

Information and *redundancy*:

fundamental concepts in schema mapping management

Information and *redundancy*:

fundamental concepts in schema mapping management

legacy DB

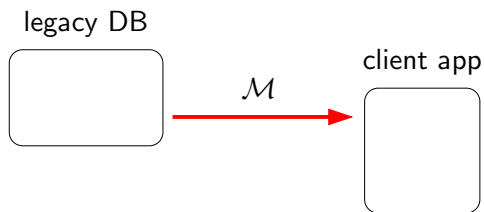


client app



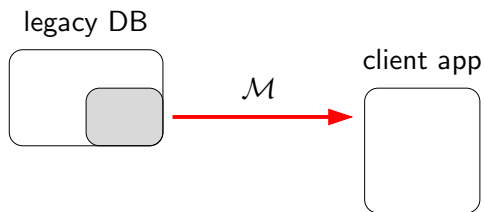
Information and *redundancy*:

fundamental concepts in schema mapping management



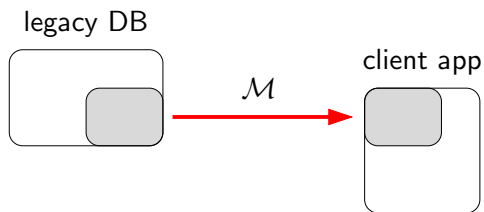
Information and *redundancy*:

fundamental concepts in schema mapping management

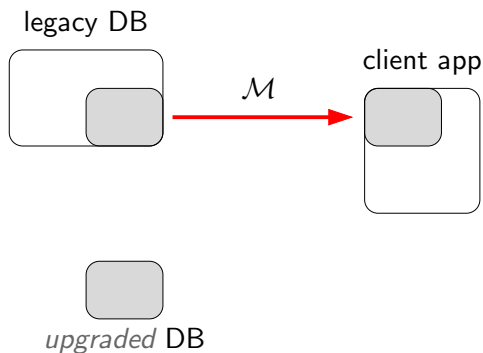


Information and *redundancy*:

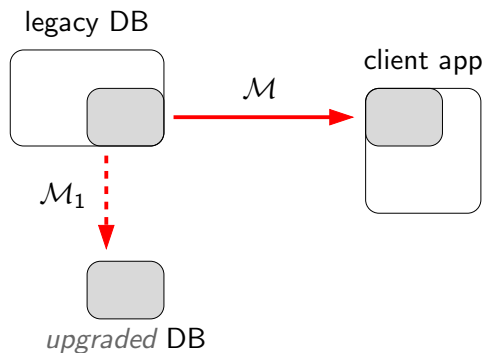
fundamental concepts in schema mapping management



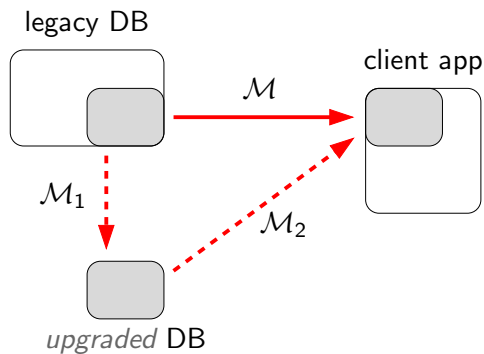
Information and redundancy:
fundamental concepts in schema mapping management



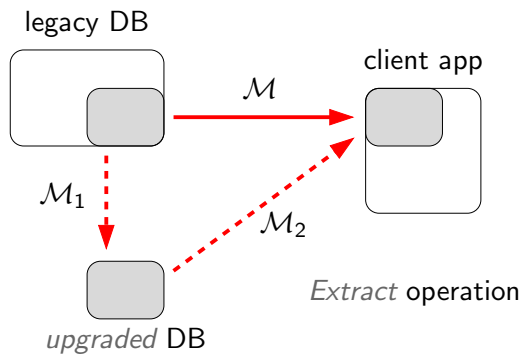
Information and redundancy:
fundamental concepts in schema mapping management



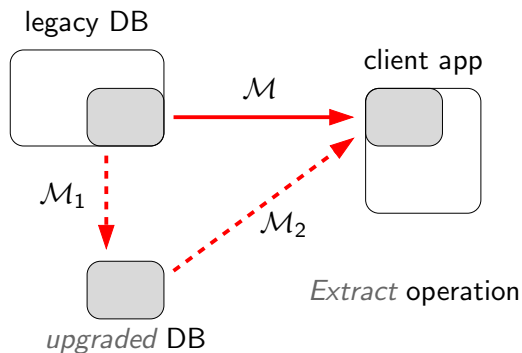
Information and redundancy:
fundamental concepts in schema mapping management



Information and redundancy:
fundamental concepts in schema mapping management

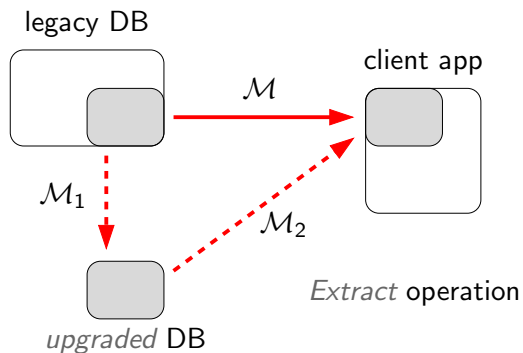


Information and redundancy:
fundamental concepts in schema mapping management



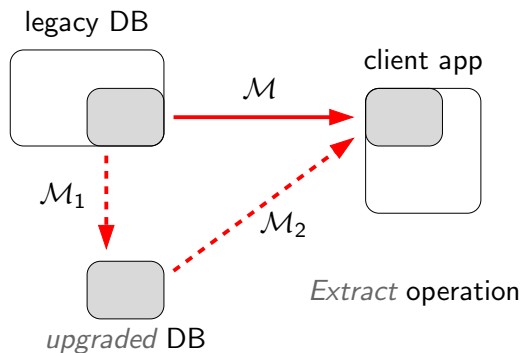
The new DB should only contain the information transferred by \mathcal{M} .

Information and redundancy:
fundamental concepts in schema mapping management



The new DB should only contain the information transferred by \mathcal{M} .

Information and redundancy:
fundamental concepts in schema mapping management



The new DB should **only** contain **the information transferred** by \mathcal{M} .

Information and *redundancy*:

fundamental concepts in schema mapping management

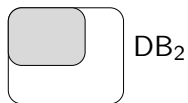
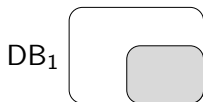
Information and *redundancy*:

fundamental concepts in schema mapping management



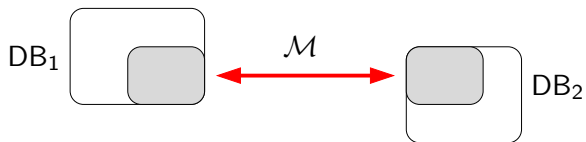
Information and *redundancy*:

fundamental concepts in schema mapping management



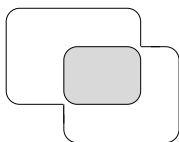
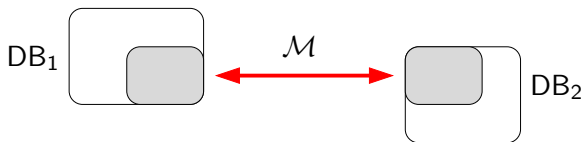
Information and *redundancy*:

fundamental concepts in schema mapping management



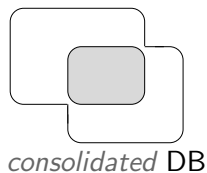
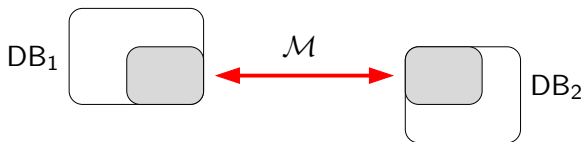
Information and *redundancy*:

fundamental concepts in schema mapping management



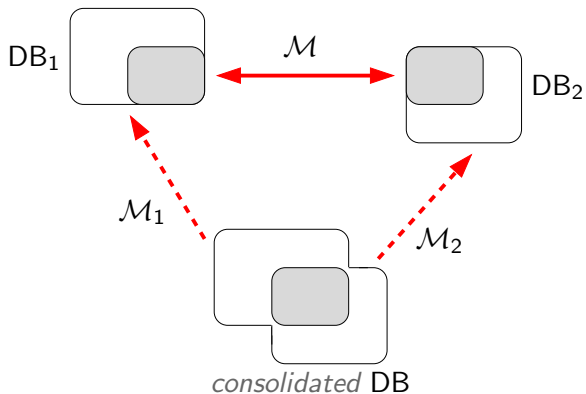
Information and redundancy:

fundamental concepts in schema mapping management

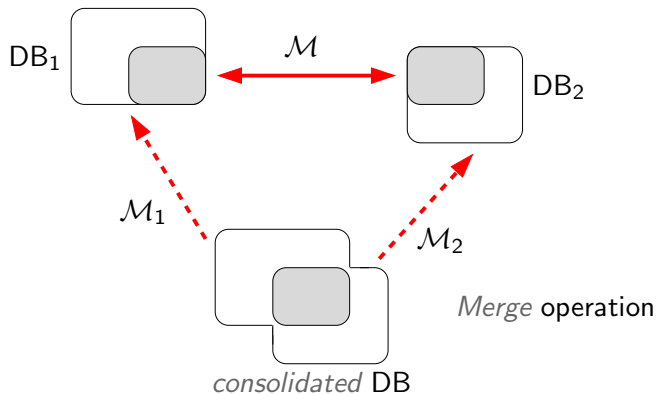


Information and redundancy:

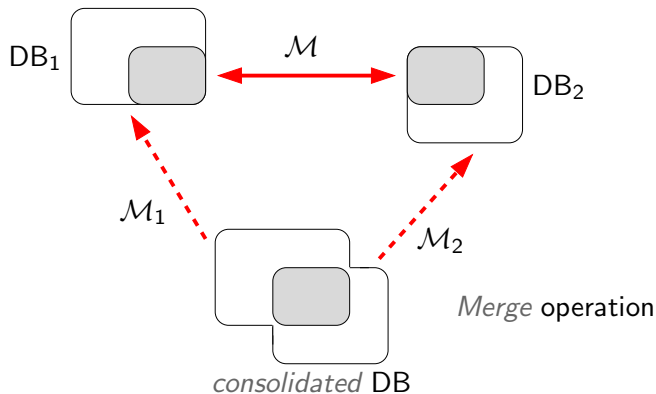
fundamental concepts in schema mapping management



Information and redundancy:
fundamental concepts in schema mapping management

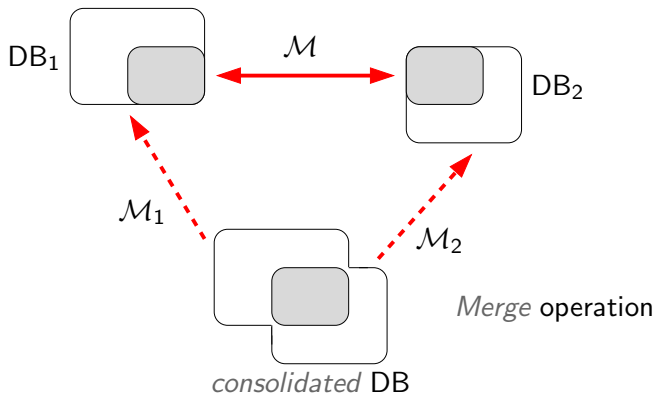


Information and redundancy:
fundamental concepts in schema mapping management



The new DB should only store the non redundant information w.r.t. \mathcal{M} .

Information and redundancy:
fundamental concepts in schema mapping management



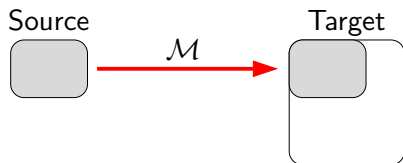
The new DB should only store the **non redundant information** w.r.t. \mathcal{M} .

Information and *redundancy*:

fundamental concepts in schema mapping management

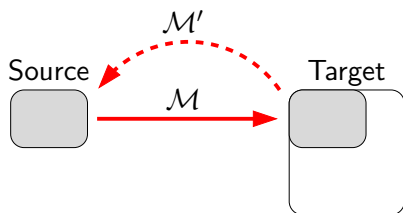
Information and *redundancy*:

fundamental concepts in schema mapping management

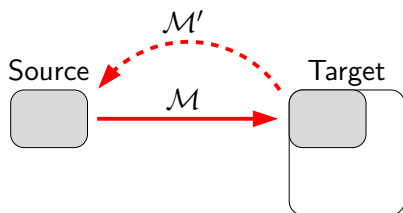


Information and *redundancy*:

fundamental concepts in schema mapping management

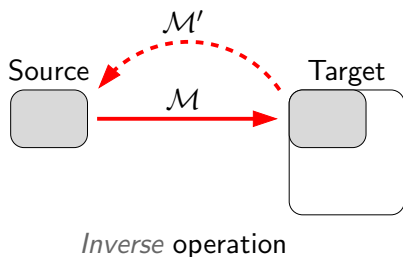


Information and redundancy:
fundamental concepts in schema mapping management



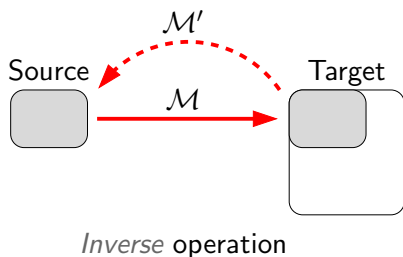
Inverse operation

Information and redundancy:
fundamental concepts in schema mapping management



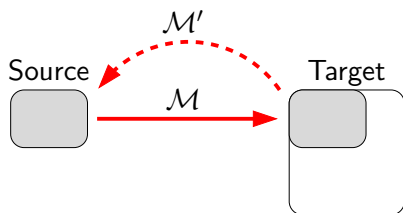
Invertibility for \mathcal{M} should coincide with no loss of information.

Information and redundancy:
fundamental concepts in schema mapping management



Invertibility for \mathcal{M} should coincide with **no loss of information**.

Information and *redundancy*:
fundamental concepts in schema mapping management



Inverse operation

Although fundamental, the notions of *information* and *redundancy* have received little attention in the schema mapping context.

Foundations of Schema Mapping Management

Marcelo Arenas, Jorge Pérez, Juan L. Reutter, Cristian Riveros

PUC Chile, U. Edinburgh, U. Oxford

We provide foundations for schema mapping management by formalizing the notions of *information* and *redundancy*.

We provide foundations for schema mapping management by formalizing the notions of *information* and *redundancy*.

Main contributions:

1. *Information* and *redundancy* in schema mappings
 - ▶ general formalization
 - ▶ characterizations and algorithmic issues

We provide foundations for schema mapping management by formalizing the notions of *information* and *redundancy*.

Main contributions:

1. *Information* and *redundancy* in schema mappings
 - ▶ general formalization
 - ▶ characterizations and algorithmic issues
2. Applications of the notions:
 - ▶ schema evolution problem
 - ▶ *Extract*, *Merge* and *Inverse* operators

Outline

Motivation

Source information

Algorithmic issues

Application: Invertibility

Target information

Application: Extract, first approach

Target and source redundancy

Application: Extract

Concluding remarks

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance
(J is called a *solution* for I under \mathcal{M}).

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance (J is called a *solution* for I under \mathcal{M}).
- ▶ the *composition of mappings*, $\mathcal{M} \circ \mathcal{M}'$, is the usual composition of binary relations.

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance (J is called a *solution* for I under \mathcal{M}).
- ▶ the *composition of mappings*, $\mathcal{M} \circ \mathcal{M}'$, is the usual composition of binary relations.

Mappings can be specified by formulas (dependencies):

$$\varphi_{\mathbf{S}}(\bar{x}) \rightarrow \psi_{\mathbf{T}}(\bar{x})$$

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance
(J is called a *solution* for I under \mathcal{M}).
- ▶ the *composition of mappings*, $\mathcal{M} \circ \mathcal{M}'$, is the usual composition of binary relations.

Mappings can be specified by formulas (dependencies):

$$\varphi_{\mathbf{S}}(\bar{x}) \rightarrow \psi_{\mathbf{T}}(\bar{x})$$

with $\varphi_{\mathbf{S}}(\bar{x})$ formula over the source and $\psi_{\mathbf{T}}(\bar{x})$ over the target.

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance (J is called a *solution* for I under \mathcal{M}).
- ▶ the *composition of mappings*, $\mathcal{M} \circ \mathcal{M}'$, is the usual composition of binary relations.

Mappings can be specified by formulas (dependencies):

$$\varphi_{\mathbf{S}}(\bar{x}) \rightarrow \psi_{\mathbf{T}}(\bar{x})$$

with $\varphi_{\mathbf{S}}(\bar{x})$ formula over the source and $\psi_{\mathbf{T}}(\bar{x})$ over the target.

- ▶ **L₁-to-L₂** dependency: $\varphi_{\mathbf{S}}(\bar{x}) \in \mathbf{L}_1$ and $\psi_{\mathbf{T}}(\bar{x}) \in \mathbf{L}_2$.

A bit of notation...

A *mapping* \mathcal{M} is a set of pairs (I, J) with

- ▶ I a source instance and J a target instance (J is called a *solution* for I under \mathcal{M}).
- ▶ the *composition of mappings*, $\mathcal{M} \circ \mathcal{M}'$, is the usual composition of binary relations.

Mappings can be specified by formulas (dependencies):

$$\varphi_{\mathbf{S}}(\bar{x}) \rightarrow \psi_{\mathbf{T}}(\bar{x})$$

with $\varphi_{\mathbf{S}}(\bar{x})$ formula over the source and $\psi_{\mathbf{T}}(\bar{x})$ over the target.

- ▶ **L₁-to-L₂** dependency: $\varphi_{\mathbf{S}}(\bar{x}) \in \mathbf{L}_1$ and $\psi_{\mathbf{T}}(\bar{x}) \in \mathbf{L}_2$.
- ▶ **CQ-to-CQ** = st-tgds.
- ▶ we are also interested in **CQ[≠]-to-CQ** and **FO-to-CQ**.

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

$\mathcal{M}_1: \text{Emp}(x, y, z) \rightarrow \exists u \text{ Person}(u, x)$

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

Intuitively:

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_1 .

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

Target₃: {Workplace(**place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

Intuitively:

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_1 .

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

Target₃: {Workplace(**place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_3 : Emp(x, y, z) \rightarrow Workplace(z)

Intuitively:

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_1 .

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

Target₃: {Workplace(**place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_3 : Emp(x, y, z) \rightarrow Workplace(z)

Intuitively:

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_1 .

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_3 .

Source information transferred by a mapping: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {Person(**ssn**, **name**) }

Target₂: {ENames(**name**), WorksIn(**name**, **place**) }

Target₃: {Workplace(**place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow $\exists u$ Person(u, x)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_3 : Emp(x, y, z) \rightarrow Workplace(z)

Intuitively:

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_1 .

\mathcal{M}_2 is *more source-informative* than \mathcal{M}_3 .

\mathcal{M}_1 and \mathcal{M}_3 are *incomparable*.

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

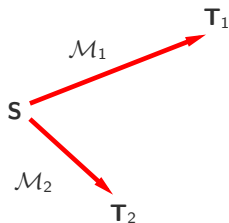
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

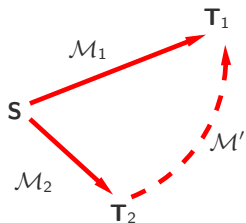
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

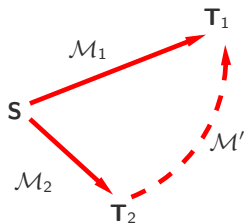
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



$$\mathcal{M}_1 : \text{Emp}(x, y, z) \rightarrow \exists u \text{ Person}(u, x)$$

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

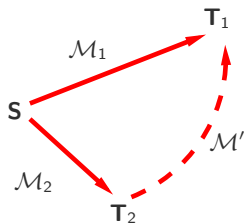
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



$\mathcal{M}_1 : \text{Emp}(x, y, z) \rightarrow \exists u \text{ Person}(u, x)$

$\mathcal{M}_2 : \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

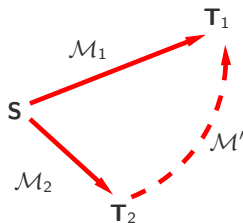
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



$$\mathcal{M}_1 : \text{Emp}(x, y, z) \rightarrow \exists u \text{ Person}(u, x)$$

$$\mathcal{M}_2 : \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$$

$$\mathcal{M}' : \text{ENames}(x) \rightarrow \exists u \text{ Person}(u, x)$$

Source information transferred by a mapping: Formalization

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the source schema.

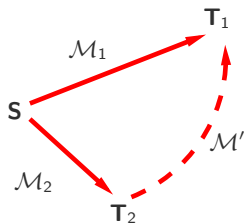
Definition

\mathcal{M}_2 is *more (or equally) source-informative than* \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_s \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1$.

\mathcal{M}_2 transfers information enough to reconstruct \mathcal{M}_1 .



$$\mathcal{M}_1 : \text{Emp}(x, y, z) \rightarrow \exists u \text{ Person}(u, x)$$

$$\mathcal{M}_2 : \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$$

$$\mathcal{M}' : \text{ENames}(x) \rightarrow \exists u \text{ Person}(u, x)$$

$$\mathcal{M}_2 \circ \mathcal{M}' = \mathcal{M}_1 \implies \mathcal{M}_1 \preceq_s \mathcal{M}_2$$

Axiomatization of \preceq_S

In the paper, we first define 4 *axioms* for an order \preceq on mappings

Axiomatization of \preceq_S

In the paper, we first define 4 *axioms* for an order \preceq on mappings

(C1) *reflexivity* : $\mathcal{M} \preceq \mathcal{M}$

(C2) *transitivity* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ and $\mathcal{M}_2 \preceq \mathcal{M}_3$, then $\mathcal{M}_1 \preceq \mathcal{M}_3$

Axiomatization of \preceq_S

In the paper, we first define 4 *axioms* for an order \preceq on mappings

(C1) *reflexivity* : $\mathcal{M} \preceq \mathcal{M}$

(C2) *transitivity* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ and $\mathcal{M}_2 \preceq \mathcal{M}_3$, then $\mathcal{M}_1 \preceq \mathcal{M}_3$

(C3) *maximum* : $\mathcal{M} \preceq \text{Id} = \{(I, I) \mid I \text{ is a source instance}\}$

Axiomatization of \preceq_S

In the paper, we first define 4 *axioms* for an order \preceq on mappings

(C1) *reflexivity* : $\mathcal{M} \preceq \mathcal{M}$

(C2) *transitivity* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ and $\mathcal{M}_2 \preceq \mathcal{M}_3$, then $\mathcal{M}_1 \preceq \mathcal{M}_3$

(C3) *maximum* : $\mathcal{M} \preceq \text{Id} = \{(I, I) \mid I \text{ is a source instance}\}$

(C4) *preservation* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ then $\mathcal{M} \circ \mathcal{M}_1 \preceq \mathcal{M} \circ \mathcal{M}_2$

Axiomatization of \preceq_S

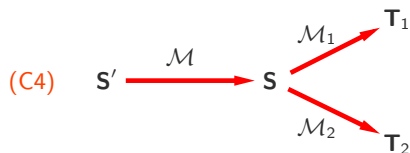
In the paper, we first define 4 *axioms* for an order \preceq on mappings

(C1) *reflexivity* : $\mathcal{M} \preceq \mathcal{M}$

(C2) *transitivity* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ and $\mathcal{M}_2 \preceq \mathcal{M}_3$, then $\mathcal{M}_1 \preceq \mathcal{M}_3$

(C3) *maximum* : $\mathcal{M} \preceq \text{Id} = \{(I, I) \mid I \text{ is a source instance}\}$

(C4) *preservation* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ then $\mathcal{M} \circ \mathcal{M}_1 \preceq \mathcal{M} \circ \mathcal{M}_2$



Axiomatization of \preceq_S

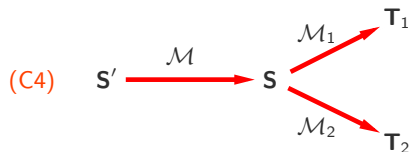
In the paper, we first define 4 *axioms* for an order \preceq on mappings

(C1) *reflexivity* : $\mathcal{M} \preceq \mathcal{M}$

(C2) *transitivity* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ and $\mathcal{M}_2 \preceq \mathcal{M}_3$, then $\mathcal{M}_1 \preceq \mathcal{M}_3$

(C3) *maximum* : $\mathcal{M} \preceq \text{Id} = \{(I, I) \mid I \text{ is a source instance}\}$

(C4) *preservation* : $\mathcal{M}_1 \preceq \mathcal{M}_2$ then $\mathcal{M} \circ \mathcal{M}_1 \preceq \mathcal{M} \circ \mathcal{M}_2$



Theorem

The order \preceq_S is the strictest relation that satisfies (C1-C4).

Towards deciding \preceq_s : target rewritability

Certain answers

Mapping \mathcal{M} , target query Q_T , source instance I :

$$\underline{\text{certain}}_{\mathcal{M}}(Q_T, I) = \bigcap_{(I, J) \in \mathcal{M}} Q_T(J)$$

Towards deciding \preceq_S : target rewritability

Certain answers

Mapping \mathcal{M} , target query Q_T , source instance I :

$$\underline{\text{certain}}_{\mathcal{M}}(Q_T, I) = \bigcap_{(I, J) \in \mathcal{M}} Q_T(J)$$

Definition

A source query Q_S is *target rewritable under \mathcal{M}* if there exists a target query Q_T such that

$$Q_S(I) = \underline{\text{certain}}_{\mathcal{M}}(Q_T, I)$$

for every source instance I .

Towards deciding \preceq_S : target rewritability

Certain answers

Mapping \mathcal{M} , target query Q_T , source instance I :

$$\underline{\text{certain}}_{\mathcal{M}}(Q_T, I) = \bigcap_{(I, J) \in \mathcal{M}} Q_T(J)$$

Definition

A source query Q_S is *target rewritable under \mathcal{M}* if there exists a target query Q_T such that

$$Q_S(I) = \underline{\text{certain}}_{\mathcal{M}}(Q_T, I)$$

for every source instance I .

- ▶ Intuitively: if Q_S is target rewritable under \mathcal{M} , then \mathcal{M} transfers all the source data retrieved by Q_S .

Source information transferred by a mapping can be characterized in terms of queries.

Theorem

Let \mathcal{M}_1 and \mathcal{M}_2 be specified by **FO-to-CQ**, then:

$\mathcal{M}_1 \preceq_s \mathcal{M}_2$ if and only if

every source query that is target rewritable under \mathcal{M}_1
is also target rewritable under \mathcal{M}_2 .

Source information transferred by a mapping can be characterized in terms of queries.

Theorem

Let \mathcal{M}_1 and \mathcal{M}_2 be specified by **FO-to-CQ**, then:

$\mathcal{M}_1 \preceq_s \mathcal{M}_2$ if and only if

every source query that is target rewritable under \mathcal{M}_1
is also target rewritable under \mathcal{M}_2 .

- ▶ The characterization is particular for **FO-to-CQ**.
For example, it does not work for **CQ-to-UCQ**.

Deciding \preceq_s

Theorem

For mappings specified by **FO-to-CQ**:

testing $\mathcal{M}_1 \preceq_s \mathcal{M}_2$ is undecidable

Deciding \preceq_s

Theorem

For mappings specified by **FO-to-CQ**:

testing $\mathcal{M}_1 \preceq_s \mathcal{M}_2$ is undecidable

Theorem

For mappings specified by **CQ[≠]-to-CQ**

testing $\mathcal{M}_1 \preceq_s \mathcal{M}_2$ is decidable

Deciding \preceq_s

Theorem

For mappings specified by **FO-to-CQ**:

testing $\mathcal{M}_1 \preceq_s \mathcal{M}_2$ is undecidable

Theorem

For mappings specified by **CQ \neq -to-CQ**

testing $\mathcal{M}_1 \preceq_s \mathcal{M}_2$ is decidable

Proof idea

For **CQ \neq -to-CQ** mappings, we prove that:

- ▶ checking target rewritability for **UCQ \neq** is decidable,
- ▶ only a finite number of queries in **UCQ \neq** need to be checked to determine if $\mathcal{M}_1 \preceq_s \mathcal{M}_2$.

Application: Invertibility can be characterized using \preceq_s .

Let $\overline{\text{Id}}$ be a mapping specified by a set of *copying*
(rules of the form $R(\bar{x}) \rightarrow \hat{R}(\bar{x})$ with R a source relation).

Application: Invertibility can be characterized using \preceq_s .

Let $\overline{\text{Id}}$ be a mapping specified by a set of *copying* (rules of the form $R(\bar{x}) \rightarrow \hat{R}(\bar{x})$ with R a source relation).

Definition [F06]: \mathcal{M}' is an *inverse* of \mathcal{M} if $\mathcal{M} \circ \mathcal{M}' = \overline{\text{Id}}$.

Application: Invertibility can be characterized using \preceq_s .

Let $\overline{\text{Id}}$ be a mapping specified by a set of *copying* (rules of the form $R(\bar{x}) \rightarrow \hat{R}(\bar{x})$ with R a source relation).

Definition [F06]: \mathcal{M}' is an *inverse* of \mathcal{M} if $\mathcal{M} \circ \mathcal{M}' = \overline{\text{Id}}$.

Theorem

Consider the class of total and closed-down on the left mappings:

- ▶ \mathcal{M} is invertible $\iff \overline{\text{Id}} \preceq_s \mathcal{M}$
- ▶ \mathcal{M} is invertible $\iff \mathcal{M}$ is \preceq_s -maximal

Application: Invertibility can be characterized using \preceq_s .

Let $\overline{\text{Id}}$ be a mapping specified by a set of *copying* (rules of the form $R(\bar{x}) \rightarrow \hat{R}(\bar{x})$ with R a source relation).

Definition [F06]: \mathcal{M}' is an *inverse* of \mathcal{M} if $\mathcal{M} \circ \mathcal{M}' = \overline{\text{Id}}$.

Theorem

Consider the class of total and closed-down on the left mappings:

- ▶ \mathcal{M} is invertible $\iff \overline{\text{Id}} \preceq_s \mathcal{M}$
- ▶ \mathcal{M} is invertible $\iff \mathcal{M}$ is \preceq_s -maximal

Invertibility do coincide with transferring the maximum amount of source information!

Application: Invertibility can be characterized using \preceq_s .

Let $\overline{\text{Id}}$ be a mapping specified by a set of *copying* (rules of the form $R(\bar{x}) \rightarrow \hat{R}(\bar{x})$ with R a source relation).

Definition [F06]: \mathcal{M}' is an *inverse* of \mathcal{M} if $\mathcal{M} \circ \mathcal{M}' = \overline{\text{Id}}$.

Theorem

Consider the class of total and closed-down on the left mappings:

- ▶ \mathcal{M} is invertible $\iff \overline{\text{Id}} \preceq_s \mathcal{M}$
- ▶ \mathcal{M} is invertible $\iff \mathcal{M}$ is \preceq_s -maximal

Invertibility do coincide with transferring the maximum amount of source information!

Corollary [FN09]

Testing invertibility for **CQ[≠]-to-CQ** mappings is decidable.

Covering target information: the *dual* definition

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the target schema.

Definition

\mathcal{M}_2 is *more (or equally) target-informative* than \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_T \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}' \circ \mathcal{M}_2 = \mathcal{M}_1$.

Covering target information: the *dual* definition

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the target schema.

Definition

\mathcal{M}_2 is *more (or equally) target-informative* than \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_T \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}' \circ \mathcal{M}_2 = \mathcal{M}_1$.

- ▶ *Universal solutions* [FKMP05]: A solution J^* for an instance I that represents the *entire space of solutions* of I under \mathcal{M} .

Covering target information: the *dual* definition

Assume that \mathcal{M}_1 and \mathcal{M}_2 share the target schema.

Definition

\mathcal{M}_2 is *more (or equally) target-informative* than \mathcal{M}_1 , denoted by

$$\mathcal{M}_1 \preceq_T \mathcal{M}_2,$$

if there exists a mapping \mathcal{M}' such that $\mathcal{M}' \circ \mathcal{M}_2 = \mathcal{M}_1$.

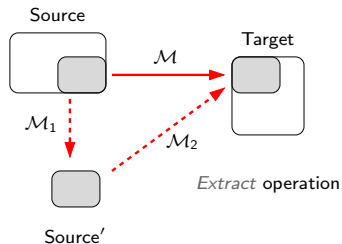
- ▶ *Universal solutions* [FKMP05]: A solution J^* for an instance I that represents the *entire space of solutions* of I under \mathcal{M} .

Theorem

Let \mathcal{M}_1 and \mathcal{M}_2 be specified by **FO-to-CQ**, then:

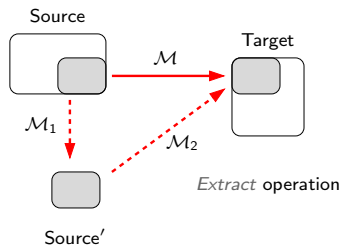
$\mathcal{M}_1 \preceq_T \mathcal{M}_2$ if and only if
every target instance that is universal solution under \mathcal{M}_1
is also universal solution under \mathcal{M}_2 .

Application 2: formalization of *Extract* (first attempt)



- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.

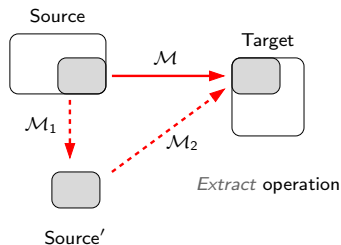
Application 2: formalization of *Extract* (first attempt)



- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.

$$(E1) \mathcal{M}_1 \equiv_s \mathcal{M}$$

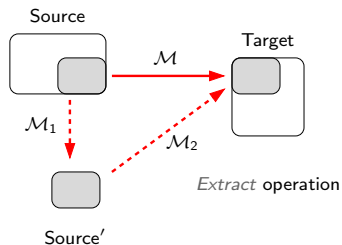
Application 2: formalization of *Extract* (first attempt)



- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.

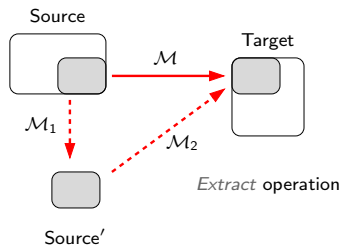
(E1) $\mathcal{M}_1 \equiv_s \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_s \mathcal{M}$ and $\mathcal{M} \preceq_s \mathcal{M}_1$).

Application 2: formalization of *Extract* (first attempt)



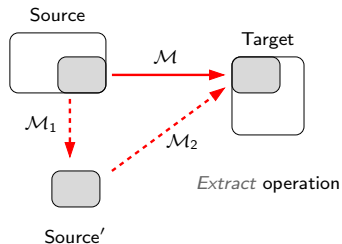
- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.
 - (E1) $\mathcal{M}_1 \equiv_S \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_S \mathcal{M}$ and $\mathcal{M} \preceq_S \mathcal{M}_1$).
 - (E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$

Application 2: formalization of *Extract* (first attempt)



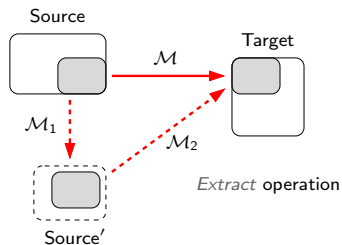
- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.
 - (E1) $\mathcal{M}_1 \equiv_S \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_S \mathcal{M}$ and $\mathcal{M} \preceq_S \mathcal{M}_1$).
 - (E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$ (i.e. $\mathcal{M}_2 \preceq_T \mathcal{M}$ and $\mathcal{M} \preceq_T \mathcal{M}_2$).

Application 2: formalization of *Extract* (first attempt)



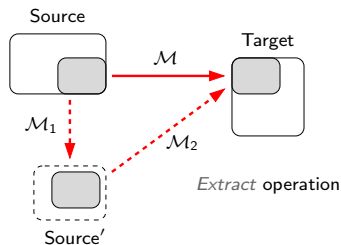
- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.
 - (E1) $\mathcal{M}_1 \equiv_S \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_S \mathcal{M}$ and $\mathcal{M} \preceq_S \mathcal{M}_1$).
 - (E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$ (i.e. $\mathcal{M}_2 \preceq_T \mathcal{M}$ and $\mathcal{M} \preceq_T \mathcal{M}_2$).
 - (E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

Application 2: formalization of *Extract* (first attempt)



- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.
 - (E1) $\mathcal{M}_1 \equiv_S \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_S \mathcal{M}$ and $\mathcal{M} \preceq_S \mathcal{M}_1$).
 - (E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$ (i.e. $\mathcal{M}_2 \preceq_T \mathcal{M}$ and $\mathcal{M} \preceq_T \mathcal{M}_2$).
 - (E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

Application 2: formalization of *Extract* (first attempt)



- ▶ We model the *extract* of \mathcal{M} as a pair $(\mathcal{M}_1, \mathcal{M}_2)$ s.t.
 - (E1) $\mathcal{M}_1 \equiv_S \mathcal{M}$ (i.e. $\mathcal{M}_1 \preceq_S \mathcal{M}$ and $\mathcal{M} \preceq_S \mathcal{M}_1$).
 - (E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$ (i.e. $\mathcal{M}_2 \preceq_T \mathcal{M}$ and $\mathcal{M} \preceq_T \mathcal{M}_2$).
 - (E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

¿How do we ensure the *optimality* of the new source schema?

Outline

Motivation

Source information

Algorithmic issues

Application: Invertibility

Target information

Application: Extract, first approach

Target and source redundancy

Application: Extract

Concluding remarks

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

$$\mathcal{M}_1: \quad \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$$

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target₂: {Worker(**name**, **working_place**) }

$\mathcal{M}_1: \quad \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target₂: {Worker(**name**, **working_place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow Worker(x, z)

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target₂: {Worker(**name**, **working_place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow Worker(x, z)

Intuitively:

▶ \mathcal{M}_1 is *target redundant*:

employee names are stored twice in the target schema.

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target₂: {Worker(**name**, **working_place**) }

$$\mathcal{M}_1: \text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$$
$$\mathcal{M}_2: \text{Emp}(x, y, z) \rightarrow \text{Worker}(x, z)$$

Intuitively:

- ▶ \mathcal{M}_1 is *target redundant*:
employee names are stored twice in the target schema.
- ▶ \mathcal{M}_2 is *not target redundant*:
all information in the target is *essential* for \mathcal{M}_2 .

Target redundancy in mappings: Intuition

Source: {Emp(**name**, **lives_in**, **works_in**) }

Target₁: {ENames(**name**), WorksIn(**name**, **place**) }

Target₂: {Worker(**name**, **working_place**) }

\mathcal{M}_1 : Emp(x, y, z) \rightarrow ENames(x) \wedge WorksIn(x, z)

\mathcal{M}_2 : Emp(x, y, z) \rightarrow Worker(x, z)

Intuitively:

- ▶ \mathcal{M}_1 is *target redundant*:
employee names are stored twice in the target schema.
- ▶ \mathcal{M}_2 is *not target redundant*:
all information in the target is *essential* for \mathcal{M}_2 .

Notice that \mathcal{M}_1 and \mathcal{M}_2 are *equally source-informative*, $\mathcal{M}_1 \equiv_s \mathcal{M}_2$

Target redundancy in mappings: Formalization

Definition

\mathcal{M} is *target redundant* if there is an instance $J^* \in \text{range}(\mathcal{M})$ such that the mapping

$$\mathcal{M}' = \{(I, J) \in \mathcal{M} \mid J \neq J^*\}$$

and \mathcal{M} are equally source-informative ($\mathcal{M} \equiv_s \mathcal{M}'$).

Target redundancy in mappings: Formalization

Definition

\mathcal{M} is *target redundant* if there is an instance $J^* \in \text{range}(\mathcal{M})$ such that the mapping

$$\mathcal{M}' = \{(I, J) \in \mathcal{M} \mid J \neq J^*\}$$

and \mathcal{M} are equally source-informative ($\mathcal{M} \equiv_s \mathcal{M}'$).

We can *lose a target instance*, and still be able to transfer the same amount of source information.

Target redundancy in mappings: Formalization

Definition

\mathcal{M} is *target redundant* if there is an instance $J^* \in \text{range}(\mathcal{M})$ such that the mapping

$$\mathcal{M}' = \{(I, J) \in \mathcal{M} \mid J \neq J^*\}$$

and \mathcal{M} are equally source-informative ($\mathcal{M} \equiv_s \mathcal{M}'$).

We can *lose a target instance*, and still be able to transfer the same amount of source information.

\mathcal{M}_1 : $\text{Emp}(x, y, z) \rightarrow \text{ENames}(x) \wedge \text{WorksIn}(x, z)$

J^* :

ENames:

name
Juan
Cristian

WorksIn:

name	place
Juan	Santiago

Source redundancy in mappings: the *dual* definition

Definition

\mathcal{M} is *source redundant* if there is an instance $I^* \in \text{dom}(\mathcal{M})$ such that the mapping

$$\mathcal{M}' = \{(I, J) \in \mathcal{M} \mid I \neq I^*\}$$

and \mathcal{M} are equally target-informative ($\mathcal{M} \equiv_T \mathcal{M}'$).

Source redundancy in mappings: the *dual* definition

Definition

\mathcal{M} is *source redundant* if there is an instance $I^* \in \text{dom}(\mathcal{M})$ such that the mapping

$$\mathcal{M}' = \{(I, J) \in \mathcal{M} \mid I \neq I^*\}$$

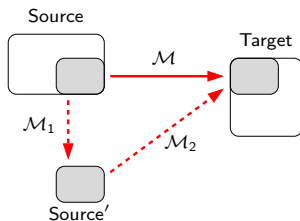
and \mathcal{M} are equally target-informative ($\mathcal{M} \equiv_T \mathcal{M}'$).

Theorem

Let \mathcal{M} be specified by **FO-to-CQ**, then:

- ▶ \mathcal{M} is target redundant iff there is a target instance that is not a universal solution under \mathcal{M} (onto mapping [FN09]).
- ▶ \mathcal{M} is source redundant iff there are two source instances with the same space of solutions under \mathcal{M} (unique solutions property [F06]).

Application 2: formalization of *Extract*



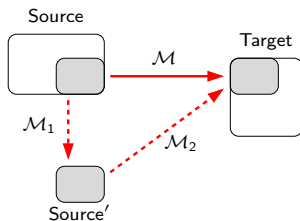
$(\mathcal{M}_1, \mathcal{M}_2)$ is an *extract* of \mathcal{M} iff:

(E1) $\mathcal{M}_1 \equiv_s \mathcal{M}$

(E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$

(E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

Application 2: formalization of *Extract*



$(\mathcal{M}_1, \mathcal{M}_2)$ is an *extract* of \mathcal{M} iff:

(E1) $\mathcal{M}_1 \equiv_s \mathcal{M}$

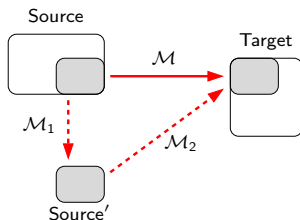
(E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$

(E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

(E4) \mathcal{M}_1 is *not target redundant*

(E5) \mathcal{M}_2 is *not source redundant*

Application 2: formalization of *Extract*



$(\mathcal{M}_1, \mathcal{M}_2)$ is an *extract* of \mathcal{M} iff:

(E1) $\mathcal{M}_1 \equiv_s \mathcal{M}$

(E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$

(E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

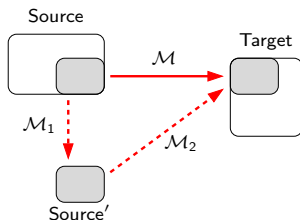
(E4) \mathcal{M}_1 is *not target redundant*

(E5) \mathcal{M}_2 is *not source redundant*

Theorem

For mappings specified by **FO-to-CQ** an *extract* always exists.

Application 2: formalization of *Extract*



$(\mathcal{M}_1, \mathcal{M}_2)$ is an *extract* of \mathcal{M} iff:

(E1) $\mathcal{M}_1 \equiv_s \mathcal{M}$

(E2) $\mathcal{M}_2 \equiv_T \mathcal{M}$

(E3) $\mathcal{M} = \mathcal{M}_1 \circ \mathcal{M}_2$

(E4) \mathcal{M}_1 is *not target redundant*

(E5) \mathcal{M}_2 is *not source redundant*

Theorem

For mappings specified by **FO-to-CQ** an *extract* always exists.

In the paper: an algorithm to compute an extract.

Information and redundancy are fundamental notions for schema mappings

In our work:

- ▶ we provide a formalization for both notions
- ▶ we study algorithmic issues, and natural characterizations
- ▶ we use these notions to re-study some schema mapping operators (schema evolution, extract, merge, inverse).

Foundations of Schema Mapping Management

Marcelo Arenas, Jorge Pérez, Juan L. Reutter, Cristian Riveros

PUC Chile, U. Edinburgh, U. Oxford

Outline

Motivation

Source information

Algorithmic issues

Application: Invertibility

Target information

Application: Extract, first approach

Target and source redundancy

Application: Extract

Concluding remarks