

# Propagating Dependencies under Schema Mappings

## – A Graph-based Approach

**Qing Wang**

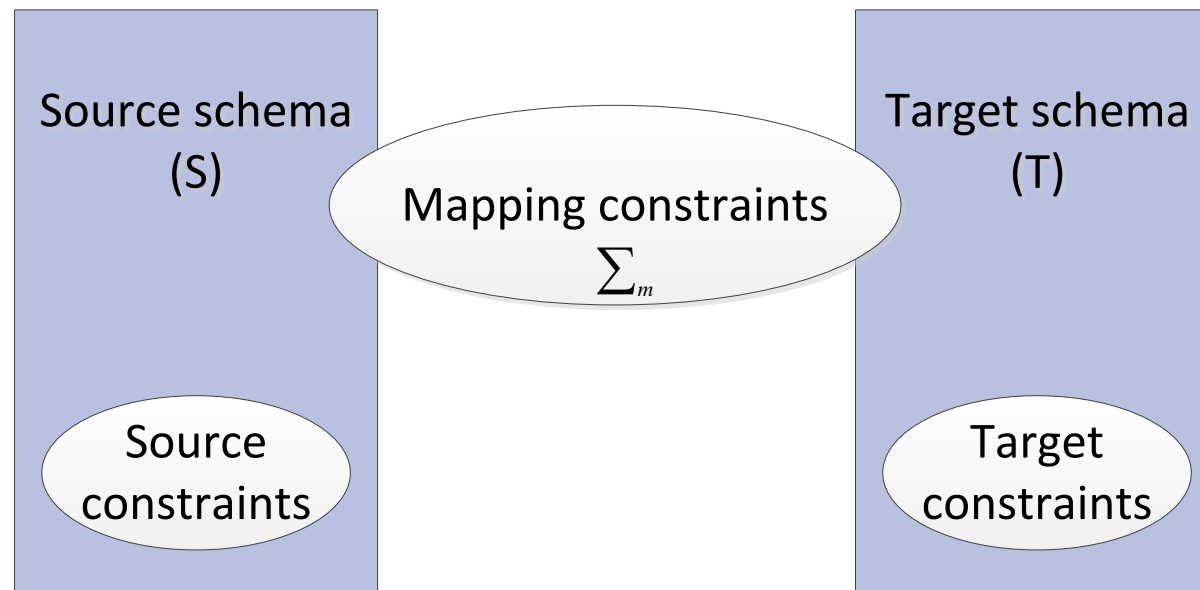
Research School of Computer Science  
Australian National University  
Canberra ACT 0200, Australia  
qing.wang@anu.edu.au

**Xi Wen**

Department of Computer Science  
Nanchang University  
Nanchang City, China  
wenxi@ncu.edu.cn

## Schema Mappings – Definition

- Schema mapping plays an important role in many database-related transformation tasks, such as data exchange, data integration and data migration.
- A *schema mapping* is a triple  $M = (S, T, \Sigma_m)$  consisting of
  - a *source schema*  $S$ ,
  - a *target schema*  $T$ , and
  - a set  $\Sigma_m$  of *mapping constraints* over  $S$  and  $T$ .



# Schema Mappings – Background

- However, designing a schema mapping is not an easy task. Generally, two lines of research exist:
  - (1) *Generate schema mappings from a visual specification* provided by users (traditionally studied);  
However, a visual specification is often ambiguous.
  - (2) *Derive schema mappings based on data examples* (attracted more interest in recent years).  
However, data examples may not be available, or could be biased.
- Existing approaches either require a manual process of tuning schema mappings or demand more data examples for improving accuracy.
- Even though, the design quality of schema mappings is still often not satisfactory.

# Schema Mappings – Questions

- Some common questions in practice:
  - (1) Can we ensure certain properties of a source database to be preserved in the desired target database through the design of a schema mapping?
  - (2) Can we determine whether or not a target constraint can be enforced on a target database before the target database is transformed from a given source database?
  - (3) If some target constraints cannot be enforced, can we efficiently identify which data in the source database need to be cleansed, or determine whether the schema mapping and target constraints need to be re-designed?
  - (4) ...

# Our Research

- **General goal** of our research:

To develop methods/tools that help check whether a schema mapping is designed meaningfully and effectively in advance, before an implementation takes place.

- **Specific task** in this paper:

Given a schema mapping, a question is:

How can we discover logical consequences among its source, target and mapping constraints?

# Schema Mappings – Motivating Example

- Source schema  $S$

- $\text{RENTCLIENT}(id, name, address)$ ,  $\text{RENTPROPERTY}(no, address, rent)$ , and  $\text{ALLCLIENT}(name, dob, gender, cid)$
- $\Sigma_s = \emptyset$

- Target schema  $T$

- $\text{PROPERTY}(no, address)$ ,  $\text{CLIENT}(id, name)$ , and  $\text{RENTAL}(id, no, rent)$
- $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

- Mapping constraints  $\Sigma_m$

- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$
- (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$
- (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$
- (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$
- (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$
- (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$

## Example – Schema Mappings

- Source schema  $S$

- $\text{RENTCLIENT}(id, name, address)$ ,  $\text{RENTPROPERTY}(no, address, rent)$ , and  $\text{ALLCLIENT}(name, dob, gender, cid)$
- $\Sigma_s = \emptyset$

- Target schema  $T$

- $\text{PROPERTY}(no, address)$ ,  $\text{CLIENT}(id, name)$ , and  $\text{RENTAL}(id, no, rent)$
- $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

- Mapping constraints  $\Sigma_m$

- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$
- (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$
- (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$
- (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$
- (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$
- (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$

*How much semantics specified by  $\Sigma_t$  and  $\Sigma_m$  is captured by  $\Sigma_s$ ?*

## Example – Schema Mappings

- Source schema  $S$

- $\text{RENTCLIENT}(id, name, address)$ ,  $\text{RENTPROPERTY}(no, address, rent)$ , and  $\text{ALLCLIENT}(name, dob, gender, cid)$
- $\Sigma_s = \emptyset$

- Target schema  $T$

- $\text{PROPERTY}(no, address)$ ,  $\text{CLIENT}(id, name)$ , and  $\text{RENTAL}(id, no, rent)$
- $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

- Mapping constraints  $\Sigma_m$

- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$
- (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$
- (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$
- (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$
- (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$
- (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$

*How much semantics specified by  $\Sigma_t$  and  $\Sigma_m$  is captured by  $\Sigma_s$ ?*

–  $\Sigma_s^\dagger = \{\text{RENTPROPERTY} : no \rightarrow rent\}$



## Example – Schema Mappings

- By comparing  $\Sigma_s$  and  $\Sigma_s^\dagger$ ,
  - $\Sigma_s = \emptyset$
  - $\Sigma_s^\dagger = \{\text{RENTPROPERTY} : no \rightarrow rent\}$

we know that every source instance that violates *RentProperty* :  $no \rightarrow rent$  must violate  $\Sigma_t$  after being transformed into a target instance.

RENTCLIENT		
<i>id</i>	<i>name</i>	<i>address</i>
c1	Tim Jenkin	5 Jicket St, Dunedin
c2	Linda Lee	36 Novar St, Dunedin
c3	Mike Carl	2 Manor St, Dunedin

RENTPROPERTY		
<i>no</i>	<i>address</i>	<i>rent</i>
1	5 Jicket St, Dunedin	350
1	5 Jicket St, Dunedin	500
2	2 Manor St, Dunedin	450

- It means that either RENTCLIENT needs to be cleansed, or target constraints need to be reconsidered.

## Example – Schema Mappings

- Source schema  $S$

- $\text{RENTCLIENT}(id, name, address)$ ,  $\text{RENTPROPERTY}(no, address, rent)$ , and  $\text{ALLCLIENT}(name, dob, gender, cid)$
- $\Sigma_s = \emptyset$

- Target schema  $T$

- $\text{PROPERTY}(no, address)$ ,  $\text{CLIENT}(id, name)$ , and  $\text{RENTAL}(id, no, rent)$
- $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

- Mapping constraints  $\Sigma_m$

- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$
- (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$
- (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$
- (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$
- (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$
- (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$

How much semantics specified by  $\Sigma_s$  and  $\Sigma_m$  is preserved by  $\Sigma_t$ ?

## Example – Schema Mappings

- Source schema  $S$

- $\text{RENTCLIENT}(id, name, address)$ ,  $\text{RENTPROPERTY}(no, address, rent)$ , and  $\text{ALLCLIENT}(name, dob, gender, cid)$
- $\Sigma_s = \emptyset$

- Target schema  $T$

- $\text{PROPERTY}(no, address)$ ,  $\text{CLIENT}(id, name)$ , and  $\text{RENTAL}(id, no, rent)$
- $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

- Mapping constraints  $\Sigma_m$

- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$
- (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$
- (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$
- (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$
- (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$
- (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$

*How much semantics specified by  $\Sigma_s$  and  $\Sigma_m$  is preserved by  $\Sigma_t$ ?*

–  $\Sigma_t^\dagger = \{\text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

## Example – Schema Mappings

- By comparing  $\Sigma_t$  and  $\Sigma_t^\dagger$ ,
  - $\Sigma_t = \{\text{RENTAL} : no \rightarrow rent, \text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$
  - $\Sigma_t^\dagger = \{\text{RENTAL}[id] \subseteq \text{CLIENT}[id], \text{RENTAL}[no] \subseteq \text{PROPERTY}[no]\}$

we know that  $\text{RENTAL}[id] \subseteq \text{CLIENT}[id]$  and  $\text{RENTAL}[no] \subseteq \text{PROPERTY}[no]$  can hold on every target instance under this schema mapping.

CLIENT	
<i>id</i>	<i>name</i>
c1	Tim Jenkin
c2	Linda Lee
c3	Mike Carl

PROPERTY	
<i>no</i>	<i>address</i>
1	5 Jicket St, Dunedin
2	2 Manor St, Dunedin

RENTAL		
<i>id</i>	<i>no</i>	<i>rent</i>
c1	1	500
c3	2	450

- Hence we only need to check whether or not  $\text{RENTAL} : no \rightarrow rent$  holds on the target instance.

## Specifying Mapping Constraints

- A **tuple-generating dependency (tgd)** is an expression of the form

$$\forall \bar{x}, \bar{y}. (\varphi(\bar{x}, \bar{y}) \Rightarrow \exists \bar{z}. \psi(\bar{x}, \bar{z})),$$

where  $\varphi$  and  $\psi$  are conjunctions of atoms, but only relational atoms occur in  $\psi$ .

- $\varphi$  is called the **premise** of  $\sigma$ , and
  - $\psi$  is called the **conclusion** of  $\sigma$ .
- Mapping constraints  $\Sigma_m$  are **bipartite tgds** over  $S$  and  $T$ , i.e., tgds with  $\varphi$  and  $\psi$  over  $S$  and  $T$ , respectively, or vice versa.
  - The co-existence of target-to-source and source-to-target tgds in  $\Sigma_m$  enables us to precisely build a high-level specification for the relationship between source and target instances.

# Propagation Graphs

- A **propagation graph**  $G = (V, E)$  consists of
  - A set  $V$  of vertices, where each  $v \in V$  is a relation schema name  $R$ , and
  - A set  $E$  of edges, where each  $R \hookrightarrow_f R' \in E$  is labelled by a function

$$f : attr(R) \mapsto attr(R').$$

- There are two types of edges: **approximate** and **exact**.

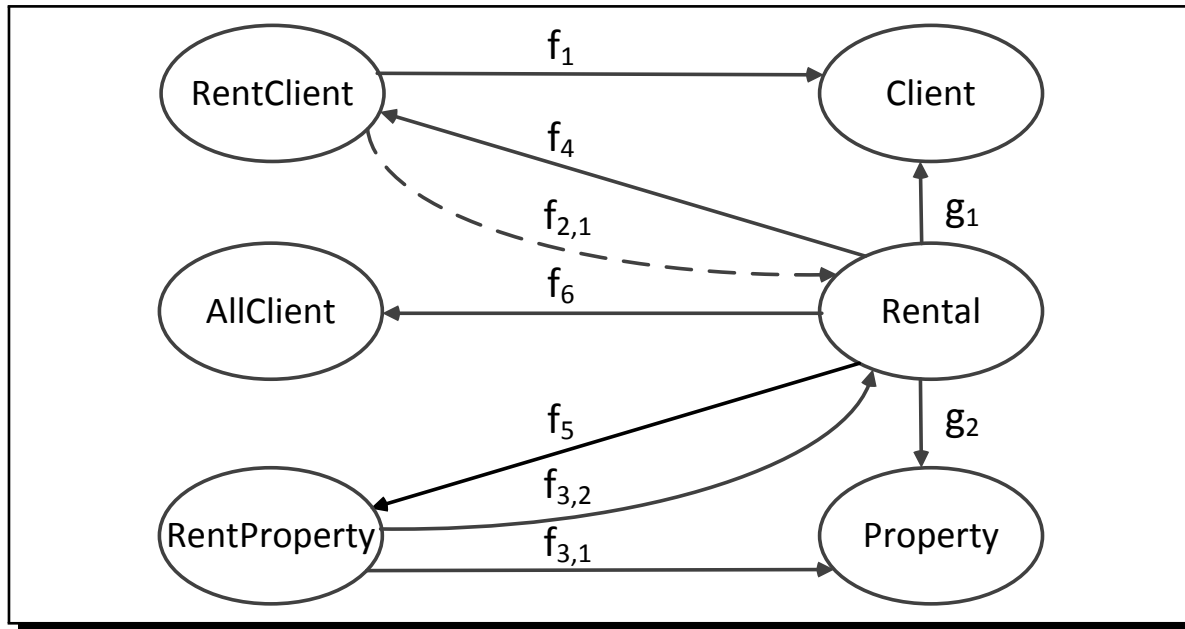
- A **propagation path** with a label  $f$  is a sequence of edges

$$\langle R_1 \hookrightarrow_{f_1} R_2, \dots, R_{n-1} \hookrightarrow_{f_{n-1}} R_n \rangle$$

such that  $f = f_{n-1} \circ \dots \circ f_1$  is a function that maps a *non-empty subset of attributes* of  $R_1$  into attributes of  $R_n$ .

## Propagation Graphs – Example

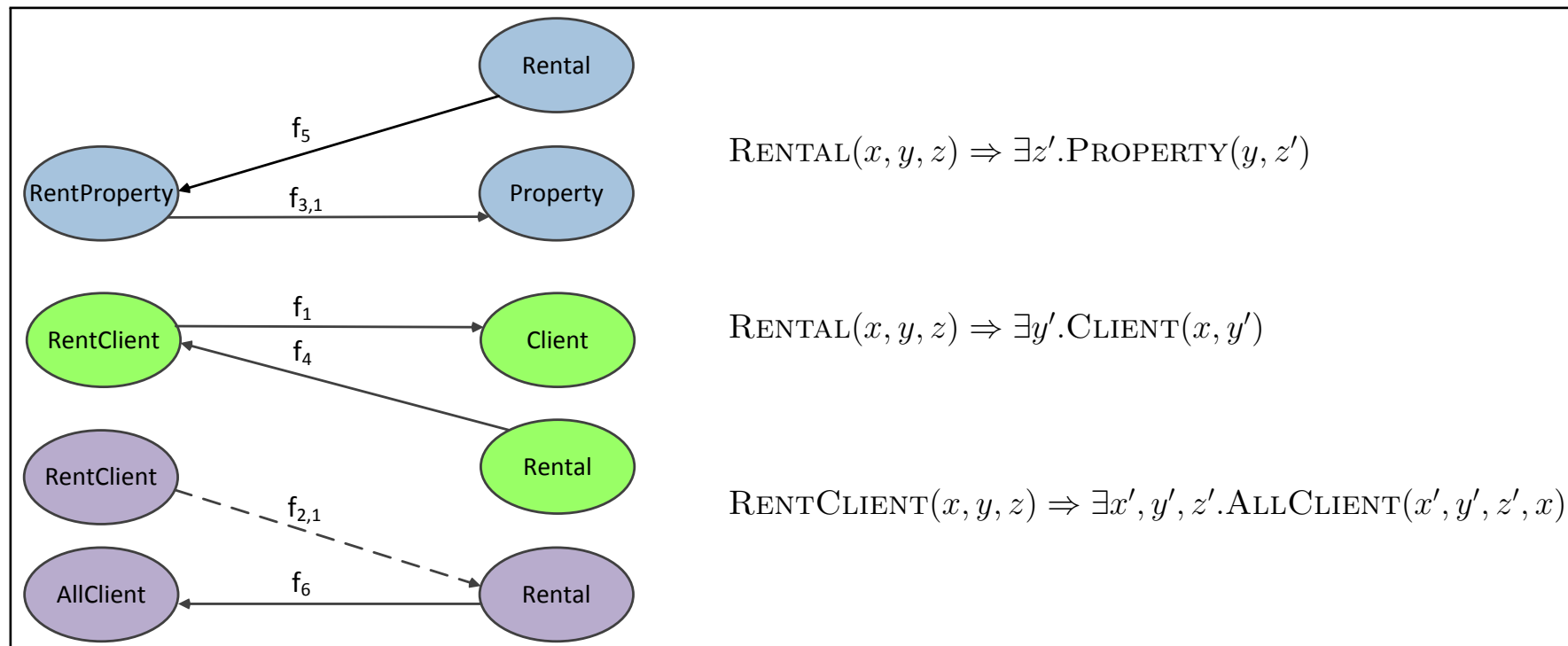
- (C1)  $\forall x, y, z. (\text{RENTCLIENT}(x, y, z) \Rightarrow \text{CLIENT}(x, y));$   
 (C2)  $\forall x, y, z, x', z'. (\text{RENTCLIENT}(x, y, z) \wedge \text{RENTPROPERTY}(x', z, z') \Rightarrow \text{RENTAL}(x, x', z'));$   
 (C3)  $\forall x, y, z. (\text{RENTPROPERTY}(x, y, z) \Rightarrow \exists x'. \text{PROPERTY}(x, y) \wedge \text{RENTAL}(x', x, z)).$   
 (C4)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists z'. \text{RENTPROPERTY}(y, z', z));$   
 (C5)  $\forall x, y, z. (\text{RENTAL}(x, y, z) \Rightarrow \exists y', z'. \text{RENTCLIENT}(x, y', z'));$   
 (C6)  $\text{RENTAL}(x, y, z) \Rightarrow \exists x', y', z'. \text{ALLCLIENT}(x', y', z', x).$   
 $\Sigma_t: \text{RENTAL}[id] \subseteq \text{CLIENT}[id] \text{ and } \text{RENTAL}[no] \subseteq \text{PROPERTY}[no].$



LABELS OF EDGES
$f_1 : 1 \mapsto 1, 2 \mapsto 2.$
$f_{2,1} : 1 \mapsto 1.$
$f_{2,2} : 1 \mapsto 2, 3 \mapsto 3.$
$f_{3,1} : 1 \mapsto 1, 2 \mapsto 2.$
$f_{3,2} : 1 \mapsto 2, 3 \mapsto 3.$
$f_4 : 1 \mapsto 1.$
$f_5 : 2 \mapsto 1, 3 \mapsto 3.$
$g_1 : 1 \mapsto 1.$
$g_2 : 2 \mapsto 1.$

## Propagating INDs

- Each propagation path from  $R_1$  to  $R_2$  over  $T$  corresponds to an IND over  $T$ .
- Depending on the types of the edges, an IND is either exact or approximate.
- Approximate INDs are prompted to users for fine-tuning the semantics of a schema mapping.





# Propagating FDs

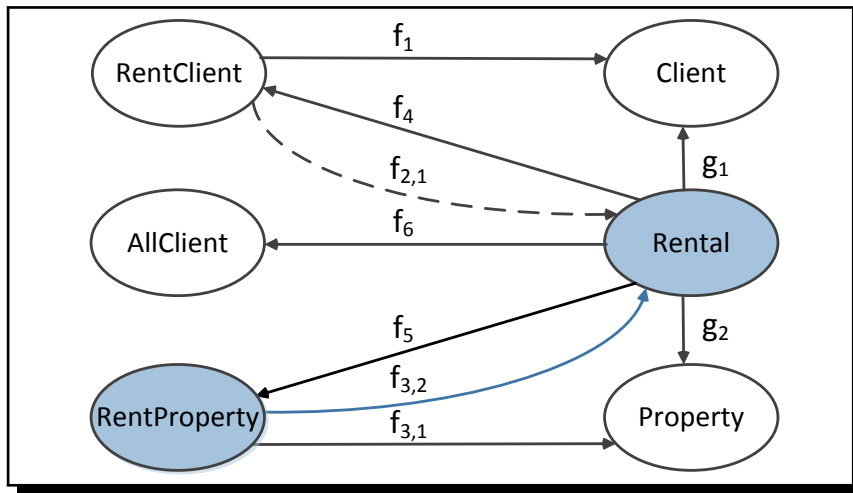
- Let  $R : X \rightarrow Y$  be a FD over  $T$ , and  $R' \in S$

- Push backward**

If there is a propagation path  $\langle R', \dots, R \rangle$  labelled by  $f$ , where  $X' = \{f^{-1}(x) | x \in X\}$  and  $Y' = \{f^{-1}(y) | y \in Y\}$ , then  $R' : X' \rightarrow Y'$  over  $S$  can be propagated.

- Push forward**

If there is a propagation path  $\langle R, \dots, R' \rangle$  labelled by  $f$ , where  $X' = \{f(x) | x \in X\}$  and  $Y' = \{f(y) | y \in Y\}$ , then  $R' : X' \rightarrow Y'$  over  $S$  can be propagated *with the condition on a subset of tuples in a relation of  $R'$* .



$$\text{RENTAL} : no \rightarrow rent \in \Sigma_t$$

$$\Downarrow$$

$$\Sigma_s^\dagger = \{\text{RENTPROPERTY} : no \rightarrow rent\}$$

# Experiments

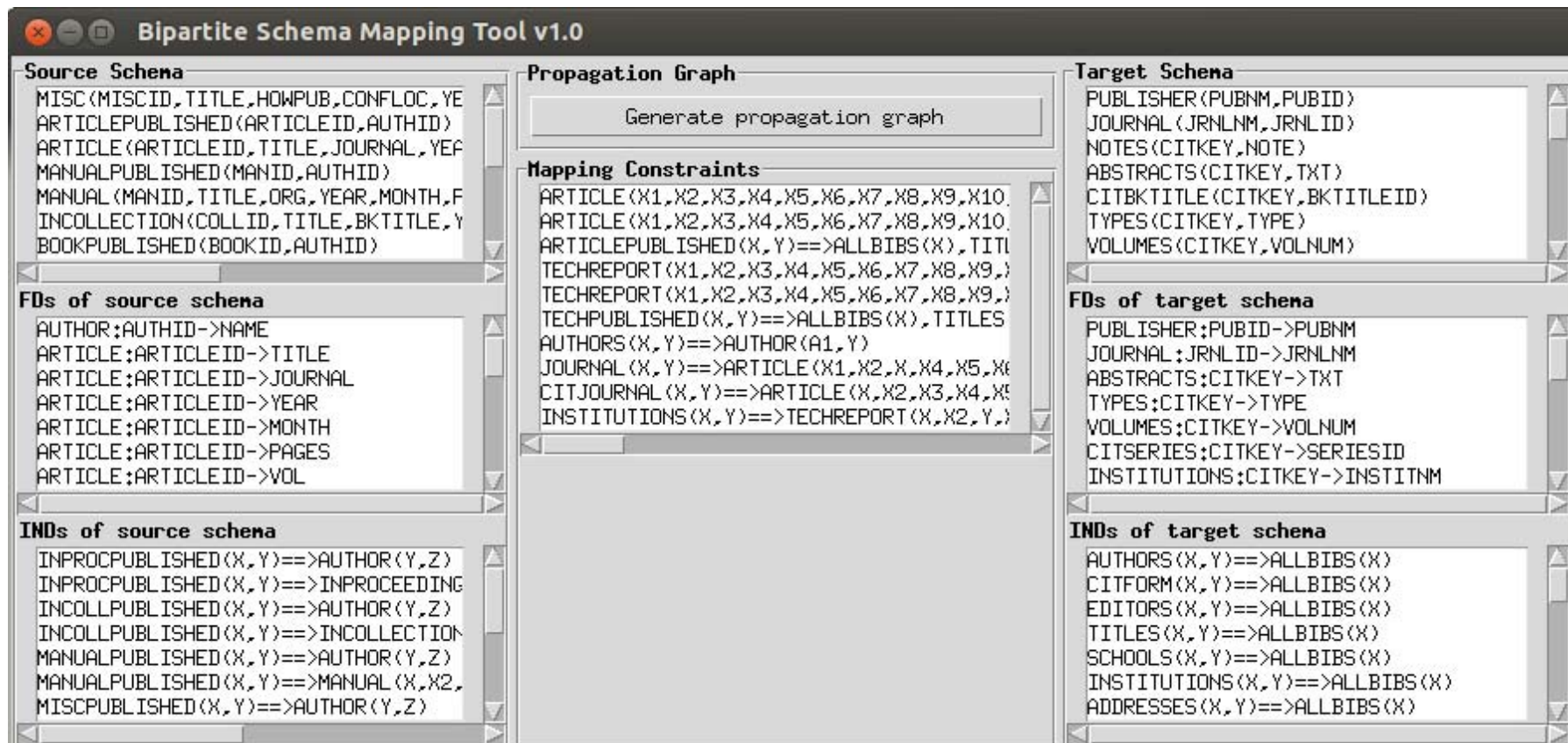
- We have developed a bipartite schema mapping (BSM) tool to help schema mapping designers in several aspects:
  - (1) Visualizing propagation graphs for any given schema mappings;
  - (2) Assessing the design quality of schema mappings;
  - (3) Facilitating the data cleaning tasks of source instances.
- We have conducted our experiments over two schema mapping data sets:
  - **RENTALAPP**: as previously described;
  - **AMALGAM**: taken from the Clio Project<sup>1</sup>.

	SOURCE SCHEMA ( $S_1$ )	TARGET SCHEMA ( $S_2$ )
No of relations	15	27
No of INDs	14	26
No of FDs	23	21
No of MCs	10	

<sup>1</sup> <http://dmlab.cs.toronto.edu/miller/amalgam/>

## Experiments – User Interface

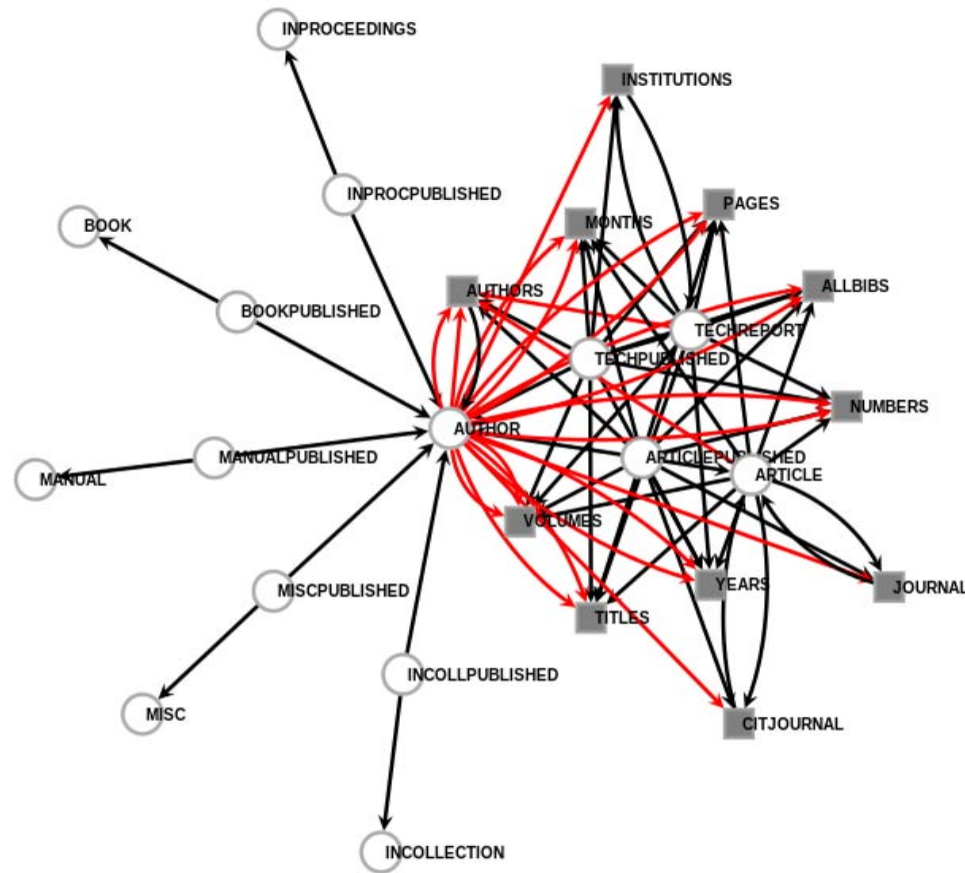
- A schema mapping over AMALGAM:  $S_1$  is the source schema and  $S_2$  is the target schema.





# Experiments – Propagation Graphs

- The propagation graph for deriving target constraints over  $S_2$  in AMALGAM.



## Conclusions

- We have developed a graphical model to represent the inter-relationships among the attributes of relation schemas.
- On that basis, we have studied the dependency propagation problem in the context of schema mappings.
- Mapping constraints of a schema mapping are permitted to be bipartite TGDs, which enables us to precisely specify the relationship between source and target databases.
- Our work can lead to a conceptual analysis tool that exploits the semantics of a schema mapping through propagation paths in propagation graphs.
- In doing so, the design quality of schema mappings can be assessed before actually implementing them.

Thanks and Questions