# Schema Mappings and Data Examples

## Balder ten Cate

UC Santa Cruz & LogicBlox

TbiLLC 2013 - Gudauri

# Relational Databases for Logicians*

*) an oversimplified picture.

- Database schema ~ a finite relational signature. E.g.,

  - { PARTICIPANT(name, email, flight-nr),  FLIGHT(flight-nr, dept-time) }

# Relational Databases for Logicians*

*) an oversimplified picture.

- Database schema ~ a finite relational signature. E.g.,

    - { PARTICIPANT(name, email, flight-nr),  FLIGHT(flight-nr, dept-time) }

- Database instance (of a given schema) ~ a finite structure.
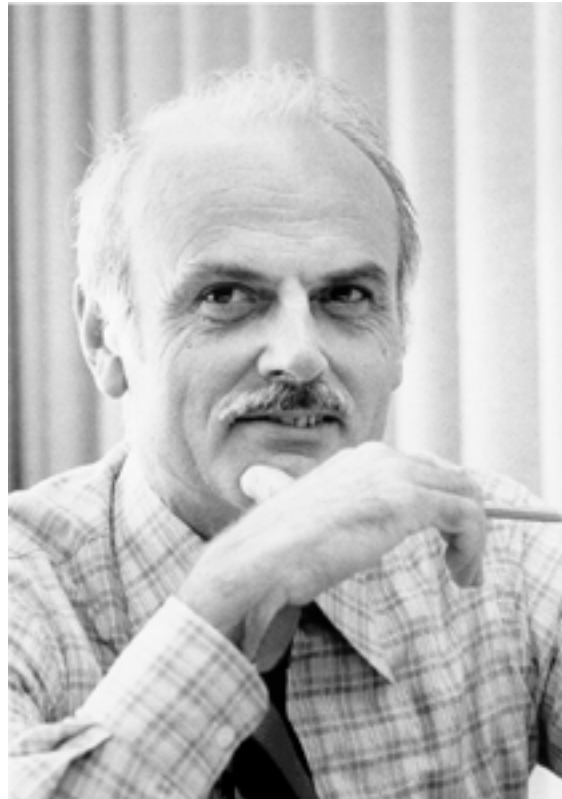
# Relational Databases for Logicians*

*) an oversimplified picture.

- **Database schema** ~ a finite relational signature. E.g.,

  - { PARTICIPANT(name, email, flight-nr),  FLIGHT(flight-nr, dept-time) }

- **Database instance** (of a given schema) ~ a finite structure.

- **Database queries** ~ logical formulas with free variables

  - $\phi(x,y) = \exists z,u$ (PARTICIPANT(x,y,z) & FLIGHT(z, 3:00AM))

# Relational Databases for Logicians*

*) an oversimplified picture.

- Database schema ~ a finite relational signature. E.g.,

  - { PARTICIPANT(name, email, flight-nr), FLIGHT(flight-nr, dept-time) }

- Database instance (of a given schema) ~ a finite structure.

- Database queries ~ logical formulas with free variables

  - $\phi(x,y) = \exists z,u$ (PARTICIPANT(x,y,z) & FLIGHT(z, 3:00AM))

- Database constraints ~ logical sentences expressing structural properties

  - $\forall x,y,z,u$ (PARTICIPANT(x,y,z) → $\exists t$ FLIGHT(z,t))
  - $\forall x,y,z$ (FLIGHT(x,y) & FLIGHT(x,z) → y=z)

Edgar F. Codd (1923-2003)

# Query Languages

- Most important query languages

  - Conjunctive Queries (CQs): $\psi(\mathbf{x}) = \exists \mathbf{y}\,(\alpha_1(\mathbf{x},\mathbf{y}) \wedge \ldots \wedge \alpha_n(\mathbf{x},\mathbf{y}))$

  - Unions of Conjunctive Queries (UCQs): disjunctions of CQs.

  - First-order Queries (~ SQL queries)

  - Datalog (the least-fixpoint extension of UCQs)

# Query Languages

- Most important query languages

  - Conjunctive Queries (CQs): $\psi(\mathbf{x}) = \exists \mathbf{y}\,(\alpha_1(\mathbf{x},\mathbf{y}) \wedge \ldots \wedge \alpha_n(\mathbf{x},\mathbf{y}))$

  - Unions of Conjunctive Queries (UCQs): disjunctions of CQs.

  - First-order Queries (~ SQL queries)

  - Datalog (the least-fixpoint extension of UCQs)

- Most database queries in practice are CQs (a.k.a. SELECT-FROM-WHERE)

# Query Languages

- Most important query languages

  - Conjunctive Queries (CQs): $\psi(\mathbf{x}) = \exists \mathbf{y}\, (\alpha_1(\mathbf{x},\mathbf{y}) \wedge \ldots \wedge \alpha_n(\mathbf{x},\mathbf{y}))$

  - Unions of Conjunctive Queries (UCQs): disjunctions of CQs.

  - First-order Queries (~ SQL queries)

  - Datalog (the least-fixpoint extension of UCQs)

- Most database queries in practice are CQs (a.k.a. SELECT-FROM-WHERE)

- UCQs form a "robustly decidable" fragment of FO logic.

  - In particular, equivalence is decidable (NP-complete).

# Excursion: decidable fragments of FO

- Unions of Conjunctive queries:

  - $\phi(\mathbf{x}) := R(\mathbf{x}) \mid x_i = x_j \mid \phi(x) \wedge \phi(x) \mid \phi(x) \vee \phi(x) \mid \exists y\, \phi(\mathbf{x}, y)$

- The modal fragment:

  - $\phi(x) := P(x) \mid \phi(x) \wedge \phi(x) \mid \neg\phi(x) \mid \exists y(Rxy \wedge \phi(y))$

# Excursion: decidable fragments of FO

- Unions of Conjunctive queries:

    - $\phi(\mathbf{x}) := R(\mathbf{x}) \mid x_i = x_j \mid \phi(x) \wedge \phi(x) \mid \phi(x) \vee \phi(x) \mid \exists y\, \phi(\mathbf{x}, y)$

- The modal fragment:

    - $\phi(x) := P(x) \mid \phi(x) \wedge \phi(x) \mid \neg\phi(x) \mid \exists y (Rxy \wedge \phi(y))$

- UNFO (Unary-Negation Fragment of FO) [tC & Segoufin 2011]

    - $\phi(\mathbf{x}) := R(\mathbf{x}) \mid x_i = x_j \mid \phi(\mathbf{x}) \wedge \phi(\mathbf{x}) \mid \phi(\mathbf{x}) \vee \phi(\mathbf{x}) \mid \neg\phi(x) \mid \exists y\, \phi(\mathbf{x}, y)$

# Excursion: decidable fragments of FO

- Unions of Conjunctive queries:

  - $\phi(\mathbf{x}) := R(\mathbf{x}) \mid x_i = x_j \mid \phi(x) \wedge \phi(x) \mid \phi(x) \vee \phi(x) \mid \exists y\, \phi(\mathbf{x}, y)$

- The modal fragment:

  - $\phi(x) := P(x) \mid \phi(x) \wedge \phi(x) \mid \neg\phi(x) \mid \exists y(Rxy \wedge \phi(y))$

- UNFO (Unary-Negation Fragment of FO) [tC & Segoufin 2011]

  - $\phi(\mathbf{x}) := R(\mathbf{x}) \mid x_i = x_j \mid \phi(\mathbf{x}) \wedge \phi(\mathbf{x}) \mid \phi(\mathbf{x}) \vee \phi(\mathbf{x}) \mid \neg\phi(x) \mid \exists y\, \phi(\mathbf{x}, y)$

- Further extension: GNFO (Guarded-Negation Fragment of FO)
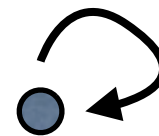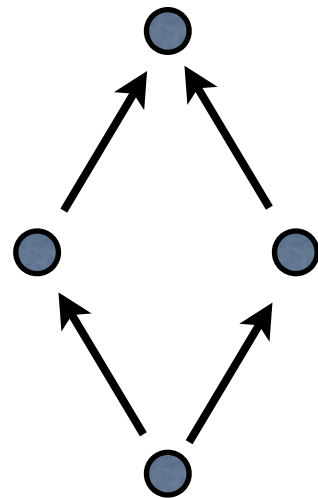  [Barany, tC & Segoufin 2011]

# Homomorphisms

- Conjunctive queries are intimately tied to homomorphisms.

# Homomorphisms

- Conjunctive queries are intimately tied to homomorphisms.

- **Definition**:

    - Let I and J be instances (i.e., finite structures) over the same schema. A homomorphism h: I → J is a map from the domain of I to the domain of J such that $(a,b,c) \in R^I$ implies $(h(a),h(b),h(c)) \in R^J$.

# Examples
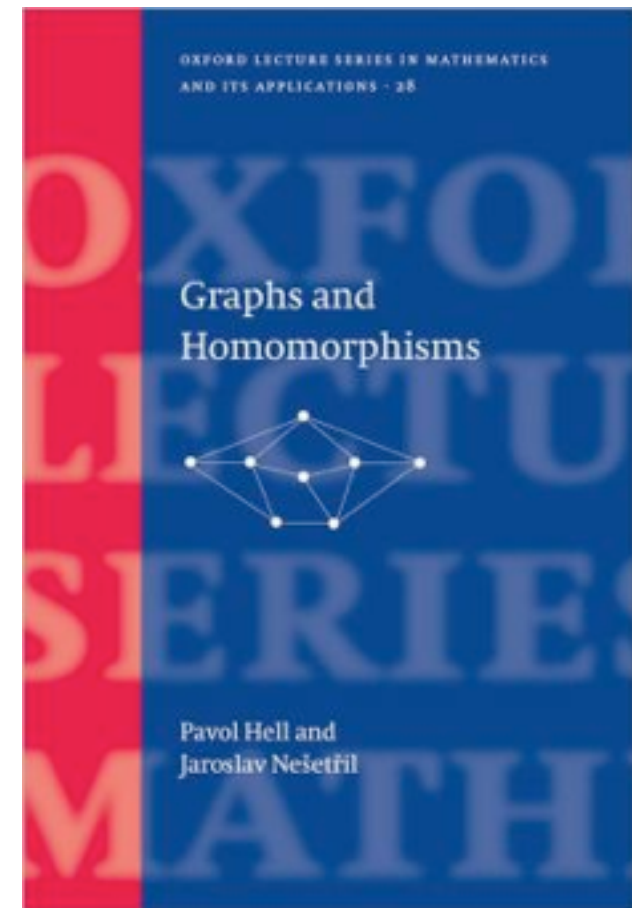
- **Def:** A query q is preserved by homomorphism if for all instances I and J and for all homomorphisms h:I $\rightarrow$ J, (a,b,c) $\in$ q(I) implies (h(a),h(b),h(c)) $\in$ q(J).

- **Def:** A query q is preserved by homomorphism if for all instances I and J and for all homomorphisms h:I → J, $(a,b,c) \in q(I)$ implies $(h(a),h(b),h(c)) \in q(J)$.

- **Thm.** A first-order query is preserved by homomorphisms if and only if it is equivalent to a union of conjunctive queries [Rossman 2005].

  - One of the few preservation theorems that hold over finite structures.

# The Homomorphism Quasi-Order

- We write $I \to J$ if there is a homomorphism h: $I \to J$.

- Fix any relational schema S and let FinStr[S] be the finite structures (i.e., instances) over S.

- (FinStr[S], $\to$) is a quasi-order (reflexive and transitive).

- Its structure has been extensively studied. We will make use of some beautiful results from this area.

# Database Constraints

- Database constraints express structural properties of relations in a schema.

    - $\forall x,y,z,u$ (PARTICIPANT(x,y,z) → ∃t FLIGHT(z,t))
    - $\forall x,y,z$ (FLIGHT(x,y) & FLIGHT(x,z) → y=z)

# Database Constraints

- Database constraints express structural properties of relations in a schema.

  - $\forall x,y,z,u \ (\text{PARTICIPANT}(x,y,z) \rightarrow \exists t \ \text{FLIGHT}(z,t))$
  - $\forall x,y,z \ (\text{FLIGHT}(x,y) \ \& \ \text{FLIGHT}(x,z) \rightarrow y=z)$

- Traditional uses of constraints:

  - Schema design,  integrity control,  query optimization

# Database Constraints

- Database constraints express structural properties of relations in a schema.

  - $\forall x,y,z,u$ (PARTICIPANT$(x,y,z) \rightarrow \exists t$ FLIGHT$(z,t)$)
  - $\forall x,y,z$ (FLIGHT$(x,y)$ & FLIGHT$(x,z) \rightarrow y=z$)

- Traditional uses of constraints:

  - Schema design,   integrity control,   query optimization

- The most well-studied language for specifying constraints:

  - Dependencies : $\forall \mathbf{x}\ (\alpha_1 \wedge \dots \wedge \alpha_n \rightarrow \exists \mathbf{y}\ (\beta_1 \wedge \dots \wedge \beta_n))$

  - Rich enough to express most database constraints in practice.

  - Unfortunately, basic tasks (e.g., entailment) are undecidable.

# Relational Databases for Logicians*

*) an oversimplified picture.

- Database schema ~ a finite relational signature. E.g.,

  - { PARTICIPANT(name, email, flight-nr), FLIGHT(flight-nr, dept-time) }

- Database instance (of a given schema) ~ a finite structure.

- Database queries ~ logical formulas with free variables

  - $\phi(x,y) = \exists z,u\ (\text{PARTICIPANT}(x,y,z)\ \&\ \text{FLIGHT}(z,\ 3{:}00\text{AM}))$

- Database constraints ~ logical sentences expressing structural properties

  - $\forall x,y,z,u\ (\text{PARTICIPANT}(x,y,z) \rightarrow \exists t\ \text{FLIGHT}(z,t))$
  - $\forall x,y,z\ (\text{FLIGHT}(x,y)\ \&\ \text{FLIGHT}(x,z) \rightarrow y{=}z)$

# The Data Interoperability Challenge

- Data-Interoperability:

    - Data may be distributed over different sources, using different schemas.

    - Applications need to access all these data.

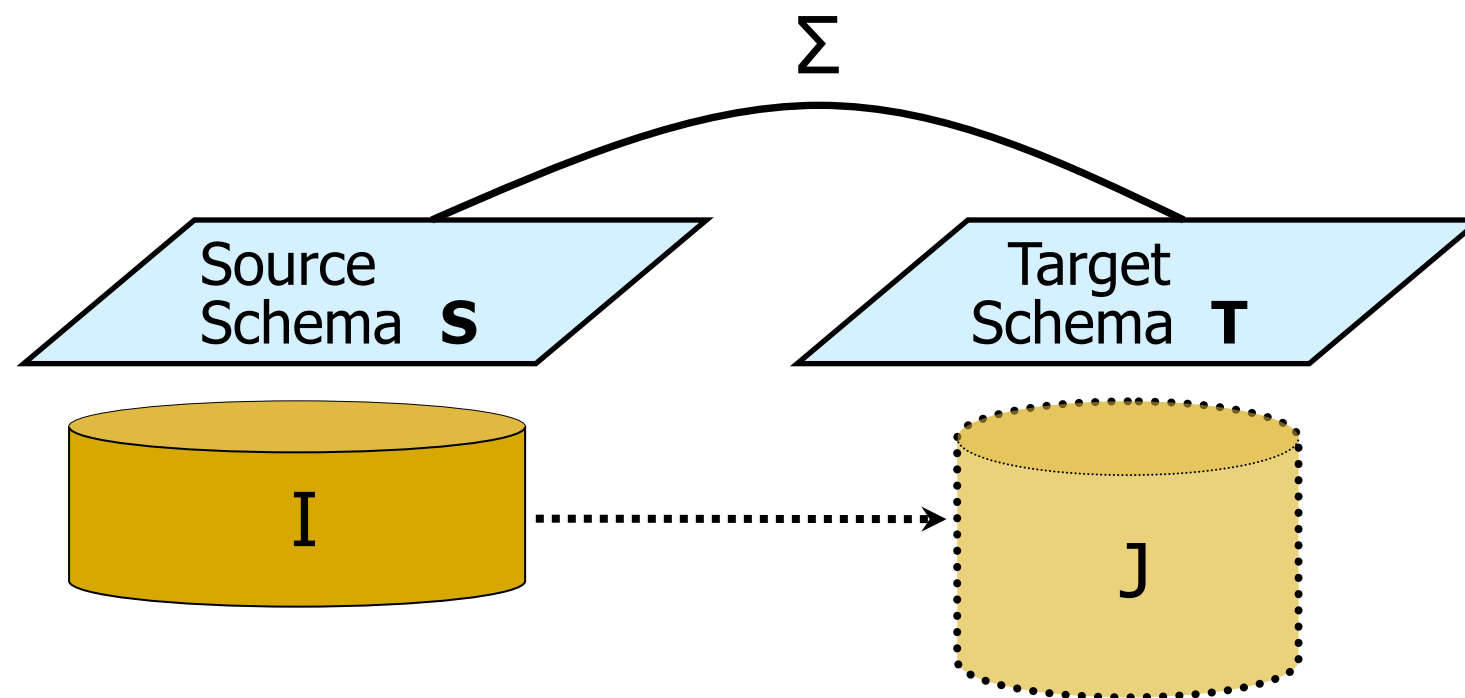# The Data Interoperability Challenge

- Data-Interoperability:

  - Data may be distributed over different sources, using different schemas.

  - Applications need to access all these data.

- How can we uniformly access and manipulate data across sources?

# The Data Interoperability Challenge

- Data-Interoperability:

  - Data may be distributed over different sources, using different schemas.

  - Applications need to access all these data.

- How can we uniformly access and manipulate data across sources?

- Two examples of data interoperability tasks:

  - Data Integration

  - Data Exchange

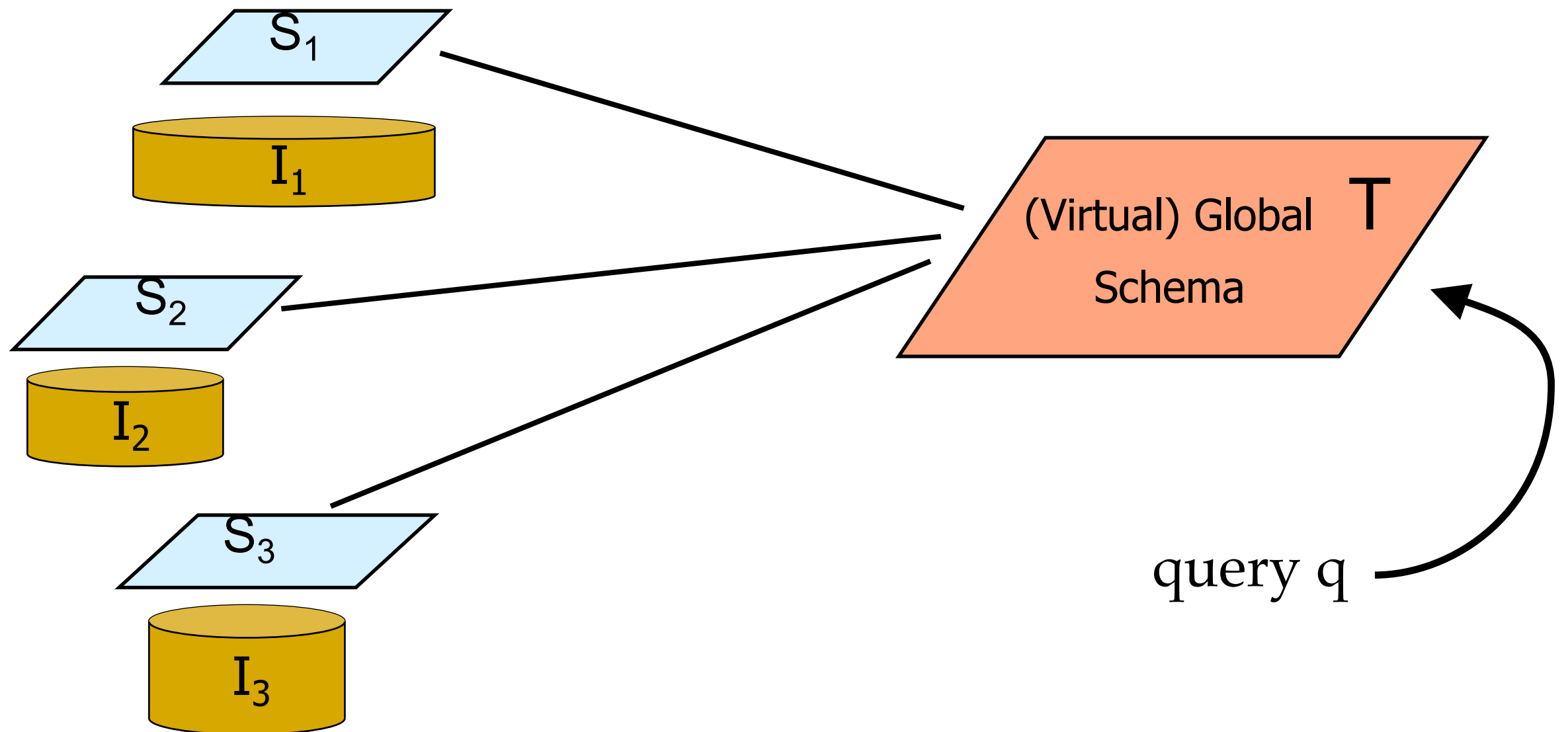# Data Exchange

Transform data structured under a source schemas into data structured under a target schema.

# Data Integration

Query heterogeneous data in different sources via a virtual global schema



$S_1$

$I_1$

$S_2$

$I_2$

$S_3$

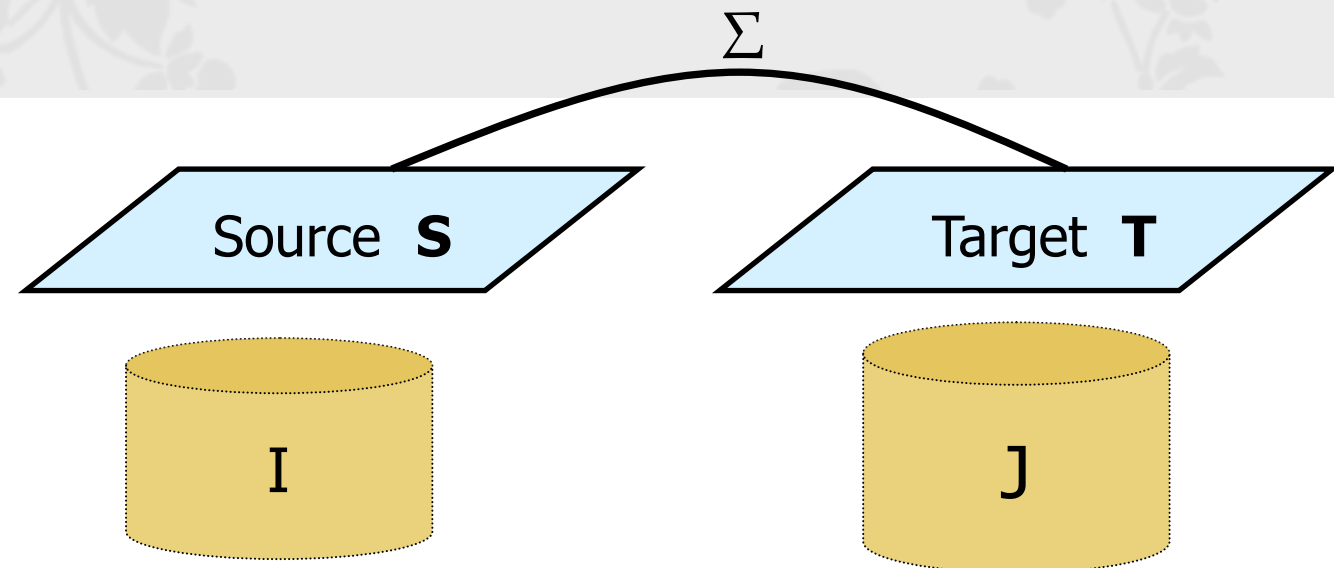$I_3$

(Virtual) Global Schema   T

query q

# Schema Mappings

- A schema mappings is a logical specification of the relationships between two database schemas.

- Schema mappings are fundamental in the formalization data interoperability tasks such as data exchange and data integration.

# Schema Mappings

- A schema mappings is a logical specification of the relationships between two database schemas.

- Schema mappings are fundamental in the formalization data interoperability tasks such as data exchange and data integration.

- Formally, a schema mapping is a triple M=(S,T,Σ), where

  - S and T are schemas (the "source schema" and the "target schema")

  - Σ is a collection of constraints involving the relations of S and T, specified in some schema mapping language (details to come). E.g., $\forall x,y,z(\text{PARTICIPANT}(x,y,z) \rightarrow \text{MAILINGLIST}(x,y))$.

# Schema Mapping Languages

- The choice of schema mapping language involves a compromise between expressive power and practical usability.

  - Allowing arbitrary FO sentences in $\Sigma$ would make the interesting problems undecidable.

- Two of the most important schema mapping specification languages:

  - **GLAV constraints**. These are dependencies $\forall \mathbf{x}\ (\phi(\mathbf{x}) \rightarrow \exists \mathbf{y}\ \psi(\mathbf{x},\mathbf{y}))$ where
    - $\phi$ is a conjunction of relational atomic formulas over the source schema
    - $\psi$ is a conjunction of relational atomic formulas over the target schema.

  - **GAV constraints**: special case of GLAV where the consequent is a single atomic formula (no existential quantification)

  - **LAV constraints:** special case of GLAV where the antecedent is a single atomic formula.
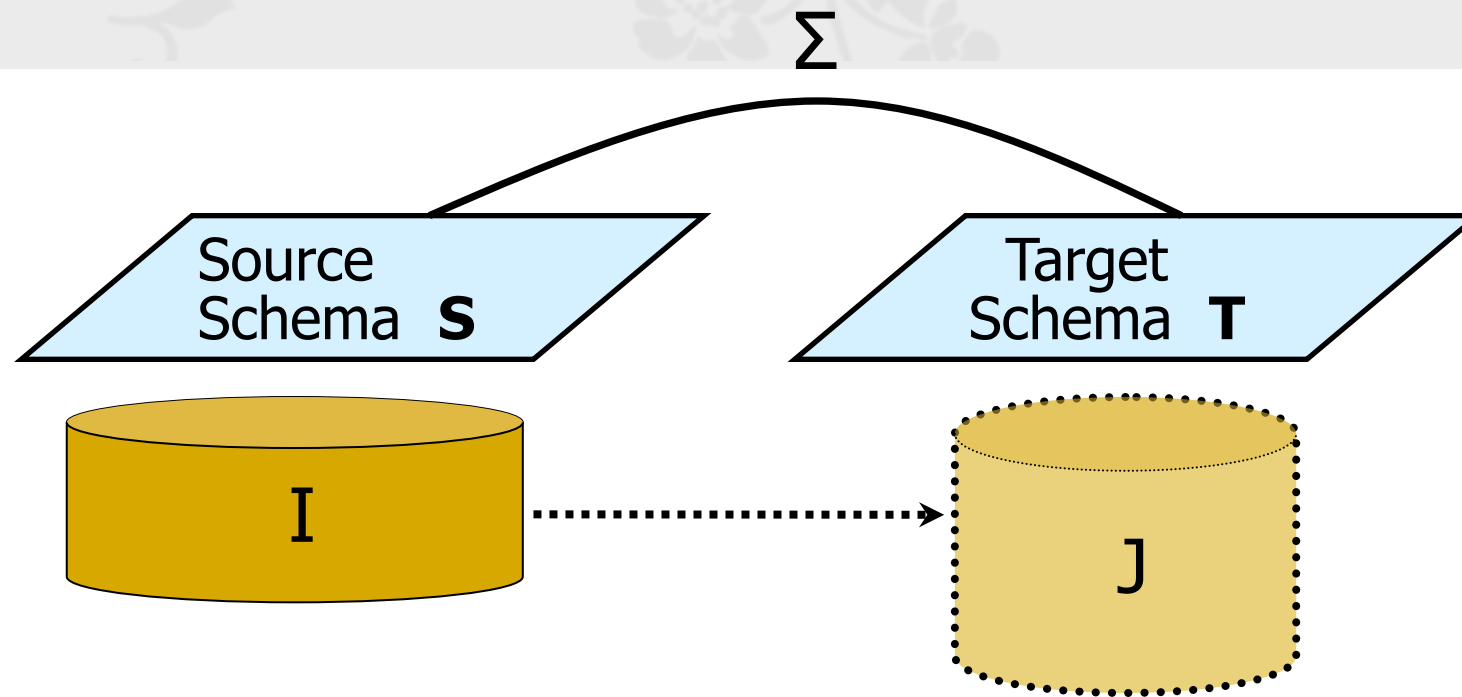
# Semantics of Schema Mappings

$\Sigma$

Source **S**         Target **T**

I          J

- $M = (S, T, \Sigma)$ schema mapping with $\Sigma$ a set of GLAV constraints.

- From a semantic point of view, M can be identified with the set of all its positive data examples.

  - Data Example: A pair (I,J) where I is a source instance and J is a target instance.
  - Positive Data Example for M: a data example (I,J) such that $(I,J) \vDash \Sigma$
  - Negative Data Example for M: a data example (I,J) such that $(I,J) \nvDash \Sigma$
  - If (I,J) is a positive data example for M, we say that J is a solution for I w.r.t. M.

  Sem(M) = { (I,J):  J is a solution for I w.r.t. M }

# Examples

- Consider the schema mapping M = ({E}, {F}, Σ), where

  - Σ = { E(x,y) → ∃z (F(x,z) ∧ F(z,y)) }

- Positive Data Examples (I,J)   (i.e., J a solution for I w.r.t. M)

  - I = { E(1,2) }        J = { F(1,1), F(1,2) }

  - I = { E(1,2) }        J = { F(1,xxx), F(xxx,2) }

  - I = { E(1,2) }        J = { F(1,xxx), F(xxx,2), F(2,3) }

- Negative Data Examples (I,J)     (i.e., J not a solution for I w.r.t. M)

  - I = { E(1,2) }        J = { F(1,3) }

  - I = { E(1,2) }        J = { F(1,3), F(4,2) }

# Data Exchange via a Schema Mapping



- **Data Exchange** via the schema mapping $\mathbf{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$:

  Given a source instance I, construct a solution J for I.

- Difficulty:
  - Typically, there are multiple solutions
  - Which one is the "best" to materialize?

# Data Exchange & Universal solutions

Fagin, Kolaitis, Miller, Popa (2003):

Identified and studied the concept of a **universal solution** in data
exchange.

- A universal solution is a most general solution.

- A universal solution "represents" the entire space of solutions.

# Universal Solutions in Data Exchange

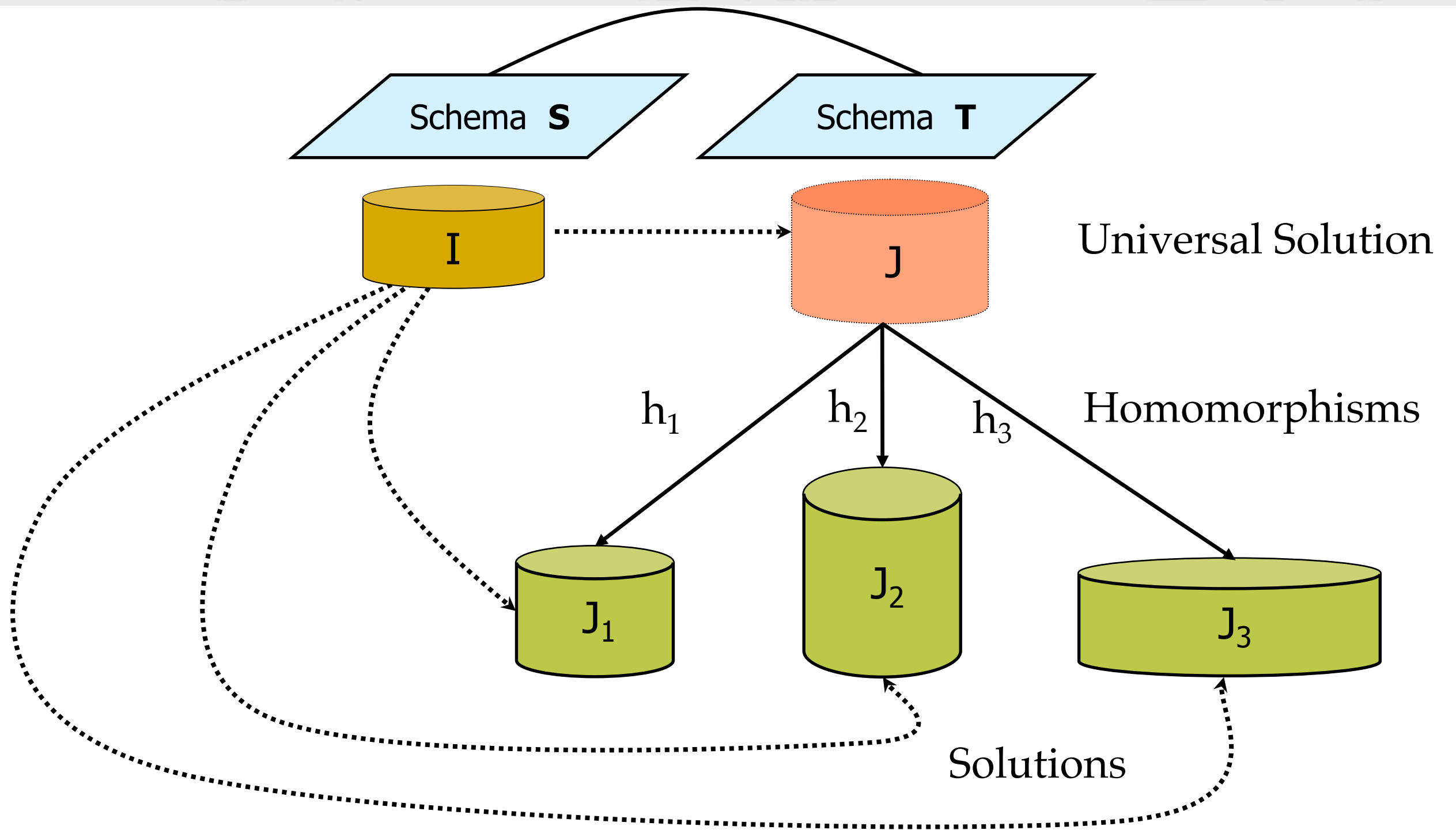Allow two types of values in instances: **constant values** and (**labelled**) **null values.**

**Definition** (FKMP):  A solution J for I is **universal** if it has homomorphisms to all other solutions for I, where the homomorphism may only change the null values.

(thus, a universal solution is a "most general" solution).

Basic result (FKMP): Universal solutions can be constructed in PTIME (data complexity) using an algorithm called the chase.
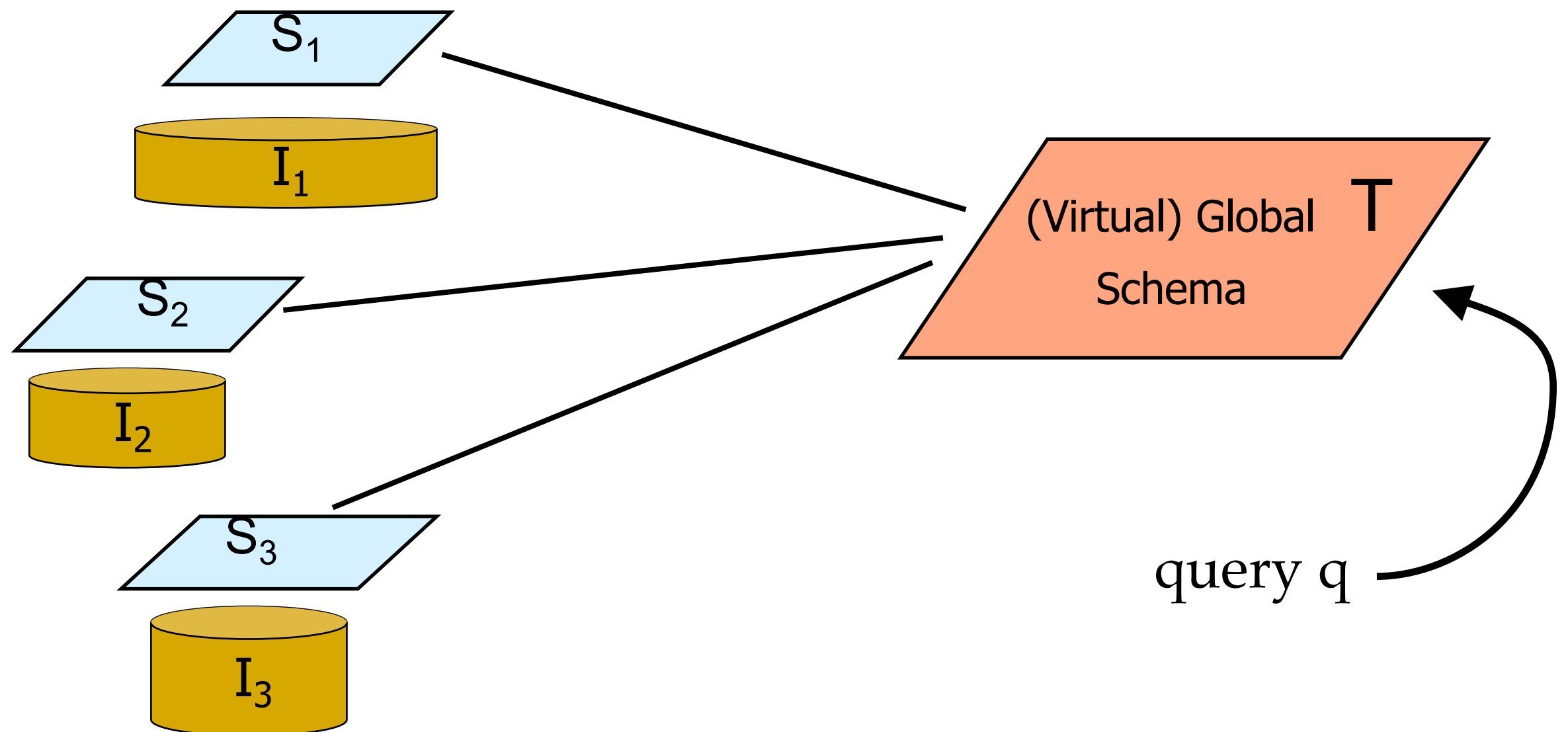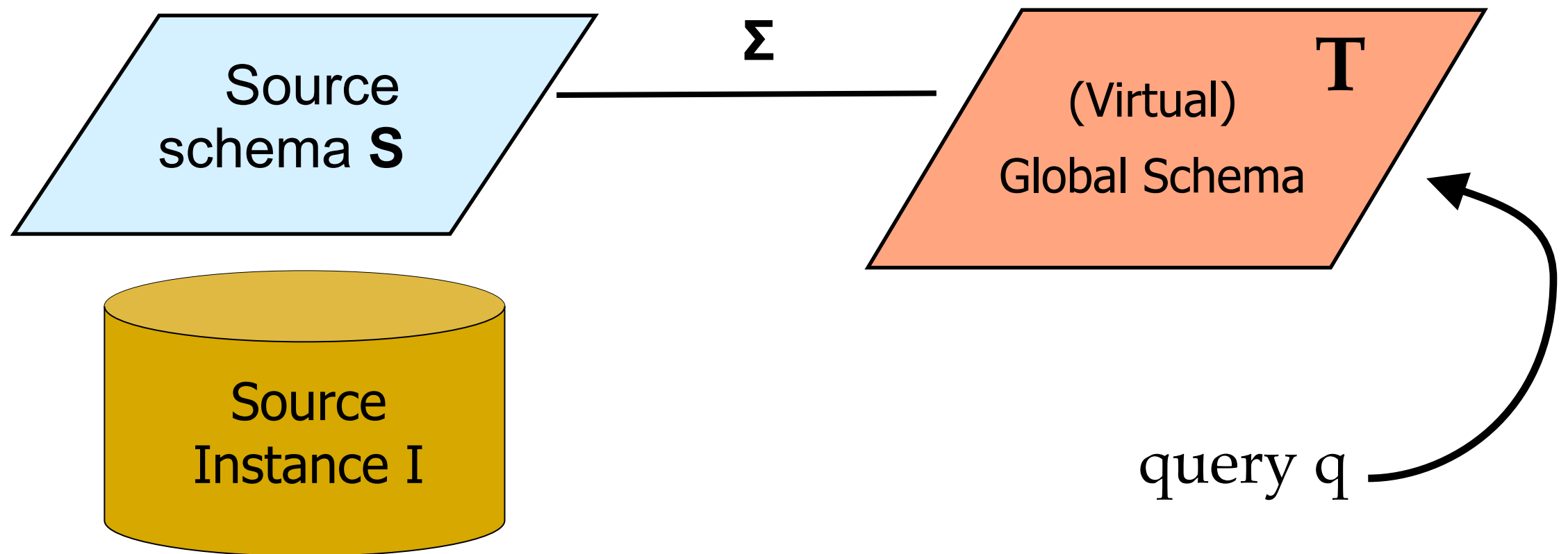
# Universal Solutions in Data Exchange

# Data Integration

Query heterogeneous data in different sources via a virtual global schema

# Data Integration

Query heterogeneous data in different sources via a virtual global schema

# Certain answers

- Let I be a source instance and let q be a target query (a query over **T**).

- **Definition**: $\text{certain}_{\mathbf{M}}(q,I) = \bigcap \{q(J) \mid J \text{ solution of I w.r.t. } \mathbf{M}\}$

  - Idea: $\text{certain}_{\mathbf{M}}(q,I)$ contains the tuples that belong to the answer of q in **all** solutions of I.

# Certain answers

- Let I be a source instance and let q be a target query (a query over **T**).

- **Definition**: $\text{certain}_{\mathbf{M}}(q,I) = \bigcap \{q(J) \mid J \text{ solution of } I \text{ w.r.t. } \mathbf{M}\}$

  - Idea: $\text{certain}_{\mathbf{M}}(q,I)$ contains the tuples that belong to the answer of q in **all** solutions of I.

- If the query is a UCQ, then $\text{certain}_{\mathbf{M}}(q,I)$ can be computed in PTIME.

  - via universal solutions or via query rewriting

# Computing certain answers

- **Theorem** (Fagin, Kolaitis, Miller, Popa 2003):

    - Let J be a universal solution of I w.r.t. **M**. Then for every UCQ q,
      $\text{certain}_{\mathbf{M}}(q,I) = q(J)_{\downarrow}$

- **Theorem** (Abiteboul, Duschka 1998 ++) :

    - For every target UCQ q, there is a source UCQ q' such that $q'(I) = \text{certain}_{\mathbf{M}}(q,I)$.
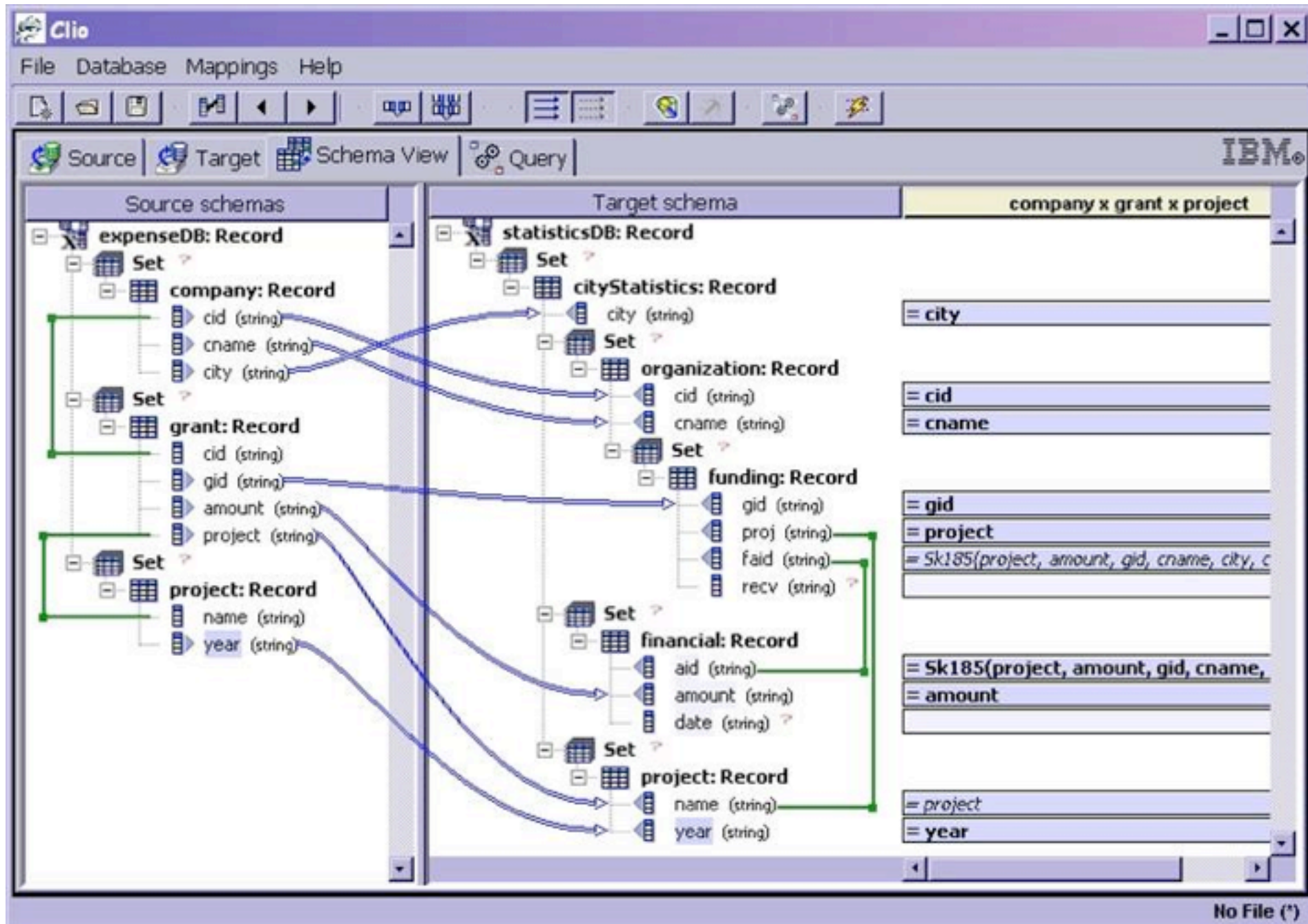
# Where to get your schema mapping

- Constructing a schema mapping is the first step in data exchange and data integration.

# Where to get your schema mapping

- Constructing a schema mapping is the first step in data exchange and data integration.

- **Common approach** (Clio, HepToX, Microsoft mapping composer):

  - derive a schema mapping from a schema matching (a collection of correspondences between attributes of the two schemas).

  - The schema matching itself is obtained semi-automatically using schema matching techniques or by interaction with a user.

  - NB: a schema matching does not uniquely determine a schema mapping.

# Data Examples

- Using data examples in schema mapping design:

  - Data examples can be used to illustrate a candidate schema mapping

  - Deriving schema mappings from examples (learning problem)

- Labeled data examples: a data example (I,J) labeled as being

  - positive -- meaning that J is a solution for I,

  - negative -- meaning that J is not a solution for I, or

  - universal -- meaning J is a universal solution for I.

# Uniquely Characterizing Data Examples

- A set E of labeled data examples uniquely characterizes a schema mapping M, within a class of schema mappings C, if

  - M fits all data examples in E.

  - every schema mapping M′ ∈ C that fits all examples in E is logically equivalent to M.

- Let M be the schema mapping specified by the GLAV constraint
  $\forall x, y \ (E(x,y) \rightarrow F(x,y))$.

  - This is both a GAV schema mapping and a LAV schema mapping.

  - The universal data example (I,J) with  I = { E(a,b) }, J = { F(a,b) }
    uniquely characterizes M w.r.t. the class of all LAV constraints.

  - There is a finite set of universal examples that uniquely characterizes
    M w.r.t. the class of all GAV constraints.

  - There is no finite set of universal examples that uniquely
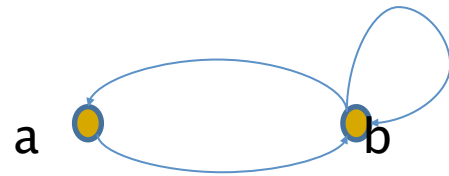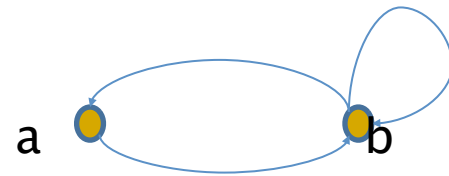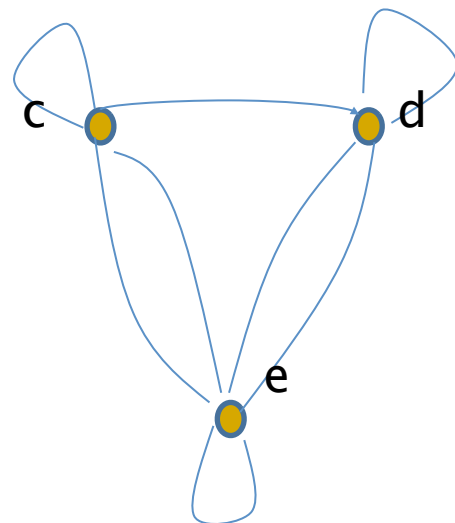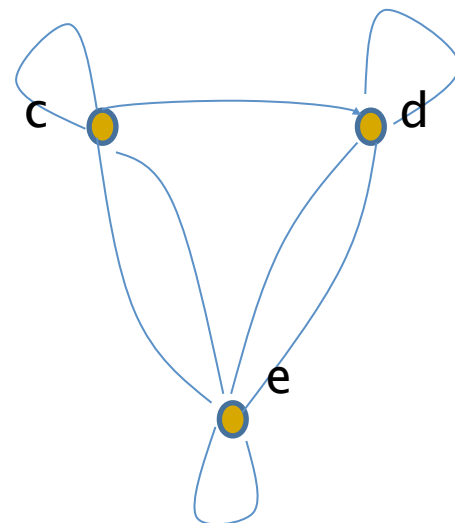    characterizes M w.r.t. the class of all GLAV constraints.

$I_1$

$J_1$

a    b

a    b

$I_2$

$J_2$

a    b

a    b

$I_3$

$J_3$

c    d

e

c    d

e

32

- **Problem**: which GAV schema mappings are uniquely characterizable, by a finite set of labeled data examples, within the class of GAV schema mappings?

- The solution was obtained through an intimate connection with dualities in the homomorphism lattice.

# More about homomorphisms

# More about homomorphisms

- Fix a schema S.

  - When we speak of structures, we will mean finite structures over S.

  - We will assume that S contains at least one non-unary relation symbol.

# More about homomorphisms

- Fix a schema S.

  - When we speak of structures, we will mean finite structures over S.

  - We will assume that S contains at least one non-unary relation symbol.

- Recall: (FinStr[S], →) is a quasi-order (reflexive and transitive).

# More about homomorphisms

- Fix a schema S.

  - When we speak of structures, we will mean finite structures over S.

  - We will assume that S contains at least one non-unary relation symbol.

- Recall: (FinStr[S], →) is a quasi-order (reflexive and transitive).

- We can construct a partially ordered set (poset) by taking the homomorphic equivalence classes.

# More about homomorphisms

- Fix a schema S.

  - When we speak of structures, we will mean finite structures over S.

  - We will assume that S contains at least one non-unary relation symbol.

- Recall: (FinStr[S], →) is a quasi-order (reflexive and transitive).

- We can construct a partially ordered set (poset) by taking the homomorphic equivalence classes.

- However, it turns out there is a nicer way to present this poset.

# The Core of a Structure

- **Definition**:

  - The core of a (finite) structure I, denoted core(I), is the smallest substructure of I that is homomorphically equivalent to I.

  - A structure I is a core if I=core(I).

# The Core of a Structure

- **Definition**:

  - The core of a (finite) structure I, denoted core(I), is the smallest substructure of I that is homomorphically equivalent to I.

  - A structure I is a core if I=core(I).

- **Theorem** [Hell and Nesetril 1992]:

  - core(I) always exists and is unique up to isomorphism

  - $I \rightleftarrows J$ iff core(I) and core(J) are isomorphic.

# The Core of a Structure

- **Definition**:

  - The core of a (finite) structure I, denoted core(I), is the smallest substructure of I that is homomorphically equivalent to I.
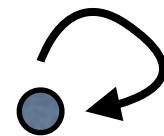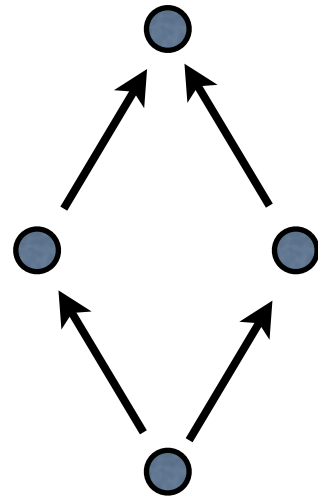
  - A structure I is a core if I=core(I).

- **Theorem** [Hell and Nesetril 1992]:

  - core(I) always exists and is unique up to isomorphism

  - I $\rightleftarrows$ J iff core(I) and core(J) are isomorphic.

- **Corollary**:

  - if I and J are cores and I $\rightleftarrows$ J then I and J are isomorphic.

  - every ~-equivalence class has a unique (up to isomorphism) smallest representative which is a core.

# Examples

# The Homomorphism Lattice.

- Let CoreStr[S] be the set of all non-isomorphic (finite) core structures over schema S. Then (CoreStr[S],→) is a poset, and in fact a lattice.

# The Homomorphism Lattice.

- Let CoreStr[S] be the set of all non-isomorphic (finite) core structures over schema S. Then (CoreStr[S],→) is a poset, and in fact a lattice.

- This lattice has been extensively studied. For example:

  - **Theorem** [Pultr and Trnkova 1980]: Every countable poset is isomorphic to a suborder of (CoreStr[S],→)

- $\rightarrow I \quad = \{ J : J \rightarrow I \}$

- $\rightarrow I = \{J : J \rightarrow I\}$

    - Example (for graphs): $\rightarrow \mathbf{K}_2$ = Class of 2-colorable graphs

- $\rightarrow I \quad = \{ J : J \rightarrow I \}$

  - Example (for graphs): $\rightarrow K_2 =$ Class of 2-colorable graphs

- $I \rightarrow \quad = \{ J : I \rightarrow J \}$

- $\rightarrow I = \{J : J \rightarrow I\}$

  – Example (for graphs): $\rightarrow K_2 = $ Class of 2-colorable graphs

- $I \rightarrow = \{J: I \rightarrow J\}$

  – Example (for graphs): $K_2 \rightarrow = $ Class of graphs with at least one edge.

- $\rightarrow I = \{ J : J \rightarrow I \}$

  - Example (for graphs): $\rightarrow K_2$ = Class of 2-colorable graphs

- $I \rightarrow = \{ J : I \rightarrow J \}$

  - Example (for graphs): $K_2 \rightarrow$ = Class of graphs with at least one edge.

- $\rightarrow I \quad = \{J : J \rightarrow I\}$

  - Example (for graphs): $\rightarrow \mathbf{K}_2$ = Class of 2-colorable graphs

- $I \rightarrow \quad = \{J: I \rightarrow J\}$

  - Example (for graphs): $\mathbf{K}_2 \rightarrow$ = Class of graphs with at least one edge.

- Note:

- $\to I \;=\; \{J : J \to I\,\}$

  – Example (for graphs): $\to \mathbf{K}_2 \;=\;$ Class of 2-colorable graphs

- $I \to \;=\; \{J : I \to J\}$

  – Example (for graphs): $\mathbf{K}_2 \to \;=\;$ Class of graphs with at least one edge.

- Note:

  – $\to I$ defines a <span style="color:darkred">downward closed</span> set in the homomorphism lattice.

38

- $\to I = \{J : J \to I\}$

  - Example (for graphs): $\to K_2$ = Class of 2-colorable graphs

- $I \to = \{J : I \to J\}$

  - Example (for graphs): $K_2 \to$ = Class of graphs with at least one edge.

- Note:

  - $\to I$ defines a downward closed set in the homomorphism lattice.

  - $I \to$ defines an upward closed set in the homomorphism lattice.

- $\to I \ = \{ J : J \to I \}$

  - Example (for graphs): $\to K_2$ = Class of 2-colorable graphs

- $I \to \ = \{ J : I \to J \}$

  - Example (for graphs): $K_2 \to$ = Class of graphs with at least one edge.

- Note:

  - $\to I$ defines a downward closed set in the homomorphism lattice.

  - $I \to$ defines an upward closed set in the homomorphism lattice.

  - $I \not\to$ defines a downward closed set in the homomorphism lattice.

# Simple Duality Pairs

- **Definition**: Let D and F be two finite structures

  - (F,D) is a duality pair if $\to D = F \nrightarrow$

  - In other words, for every structure I, $I \to D$ if and only if $F \nrightarrow I$.

  - In this case, we say that F is an obstruction for D.

# Simple Duality Pairs

- **Definition**:  Let D and F be two finite structures

  - (F,D) is a duality pair if  $\to$ D $=$ F $\not\to$

  - In other words, for every structure I, I $\to$ D if and only if F $\not\to$ I.

  - In this case, we say that F is an obstruction for D.

- **Example**:

  - For graphs,  $(\mathbf{K}_2, \mathbf{K}_1)$ is a duality pair

- **Gallai-Hasse-Roy-Vitaver Theorem** (~1965) for directed graphs:

  - Let $T_k$ be the linear order with k elements, $P_{k+1}$ be the path with k+1 elements. Then $(P_{k+1}, T_k)$ is a duality pair, since for every directed graphs H, $H \rightarrow T_k$ if and only if $P_{k+1} \nrightarrow H$.
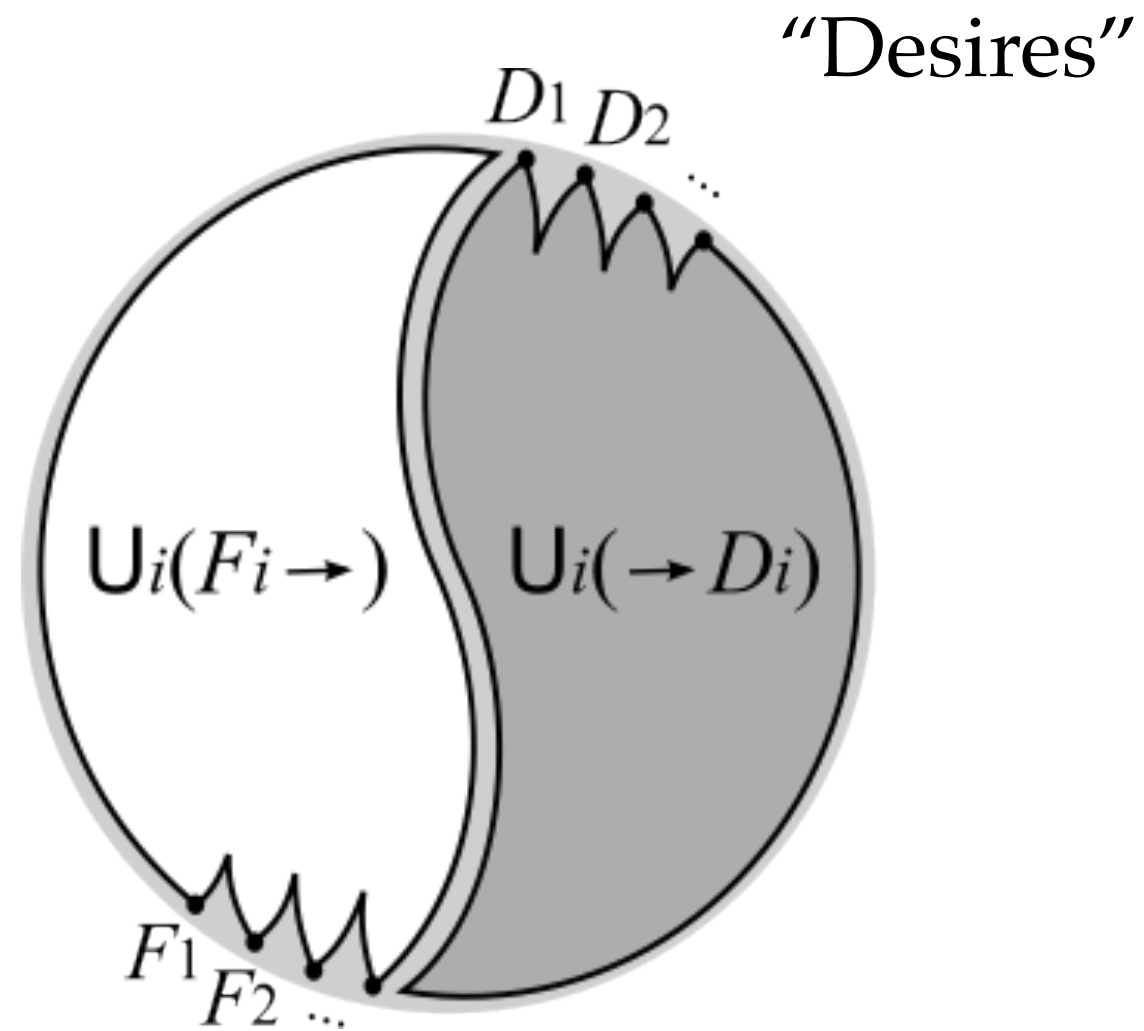
# Duality Pairs

- **Theorem** (König 1936): A graph is 2-colorable if and only if it contains no cycle of odd length. In symbols, $\rightarrow \mathbf{K}_2 = \bigcap_{i \geq 0} (\mathbf{C}_{2i+1} \nrightarrow)$.

# Duality Pairs

- **Theorem** (König 1936): A graph is 2-colorable if and only if it contains no cycle of odd length. In symbols, $\to \mathbf{K}_2 = \bigcap_{i \geq 0} (\mathbf{C}_{2i+1} \not\to)$.

- **Definition**: Let $F$ and $D$ be two sets of structures. We say that $(F, D)$ is a duality pair if $\bigcup_{D \in D} (\to D) = \bigcap_{F \in F} (F \not\to)$.

  - In other words, for every structure I, tfae:

    - There is a structure D in $D$ such that $I \to D$.

    - For every structure F in $F$, we have $F \not\to I$.

  - In this case, we say that F is an obstruction set for D.
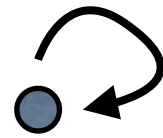
**Duality Pair** (*F*,*D*),where

$F = \{F_1, F_2, \ldots\}$
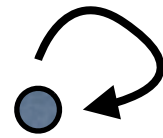
$D = \{D_1, D_2, \ldots\}$

"Desires"



"Frustrations"

# Example



- Let F be the one-element cycle.

# Example



- Let F be the one-element cycle.

- **Question**: Is {F} an obstruction set for a finite set of structures?

    - I.e., is there a duality pair of the form ({F},D) ?
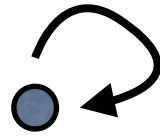
# Example

- Let F be the one-element cycle.

- **Question**: Is {F} an obstruction set for a finite set of structures?

  - I.e., is there a duality pair of the form ({F},D) ?

- No. This has to do with the fact that F contains a cycle.

# Acyclicity

- The incidence graph inc(A) of a structure A is the bipartite graph with

  - nodes: the elements of A and the atomic facts (e.g., $R(a_1,...,a_n)$) of A

  - edges between elements and facts in which they occur

# Acyclicity

- The incidence graph inc(A) of a structure A is the bipartite graph with

  - nodes: the elements of A and the atomic facts (e.g., $R(a_1,...,a_n)$) of A

  - edges between elements and facts in which they occur

- The structure A is acyclic if

  - Inc(A) is acyclic, and

  - No element occurs twice in the the same fact.

# Characterization of Obstruction Sets

- **Theorem** (Foniok, Nešetřil, and Tardif 2008):

  - Let *F* be a finite set of homomorphically incomparable core structures. Tfae:

    - *F* is an obstruction set of some finite set *D* of structures.

    - Each structure in *F* is acyclic.

  - Moreover, there is an algorithm that, given such a set *F* consisting of acyclic structures, computes a finite set *D* of structures such that (*F*, *D*) is a duality pair.

# Characterization of Obstruction Sets

- **Theorem** (Foniok, Nešetřil, and Tardif 2008):

  - Let $F$ be a finite set of homomorphically incomparable core structures. Tfae:

    - $F$ is an obstruction set of some finite set $D$ of structures.

    - Each structure in $F$ is acyclic.

  - Moreover, there is an algorithm that, given such a set $F$ consisting of acyclic structures, computes a finite set $D$ of structures such that $(F, D)$ is a duality pair.

- In particular, if F is the one-element cycle, then {F} is not an obstruction set of any finite set of structures.

# Structures with Constant Symbols

# Structures with Constant Symbols

- The preceding theorem extends to structures with constant symbols when acyclicity is replaced by c-acyclicity.

# Structures with Constant Symbols

- The preceding theorem extends to structures with constant symbols when acyclicity is replaced by c-acyclicity.

- A structure with constant symbols is c-acyclic if

  - Every cycle in Inc(A) contains an element named by a constant symbol, and

  - Only elements named by constant symbols may occur twice in the same fact.

# Back to Schema Mappings

- The canonical structure of a GAV constraint

$$\forall \mathbf{x}\ (\varphi_1(\mathbf{x}) \land ... \land \varphi_\kappa(\mathbf{x}) \rightarrow R(x_{i1},\ldots,x_{im}))$$

is the structure with

- domain: the variables in **x** themselves

- atomic facts: $\varphi_1(x), \ldots, \varphi_\kappa(x)$

- constant symbols $c_1,\ldots,c_m$ denoting $x_{i1},\ldots,x_{im}$

# Back to Schema Mappings

- The canonical structure of a GAV constraint

$$\forall \mathbf{x}\ (\varphi_1(\mathbf{x}) \wedge \ldots \wedge \varphi_\kappa(\mathbf{x}) \rightarrow R(x_{i1}, \ldots, x_{im}))$$

  is the structure with

  - domain: the variables in $\mathbf{x}$ themselves

  - atomic facts: $\varphi_1(x), \ldots, \varphi_\kappa(x)$

  - constant symbols $c_1, \ldots, c_m$ denoting $x_{i1}, \ldots, x_{im}$

- **Example**: $\forall xyz\ (E(x,y) \wedge E(y,z) \rightarrow R(x,z))$ has canonical structure
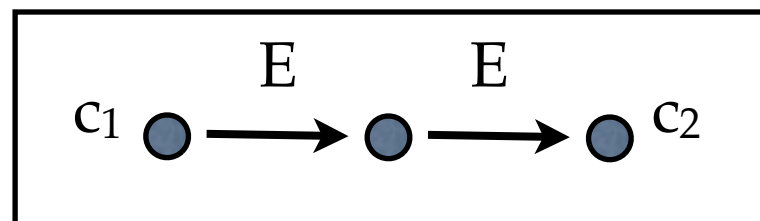
# Back to Schema Mappings

- The canonical structure of a GAV constraint

$$\forall \mathbf{x} \, (\varphi_1(\mathbf{x}) \wedge \ldots \wedge \varphi_\kappa(\mathbf{x}) \to R(x_{i1}, \ldots, x_{im}))$$

  is the structure with

  - domain: the variables in $\mathbf{x}$ themselves

  - atomic facts: $\varphi 1(x), \ldots, \varphi \kappa(x)$

  - constant symbols $c_1, \ldots, c_m$ denoting $x_{i1}, \ldots, x_{im}$

- **Example**: $\forall xyz \, (E(x,y) \wedge E(y,z) \to R(x,z))$ has canonical structure

- **Theorem**: Let M = (S, T, Σ) be a GAV schema mapping. Tfae:

  - M is uniquely characterizable within the class of all GAV constraints.

  - For every target relation symbol R, the set of the canonical structures of the GAV constraints in Σ with R as their consequent is the obstruction set of some finite set D of structures.

- **Theorem**: Let M = (S, T, Σ) be a GAV schema mapping. Tfae:

    - M is uniquely characterizable within the class of all GAV constraints.

    - For every target relation symbol R, the set of the canonical structures of the GAV constraints in Σ with R as their consequent is the obstruction set of some finite set D of structures.

- **Corollary**: testing unique characterizability is NP-complete, and one can effectively construct a uniquely characterizing finite set of data examples if it exists.

# Summary

- **Schema mappings**: a fundamental building block in the study of data-interoperability problems.

- **Homomorphism dualities**: a powerful tool from graph theory (with many applications in constraint satisfaction as well)

# Main References

- Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, Lucian Popa (2003) Data Exchange: Semantics and Query Answering. ICDT 2003: 207-224

- Bogdan Alexe, Balder ten Cate, Phokion G. Kolaitis, and Wang Chiew Tan (2011). Characterizing schema mappings via data examples. ACM Trans. Database Syst., 36(4):23

- Pavol Hell and Jaroslav Nesetril (2004). Graphs and homomorphisms.