# Reflections on Schema Mappings, Data Exchange, and Metadata Management

## Phokion G. Kolaitis

UC Santa Cruz & IBM Research - Almaden





# A Very Brief History

- Data exchange is the oldest database problem. Bernstein - 2003
- Yet, data exchange was not formalized until around 2000.
- Data exchange was formalized using schema mappings.
  - Schema Mappings as Query Discovery Miller, Haas, Hernández - VLDB 2000
  - Clio: A Semi-Automatic Tool for Schema Mapping Hernández, Miller, Haas - SIGMOD 2001
  - Data Exchange: Semantics and Query Answering Fagin, K ..., Miller, Popa - ICDT 2003
- Extensive study of data exchange and schema mappings during the past 15 years.

## Roadmap

### Part I: Schema Mappings and Data Exchange

Algorithmic and structural properties.

### Part II: Operations on Schema Mappings

Composing and inverting schema mappings.

### Part III: Understanding and Deriving Schema Mappings

Data examples to understand/derive schema mappings.

# Data Exchange



Transform data structured under a source schema **S** into data structured a target schema **T** 

# Formalizing Data Exchange via Schema Mappings

Definition: A schema mapping is a triple  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W})$  with

- S is a source schema, T is a target schema,
- ➤ W is a set of pairs (I, J) with I a source instance and J a target instance.

Syntactically, a schema mapping is a triple  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$  with

- S is a source schema, T is a target schema,
- Σ is a set of constraints in some logical formalism expressing the relationship between S and T.

 $(I,J)\vDash \Sigma \quad \text{if and only if} \quad (I,J)\in \mathcal{W}.$ 

Definition: A solution for a source instance I with respect to  $\mathcal{M}$  is a target instance J such that  $(I, J) \in \mathcal{W}$  (or,  $(I, J) \models \Sigma$ ).

# Algorithmic Problems for Schema Mappings

Let  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W})$  be a fixed schema mapping.

Existence of Solutions Problem:

Given a source instance I,

- determine whether or not a solution for I w.r.t. *M* exists;
- if so, compute a "good" solution for I w.r.t.  $\mathcal{M}$ .

### Certain Answers Problem:

Let q be a fixed query over the target schema **T**.

Given a source instance I,

compute the certain answers of q on I w.r.t. M

 $cert(q, I, M) = \bigcap \{q(J) : J \text{ is a solution for } I \text{ w.r.t. } M\}.$ 

# Schema-Mapping Languages

Question: What is a "good" schema-mapping language? Ideally, a "good" schema-mapping language should have

- sufficient expressive power to express interesting data-transformation tasks;
- tractable algorithmic behavior.

Fact: The existence-of-solutions problem is undecidable for schema mappings specified by first-order sentences.

Reduction from the finite validity problem for FO.

## Basic Tasks for a Schema-Mapping Language

- ► Copy (Nicknaming):  $\forall x_1, \dots, x_n (P(x_1, \dots, x_n) \rightarrow R(x_1, \dots, x_n))$
- ► Projection:  $\forall x, y, z(P(x, y, z) \rightarrow R(x, y))$
- Column Augmentation:  $\forall x, y(P(x, y) \rightarrow \exists z R(x, y, z))$
- Decomposition:  $\forall x, y, z(P(x, y, z) \rightarrow R(x, y) \land T(y, z))$
- Join:

 $\forall x,y,z(E(x,z) \wedge F(z,y) \rightarrow R(x,y,z))$ 

• Combinations of the above, e.g.,  $\forall x, y, z(E(x, z) \land F(z, y) \rightarrow \exists wT(x, y, z, w)))$ 

# Global-and-Local-As-View Constraints

Definition: A Global-and-Local-as-View (GLAV) constraint is a first-order sentence of the form

 $\forall \mathbf{x}(\varphi(\mathbf{x}) \rightarrow \exists \mathbf{y}\psi(\mathbf{x}, \mathbf{y})), \text{ where }$ 

- $\varphi(\mathbf{x})$  is a conjunction of atoms over the source;
- $\psi(\mathbf{x}, \mathbf{y})$ ) is a conjunction of atoms over the target.

### Example:

 $\forall c, i, s(\mathsf{ENROLLS}(s, c) \land \mathsf{TEACHES}(i, c) \rightarrow \exists g \; \mathsf{GRADES}(s, c, g))$ 

### Fact:

Each basic task (projection, decomposition, join, ...) can be expressed by a GLAV constraint.

GLAV, GAV, LAV Constraints and Mappings

- ► Global-and-Local-as-View (GLAV) constraint  $\forall \mathbf{x}(\varphi(\mathbf{x}) \rightarrow \exists \mathbf{y}\psi(\mathbf{x}, \mathbf{y}))$
- Global-As-View (GAV) constraint

 $\forall \mathbf{x}(\varphi(\mathbf{x}) \rightarrow R(\mathbf{x}))$ , where *R* is a target relation.

- Copy, Projection, Join
- ► Local-As-View (LAV) constraint  $\forall \mathbf{x}(P(\mathbf{x}) \rightarrow \exists \mathbf{y}\psi(\mathbf{x}, \mathbf{y})), \text{ where } P \text{ is a source relation.}$ 
  - Copy, Column Augmentation, Decomposition
- A GLAV mapping is a schema mapping  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ , where  $\Sigma$  is a finite set of GLAV constraints.
- Similarly, for the notions of a GAV mapping and a LAV mapping.

# Structural Properties of GLAV Mappings

Theorem (Fagin, K ..., Miller, Popa - 2003)

Let  $\mathcal{M}$  =  $(\boldsymbol{S},\boldsymbol{T},\boldsymbol{\Sigma})$  be a GLAV mapping.

•  $\ensuremath{\mathcal{M}}$  admits universal solutions

For every source instance I, there is a solution  $J^*$  for I such that for every solution J for I, there is a homomorphism from  $J^*$  to J that is the identity on elements from I.

 Moreover, such a J\* can be computed in polynomial time in the size of I via the chase procedure.

## • $\mathcal{M}$ allows for CQ-rewriting

For every conjunctive query q over **T**, there is a union q' of conjunctive queries over **S** such that

 $\operatorname{cert}(q, \mathrm{I}, \mathcal{M}) = q'(\mathrm{I}).$ 

In particular, the certain answers of q can be computed in polynomial time in the size of I.

## Universal Solutions in Data Exchange



12/40

Structural Properties of GLAV, GAV, & LAV Mappings

Proposition:

 $\bullet$  Every GLAV mapping  $\mathcal M$  is closed under target homomorphisms

If J is a solution for I and there is a homomorphism from J to J' that is the identity on elements of I, then J' is a solution for I.

• Every GAV mapping  $\mathcal{M}$  is closed under target intersections If J and J' are solutions for I, then  $J \cap J'$  is a solution for I.

• Every LAV mapping  $\mathcal{M}$  is closed under unions If J is a solution for I and if J' is a solution for I', then  $J \cup J'$  is a solution for  $I \cup I'$ .

# Characterizing GAV and LAV Mappings

## Theorem (ten Cate and K ... - 2009)

Let  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W})$  be a schema mapping. Then the following statements are equivalent:

- M is logically equivalent to a GAV mapping (respectively, M is logically to a LAV mapping).
- 2.  $\mathcal{M}$  has the following properties:
  - *M* admits universal solutions;
  - *M* allows for CQ-rewriting;
  - *M* is closed under target homomorphisms;
  - *M* is closed under target intersections (respectively, *M* is closed under unions).

# Characterizing GLAV Mappings

## Theorem (ten Cate and K ... - 2009)

Let  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W})$  be a schema mapping. Then the following statements are equivalent:

- 1.  $\mathcal{M}$  is logically equivalent to a GLAV mapping.
- 2.  ${\mathcal M}$  has the following properties:
  - *M* admits universal solutions;
  - *M* allows for CQ-rewriting;
  - *M* is closed under target homomorphisms;
  - $\mathcal{M}$  is *n*-modular, for some  $n \ge 1$ .

## Definition

A schema mapping  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W})$  is *n*-modular if whenever  $(I, J) \notin \mathcal{W}$ , there is some  $I' \subseteq I$  such that  $|I'| \leq n$  and  $(I', J) \notin \mathcal{W}$ .

# Structural Characterizations: Summary

Type of Schema-Mapping	Characterizing Properties
GAV Mapping	admits universal solutions
	allows for CQ-rewriting
	closed under target homomorphisms
	closed under target intersections
LAV Mapping	admits universal solutions
	allows for CQ-rewriting
	closed under target homomorphisms
	closed under unions
GLAV Mapping	admits universal solutions
	allows for CQ-rewriting
	closed under target homomorphisms
	<i>n</i> -modular, for some $n \ge 1$

# Managing Schema Mappings via Operators

- Schema mappings can be quite complex.
- Methods and tools are needed to automate or semi-automate schema-mapping management.
- Metadata Management Framework Bernstein 2003

Based on generic schema-mapping operators:

- Match operator
- Merge operator
- Composition operator
- Inverse operator.
- Extensive study of the Composition operator and the Inverse operator.

# **Composing Schema Mappings**

Problem:

- Given M<sub>12</sub> = (S<sub>1</sub>, S<sub>2</sub>, Σ<sub>12</sub>) and M<sub>23</sub> = (S<sub>2</sub>, S<sub>3</sub>, Σ<sub>23</sub>), derive a schema mapping M<sub>13</sub> = (S<sub>1</sub>, S<sub>3</sub>, Σ<sub>13</sub>) that is "equivalent" to the sequential application of M<sub>12</sub> and M<sub>23</sub>.
- +  $\mathcal{M}_{13}$  is a composition of  $\mathcal{M}_{12}$  and  $\mathcal{M}_{23}$  , denoted

$$\mathcal{M}_{13} = \mathcal{M}_{12} \circ \mathcal{M}_{23}.$$

But, what does it mean to say that M<sub>13</sub> is "equivalent" to the composition of M<sub>12</sub> and M<sub>23</sub>?



# Semantics of Composition

- Metadata Model Management Bernstein 2003
  - Composition is one of the fundamental operators
  - However, no precise semantics is given.
- Composing Mappings among Data Sources Madhavan and Halevy - 2003
  - First to propose a semantics for composition
  - Notion of composition relative to a class of queries
  - CQ-composition: relative to the class of conjunctive queries
- Set-theoretic semantics of composition
  Fagin, K..., Popa, Tan 2004, Melnik 2004

# Set-theoretic Semantics of Composition

Recall that

- a syntactically specified M = (S,T,Σ) is identified with
- ► the semantically specified  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \mathcal{W}(\mathcal{M}))$ , where  $\mathcal{W}(\mathcal{M}) = \{(I, J) : (I, J) \models \Sigma\}.$

 $\begin{array}{l} \text{Definition} \ (\text{FKPT - 2004}, \text{Melnik - 2004}) \\ \text{A schema mapping} \ \mathcal{M}_{13} = (\textbf{S}_1, \textbf{S}_3, \boldsymbol{\Sigma}_{13}) \ \text{is the composition} \\ \mathcal{M}_{12} \circ \mathcal{M}_{23} \ \text{of} \ \mathcal{M}_{12} = (\textbf{S}_1, \textbf{S}_2, \boldsymbol{\Sigma}_{12}) \ \text{and} \ \mathcal{M}_{23} = (\textbf{S}_2, \textbf{S}_3, \boldsymbol{\Sigma}_{23}) \ \text{if} \\ \mathcal{W}(\mathcal{M}_{13}) = \mathcal{W}(\mathcal{M}_{12}) \circ \mathcal{W}(\mathcal{M}_{23}), \end{array}$ 

i.e.,

•  $(I_1, I_3) \in \mathcal{W}(\mathcal{M}_{13})$ 

if and only if

• there is some  $I_2$  such that  $(I_1, I_2) \in \mathcal{W}(\mathcal{M}_{12})$  and  $(I_2, I_3) \in \mathcal{W}(\mathcal{M}_{23})$ .

# The Language of Composition

### Questions:

- Is the language of GLAV constraints closed under composition?
   In other words:
- If  $M_{12}$  and  $M_{23}$  are GLAV mappings, is  $M_{12} \circ M_{23}$  a GLAV mapping as well?
- If not, what is the *right* language for composing schema mappings?

# The Language of Composition

## Theorem (Fagin, K ..., Popa, Tan - 2004)

- GAV mappings are closed under composition.
- GLAV mappings are not closed under composition.
- In fact, there are GLAV mappings M<sub>12</sub> and M<sub>23</sub> whose composition M<sub>12</sub> ∘ M<sub>23</sub> is not expressible even in least fixed-point logic LFP.

### Question:

What is the *right* language for composing GLAV mappings?

# Towards the "Right" Language for Composition

Motivating Example:

- ►  $\mathcal{M}_{12}$ :  $\forall e(\operatorname{Emp}(e) \to \exists m \operatorname{Rep}(e, m))$
- ►  $\mathcal{M}_{23}$ :  $\forall e \forall m(\operatorname{Rep}(e, m) \to \operatorname{Mgr}(e, m))$  $\forall e(\operatorname{Rep}(e, e) \to \operatorname{SelfMgr}(e))$

## Theorem:

- The composition  $\mathcal{M}_{12} \circ \mathcal{M}_{23}$  is not definable by any set (finite or infinite) of GLAV constraints.
- The composition  $\mathcal{M}_{12} \circ \mathcal{M}_{23}$  is definable by the following Second-Order GLAV constraint:

 $\exists f(\forall e(\operatorname{Emp}(e) \to \operatorname{Mgr}(e, f(e)) \land \\ \forall e(\operatorname{Emp}(e) \land (e = f(e) \to \operatorname{SelfMgr}(e)))))$ 

# Second Order GLAV Constraints

Definition: Let **S** be a source schema and **T** a target schema. A Second-Order GLAV constraint (SO GLAV) is a formula of the form:

$$\exists f_1 \cdots \exists f_n((\forall \mathbf{x}_1(\varphi_1(\mathbf{x}_1) \to \psi_1(\mathbf{x}_1))) \land \cdots \land (\forall \mathbf{x}_n(\varphi_n(\mathbf{x}_1) \to \psi_n(\mathbf{x}_n)))),$$

where

- Each  $f_i$  is a function symbol.
- Each φ<sub>i</sub> is a conjunction of atoms from S and equalities of terms.
- Each  $\psi_i$  is a conjunction of atoms from **T**.

Example:  $\exists f(\forall e(\operatorname{Emp}(e) \to \operatorname{Mgr}(e, f(e))) \land \forall e(\operatorname{Emp}(e) \land (e = f(e) \to \operatorname{SelfMgr}(e)))))$ 

Data Exchange via Second-Order GLAV Constraints

Theorem (Fagin, K ..., Popa, Tan - 2004)

- SO GLAV mappings are closed under composition.
- The chase procedure can be extended to SO GLAV mappings; in particular, it produces universal solutions in polynomial time.
- Every SO GLAV mapping is the composition of finitely many GLAV mappings (in fact, just two).

Conclusion: SO GLAV constraints are the right language for the composition of GLAV mappings.

Note: SO GLAV constraints and the composition algorithm are used in the IBM InfoSphere Information Server.

# More on SO GLAV Mappings

Definition: A plain SO GLAV constraint is a SO GLAV constraint with no equalities = and no nested function terms.

## Theorem (Arenas, Pérez, Reutter, Riveros - 2013)

Plain SO GLAV constraints are the right language for the CQ-composition of GLAV mappings.

### Open Problem 1:

- Is there a structural characterization of SO GLAV mappings?
- Is there a structural characterization of plain SO GLAV mappings?

# Inverting Schema Mappings



Problem: Given a schema mapping **M**, find a schema mapping **M**\* that "undoes" what **M** did.

# Exact Inverses of Schema Mappings

Definition: Fagin - 2006  $\mathcal{M}^*$  is an inverse of  $\mathcal{M}$  if  $\mathcal{M} \circ \mathcal{M}^* = Id$ , where Id is the identity schema mapping specified by copy constraints.

Note: Schema mappings may entail inherent information loss.

Union Schema Mapping

 $\forall \mathbf{x} (P(\mathbf{x}) \to Q(\mathbf{x})) \\ \forall \mathbf{x} (R(\mathbf{x}) \to Q(\mathbf{x}))$ 

Fact: Inverses of GLAV mappings rarely exist.

# Approximate Inverses of Schema Mappings

Several different approaches, including:

- Quasi-inverse
  Fagin, K ..., Popa, Tan 2007
- Maximum Recovery Arenas, Pérez, Riveros - 2008
- Chase Inverse Fagin, K ..., Popa, Tan - 2011

Note:

- Maximum recoveries have better properties than other notions of approximate inverses do.
- In particular, every plain SO GLAV mapping has a maximum recovery (hence, so does every GLAV mapping).

# Combining Composition and Inversion

## Fact:

- The language for expressing maximum recoveries involves disjunctive constraints.
- No definitive notion of an inverse has emerged.

## Open Problem 2:

 Find a useful notion of inverse and a language for expressing it, so that the language is closed under compositions and inversions of GLAV mappings.

# **Evolution of Schema Mappings**



**Fact:** Schema evolution can be analyzed using the composition operator and the inverse operator.

## Schema Mappings Can Be Complex

Map 2:

	for sm2xA in SA dummy COUNTRY 4			
	viete tm2x0 in S27 dimmy country 10, tm2x1 in S27 dimmy organize 12			
	whor	ses (m2x0 country membershiptmext) organization id		
	entiof	where chizko.country.membership=thizki.organization.it,		
Satisi Smizko.country.angA=tmizko.country.arga, Smizko.country.arga		SHIZAO, COUNTRY, ANEA-LHIZAO, COUNTRY, ATBA, SHIZAO, COUNTRY, CARTIAL-LHIZAO, COUNTRY, CARTIAL,		
		Silizad.country.courty.rut, Silizad.country.ruter.sec.country.name,		
	sm2x0.COUNTRY.POPULATION=Tm2x0.Country.population,(			
	мар з:			
		for smsx0 in S0.dummy GEO RIVE 23, Sm3x1 in S0.dummy RIVER 24,		
		sm3x2 in S0.dummy_PROVINCE_5		
		where sm3x0.GEO_RIVER.RIVER=sm3x1.RIVER.NAME, sm3x2.PROVINCE.NAME=sm3x0.GEO_RIVER.PROVINCE,		
		sm3x2.PROVINCE.COUNTRY=sm2x0.COUNTRY.CODE,		
		exists tm3x0 in S27.dummy_river_24, tm3x1 in tm3x0.river.dummy_located_23,		
		tm3x4 in S27.dummy_country_10, tm3x5 in tm3x4.country.dummy_province_9,		
	tm3x6 in S27.dummy_organiza_13			
		where tm3x4.country.membership=tm3x6.organization.id, tm3x5.province.id=tm3x1.located.province,		
		tm2x0.country.id=tm3x1.located.country,		
		<pre>satisf sm2x0.COUNTRY.AREA=tm3x4.country.area, sm2x0.COUNTRY.CAPITAL=tm3x4.country.capital,</pre>		
		<pre>sm2x0.COUNTRY.CODE=tm3x4.country.id, sm2x0.COUNTRY.NAME=tm3x4.country.name,</pre>		
		sm2x0.COUNTRY.POPULATION=tm3x4.country.population, sm3x1.RIVER.LENGTH=tm3x0.river.length,		
		sm3x0.GE0 RIVER.COUNTRY=tm3x1.located.country, sm3x0.GE0 RIVER.PROVINCE=tm3x1.located.province,		
		sm3x1.RIVER.NAME=tm3x0.river.name ),(		
	Map 4:			
		for sm4x0 in S0.dummy GEO ISLA 25, sm4x1 in S0.dummy ISLAND 26,		
		sm4x2 in S0.dummy PROVINCE 5		
		where sm4x0.GEO ISLAND.ISLAND=sm4x1.ISLAND.NAME. sm4x2.PROVINCE.NAME=sm4x0.GEO ISLAND.PROVINCE.		
		sm4x2.PROVINCE.COUNTRY=sm2x0.COUNTRY.CODE,		
		exists tm4x0 in S27.dummy island 26, tm4x1 in tm4x0.island.dummy located 25,		
		tm4x4 in S27.dummy country 10, tm4x5 in tm4x4.country.dummy province 9,		
	tm4x6 in S27.dummy organiza 13			
	where tm4x4.country.membership=tm4x6.organization.id, tm4x5.province.id=tm4x1.located.province.			
		tm2x0.country.id=tm4x1.located.country.		
		satisf sm2x0.COUNTRY.AREA=tm4x4.country.area, sm2x0.COUNTRY.CAPITAL=tm4x4.country.capital.		
		sm2x0.COUNTRY.CODE=tm4x4.country.id, sm2x0.COUNTRY.NAME=tm4x4.country.name.		
		sm2x0.COUNTRY.POPULATION=tm4x4.country.population.sm4x1.ISLAND.AREA=tm4x0.island.area.		
		sm4x1_ISLAND.COORDINATESLAT=tm4x0_island_latitude.sm4x0_GE0_ISLAND.COUNTRY=tm4x1.located.country.		
		sm4x0.GE0_ISLAND.PROVINCE=tm4x1.located.province.sm4x1.ISLAND.COORDINATESLONG=tm4x0.island.longitude.		
		sm4x1.TSFAND.NAME=tm4x0.island.name).(		
	Map 5:			
		for sm5x0 in S0.dummy GEO SEA 19. sm5x1 in S0.dummy SEA 20.		
	sm5x2 in S0.dummy PROVINCE 5			
	where sm5x2.PROVINCE.NAME=sm5x0.GE0 SEA.PROVINCE. sm5x0.GE0 SEA.SEA=sm5x1.SEA.NAME.			
		Sm5x2 PROVINCE COUNTRY_Sm2x0 COUNTRY_CODE		
		exists tm5x0 in S27.dummy sea 19. tm5x1 in tm5x0.sea.dummy located 18.		
		tm5x4 in S27 dummy country 10, tm5x5 in tm5x4 country dummy province 9		
		tm5v6 in S27 dummy organiza 13		
		where tm5x4.country.membership=tm5x6.organization.id. tm5x5.province.id=tm5x1.located.province.		
	tm2x0 country idetm5x1 located country			
		entice country, inclusions a concern, entited country, entited entities and a country contraction of the country		
		autar ameno.cookinti.hteh-emox4.country.area, ameno.cookinti.ckFiikE-tmox4.country.capitat,		

32/40

# Understanding and Deriving Schema Mappings

#### Idea:

Use data examples to understand/derive schema mappings

## Theme I: From Syntax to Semantics

Understand schema mappings using data examples.

## Theme II: From Semantics to Syntax:

Derive schema mappings using data examples.

## Data Examples and Universal Examples

Definition: Let  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$  be a schema mapping.

- A data example for *M* is a pair (I, J) such that J is a solution for I w.r.t. *M* (i.e., (I, J) ⊨ Σ).
- ► A universal example for  $\mathcal{M}$  is a pair (I,J) such that J is a universal solution for I w.r.t.  $\mathcal{M}$ .

Note: The space of data examples and the space of universal examples are typically infinite.

Question: Can a schema mapping  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$  be captured by finitely many data examples or by finitely many universal examples?

# From Syntax to Semantics: Unique Characterizations

Definition: Let  $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$  be a schema mapping, **U** set of universal examples for  $\mathcal{M}$ , and **C** a class of GLAV constraints.

U uniquely characterizes *M* w.r.t. C if for every schema mapping *M*' = (S, T, Σ') such that Σ' ⊆ C and U is a set of universal examples for *M*', we have that Σ ≡ Σ'.

### Note:

 Unique characterizability via finitely many positive/negative examples implies unique characterizability via finitely many universal examples, but not vice-versa. From Syntax to Semantics: Unique Characterizations

Theorem (Alexe, ten Cate, K ..., Tan - 2011)

- Every LAV mapping is uniquely characterizable by a finite set of universal examples w.r.t. to LAV constraints.
- Criterion for a GAV mapping to be uniquely characterizable by a finite set of universal examples w.r.t. GAV constraints.
- The associated decision problem for GAV mappings is NP-complete.

## Open Problem 3:

- Find criteria for a GLAV mapping to by uniquely characterizable by a finite set of universal examples w.r.t. GLAV constraints.
- What is the exact complexity of the associated decision problem for GLAV mappings?

# From Semantics to Syntax: Derivations

Interactive Derivation and Refinement of Schema Mappings

- Fitting Algorithm EIRENE System Alexe, ten Cate, K ..., Tan - 2011
- Interacting Mapping Specification with Exemplar Tuples -IMS System Bonifati, Comignani, Coquery, Thion - 2017

Deriving Optimal Schema Mappings from Data Examples

- Cost Model for Optimal Repairs of Schema Mappings Gottlob and Senellart - 2010
- Cost Model for Mapping Selection via Data Examples Kimmig, Memory, Miller, Getoor - 2017

From Semantics to Syntax: Learning

## Theorem (ten Cate, Dalmau, K ... - 2012)

GAV mappings are efficiently learnable in Angluin's model with membership and equivalence queries.

Active Learning of GAV Mappings ten Cate, K ..., Qian, Tan - PODS 2018

 Algorithm uses conformance testing as a substitute for the equivalence oracle.

Open Problem 4:

- Learnability of LAV mappings.
- Learnability of GLAV mappings.

## Topics Not Covered - Partial List

- Richer schema mappings:
  - Schema mappings with target constraints
  - Schema mappings with arithmetic constraints
  - Schema mappings with bi-directional constraints.
- Alternative notions of solutions Libkin 2006
- Alternative notions of certain answers
- Beyond conjunctive queries (non-monotonic, aggregate)
- XML Data Exchange Arenas and Libkin 2005
- Benchmarks
  - STBenchmark Alexe, Tan, Velegrakis 2008
  - iBench Arocena, Glavic, Ciucanu, Miller 2015
  - Benchmarking the Chase Benedikt et al. 2017.

# Synopsis, Challenges, and Outlook

## Synopsis

- Mature body of research during the past 15 years.
- A case study of logic *in* computer science, but also of logic *from* computer science.
- Theory and practice have informed each other.

## Challenges and Outlook

- Several key technical problems remain open.
- Little penetration of advanced technical findings to practice.
- Benchmarks and open-source systems for data exchange need to be developed.