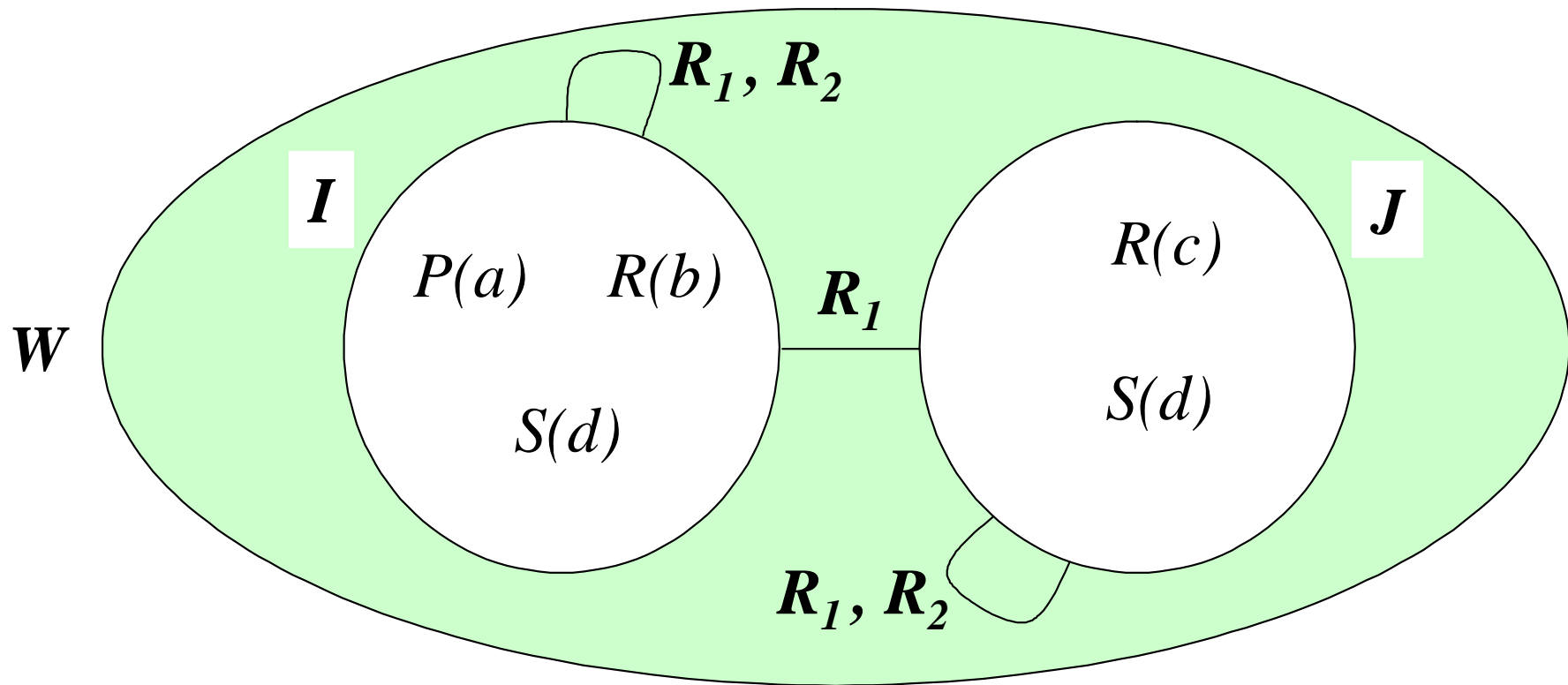
$\langle I, W, \mathcal{R}_1, \mathcal{R}_2 \rangle$  and  $\langle J, W, \mathcal{R}_1, \mathcal{R}_2 \rangle$  are **epistemic interpretations**

## Epistemic logic: basic notions

- a FOL formula constituted by an atom  $a(\vec{x})$  is satisfied in  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  by the tuples  $\vec{t}$  of constants such that  $a(\vec{t})$  is true in  $\mathcal{I}$
- an atom of the form  $\mathbf{K}_i \varphi(\vec{x})$  is satisfied in  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  by the tuples  $\vec{t}$  of constants such that  $\varphi(\vec{t})$  is satisfied in all epistemic interpretations  $\langle \mathcal{J}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  with  $(\mathcal{I}, \mathcal{J}) \in \mathcal{R}_i$

An **epistemic model** of an epistemic logic theory  $\{\varphi_1, \dots, \varphi_t\}$  is an epistemic interpretation that satisfies every formula  $\varphi_i \in \{\varphi_1, \dots, \varphi_t\}$ .

# Epistemic logic: example 1

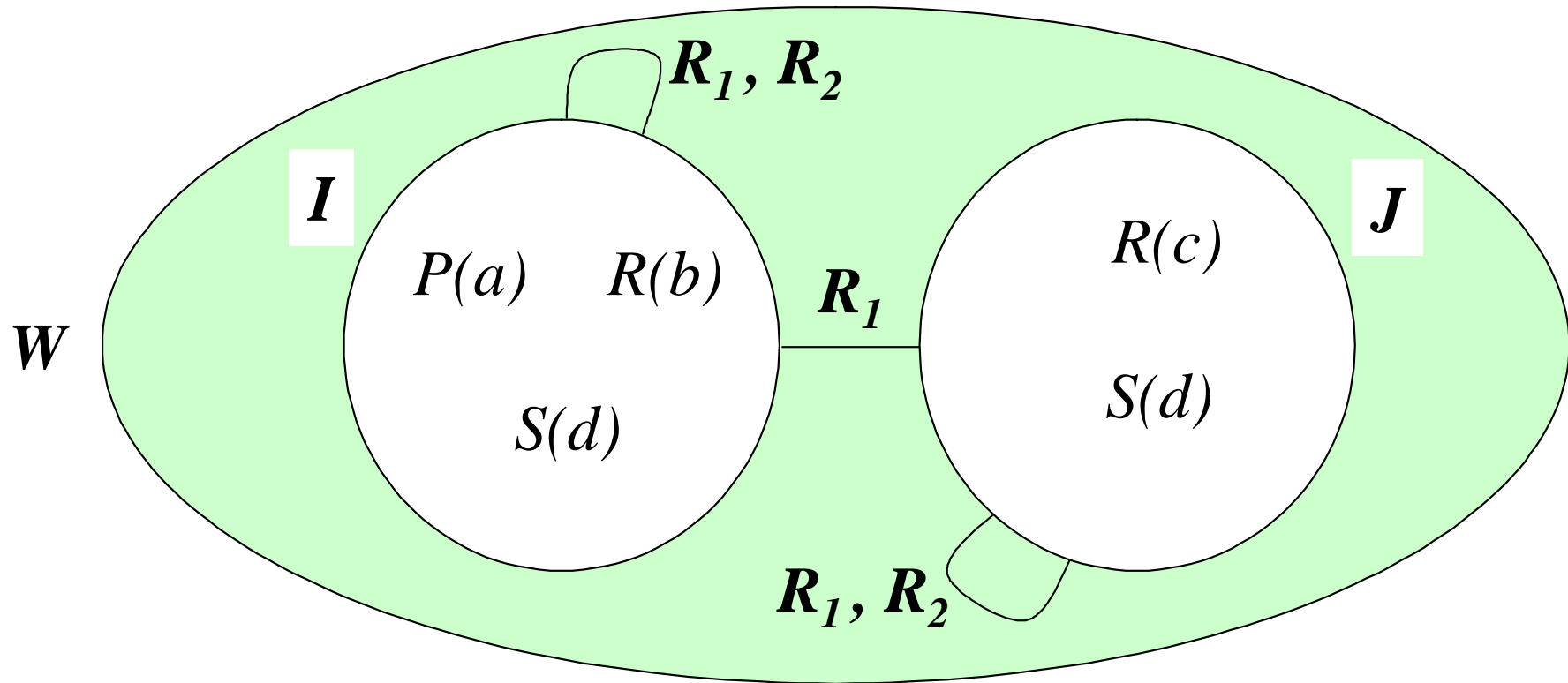


$$\langle I, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \models P(a)$$

$$\langle J, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \not\models P(a)$$

$$\langle I, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \not\models \mathbf{K}_1 P(a)$$

## Epistemic logic: example 2

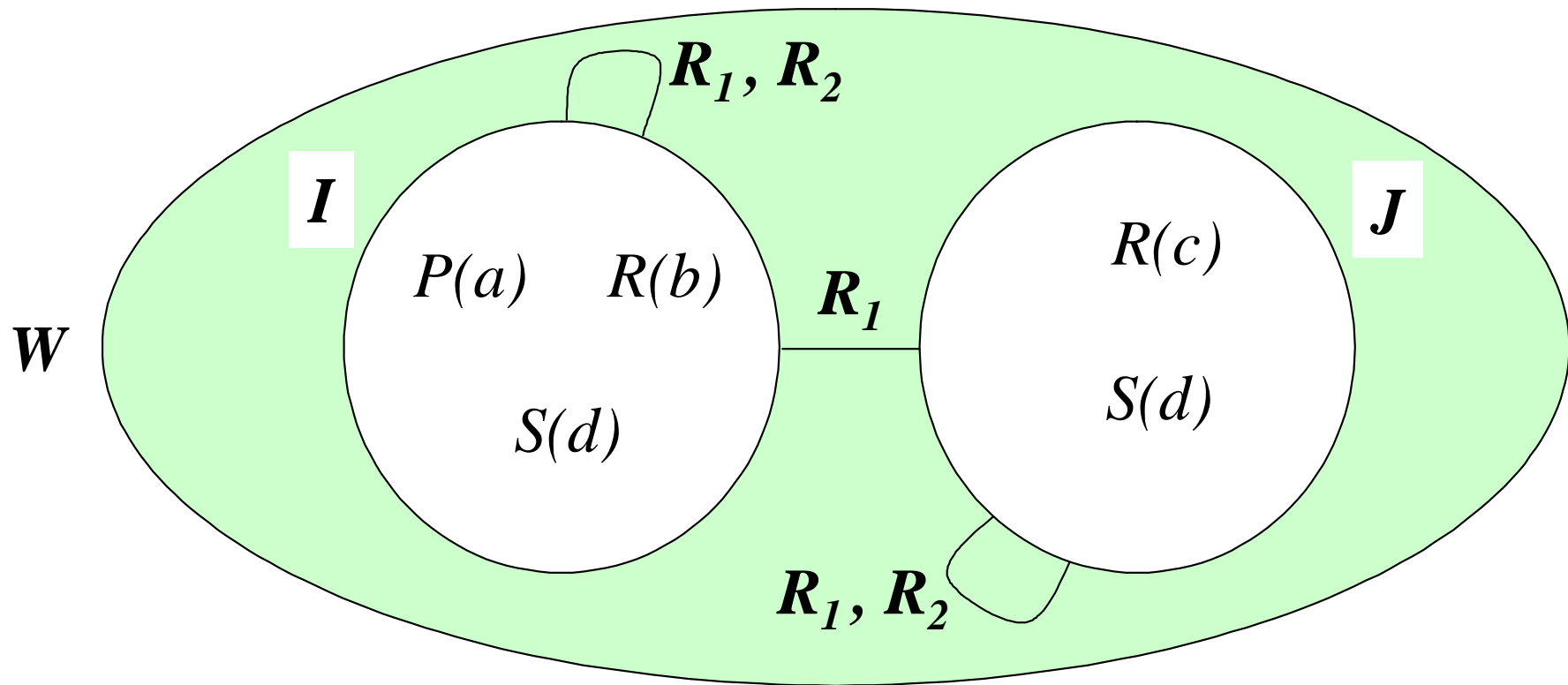


$$\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \mathcal{R}_2 \rangle \models \mathbf{K}_1 (R(b) \vee R(c))$$

$$\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \mathcal{R}_2 \rangle \not\models (\mathbf{K}_1 R(b)) \vee (\mathbf{K}_1 R(c))$$

$$\langle \mathcal{J}, \mathcal{W}, \mathcal{R}_1, \mathcal{R}_2 \rangle \models \mathbf{K}_2 \mathbf{K}_1 S(d)$$

## Epistemic logic: example 3



$$\langle I, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \models \mathbf{K}_1 (\exists x R(x))$$

$$\langle I, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \not\models \exists x (\mathbf{K}_1 R(x))$$

$$\langle J, W, \mathcal{R}_1, \mathcal{R}_2 \rangle \models \exists x (\mathbf{K}_2 S(x))$$

## Epistemic semantics for P2P mappings: basic idea

We formalize a P2P system  $\Pi$  in terms of the epistemic logic theory  $E_\Pi$ :

- the alphabet  $\mathcal{A}_\Pi$  is the disjoint union of the alphabets of the various peer theories  $T_P$ , one for each peer  $P$  in  $\Pi$
- For each peer theory  $T_{P_i}$ , for each formula  $\varphi \in T_{P_i}$ , the following formula is in  $E_\Pi$

$$\mathbf{K}_i \varphi$$

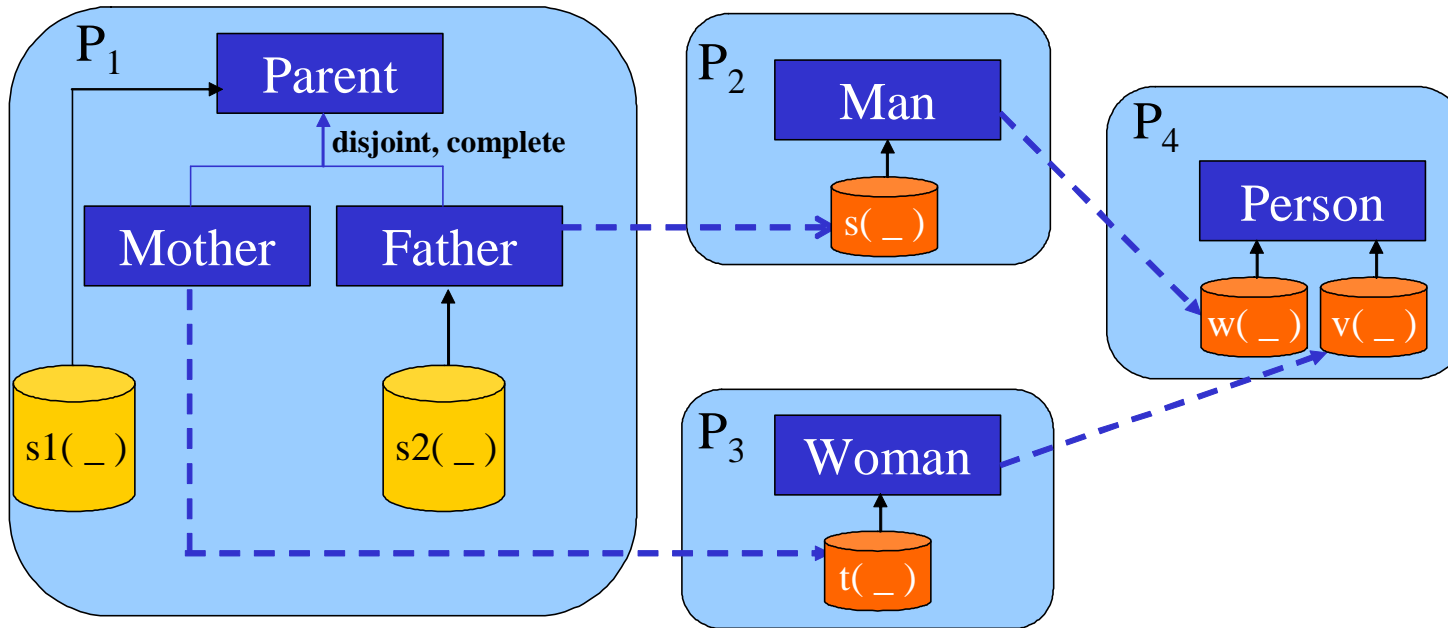
- for each P2P mapping assertion from peer  $P_i$  to peer  $P_j$

$$\{\vec{x} \mid \exists \vec{y} \varphi(\vec{x}, \vec{y})\}_i \rightsquigarrow \{\vec{x} \mid s(\vec{x})\}_j$$

in the peers of  $\Pi$ , the following formula is in  $E_\Pi$

$$\forall \vec{x} ((\mathbf{K}_i \exists \vec{y} \varphi_1(\vec{x}, \vec{y})) \equiv \mathbf{K}_j s(\vec{x}))$$

# Epistemic semantics for P2P mappings: example of theory



$\mathbf{K}_1 \forall x((Parent(x) \equiv Mother(x) \vee Father(x)) \wedge (Mother(x) \supset \neg Father(x)))$

$\mathbf{K}_1 \forall x(s1(x) \supset Parent(x))$

$\mathbf{K}_1 \forall x(s2(x) \supset Father(x))$

$\mathbf{K}_2 \forall x(s(x) \supset Man(x))$

$\mathbf{K}_3 \forall x(t(x) \supset Woman(x))$

$\mathbf{K}_4 \forall x(w(x) \supset Person(x))$

$\mathbf{K}_4 \forall x(v(x) \supset Person(x))$

$\forall x(\mathbf{K}_1 Mother(x) \equiv \mathbf{K}_3 t(x))$

$\forall x(\mathbf{K}_1 Father(x) \equiv \mathbf{K}_2 s(x))$

$\forall x(\mathbf{K}_2 Man(x) \equiv \mathbf{K}_4 w(x))$

$\forall x(\mathbf{K}_3 Woman(x) \equiv \mathbf{K}_4 v(x))$

## Epistemic semantics for P2P mappings: basic idea

An **epistemic model of  $\Pi$  relative to the source database  $\mathcal{D}$**  is an epistemic model  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  of  $E_\Pi$  such that, in every  $\mathcal{J} \in \mathcal{W}$ , the extension of all local sources in all peers is the one sanctioned by  $\mathcal{D}$ .

If  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  is an epistemic model of  $\Pi$  relative to  $\mathcal{D}$ , then

$\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  satisfies all formulas of  $E_\Pi$ , and therefore

- for each  $i$ ,  $\mathcal{I}$  is a model of  $T_{P_i}$  relative to  $\mathcal{D}$ , which means that  $\mathcal{I}$  satisfies all constraints in the various peer schemas, and
- $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  satisfies all formulas corresponding to the P2P mapping assertions in the peers of  $\Pi$

## Epistemic semantics for P2P mappings: basic idea

Note that  $\langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle$  satisfying the P2P mapping assertion  $cq_h \rightsquigarrow s_k$  means that

for every tuple  $\vec{t}$  of constants in  $\Gamma$ ,

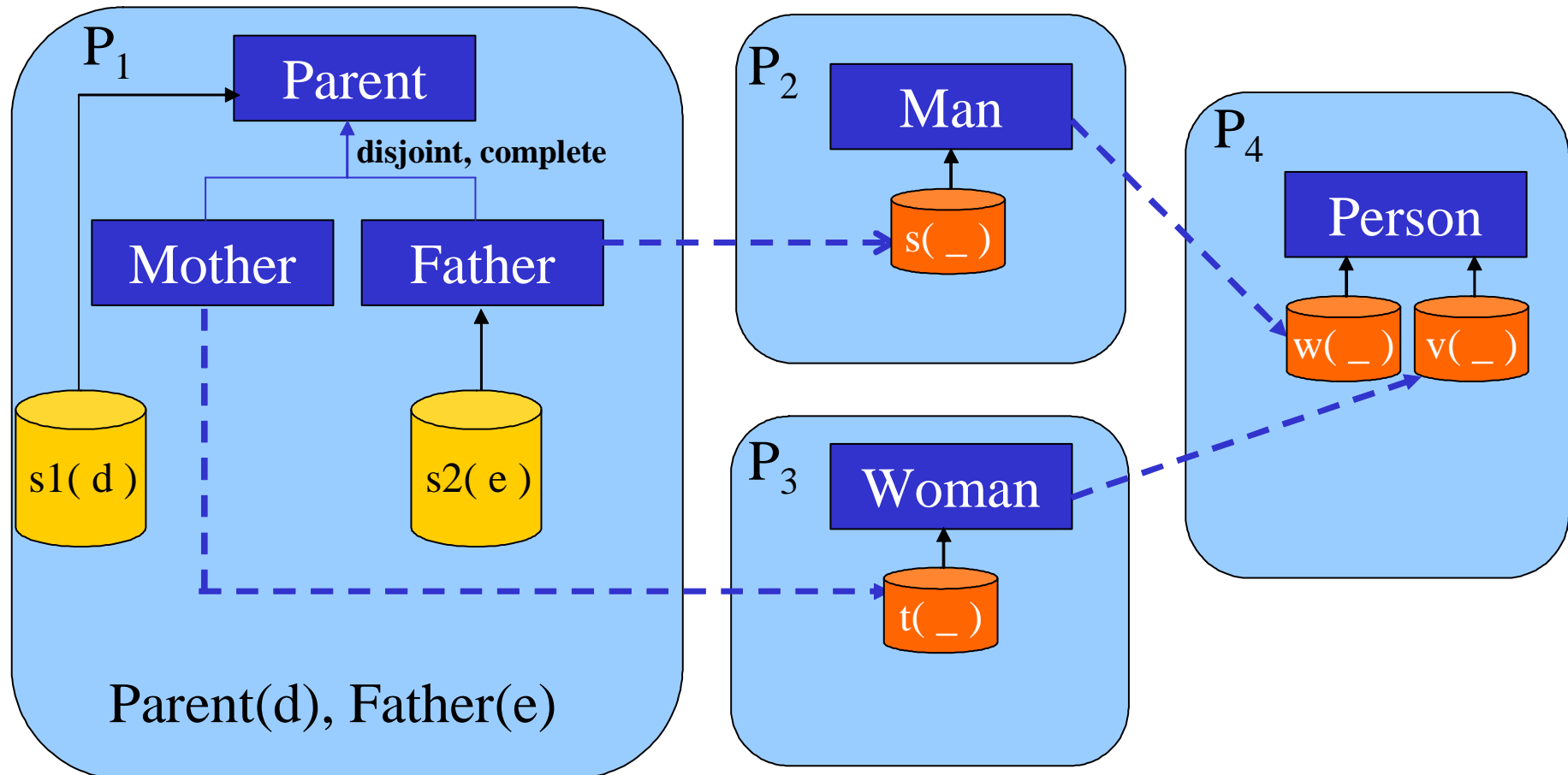
if  $\vec{t} \in cq^{\mathcal{J}}$  for every FOL model  $\mathcal{J}$  such that  $(\mathcal{I}, \mathcal{J}) \in \mathcal{R}_h$ ,

then  $\vec{t} \in s^{\mathcal{J}}$  for every FOL model  $\mathcal{J}$  such that  $(\mathcal{I}, \mathcal{J}) \in \mathcal{R}_k$

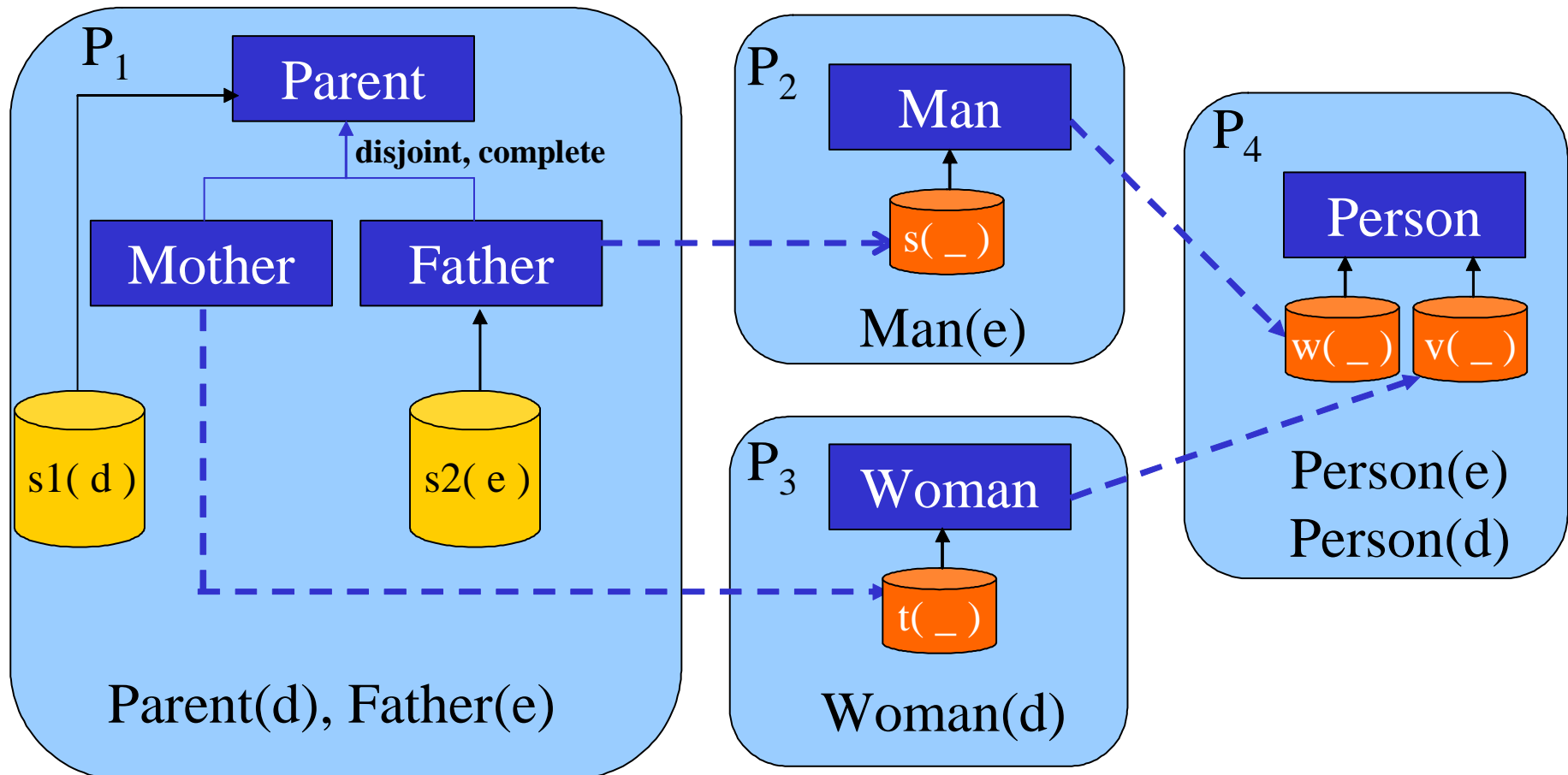
Given a query  $Q$  of arity  $n$  posed to a peer  $P_i$  of  $\Pi$ , and a source database  $\mathcal{D}$ , the certain answers to  $Q$  relative to  $\mathcal{D}$  under the epistemic semantics are

$$ans_k(Q, \Pi, \mathcal{D}) = \{ \vec{t} \in \Gamma^n \mid \vec{t} \in Q^{\mathcal{I}}, \text{ for every epistemic model } \langle \mathcal{I}, \mathcal{W}, \mathcal{R}_1, \dots, \mathcal{R}_n \rangle \text{ of } \Pi \text{ relative to } \mathcal{D} \}$$

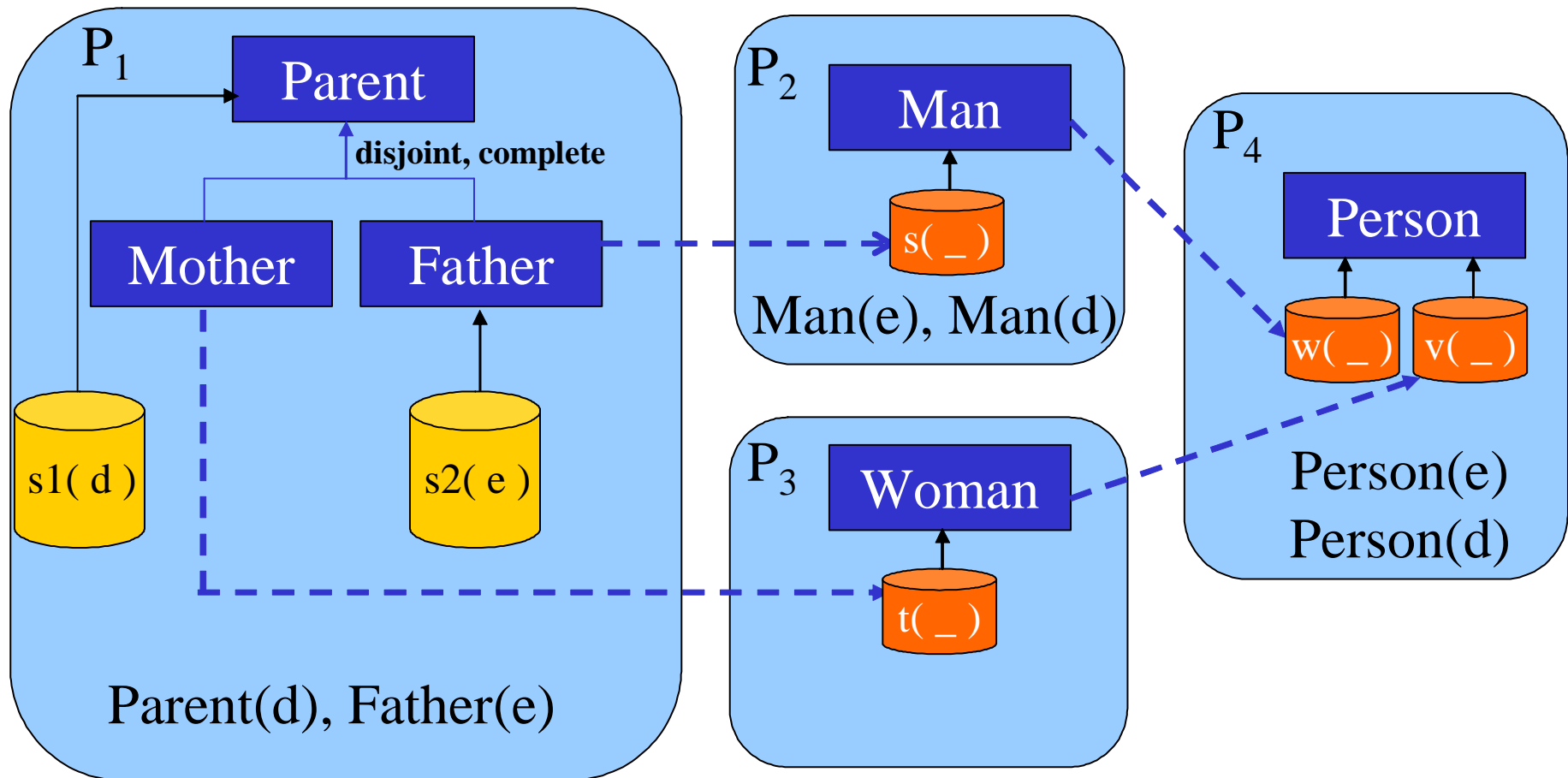
# Semantics of P2P mappings: example



# FOL semantics of P2P mappings: model 1

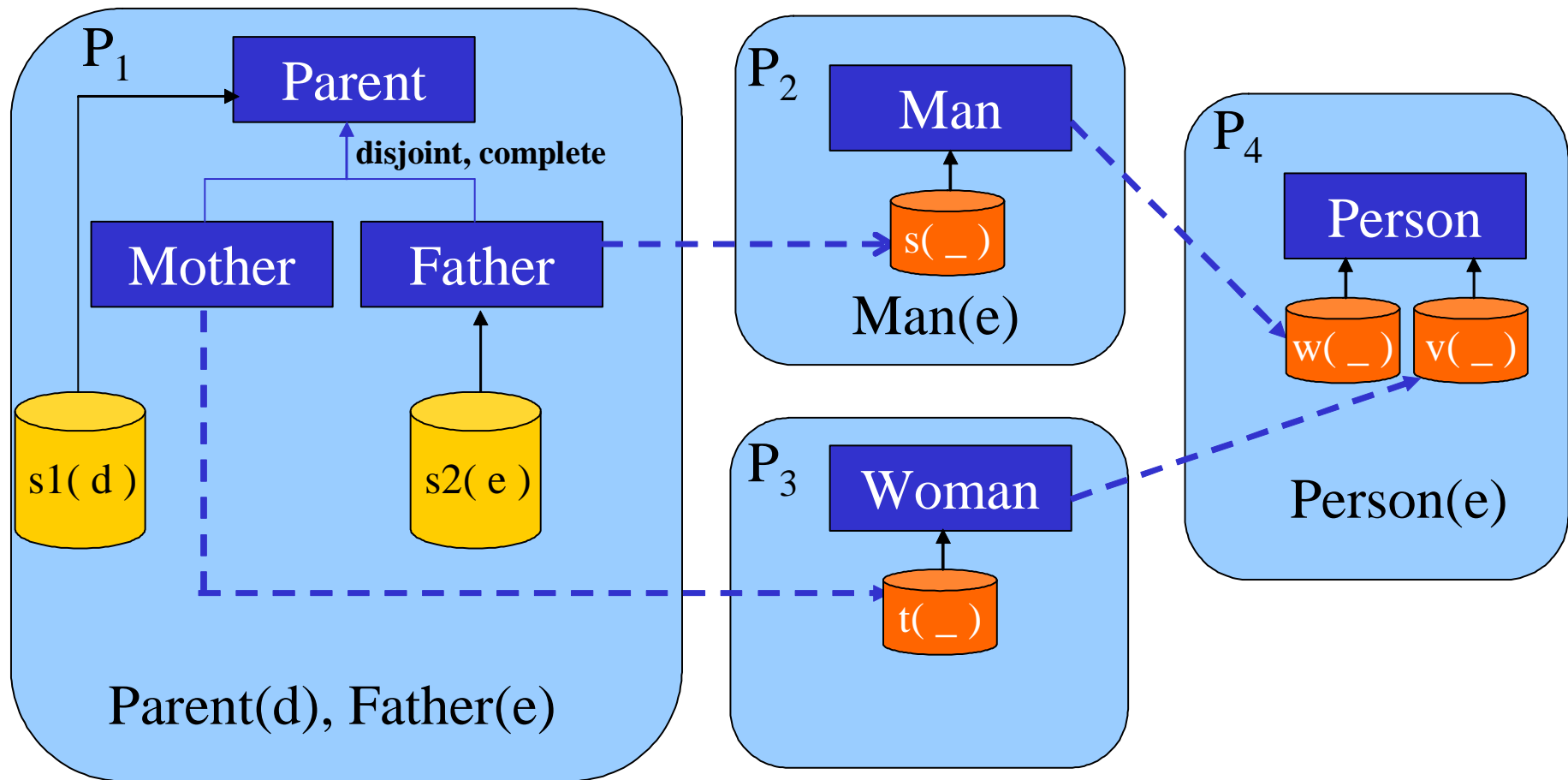


## FOL semantics of P2P mappings: model 2



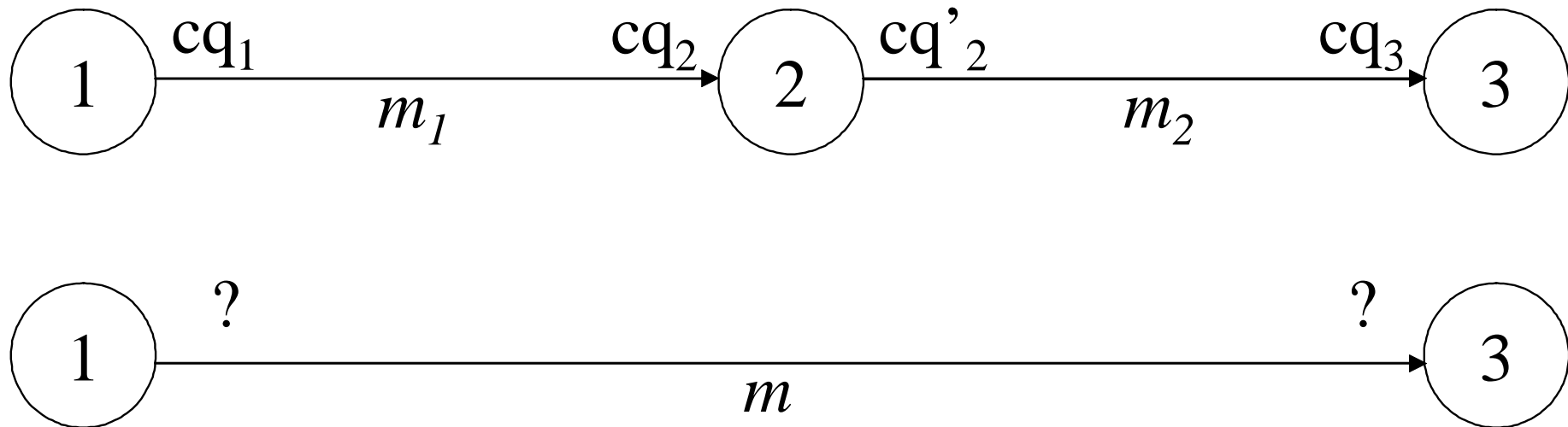
According to the FOL semantics,  $\text{Person}(d)$  is true in all cases, and therefore **is a** certain answer to  $\{x \mid \text{Person}(x)\}$

# Epistemic semantics of P2P mappings



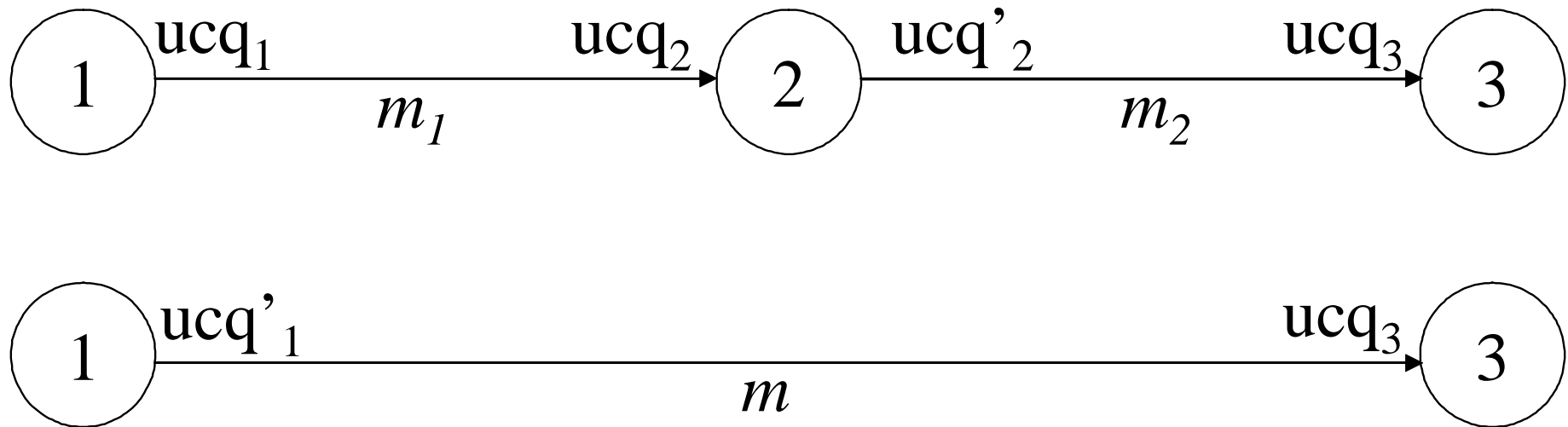
According to the epistemic semantics,  $\text{Person}(d)$  is **not a** certain answer to  $\{x \mid \text{Person}(x)\}$ , since neither  $\mathbf{K}_2 \text{Man}(d)$  nor  $\mathbf{K}_3 \text{Woman}(d)$  is valid, and therefore,  $\mathbf{K}_4 \text{Person}(d)$  is not valid too.

## Epistemic semantics of P2P mappings: composition



In general, in FOL,  $m$  is not captured by any finite set of GLAV conjunctive mapping assertions.

## Epistemic semantics of P2P mappings: composition



In Hyper,  $m$  is always captured by a finite set of GLAV mapping assertions.

## Answering queries in P2P data integration

- **Distributed query answering**
  - the query is posed to one peer in the system
  - each peer executes the same algorithm, and in doing so exchanges information only with the peers it is connected to
- **Step-by-step algorithm**
  - the query is posed to one peer in the system
  - each peer answers extensionally by taking into account its own data, and then answers intensionally by directing the client to other peers

In both cases, two important issues are

- How to reformulate a query expressed over a peer schema in terms of the local and external sources
- **Loop detection**

## Query answering in P2P systems under epistemic semantics

- We are interested in an algorithm for **distributed query answering**, with no central coordination or centralized data structures
- We assume that peers accept queries in a query language  $\mathcal{L}$  (subsuming at least conjunctive queries)
- We require that each peer, given a query  $Q$  in  $\mathcal{L}$ , is able to compute a **Datalog query**  $Q'$  that is a **perfect rewriting** of  $Q$

## Perfect rewriting

- Given peer  $P$ , and query  $Q$  in  $\mathcal{L}$ , a query  $Q_1$  is a **perfect rewriting** of  $Q$  in  $P$  if
  - $Q_1$  is expressed over the source alphabet of  $P$
  - for each source database  $D$  for  $P$  (i.e., over the local and the external sources), we have that

$$Q_1^D = \text{ans}(Q, P, D)$$

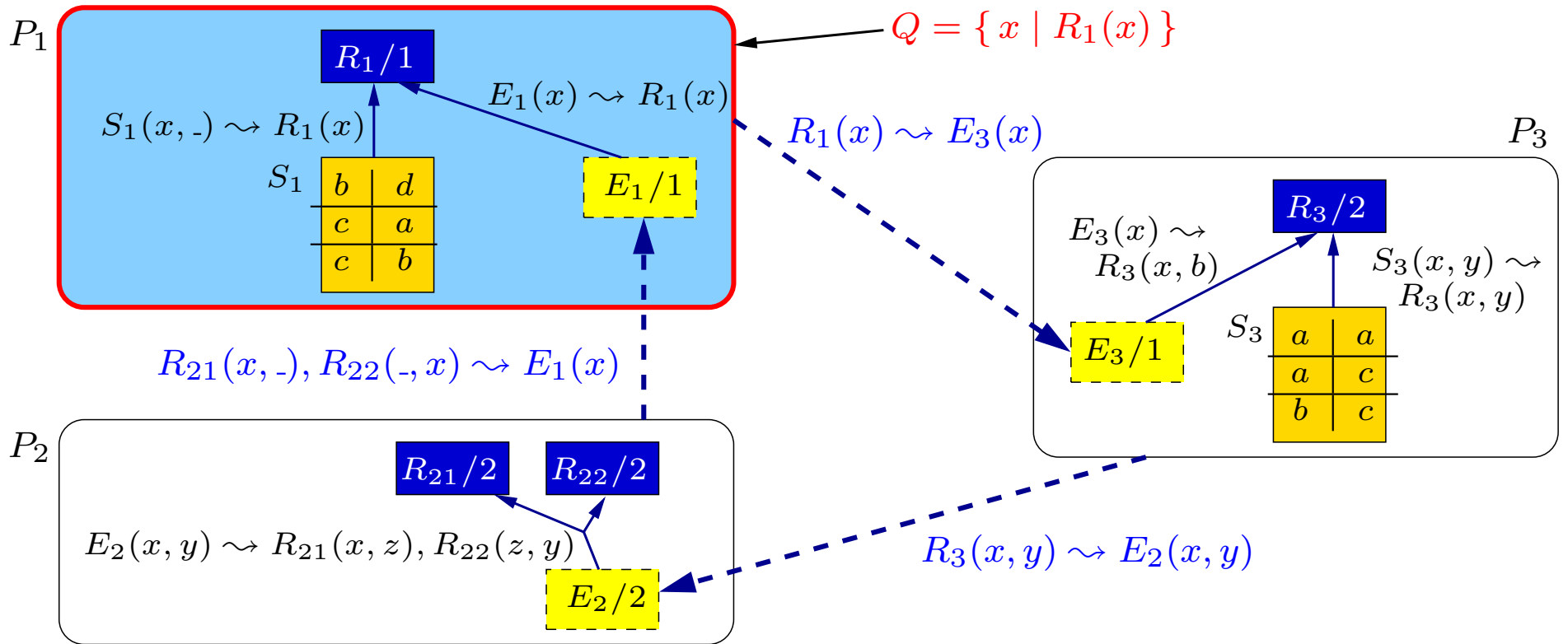
- Perfect rewritings exists in several settings, such as conjunctive client queries, and
  - conjunctive GLAV mappings
  - conjunctive GLAV mappings, plus key and foreign key constraints
  - conjunctive GLAV mappings, plus inclusion dependencies

## Query answering: distributed algorithm

[Calvanese & al PODS'04] presents a distributed query answering algorithm

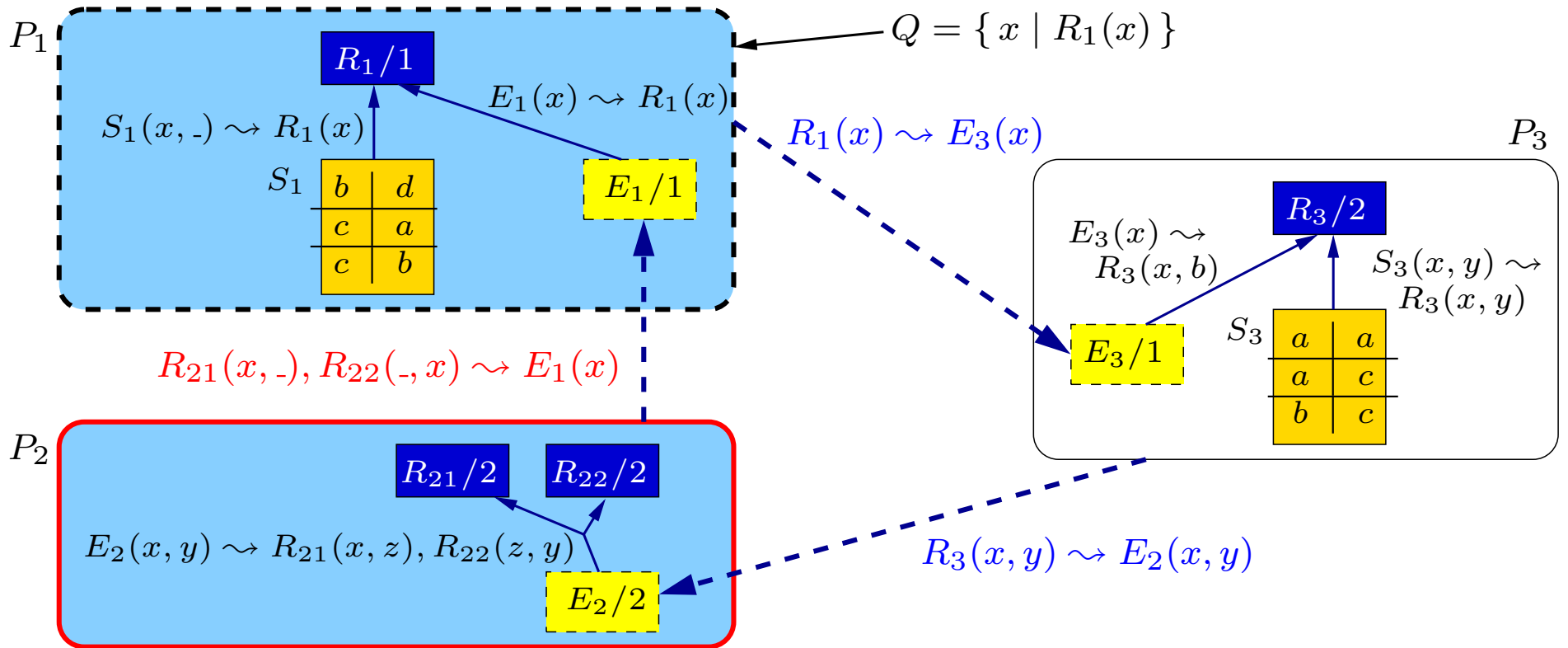
- Each peer rewrites the queries that are requested to it in terms of the local and external sources
- A reference to an external source triggers a request to the peer to which the external source is connected
- Answers to such requests consist of a Datalog program with two parts:
  - an extensional part, which is a set of facts (about local source relations received from other peers)
  - an intensional part, which is a set of Datalog rules
- The final Datalog program is executed at the initiating peer
- Infinite looping is avoided by:
  - associating to each client query a unique (global) transaction id
  - avoiding requests that have already been made for the same transaction id

# Query answering technique: example



$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

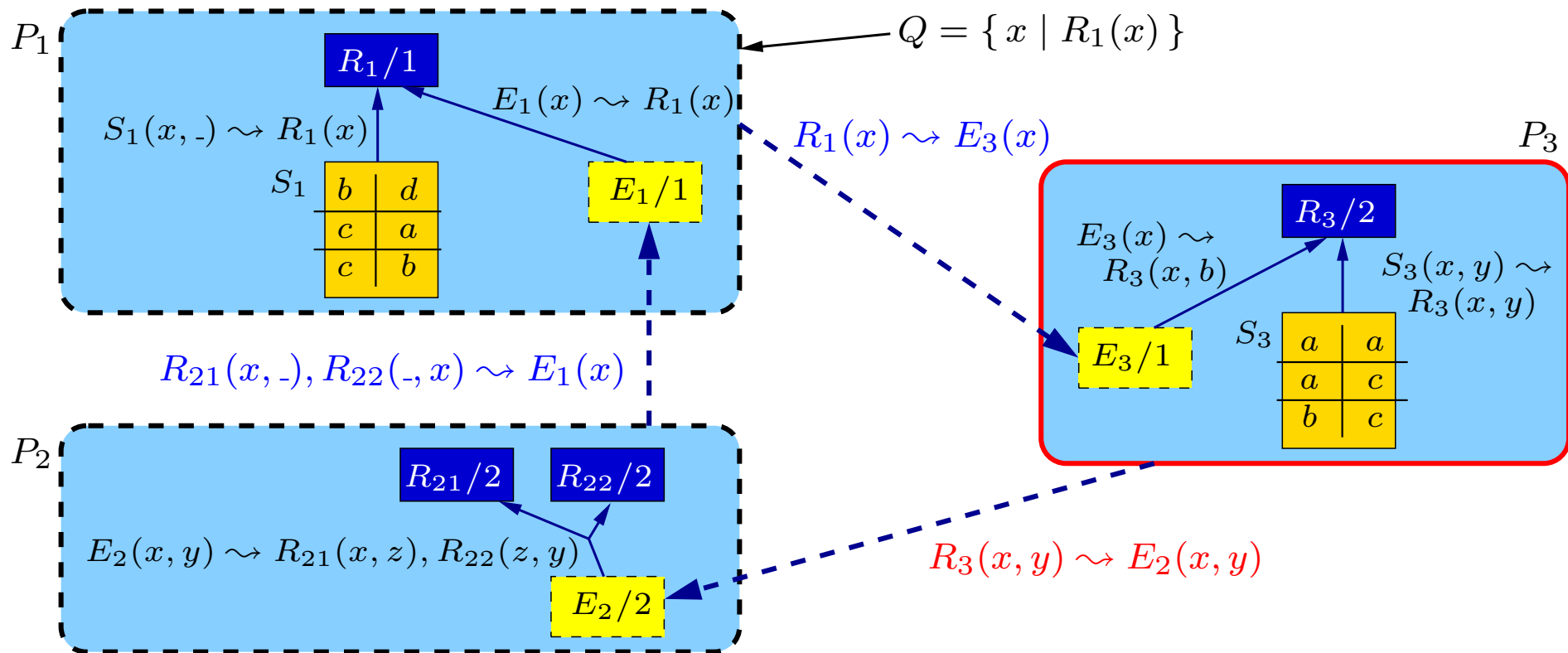
# Query answering technique: example



$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

# Query answering technique: example

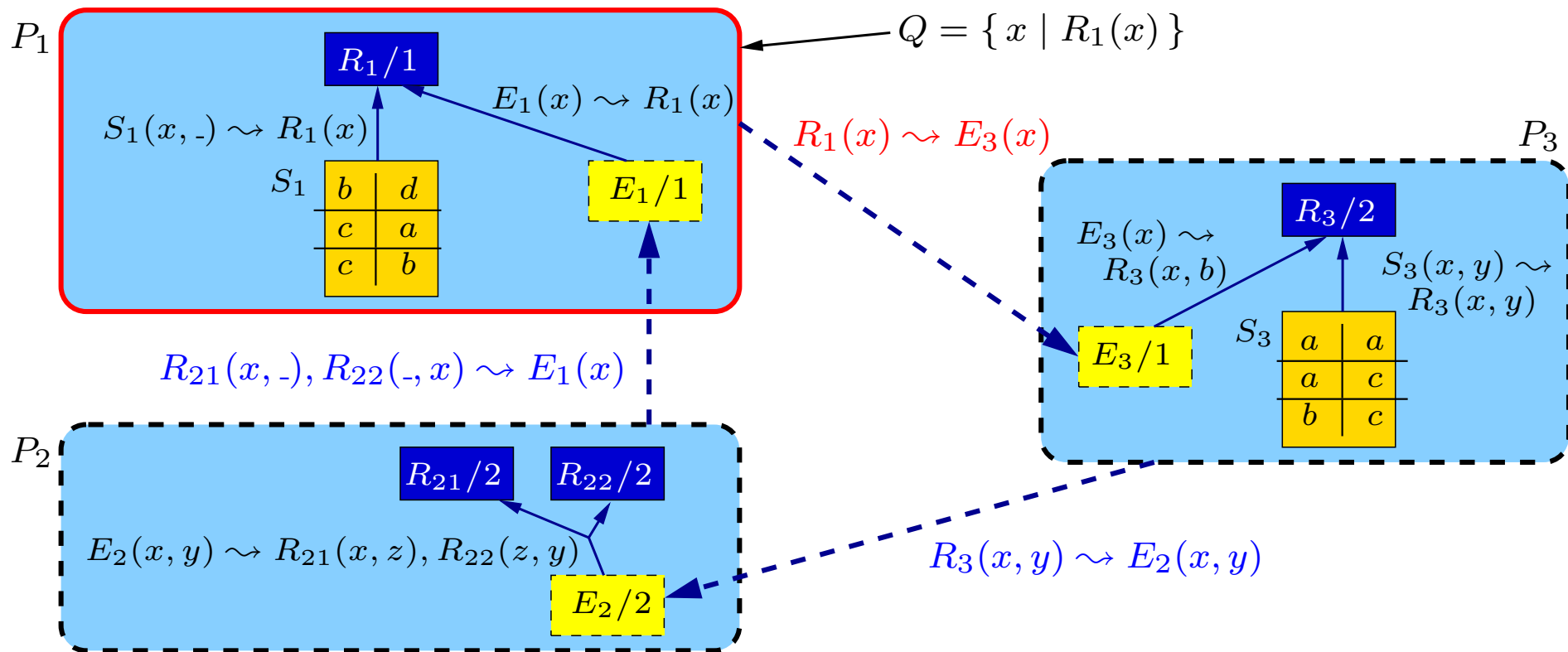


$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

# Query answering technique: example



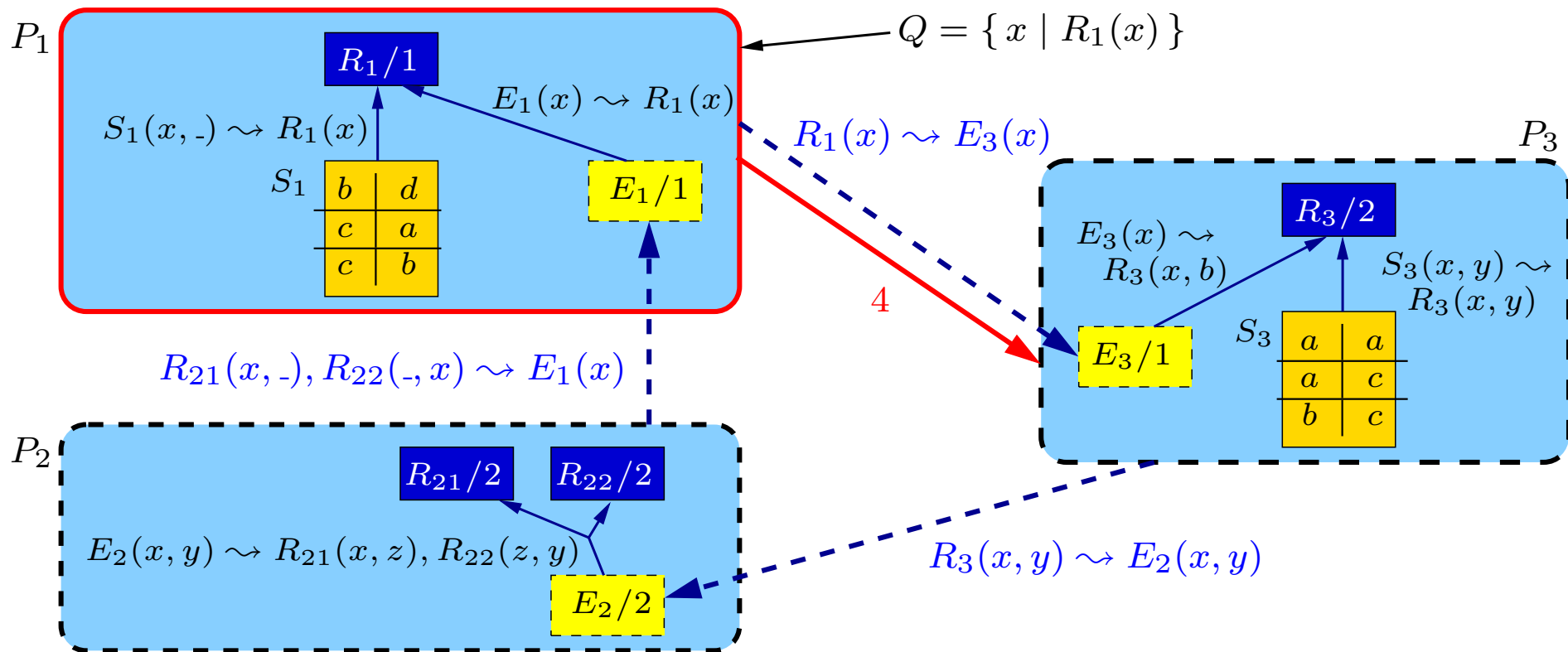
$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$4 \begin{cases} E_3(x) \leftarrow S_1(x, -) \\ E_3(x) \leftarrow E_1(x) \end{cases}$$

# Query answering technique: example



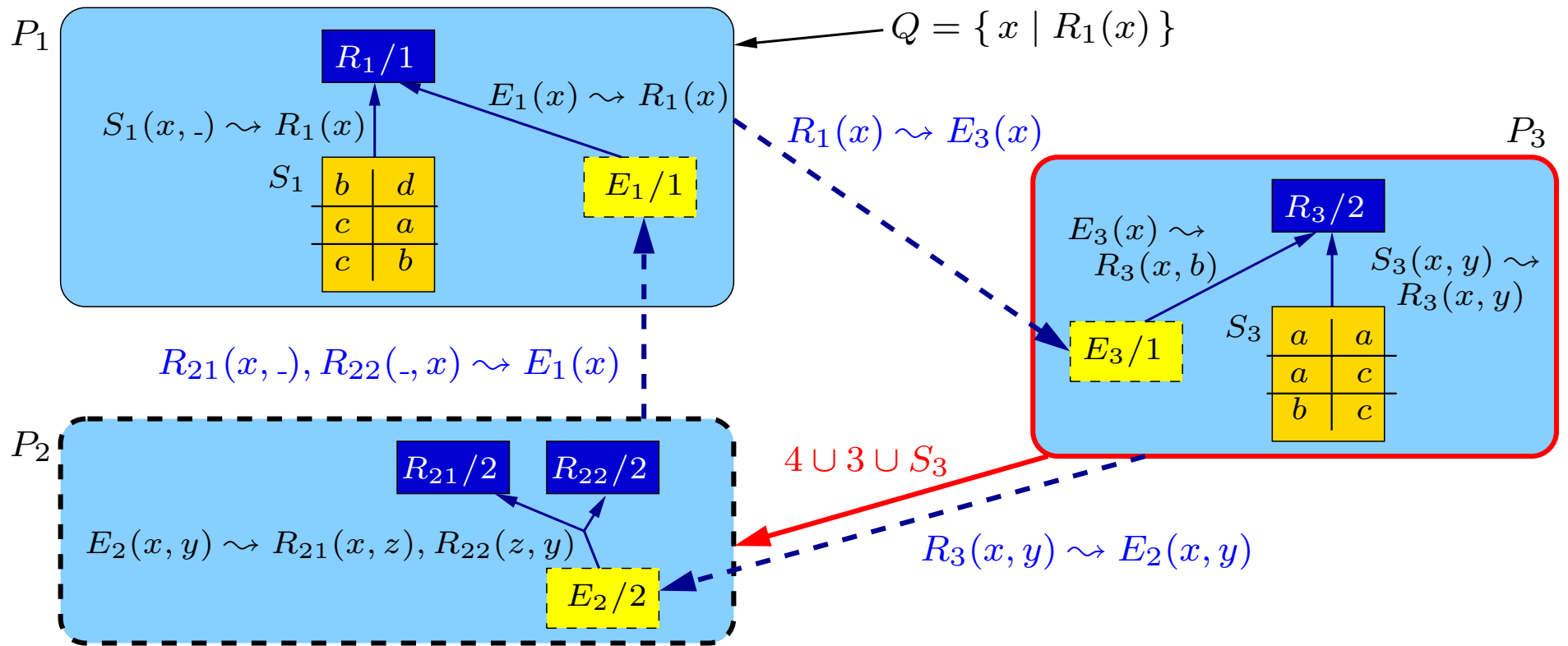
$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$4 \begin{cases} E_3(x) \leftarrow S_1(x, -) \\ E_3(x) \leftarrow E_1(x) \end{cases}$$

# Query answering technique: example



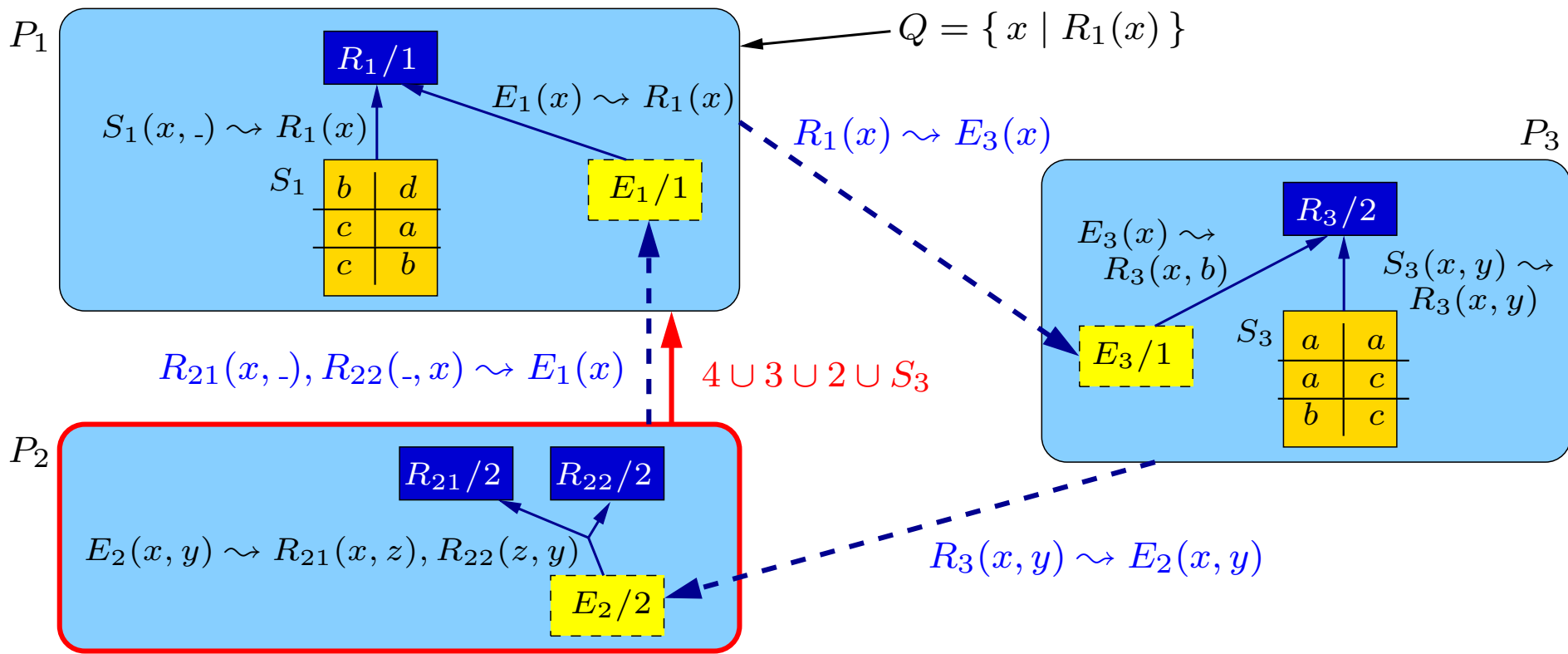
$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$4 \begin{cases} E_3(x) \leftarrow S_1(x, -) \\ E_3(x) \leftarrow E_1(x) \end{cases}$$

# Query answering technique: example



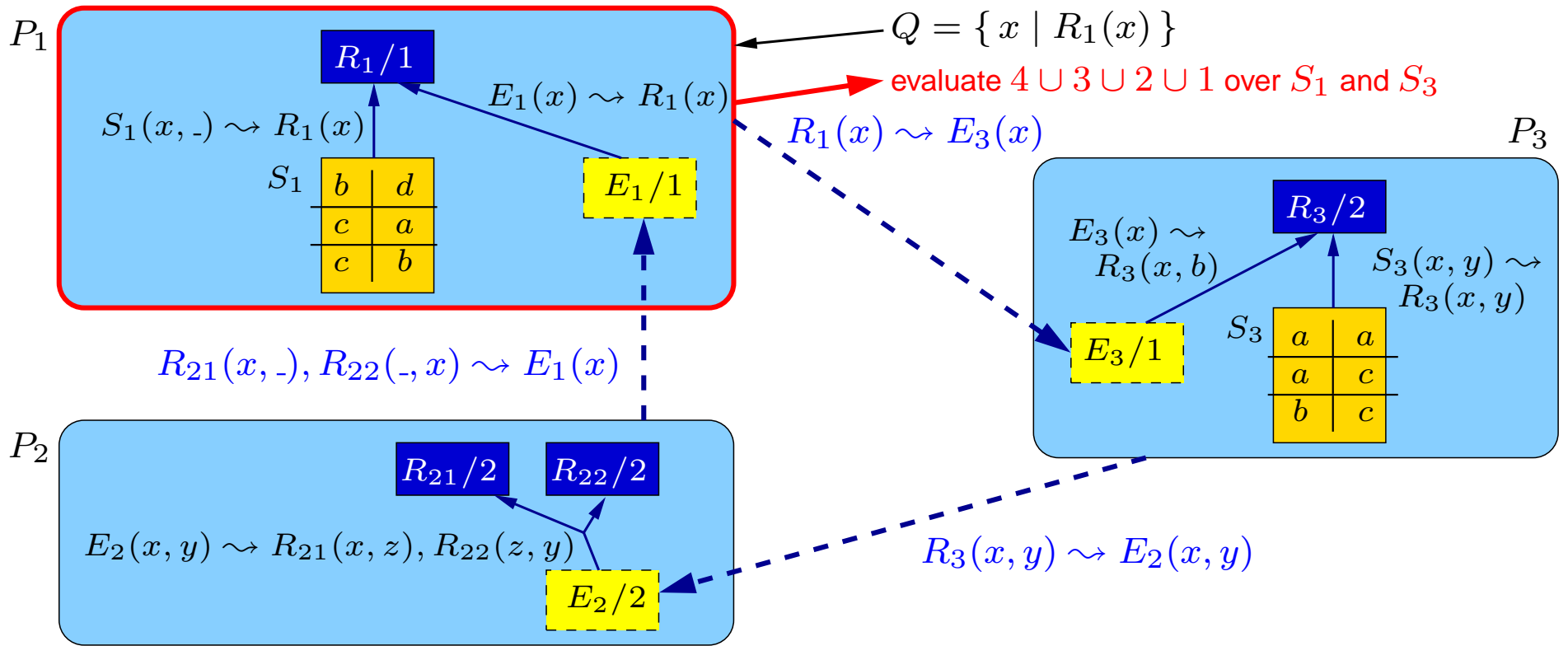
$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

$$4 \begin{cases} E_3(x) \leftarrow S_1(x, -) \\ E_3(x) \leftarrow E_1(x) \end{cases}$$

# Query answering technique: example



$$1 \begin{cases} Q(x) \leftarrow S_1(x, -) \\ Q(x) \leftarrow E_1(x) \end{cases}$$

$$2 \begin{cases} E_1(x) \leftarrow E_2(x, -), E_2(-, x) \end{cases}$$

$$3 \begin{cases} E_2(x, y) \leftarrow S_3(x, y) \\ E_2(x, y) \leftarrow E_3(x), y = b \end{cases}$$

$$4 \begin{cases} E_3(x) \leftarrow S_1(x, -) \\ E_3(x) \leftarrow E_1(x) \end{cases}$$

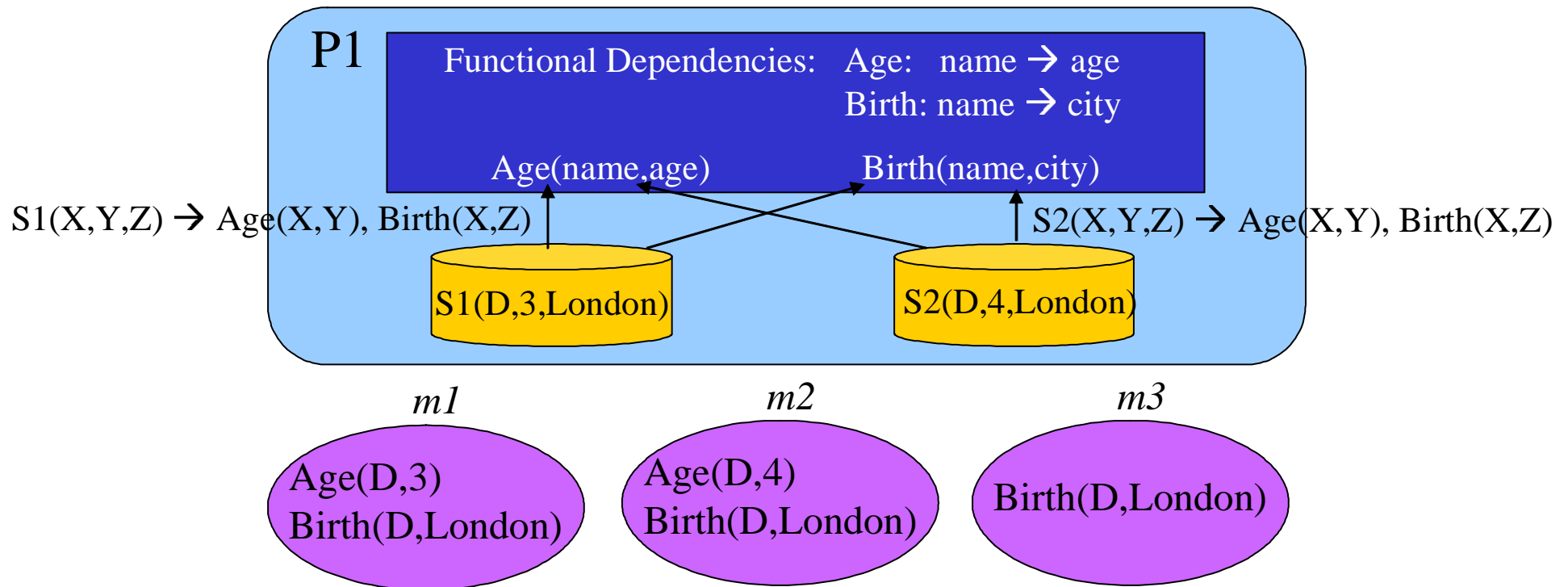
## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- Conclusions

## Dealing with inconsistencies in P2P data integration

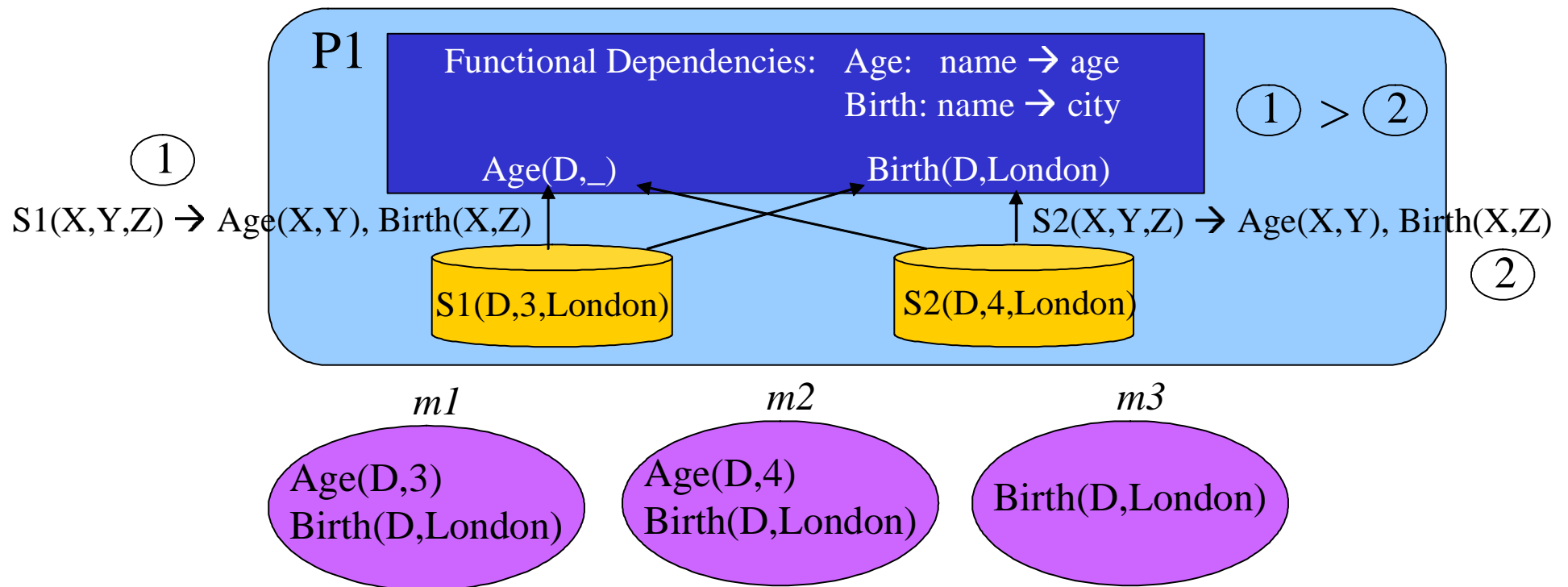
- The whole system blows up with a single inconsistency in one peer:  
unacceptable quality of query answering
- Extraction programs are supposed to perform data cleaning, but still may miss the resolution of some inconsistencies
- We resort to **quality information for dealing with inconsistencies**: we conceive information on source and peer qualities as a mechanism for deciding among possible inconsistency resolutions
- If quality information do not suffice to decide, we reason disjunctively
- Our main goal is to come up with a **well-defined semantics**, which extends the one based on epistemic logic with new (non-monotonic) features

## Dealing with inconsistencies in one peer



- A model  $m$  is preferred to model  $n$  if  $n$  misses some data from the sources that  $m$  does not miss (In the figure, both  $m1$  and  $m2$  are better than  $m3$ )
- The models of a peer are the **most preferred models**
- We are using the notion of repair [Fagin&al 1983, Arenas&al 1999, Arenas&al 2000, Chomicki&al 2002, Cali&al 2003]

# Dealing with quality-based preferences in one peer



- A model  $m$  is preferred to model  $n$  if  $n$  misses some data from the sources that  $m$  does not miss, or if  $m$  respects the preferences more than  $n$
- In the figure,  $m1$  is better than both  $m2$  and  $m3$
- The models of a peer are the **most preferred models**

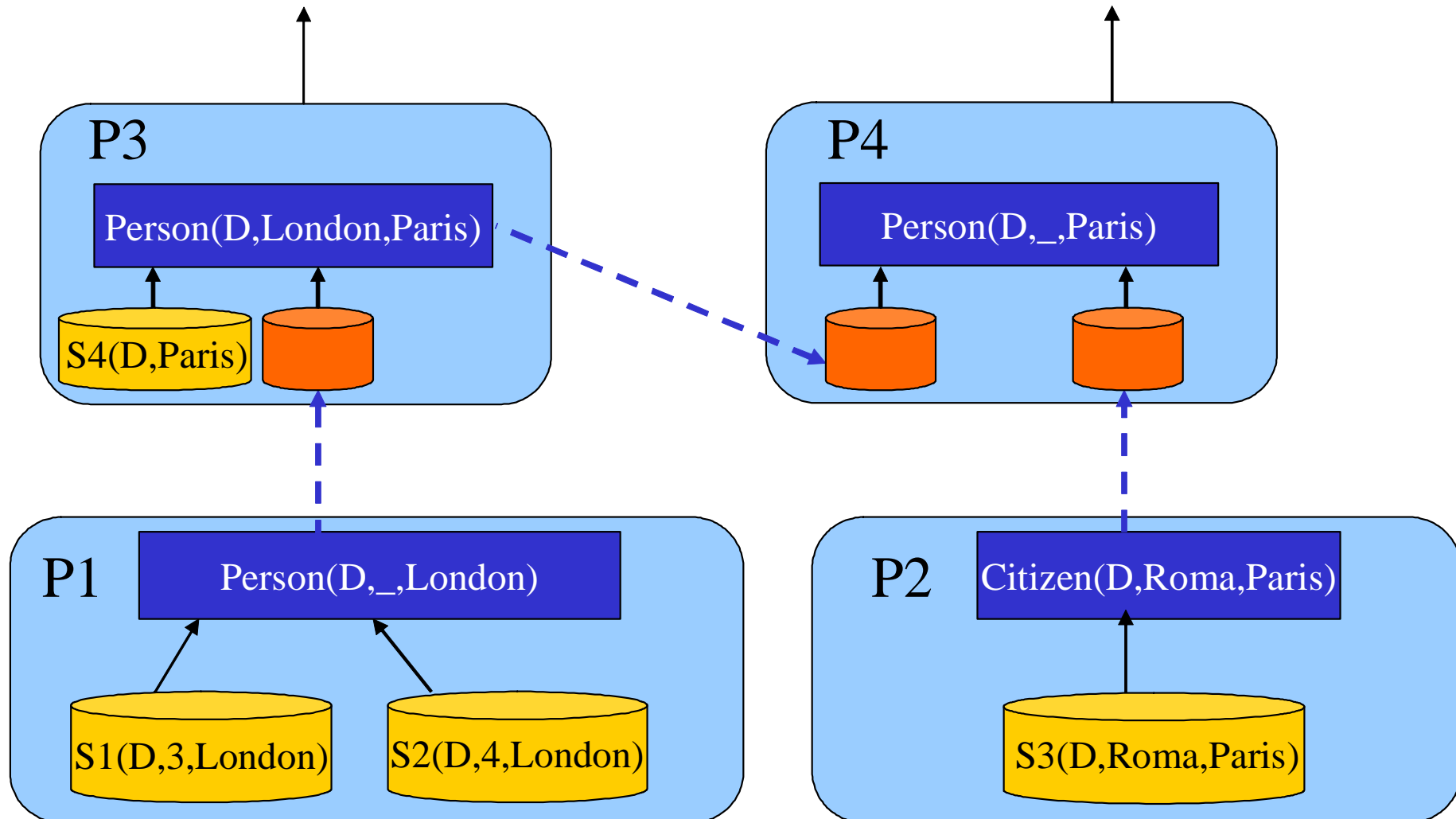
## Dealing with inconsistencies in P2P data integration

- To generalize the idea to the case of multiple peers, we have to be able to compare epistemic models
- The models of the P2P data integration system are the **most preferred epistemic models**

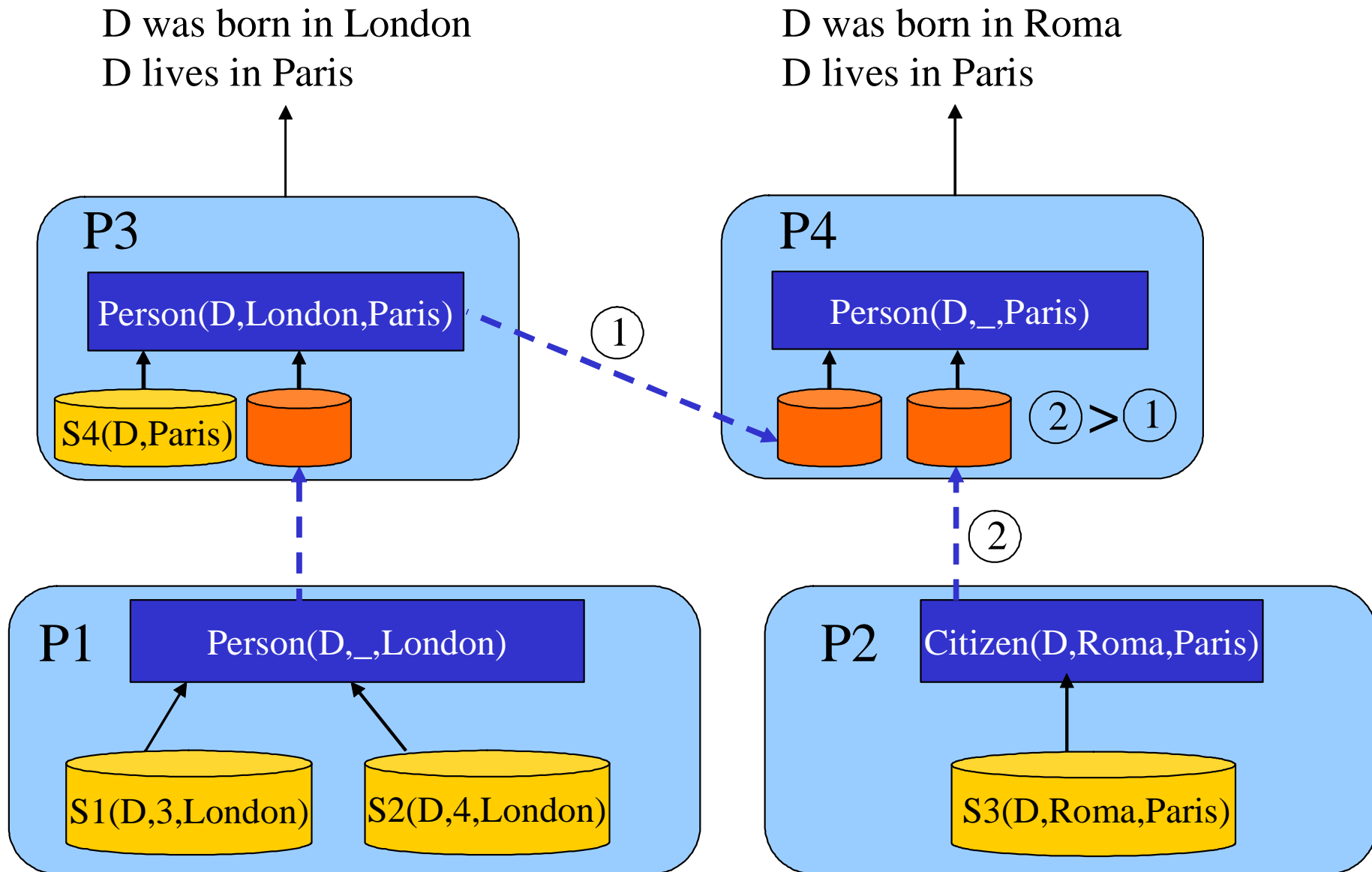
# Dealing with inconsistencies in P2P data integration

D was born in London  
D lives in Paris

Don't know where D was born  
D lives in Paris



# Dealing with inconsistencies and quality in P2P data integration



## Outline

- Data integration architectures
- Centralized data integration
- P2P data integration
- Hyper: epistemic semantics for P2P data integration
- Dealing with inconsistencies in Hyper
- **Conclusions**

# Conclusions

Many open problems and issues, including

- Object identification and mapping
- How to obtain information on mapping, quality and preferences
- Algorithm and complexity in the extended epistemic semantics
- Global schema (or target schema, or peer schemas) expressed in terms of semi-structured data (with constraints)
- Limitations in accessing the sources
- Privacy-based restrictions on peer answers
- Optimization and dynamicity during query processing
- Ongoing experiments within
  - **Hyper**, an IBM Shared University Research program, see <http://www.dis.uniroma1.it/~lenzerini/progetti/hyper/>
  - the EU project **Sewasie**