

A Collective, Probabilistic Approach to Schema Mapping

Angelika Kimmig, Alex Memory, Renée Miller, Lise Getoor

ILP 2017 (published at ICDE 2017)

KU LEUVEN

CARDIFF
UNIVERSITY
PRIFYSGOL
CAERDYDD



Context: Data Exchange & Data Integration

source

emp

id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS
4	Carl	X
5	Tom	Y7

proj

topic	mgr	lead
BigData	1	2
ML	1	1
eGov	4	5
DM	5	5

target

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111
MT	Justin	444
Deep	Ann	333

leader

name
Alice
Bob
Jim
Ann
Igor

org

oid	name
111	SAP
222	MS
333	Z
444	HC

Context: Data Exchange & Data Integration

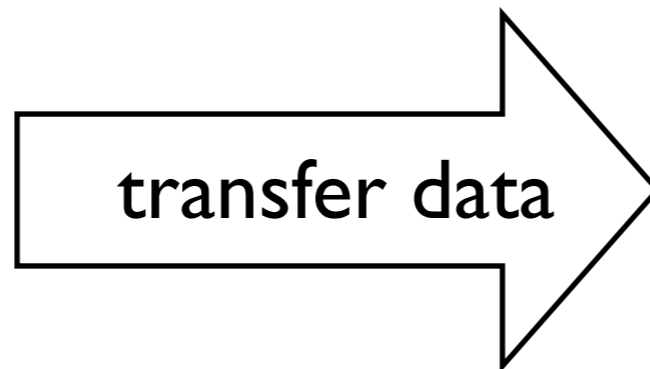
source

emp

id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS
4	Carl	X
5	Tom	Y7

proj

topic	mgr	lead
BigData	1	2
ML	1	1
eGov	4	5
DM	5	5



target

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111
MT	Justin	444
Deep	Ann	333

leader

name
Alice
Bob
Jim
Ann
Igor

org

oid	name
111	SAP
222	MS
333	Z
444	HC

Context: Data Exchange & Data Integration

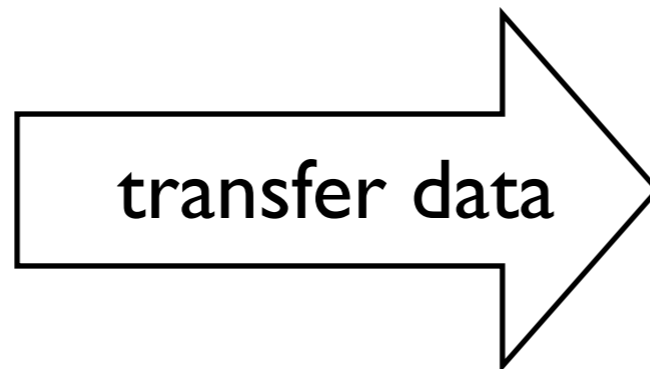
source

emp

id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS
4	Carl	X
5	Tom	Y7

proj

topic	mgr	lead
BigData	1	2
ML	1	1
eGov	4	5
DM	5	5



target

task

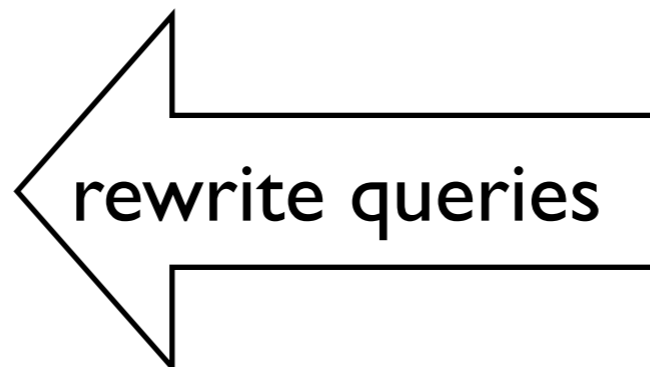
title	supervisor	oid
BigData	Alice	111
ML	Alice	111
MT	Justin	444
Deep	Ann	333

leader

name
Alice
Bob
Jim
Ann
Igor

org

oid	name
111	SAP
222	MS
333	Z
444	HC



Context: Data Exchange & Data Integration

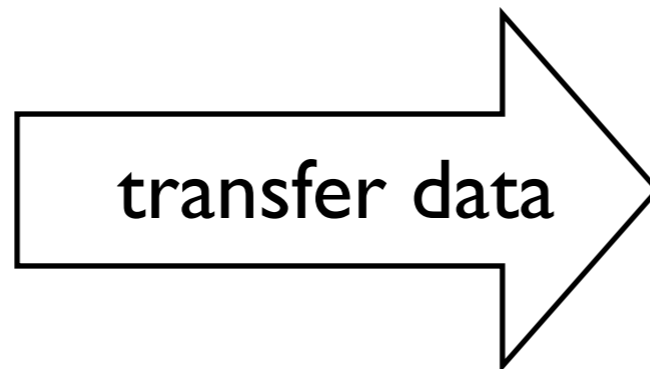
source

emp

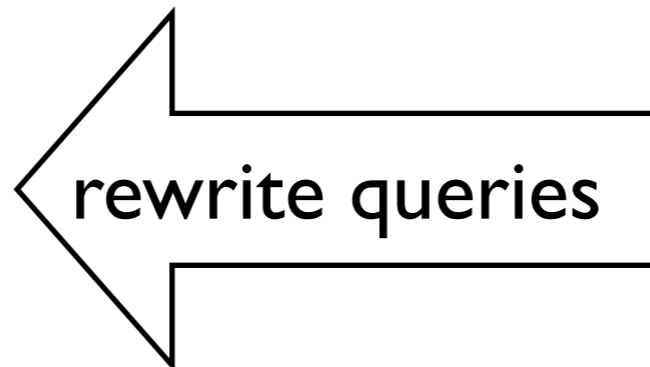
id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS
4	Carl	X
5	Tom	Y7

proj

topic	mgr	lead
BigData	1	2
ML	1	1
eGov	4	5
DM	5	5



need
**schema
mapping**



target

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111
MT	Justin	444
Deep	Ann	333

leader

name
Alice
Bob
Jim
Ann
Igor

org

oid	name
111	SAP
222	MS
333	Z
444	HC

Context: Data Exchange & Data Integration

source

target

emp

task

id	name	company
1	Alice	SAP

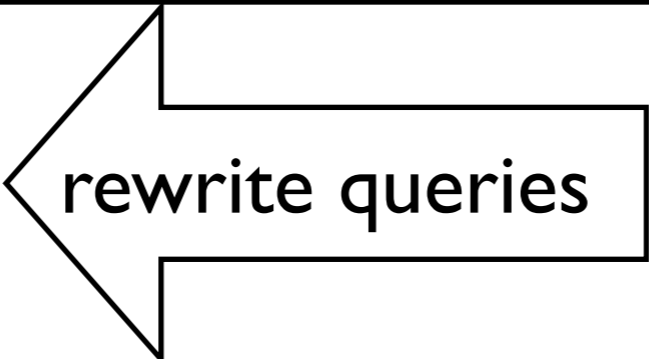


title	supervisor	oid
BigData	Alice	111
ML	Alice	111

$emp(I,N,C) \rightarrow leader(N)$
 $proj(T,M,L) \rightarrow \exists S. \exists O. task(T,S,O)$
 $proj(T,M,L) \ \& \ emp(L,N,C) \rightarrow \exists O. task(T,N,O) \ \& \ org(O,C)$

proj

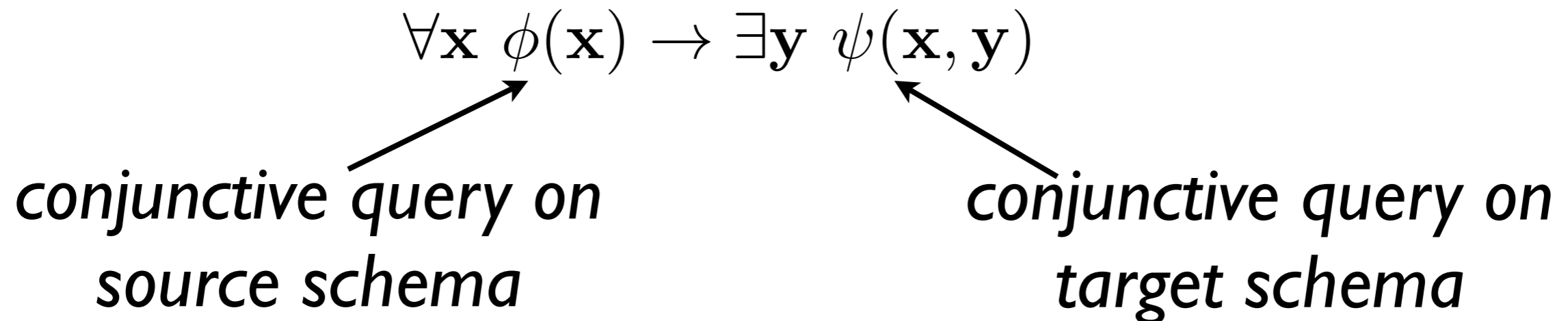
topic	mgr	lead
BigData	1	2
ML	1	1
eGov	4	5
DM	5	5



name	oid	org
Bob	222	MS
Jim	333	Z
Ann	444	HC
Igor		

Schema Mapping

- **st tgd** = source-target tuple generating dependency
[Fagin et al, 05; ten Cate & Kolaitis, 10]
- first order rule



- a **schema mapping** is a set of st tgds

Goal: learn schema mapping from data

proj

topic	mgr	lead
BigData	1	2
ML	1	1

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111

emp

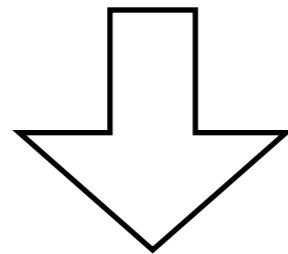
id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS

leader

name
Alice
Bob

org

oid	name
111	SAP
222	MS



$\text{proj}(T,M,L) \ \& \ \text{emp}(L,N,C) \rightarrow \text{leader}(N)$

$\text{emp}(I,N,C) \rightarrow \exists O. \text{org}(O,C)$

$\text{proj}(T,M,L) \ \& \ \text{emp}(M,N,C) \rightarrow \exists O. \text{task}(T,N,O) \ \& \ \text{org}(O,C)$

Goal: learn schema mapping from data

proj

topic	mgr	lead
BigData	1	2
ML	1	1

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111

emp

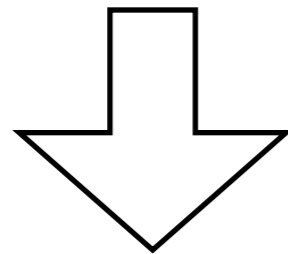
id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS

leader

name
Alice
Bob

org

oid	name
111	SAP
222	MS



$\text{proj}(T,M,L) \ \& \ \text{emp}(L,N,C) \ \rightarrow \ \text{leader}(N)$

$\text{emp}(I,N,C) \ \rightarrow \ \exists \ O. \ \text{org}(O,C)$

$\text{proj}(T,M,L) \ \& \ \text{emp}(M,N,C) \ \rightarrow \ \exists \ O. \ \text{task}(T,N,O) \ \& \ \text{org}(O,C)$

Challenges:

- ambiguous metadata
- imperfect data
- existentials / nulls

Here:

- **Given** data example (I,J), candidate set C
- **Find** optimal $M \subseteq C$ for (I,J)

I proj

topic	mgr	lead
BigData	1	2
ML	1	1

emp

id	name	company
1	Alice	SAP
2	Bob	IBM
3	Pat	MS

J task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111

leader

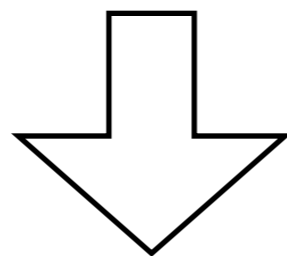
name
Alice
Bob

org

oid	name
111	SAP
222	MS

C

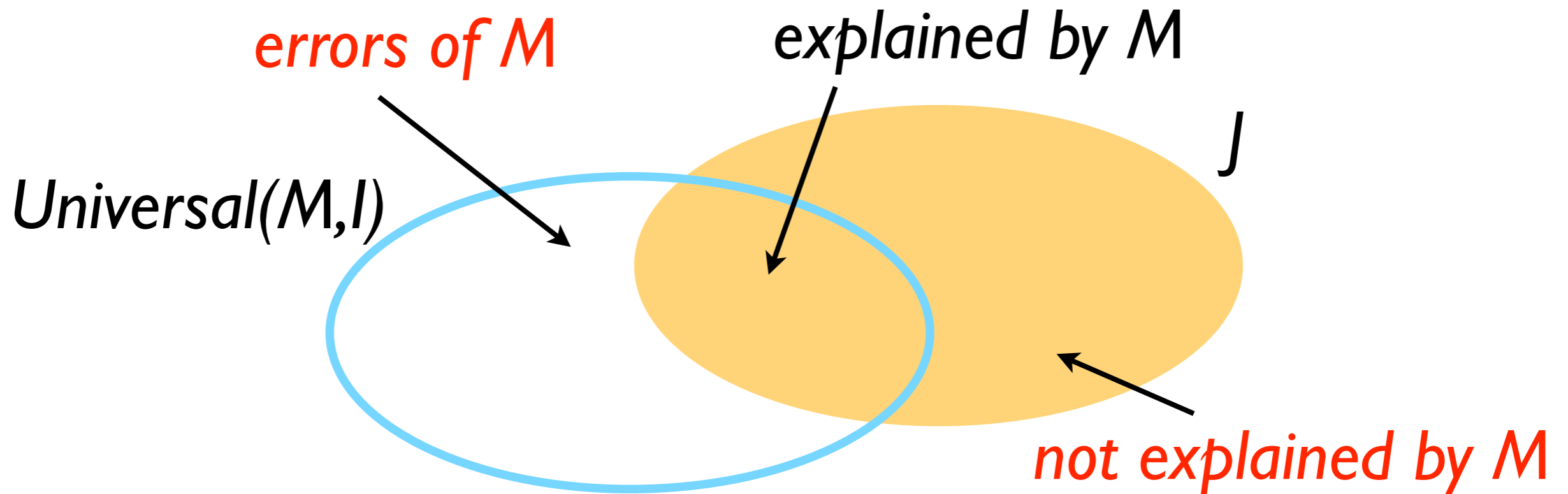
$\text{emp}(I,N,C) \rightarrow \text{leader}(N)$
 $\text{proj}(T,M,L) \rightarrow \exists S. \exists O. \text{task}(T,S,O)$
 $\text{proj}(T,M,L) \ \& \ \text{emp}(L,N,C) \rightarrow \text{leader}(N)$
 $\text{proj}(T,M,L) \ \& \ \text{emp}(M,N,C)$
 $\rightarrow \exists O. \text{task}(T,N,O) \ \& \ \text{org}(O,C)$
 $\text{proj}(T,M,L) \ \& \ \text{emp}(L,N,C)$
 $\rightarrow \exists O. \text{task}(T,N,O) \ \& \ \text{org}(O,C)$
 $\text{emp}(I,N,C) \rightarrow \exists O. \text{org}(O,C)$
 $\text{emp}(I,N,C) \rightarrow \exists O. \text{org}(O,N)$
...



M

$\text{proj}(T,M,L) \ \& \ \text{emp}(L,N,C) \rightarrow \text{leader}(N)$
 $\text{emp}(I,N,C) \rightarrow \exists O. \text{org}(O,C)$
 $\text{proj}(T,M,L) \ \& \ \text{emp}(M,N,C) \rightarrow \exists O. \text{task}(T,N,O) \ \& \ \text{org}(O,C)$

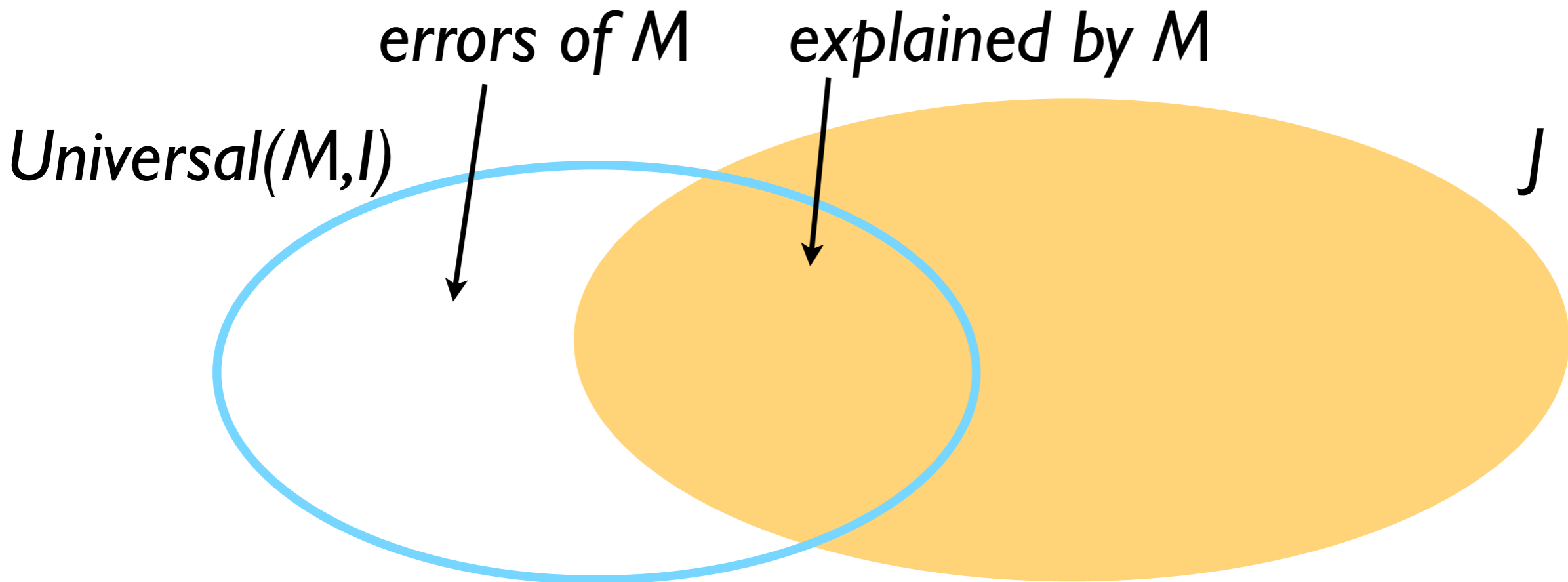
Full st tgds (no \exists)



goal: find small M that maximizes intersection

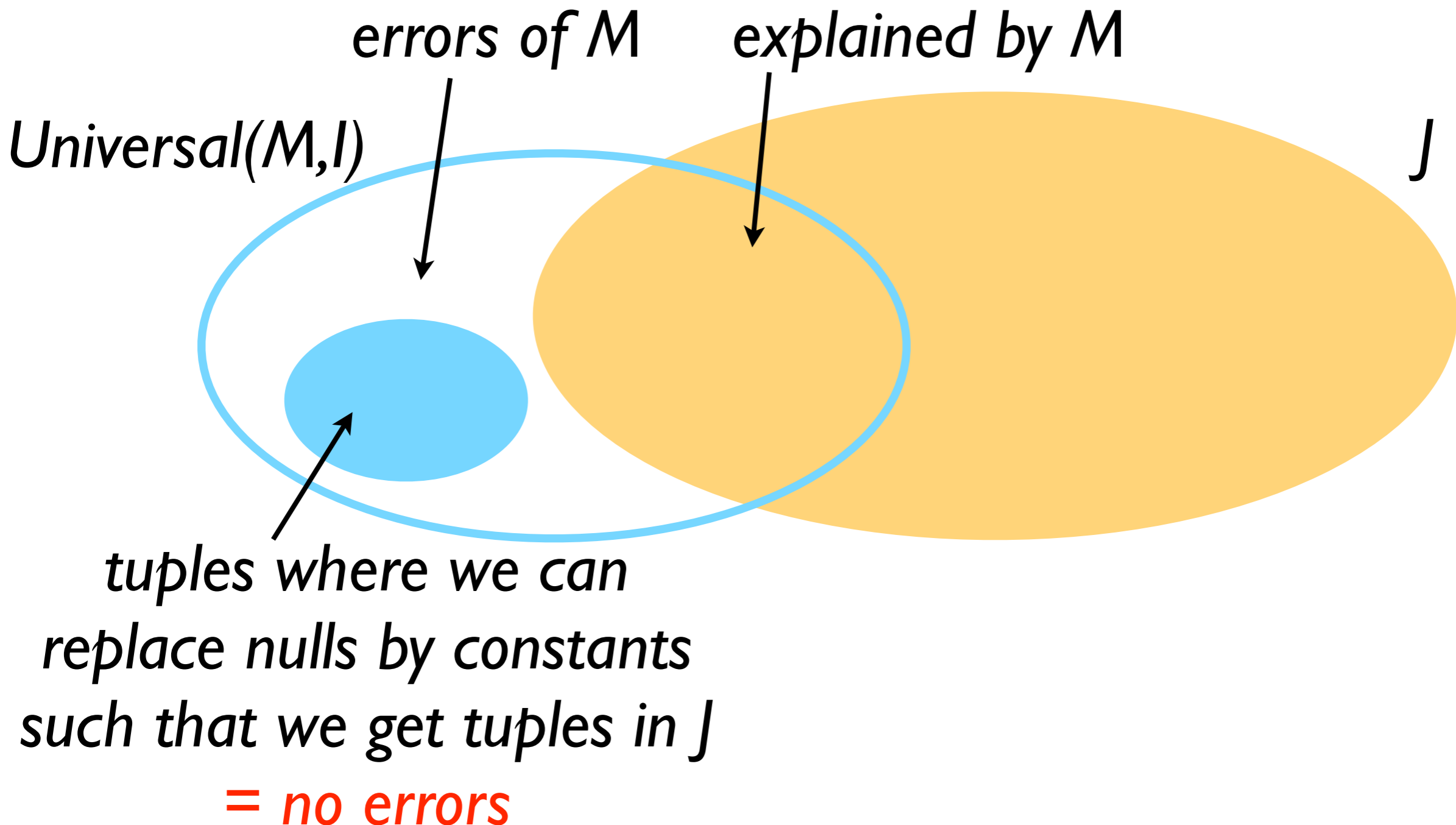
Arbitrary st tgds:

replace containment with homomorphism checks



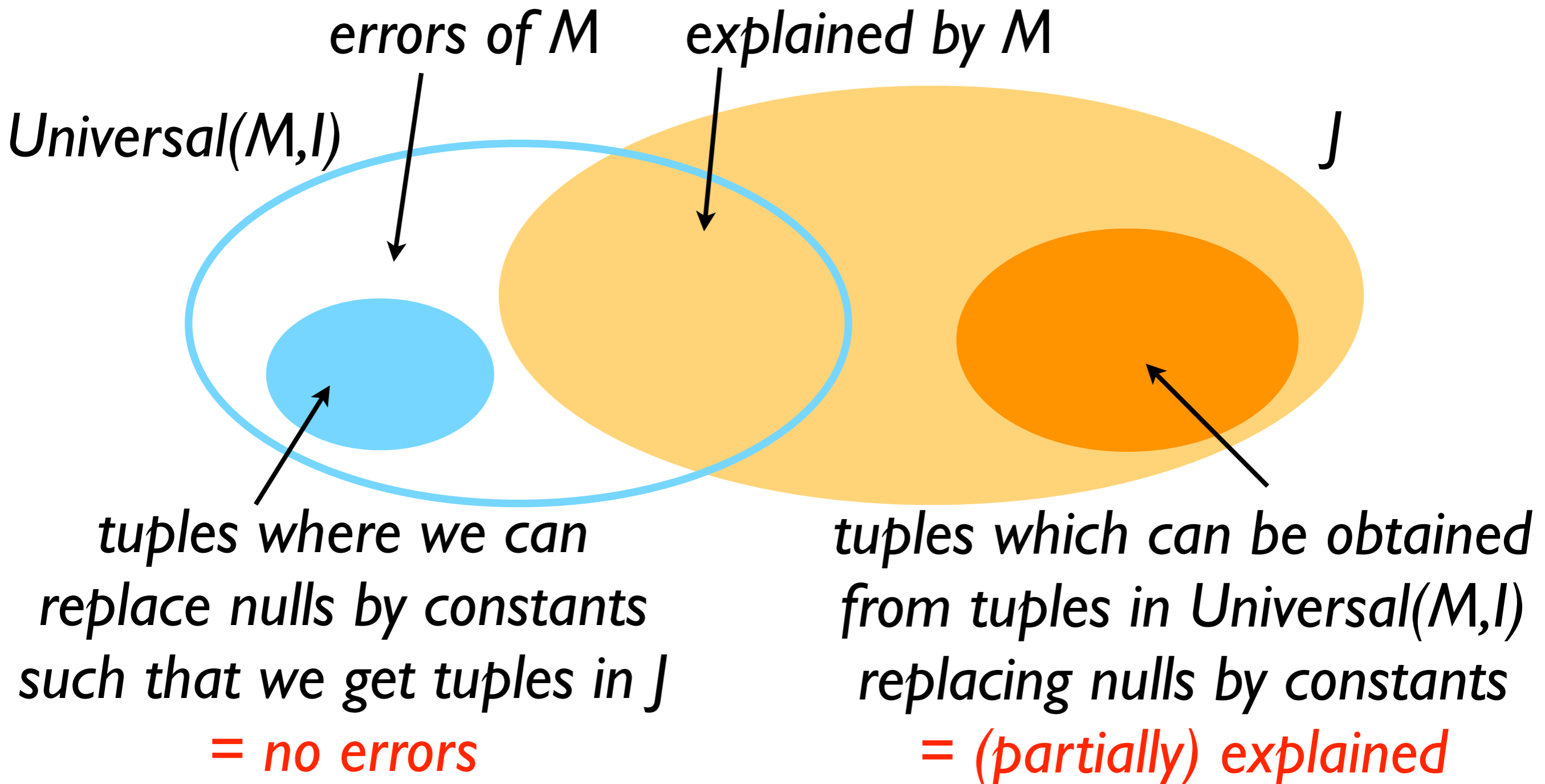
Arbitrary st tgds:

replace containment with homomorphism checks



Arbitrary st tgds:

replace containment with homomorphism checks



Example

task

title	supervisor	oid
BigData	Bob	null1
ML	Alice	null2

leader

name

org

oid	name
null1	IBM
null2	SAP

task

title	supervisor	oid
BigData	Alice	111
ML	Alice	111

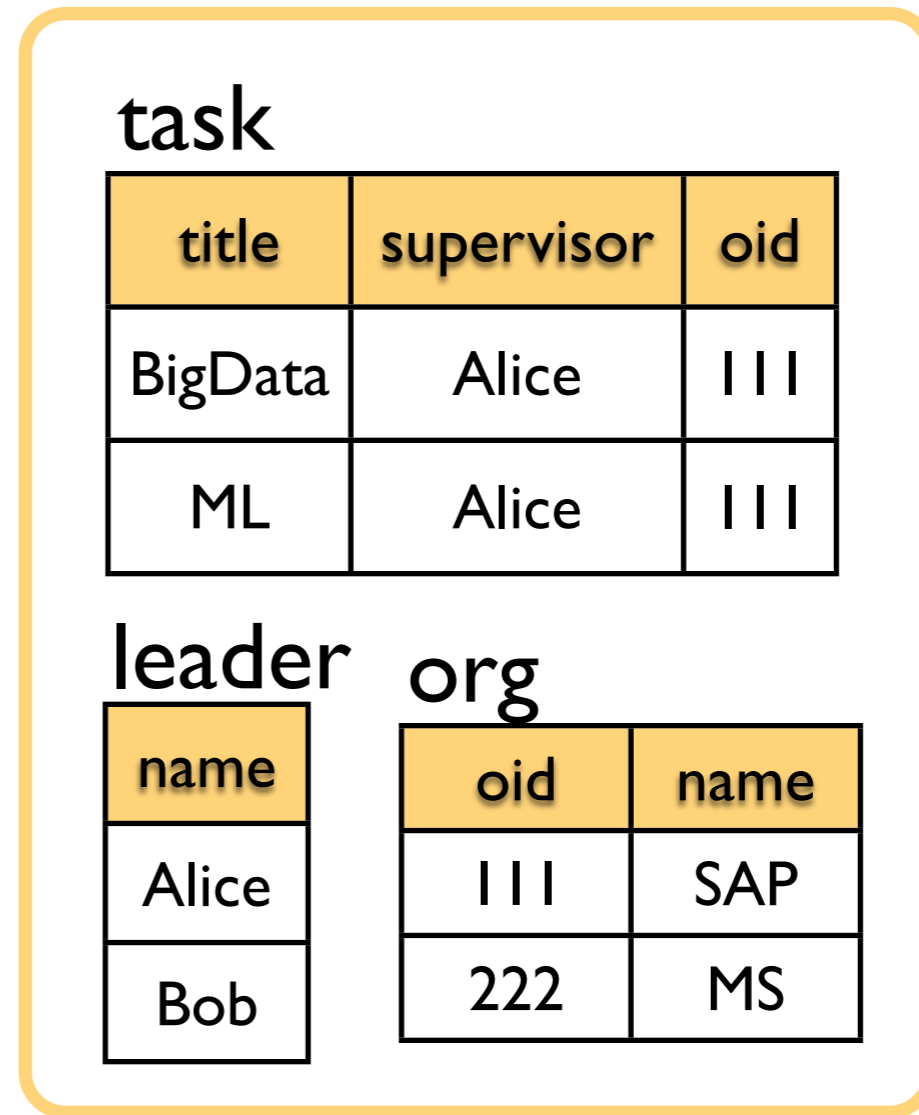
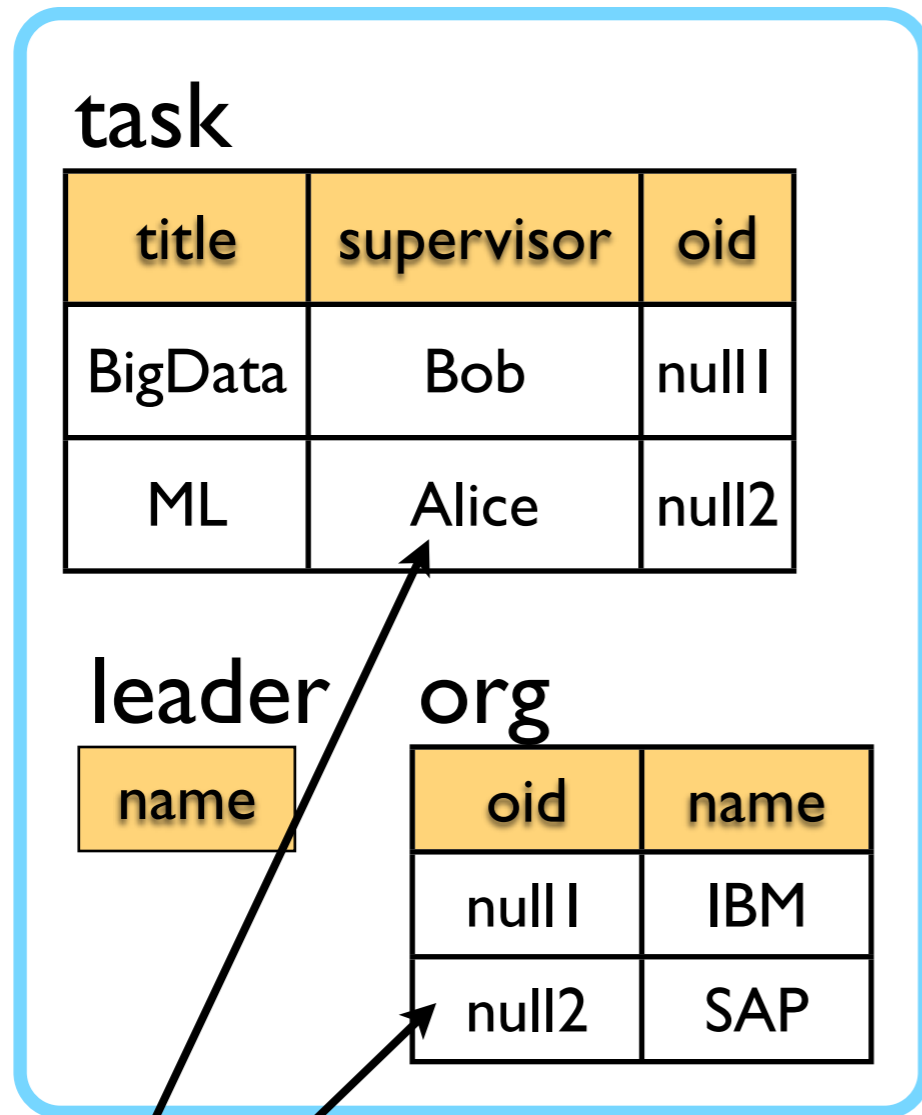
leader

name
Alice
Bob

org

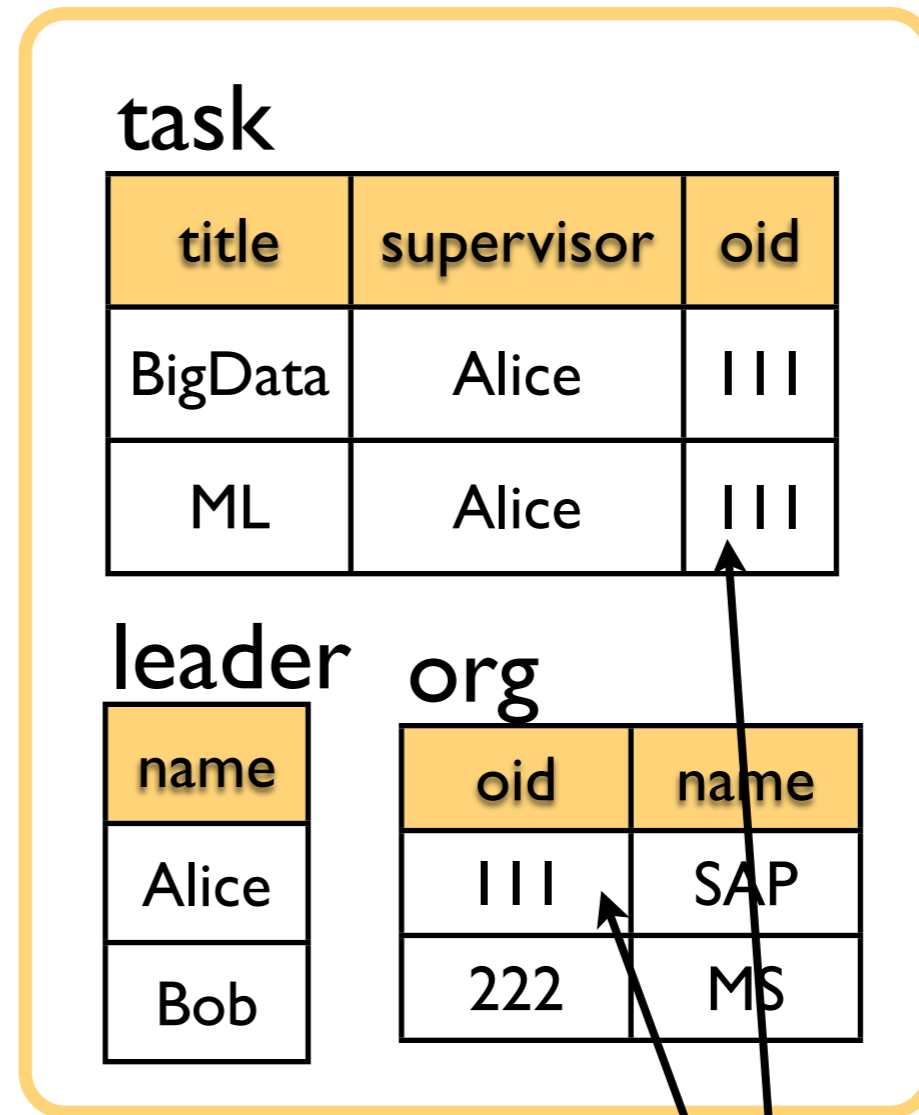
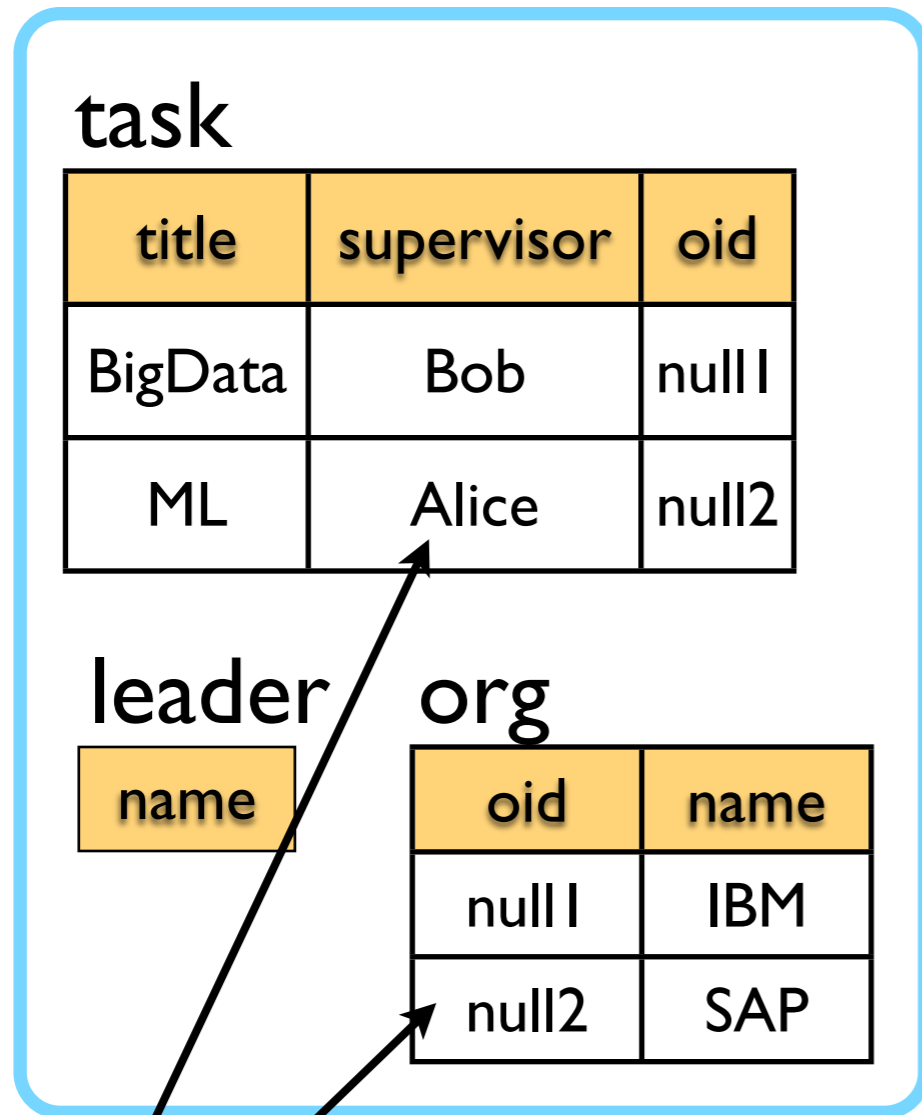
oid	name
111	SAP
222	MS

Example



no errors: replace null2 by 111...

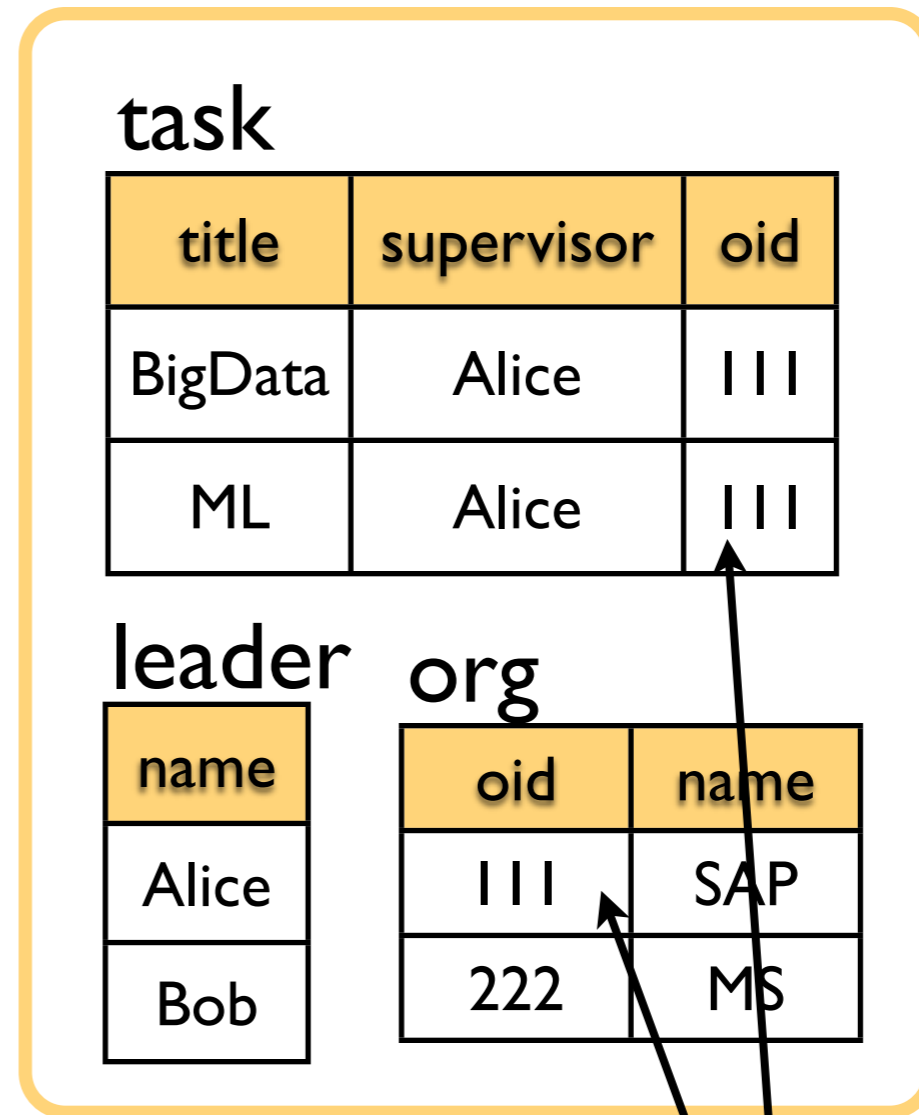
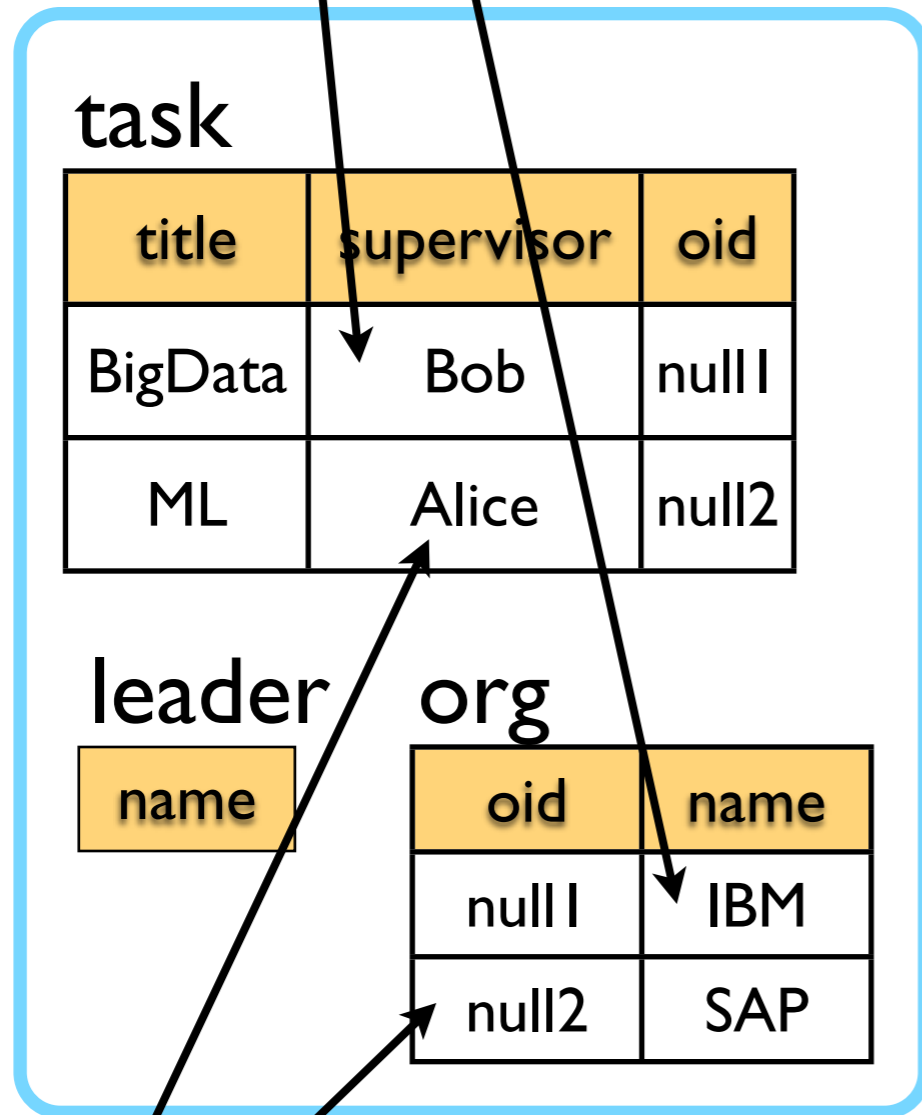
Example



no errors: replace null2 by 111... .. to explain these

Example

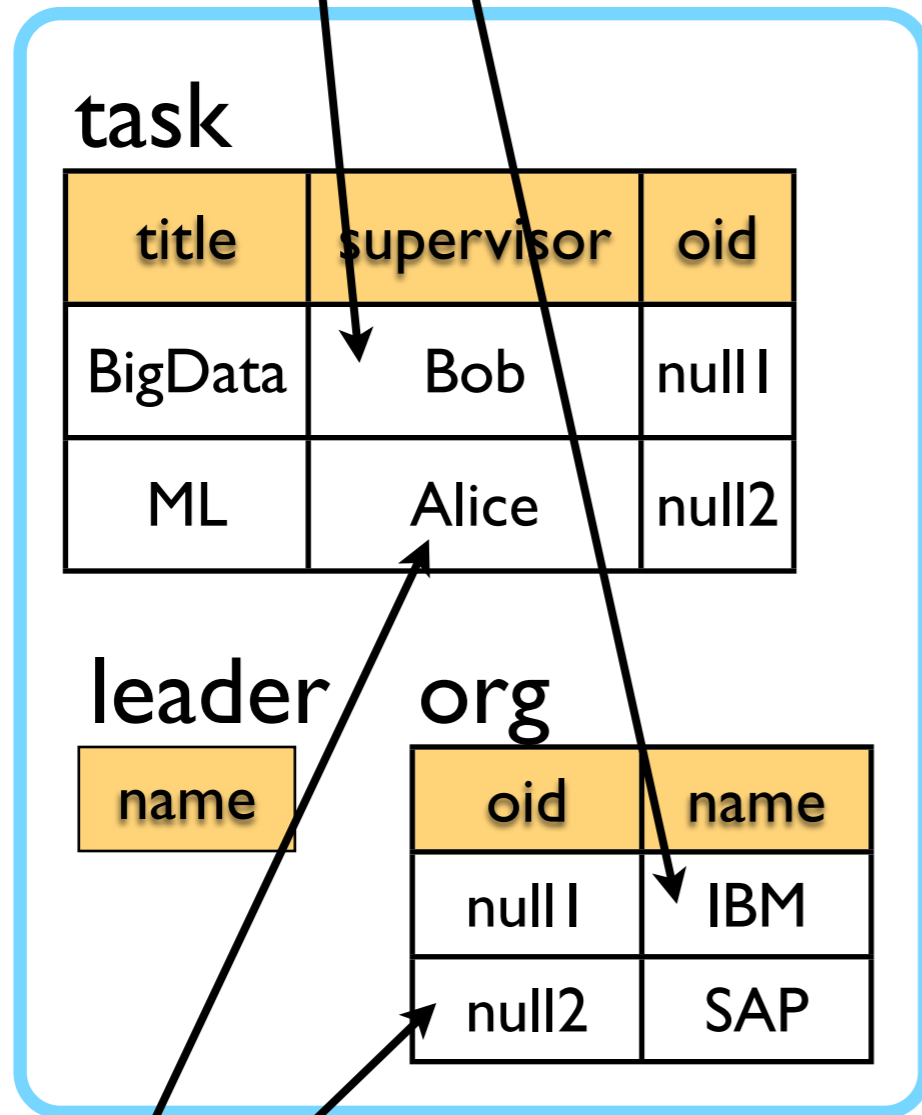
errors: can't replace null1 to get tuples in J



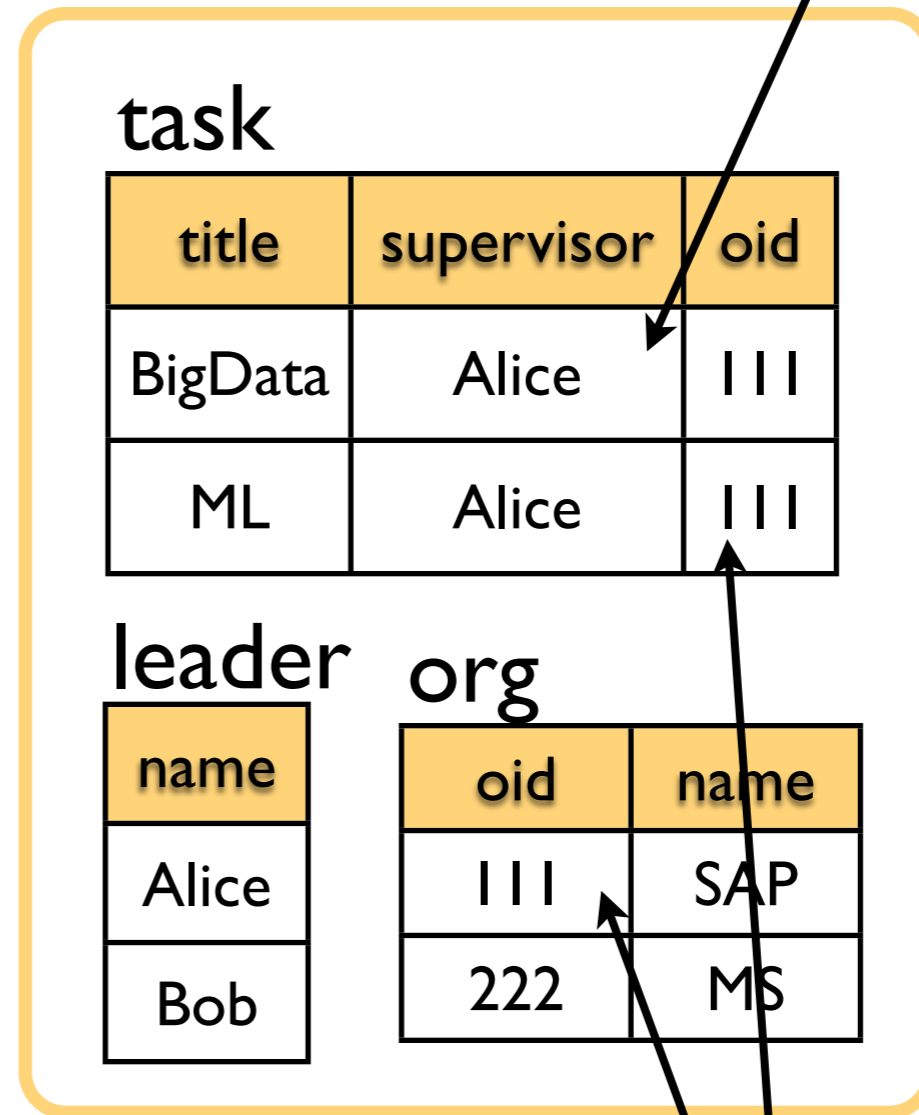
no errors: replace null2 by 111... .. to explain these

Example

errors: can't replace null1 to get tuples in J



not explained by any tuple in $Universal(M,I)$



no errors: replace null2 by 111... .. to explain these

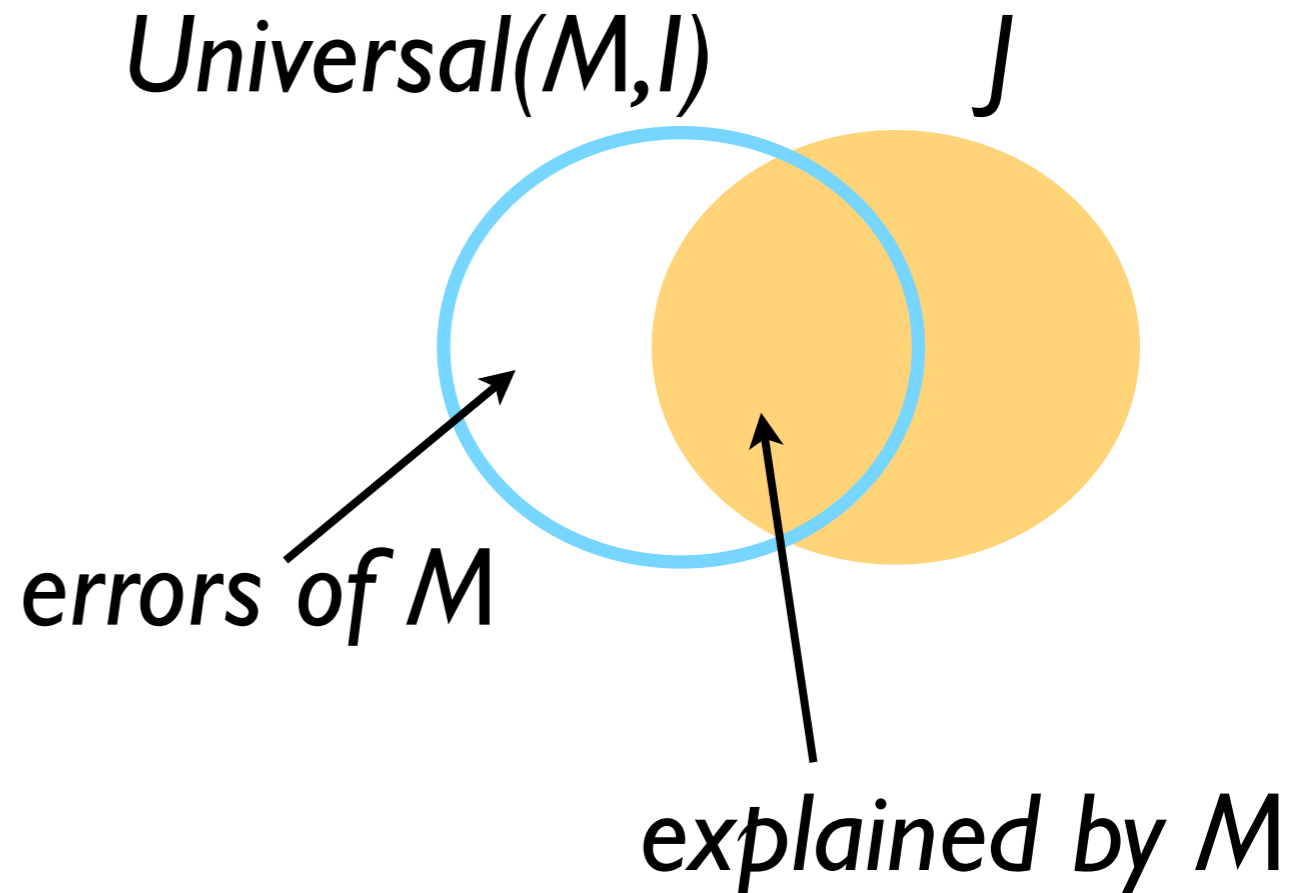
Task

- **Given**

- source schema S , target schema T
- data example (I, J)
- set C of candidate st tgds

- **Find** an optimal mapping M , i.e.,

$$\arg \min_{M \subseteq C} \left[\text{size}(M) + \sum_{t \in J} (1 - \text{explains}(M, t)) + \sum_{t \in \text{Universal}(C, I) - J} \text{error}(M, t) \right]$$



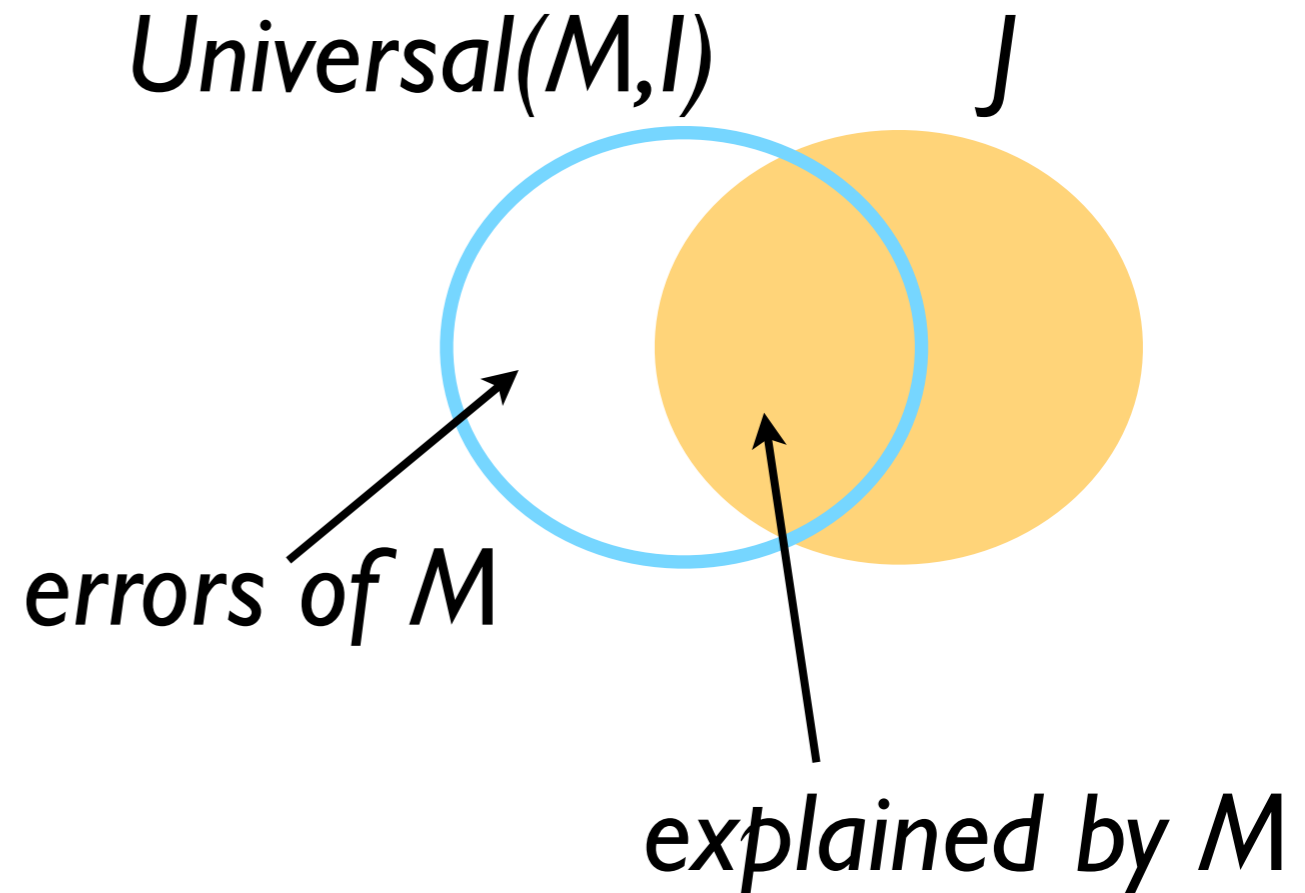
Task

- **Given**

- source schema S , target schema T
- data example (I, J)
- set C of candidate st tgds

- **Find** an optimal mapping M , i.e.,

$$\arg \min_{M \subseteq C} \left[\text{size}(M) + \sum_{t \in J} (1 - \text{explains}(M, t)) + \sum_{t \in \text{Universal}(C, I) - J} \text{error}(M, t) \right]$$



**NP-hard even
for full st tgds**

Probabilistic Soft Logic (PSL)

- declarative **language** to specify probabilistic models over logical atoms / relational domains
- PSL program = set of **weighted first order rules**
$$w : b_1(\vec{X}) \wedge \dots \wedge b_n(\vec{X}) \rightarrow h_1(\vec{X}) \vee \dots \vee h_m(\vec{X})$$
- **MPE inference** = finding most likely model
- **efficient** approximate solver with **guarantee** on solution quality

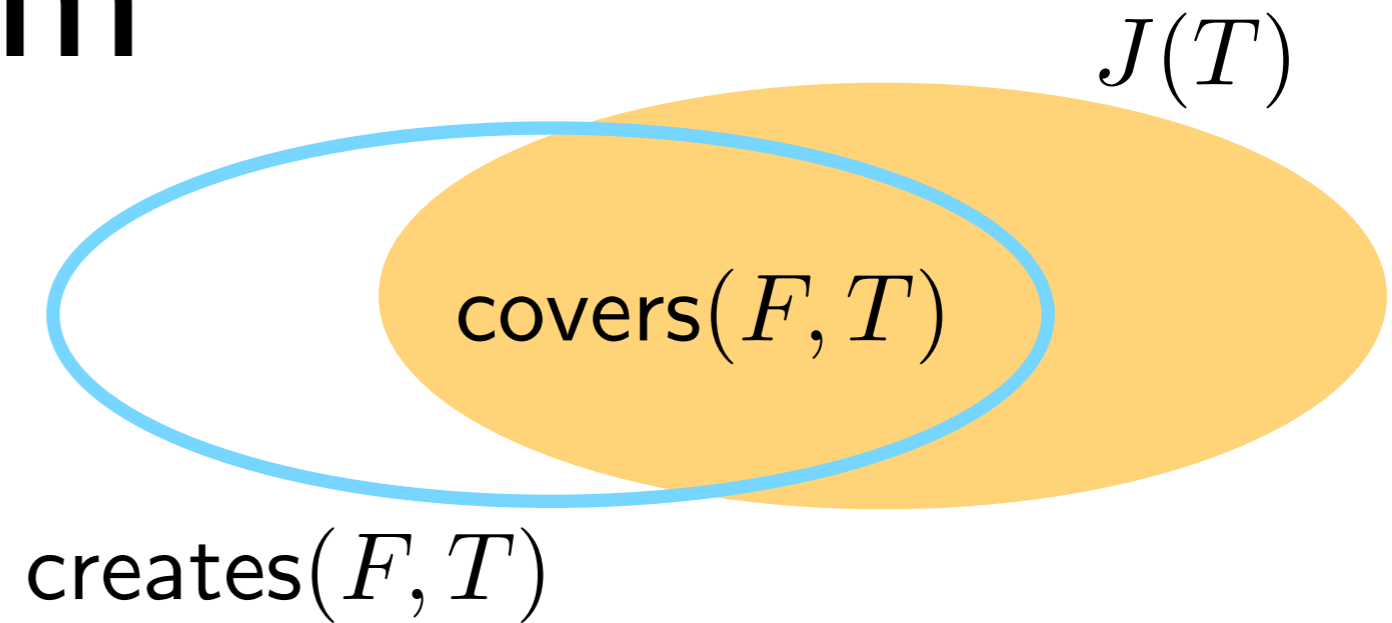
Our PSL program

given

$J(T)$

$\text{covers}(F, T)$

$\text{creates}(F, T)$



Our PSL program

given

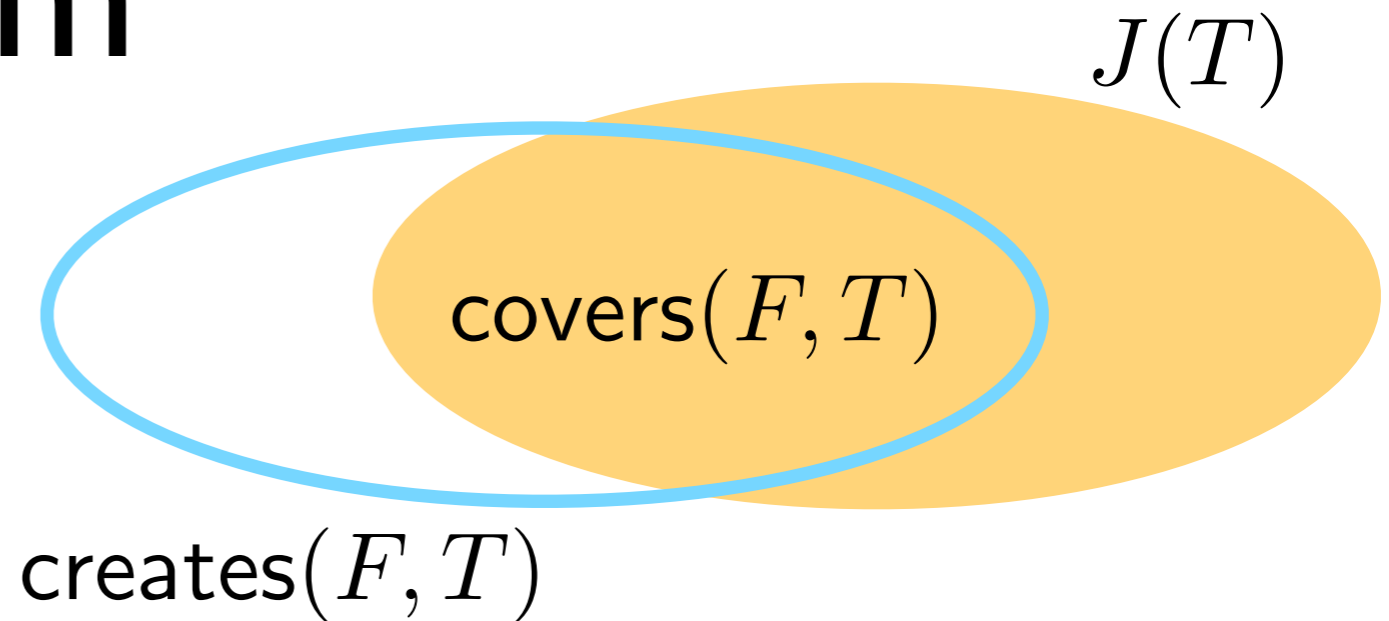
$J(T)$

$\text{covers}(F, T)$

$\text{creates}(F, T)$

find optimal M :

$$\text{in}(F) \Leftrightarrow F \in M$$



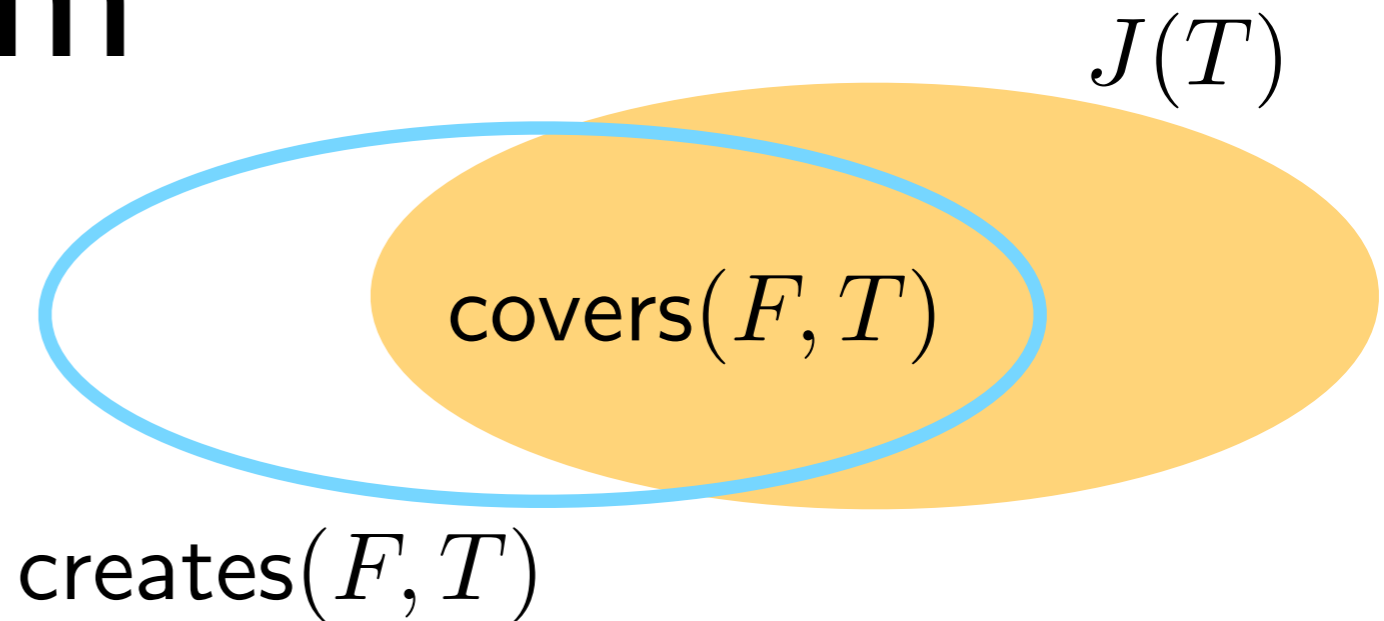
Our PSL program

given

$J(T)$

$\text{covers}(F, T)$

$\text{creates}(F, T)$



find optimal M :

$$\text{in}(F) \Leftrightarrow F \in M$$

minimize #errors:

$$1 : \text{in}(F) \wedge \text{creates}(F, T) \rightarrow J(T)$$

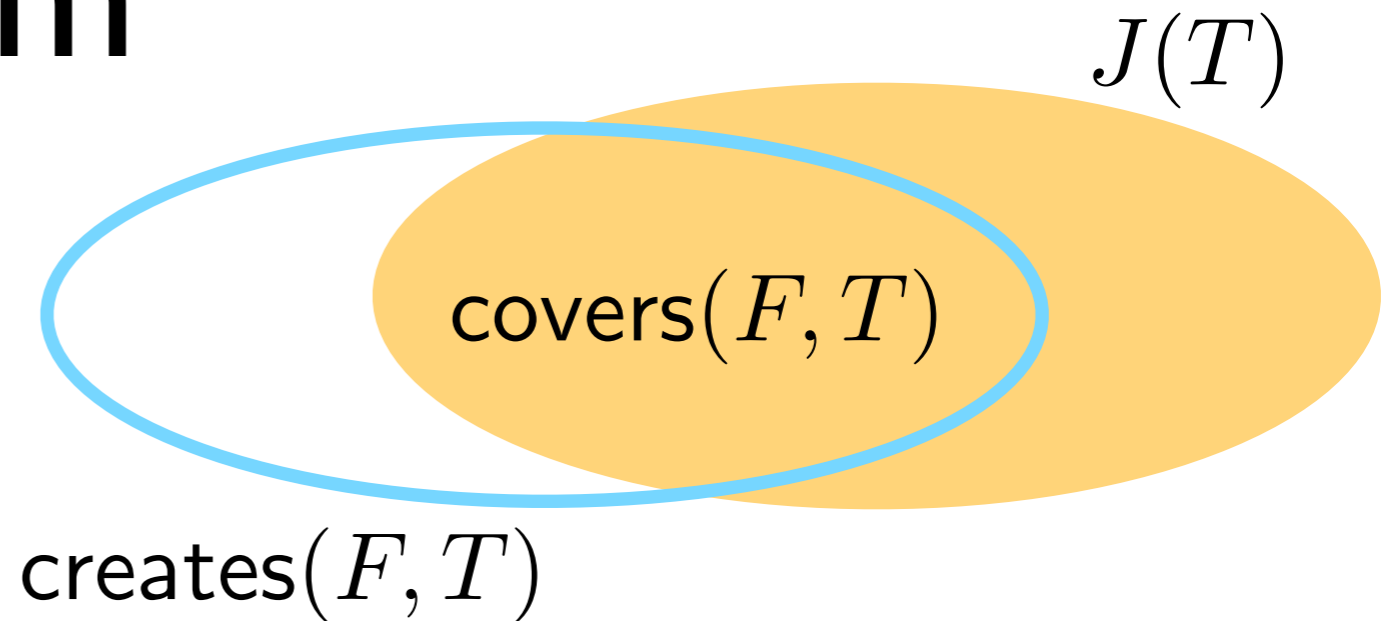
Our PSL program

given

$J(T)$

$\text{covers}(F, T)$

$\text{creates}(F, T)$



find optimal M :

$$\text{in}(F) \Leftrightarrow F \in M$$

minimize #errors:

$$1 : \text{in}(F) \wedge \text{creates}(F, T) \rightarrow J(T)$$

minimize #unexplained:

$$1 : J(T) \rightarrow \exists F. \text{covers}(F, T) \wedge \text{in}(F)$$

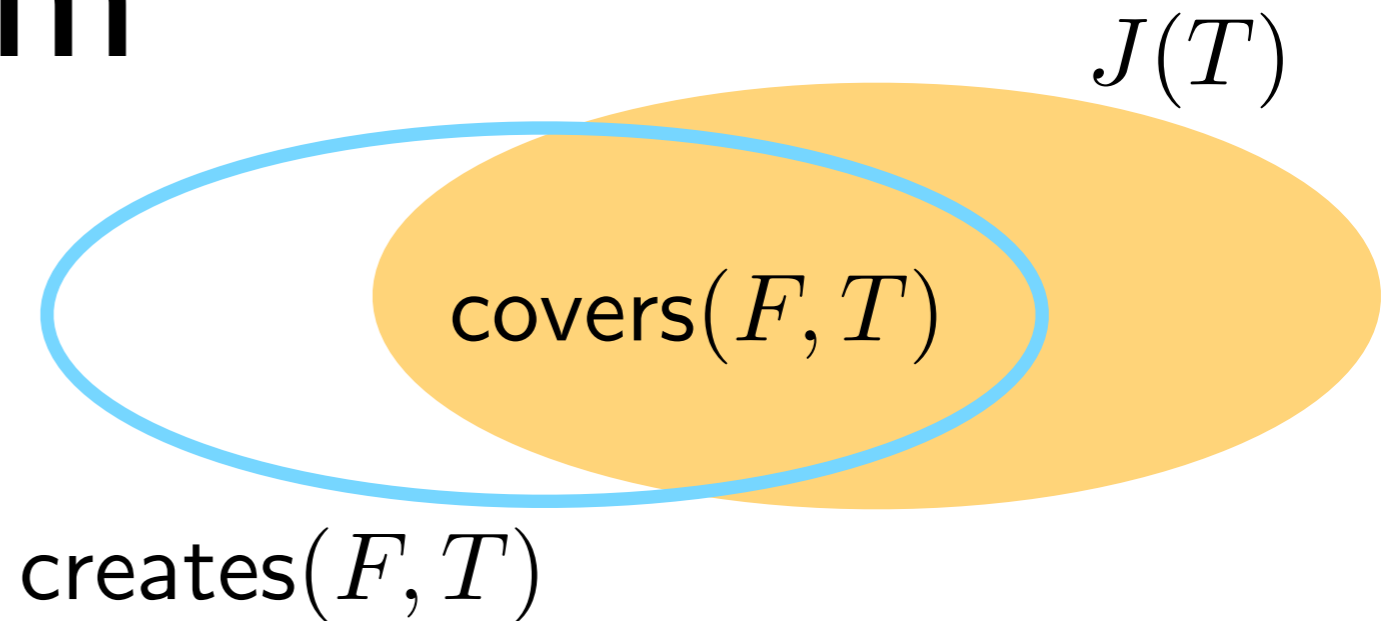
Our PSL program

given

$J(T)$

$\text{covers}(F, T)$

$\text{creates}(F, T)$



find optimal M :

$$\text{in}(F) \Leftrightarrow F \in M$$

minimize #errors:

$$1 : \text{in}(F) \wedge \text{creates}(F, T) \rightarrow J(T)$$

minimize #unexplained:

$$1 : J(T) \rightarrow \exists F. \text{covers}(F, T) \wedge \text{in}(F)$$

minimize size of M :

$$\text{size}(F) : \text{in}(F) \rightarrow \perp$$

New PSL construct: prioritized disjunction rules

$$1 : J(T) \rightarrow \exists F. \text{covers}(F, T) \wedge \text{in}(F)$$

to be inferred

observed **priority** between
0 (low) and 1 (high)

automatically transformed into set of standard
PSL rules expressing a preference for inferred
atoms with higher priority

Experimental evaluation

- scenarios generated using iBench [Arocena et al, 15]
- candidate st tgds generated using Clio [Fagin et al, 09]
- E1: increasingly noisy candidates, perfect data
 - metadata-only baseline suffers, we get perfect mappings
- E2: ambiguous set of candidates, increasingly noisy data
 - high quality mappings found for up to 25% unexpected and 10% missing target tuples in J

Our Contributions

Given

- metadata
- data example
- candidate st tgds

Find

- small set of st tgds
- minimally invalid
- maximally explaining

Our Contributions

Given

- metadata
- data example
- candidate st tgds



Collective Mapping Discovery (CMD)

declarative probabilistic model
+
MPE inference

using Probabilistic Soft Logic (PSL)



Find

- small set of st tgds
- minimally invalid
- maximally explaining

Our Contributions

Given

- metadata
- data example
- candidate st tgds



Collective Mapping Discovery (CMD)

declarative probabilistic model

+

MPE inference

using Probabilistic Soft Logic (PSL)



Find

- small set of st tgds
- minimally invalid
- maximally explaining

- supports arbitrary st tgds
- jointly reasons about metadata and data
- handles noisy input
- efficient solver with quality guarantee
- declarative, extensible definition of optimization task

Our Contributions

Given

- metadata
- data example
- candidate st tgds



Collective Mapping Discovery (CMD)

declarative probabilistic model
+
MPE inference

using Probabilistic Soft Logic (PSL)



Find

- small set of st tgds
- minimally invalid
- maximally explaining

- supports arbitrary st tgds
- jointly reasons about metadata and data
- handles noisy input
- efficient solver with quality guarantee
- declarative, extensible definition of optimization task

Thanks!