

Automatic Extraction of Ontologies Wrapping Relational Data Sources

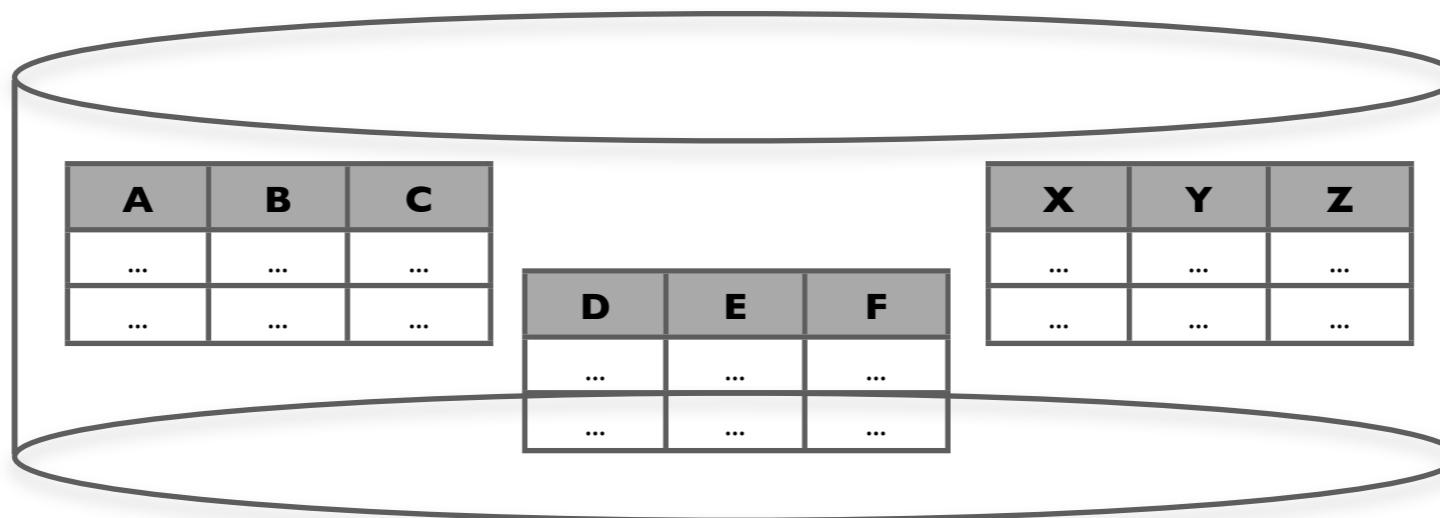
Lina Lubyte, Sergio Tessaris

KRDB Research Centre for Knowledge and Data
Free University of Bozen-Bolzano

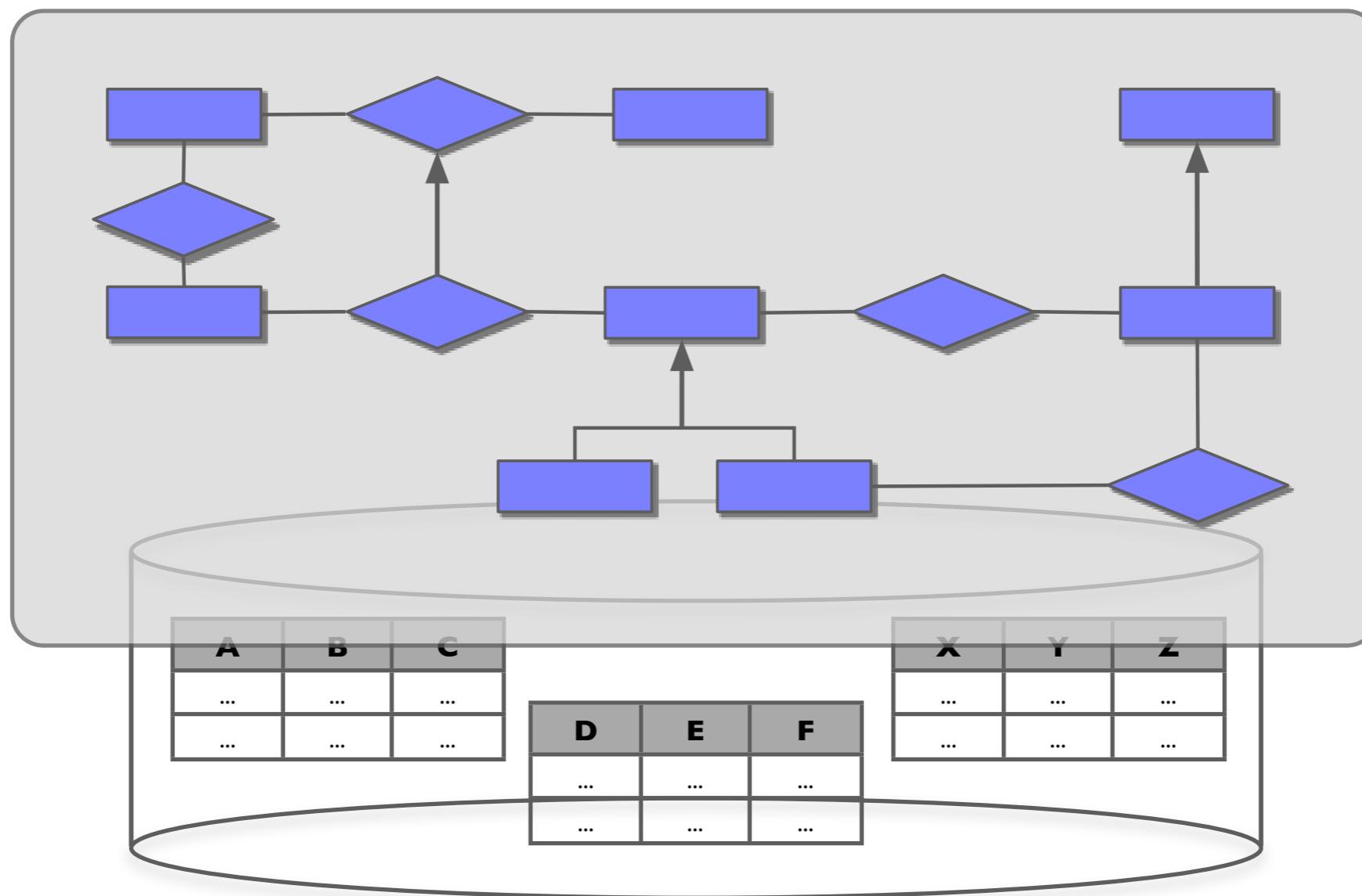
September 1, 2009

Semantic data access

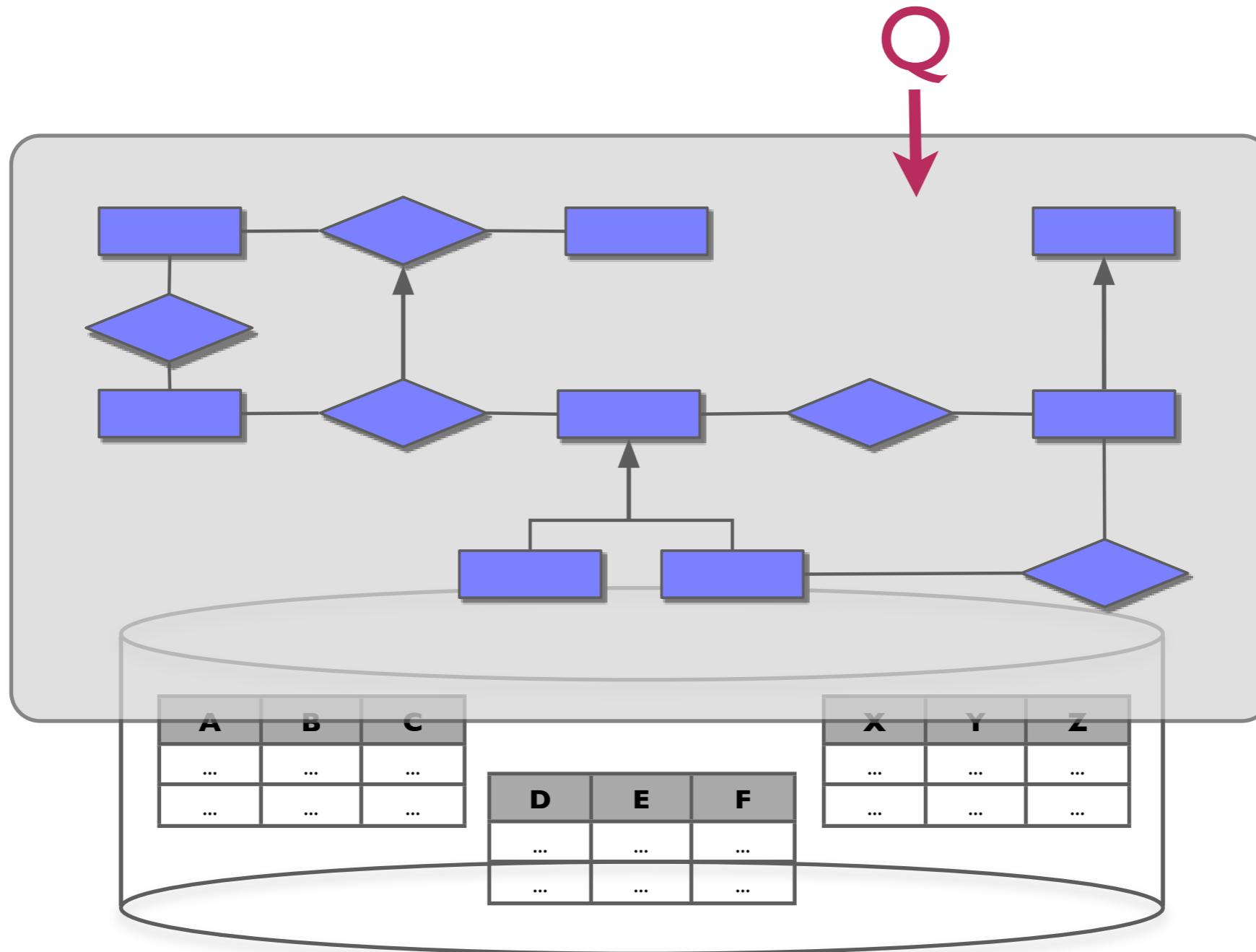
Semantic data access



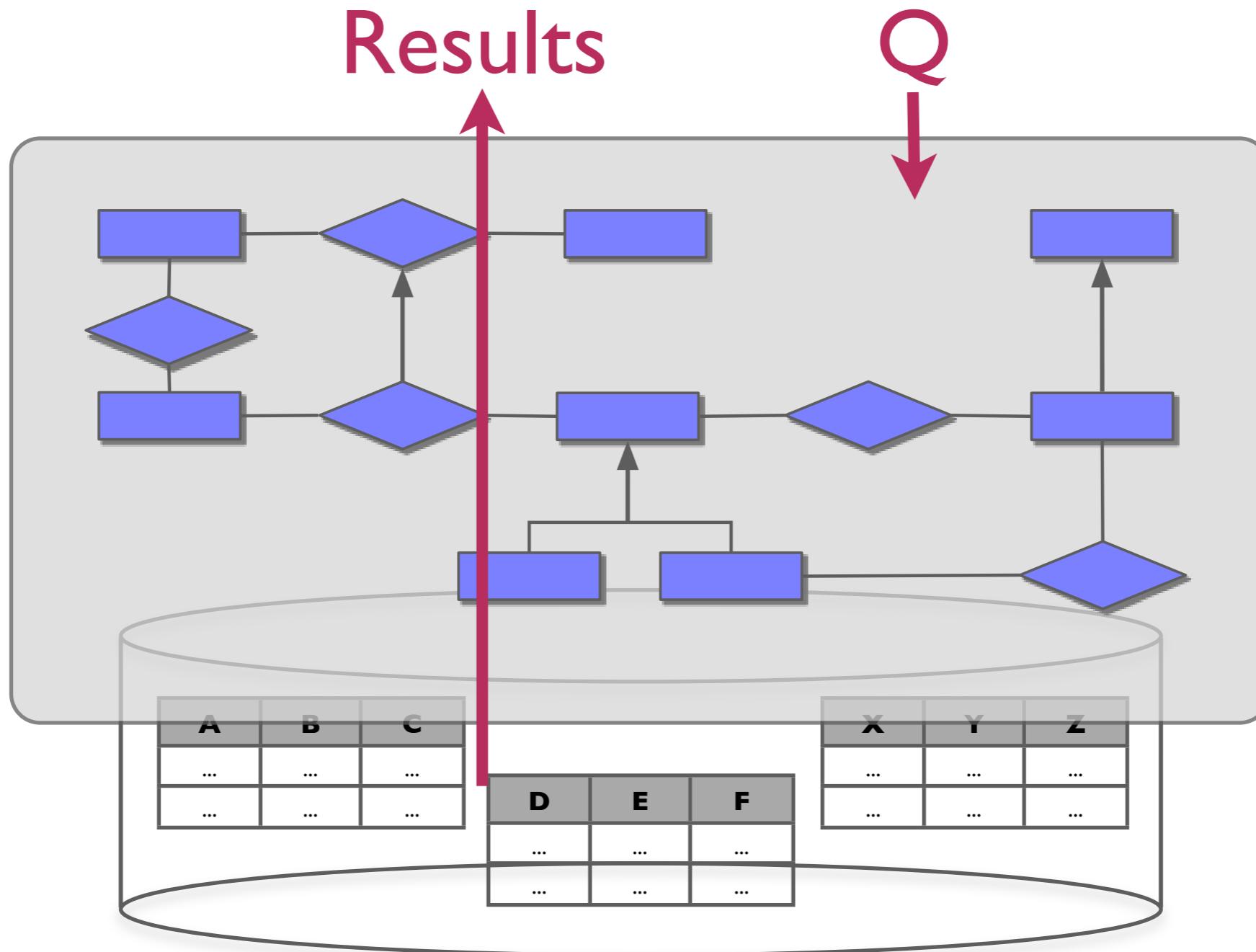
Semantic data access



Semantic data access

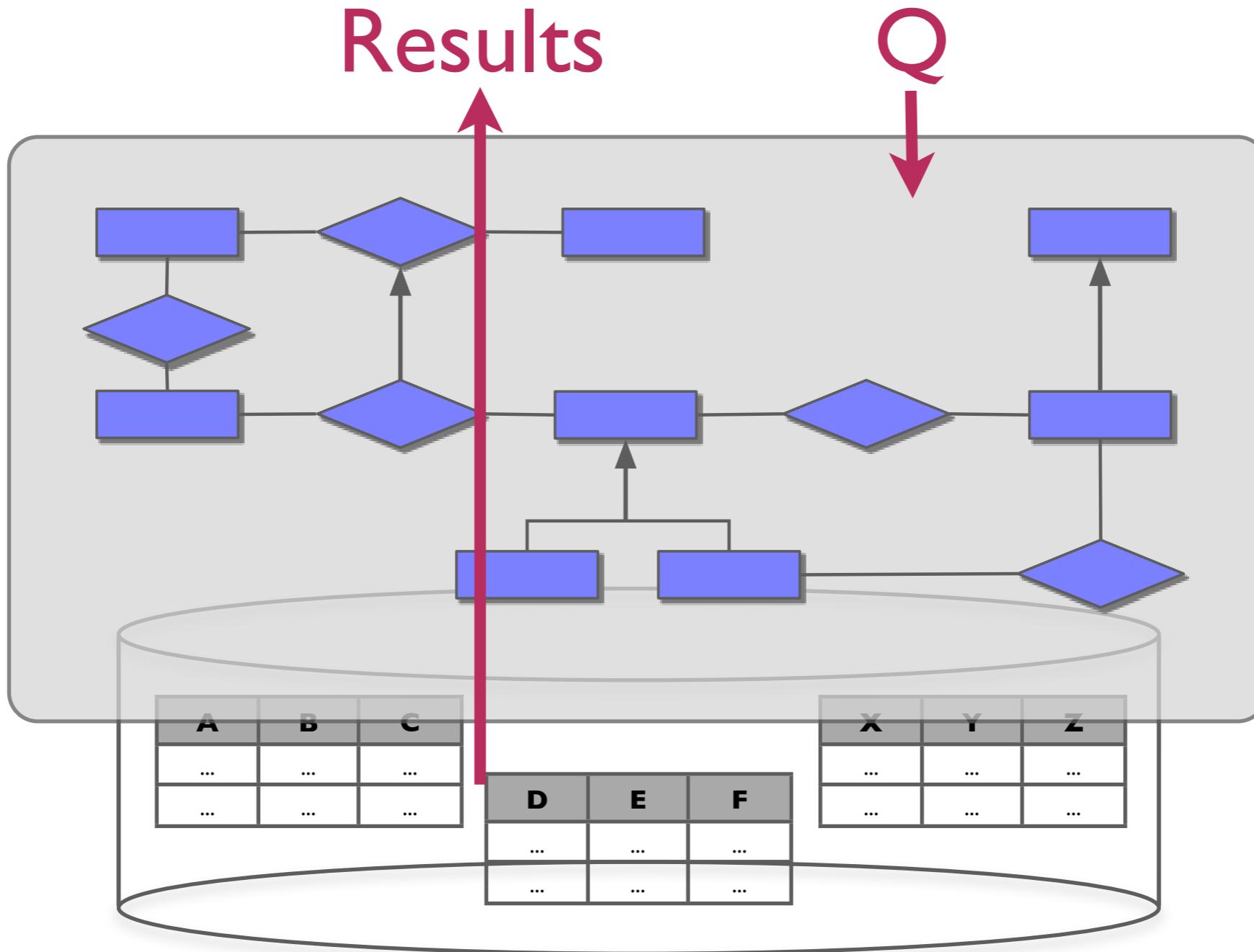


Semantic data access

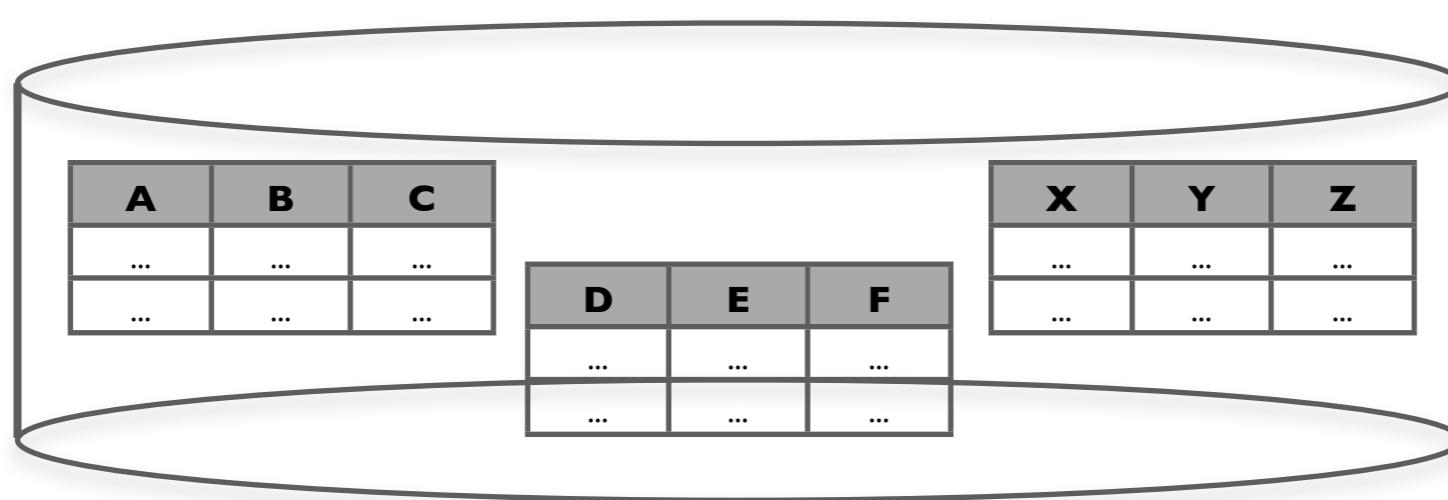


Semantic data access

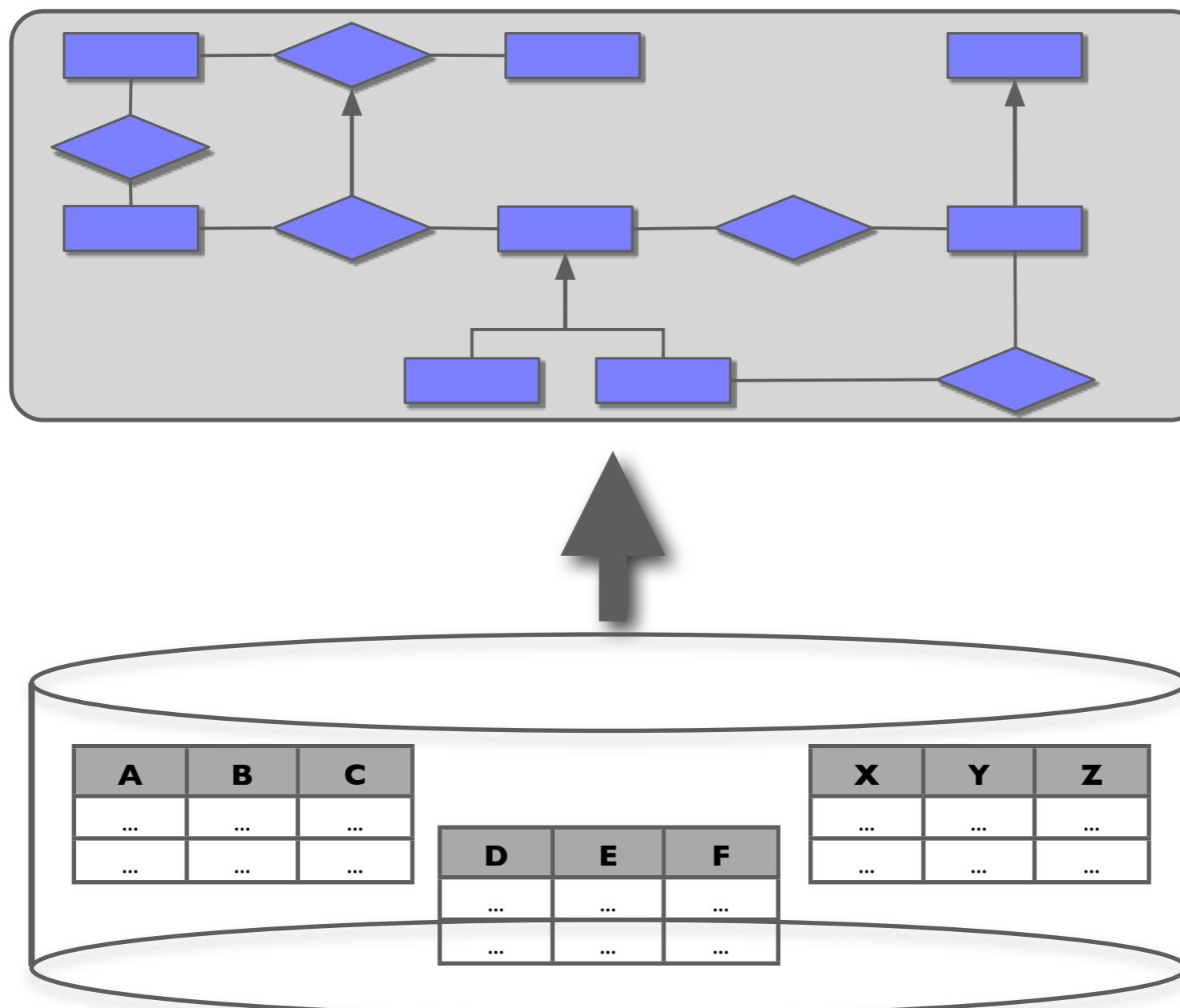
so far done mainly manually



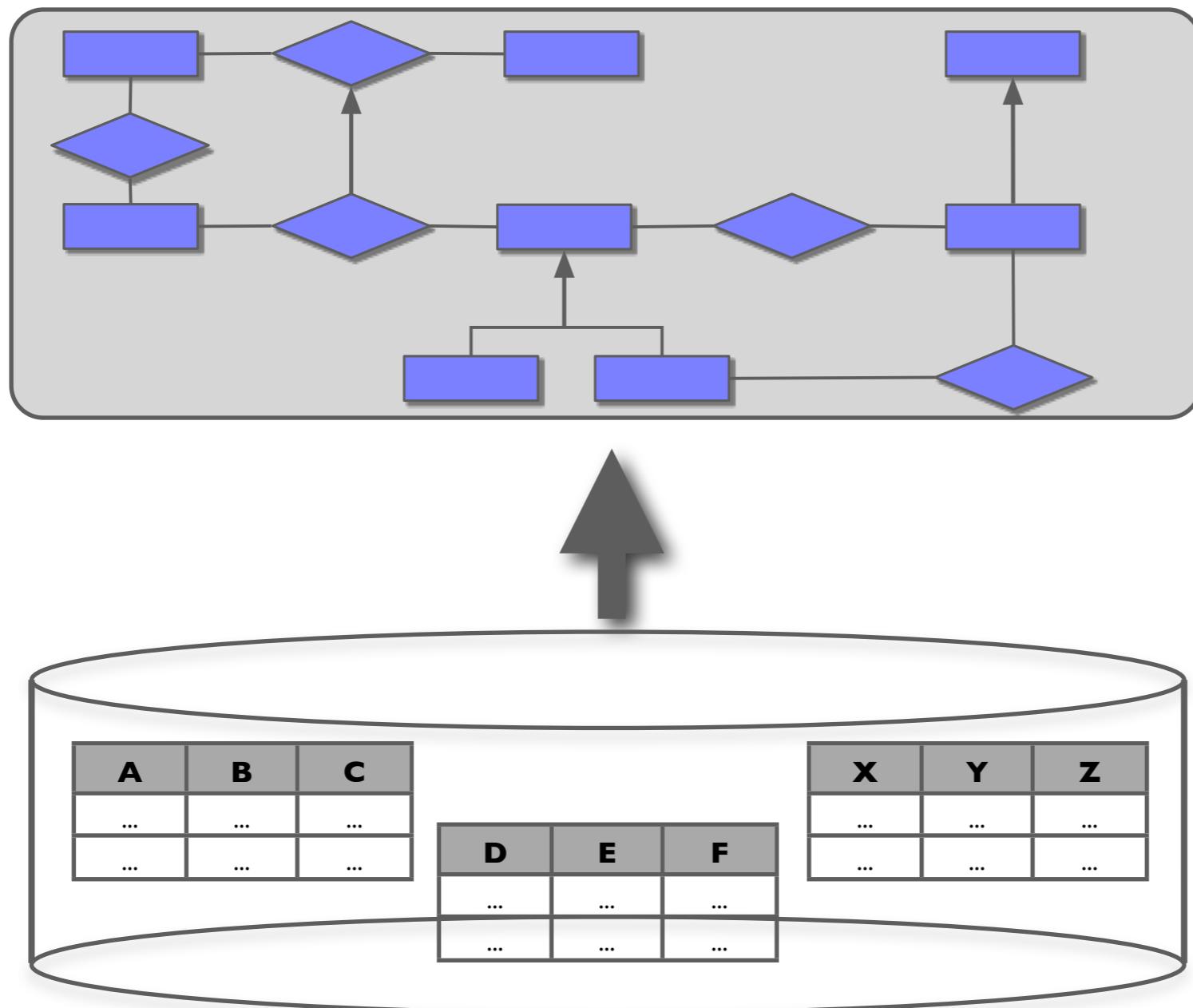
Extracting a wrapping ontology



Extracting a wrapping ontology

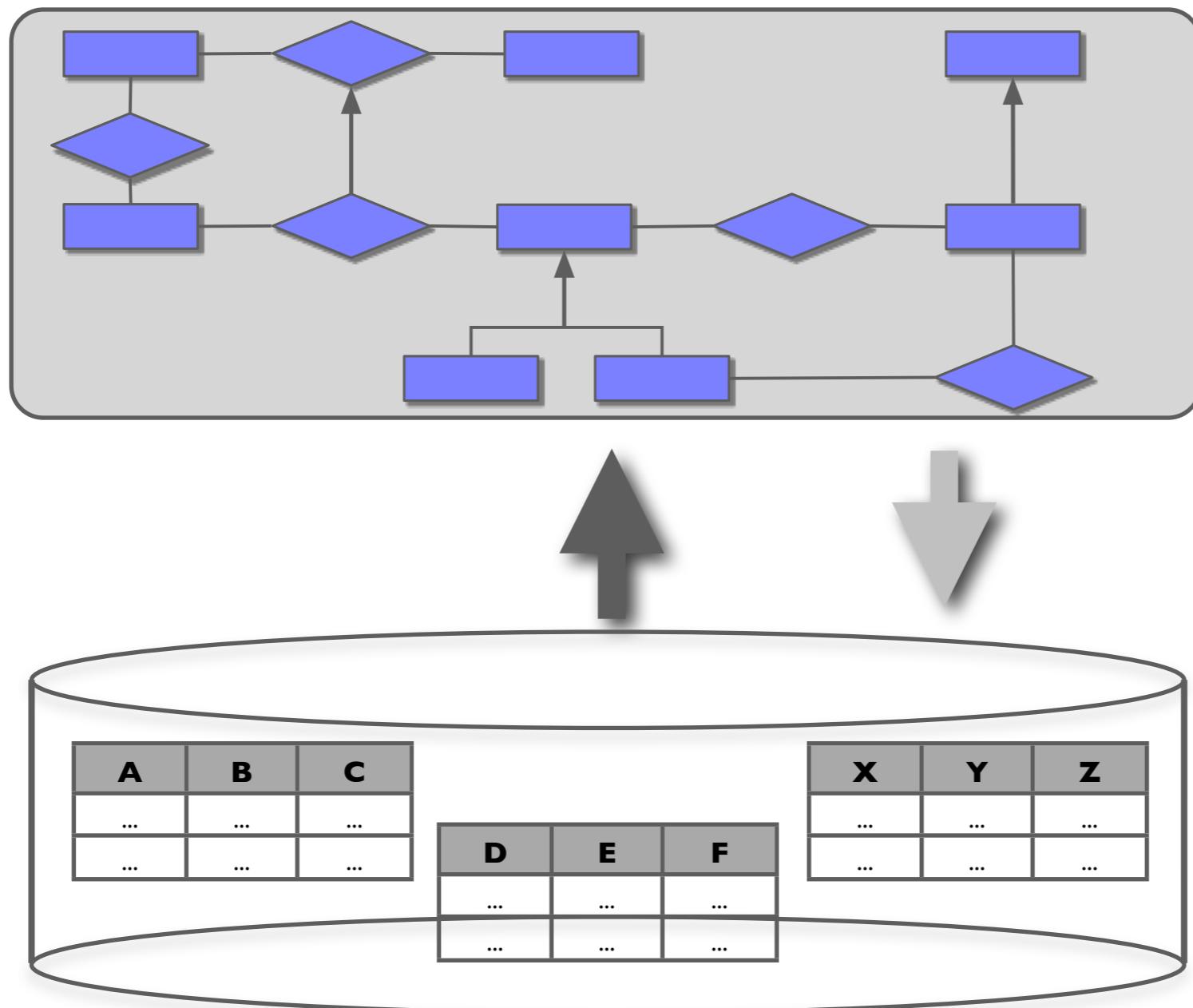


Extracting a wrapping ontology



preserve the
semantics of the
constraints in the DB

Extracting a wrapping ontology



preserve the
semantics of the
constraints in the DB

maintain the link with
the DB

Outline

- Formal Framework
- Ontology Extraction
- Implementation and Case Study
- Related Work
- Current Work

Relational source

- A **relational source** $\text{DB} = \{\Psi, \Sigma\}$, with Ψ a **relational schema** and Σ set of **integrity constraints**:
 - **nulls-not-allowed**: $\text{nonnull}(r, A)$
 - **unique**: $\text{unique}(r, A)$
 - **key**: $\text{key}(r, A)$, if $\text{nonnull}(r, A)$ and $\text{unique}(r, A)$
 - **inclusion**: $r_1[s_1] \subseteq r_2[s_2]$
 - **foreign key**: if $r_1[s_1] \subseteq r_2[s_2]$ and $\text{key}(r_2, s_2)$
 - **exclusion**: $(r_1[s_1] \cap r_2[s_2]) = \emptyset$
 - **covering**: $(r_1[s_1] \cup \dots \cup r_m[s_m]) \subseteq r_0[s_0]$

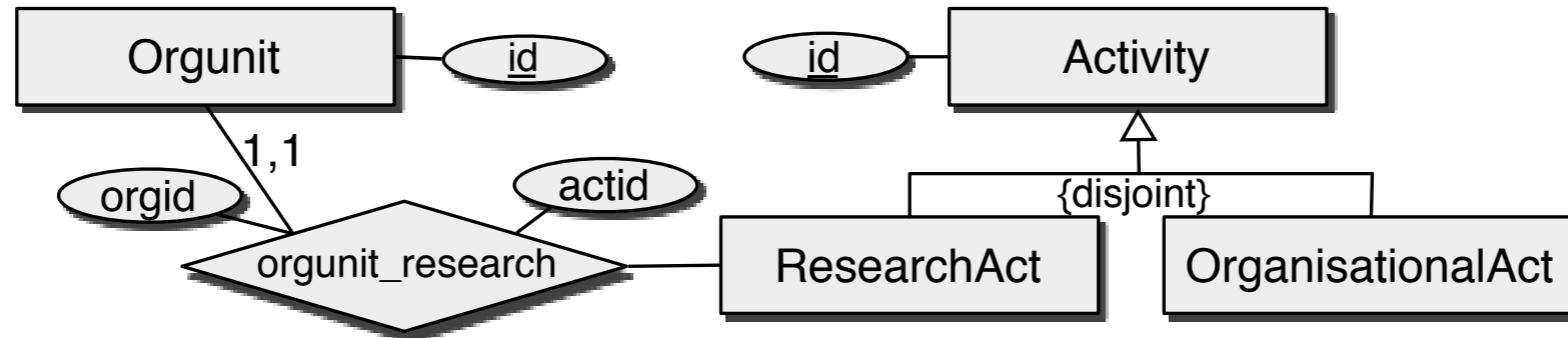
Ontology language

- DLR-DB system $S = \{R, K\}$, with R a relational schema and K set of axioms involving names from R :

$R[s] \sqsubseteq R'[s']$	$\pi_s R^{\mathcal{D}} \subseteq \pi_{s'} R'^{\mathcal{D}}$	Inclusion
$R[s] \text{ disj } R'[s']$	$\pi_s R^{\mathcal{D}} \cap \pi_{s'} R'^{\mathcal{D}} = \emptyset$	Disjointness
$\text{key}(R[s_1, \dots, s_k])$	for all $\phi_1, \phi_2 \in R^{\mathcal{D}}$ with $\phi_1 \neq \phi_2$, we have $\phi_1(s_i) \neq \phi_2(s_i)$ for some s_i , $1 \leq i \leq k$	Key
$R_1[s_1], \dots, R_k[s_k]$ cover $R[s]$	$\pi_s R^{\mathcal{D}} \subseteq \bigcup_{i=1 \dots k} \pi_{s_i} R_i^{\mathcal{D}}$	Covering

- able to express common constraints of ER diagrams
- can be reduced to DL DLR-Lite
 - use of efficient (LogSpace) query answering technique
 - compatible with OWL

Ontology language: example



$\text{orgunit_research}[\text{orgid}] \sqsubseteq \text{Orgunit}[\text{id}]$

$\text{orgunit_research}[\text{actid}] \sqsubseteq \text{ResearchAct}[\text{id}]$

$\text{key}(\text{Orgunit}[\text{id}])$

$\text{key}(\text{Activity}[\text{id}])$

$\text{key}(\text{ResearchAct}[\text{id}])$

$\text{key}(\text{OrganisationalAct}[\text{id}])$

$\text{key}(\text{orgunit_research}[\text{orgid}])$

$\text{Orgunit}[\text{id}] \sqsubseteq \text{orgunit_research}[\text{orgid}]$

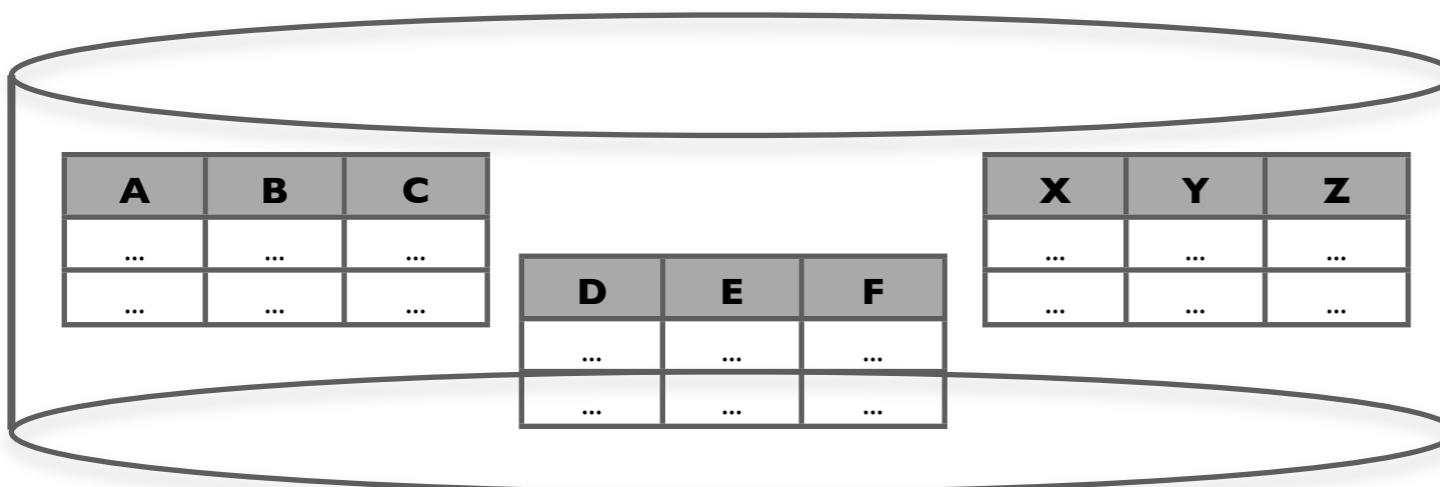
$\text{ResearchAct}[\text{id}] \sqsubseteq \text{Activity}[\text{id}]$

$\text{OrganisationalAct}[\text{id}] \sqsubseteq \text{Activity}[\text{id}]$

$\text{ResearchAct}[\text{id}] \text{ disj } \text{OrganisationalAct}[\text{id}]$

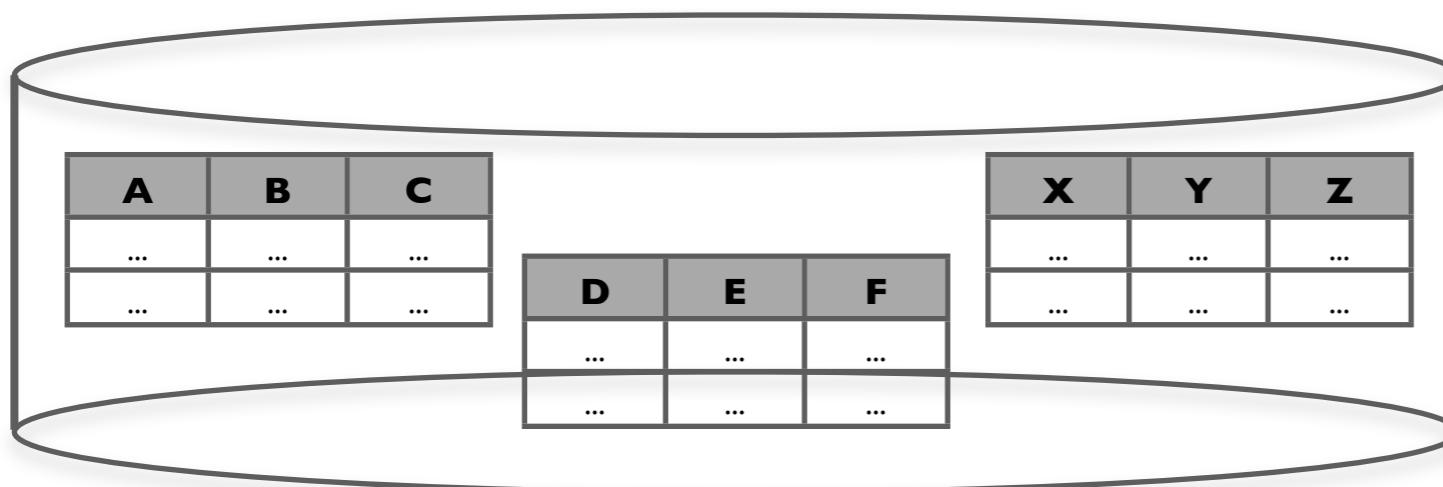
Ontology extraction: intuition

Basic idea: **reverse** the standard database modeling process



Ontology extraction: intuition

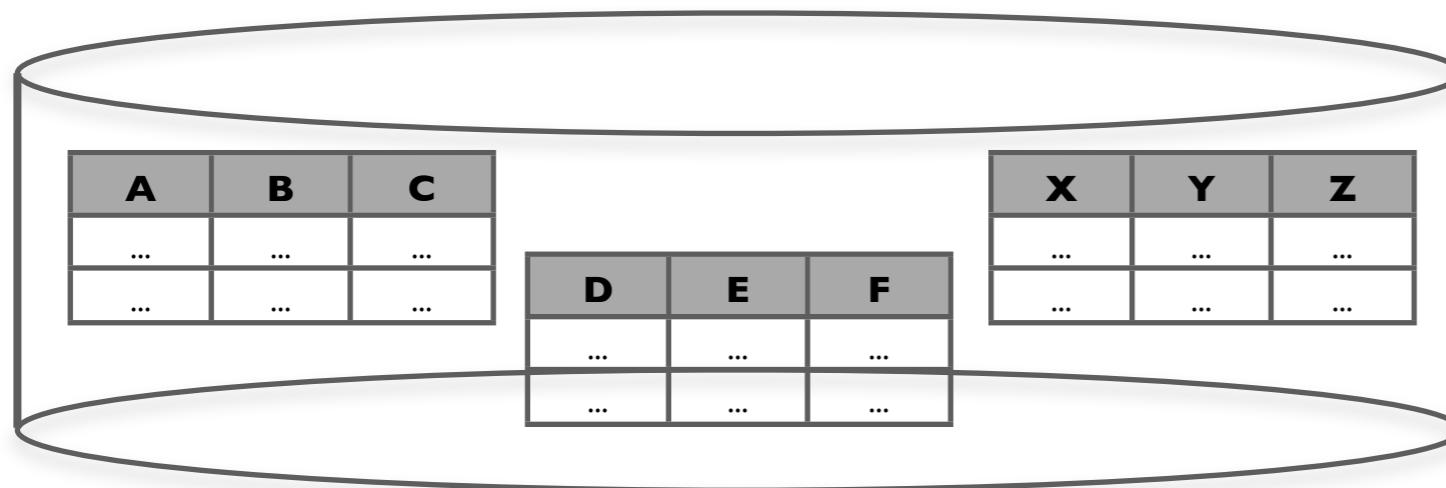
Basic idea: **reverse** the standard database modeling process



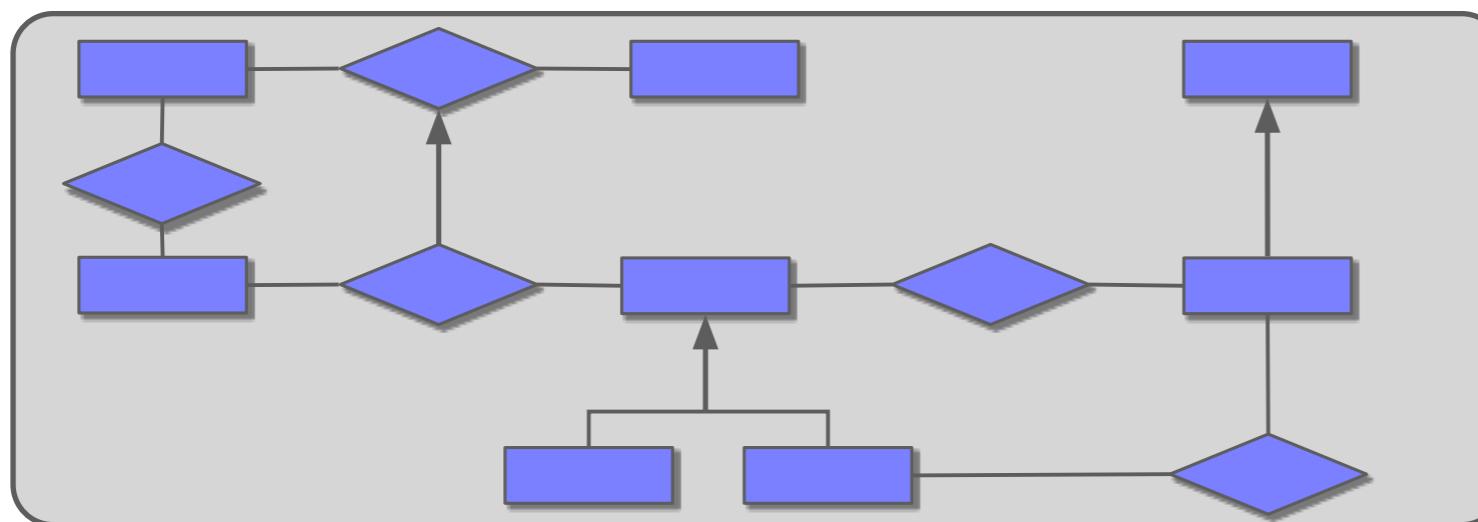
Step 1: classify relations by
analyzing the integrity
constraints in Σ

Ontology extraction: intuition

Basic idea: **reverse** the standard database modeling process



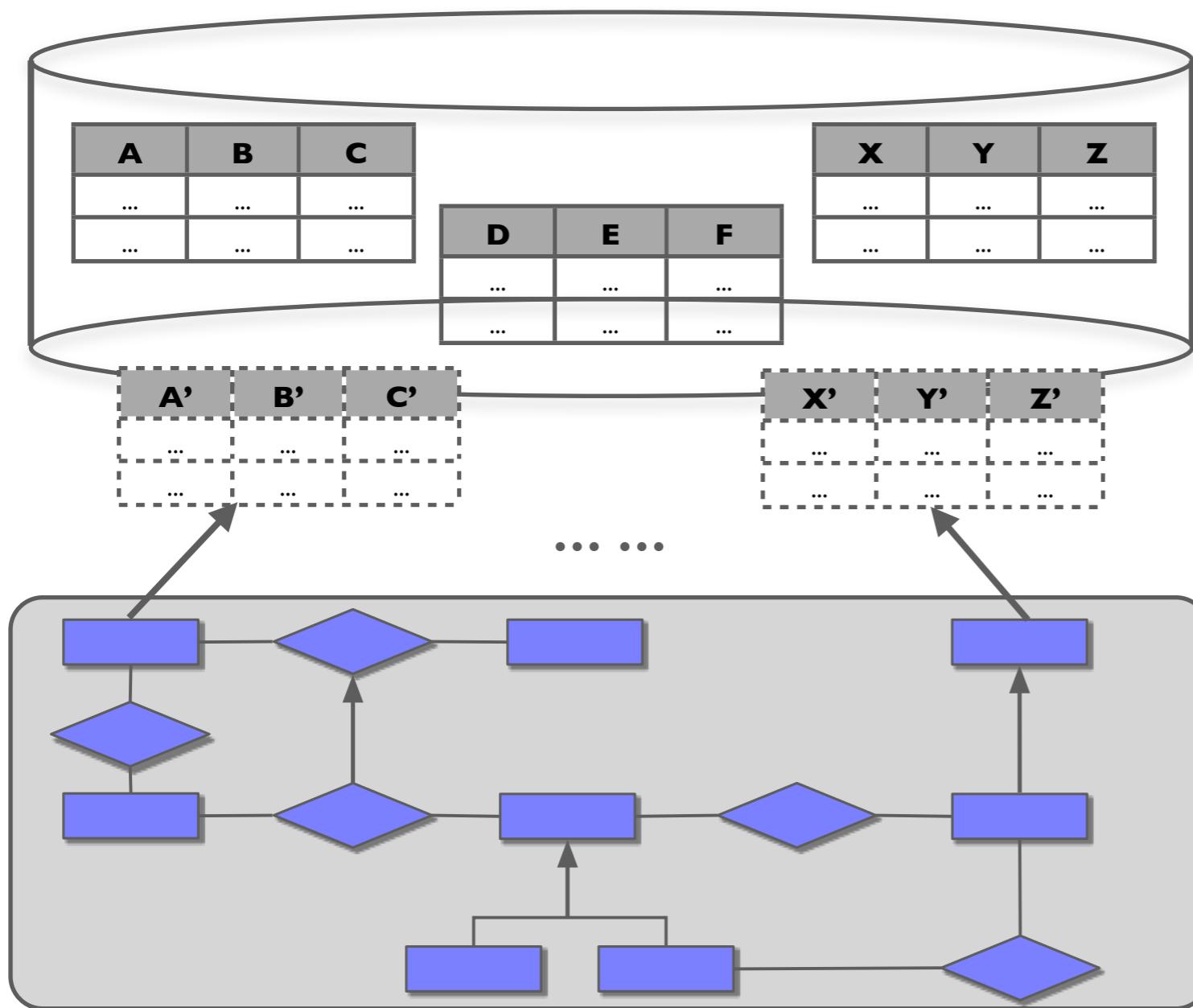
Step 1: classify relations by analyzing the integrity constraints in Σ



Step 2: build the ontology by means of DLR-DB axioms K and generate view in R for each ontology element

Ontology extraction: intuition

Basic idea: **reverse** the standard database modeling process



Step 1: classify relations by analyzing the integrity constraints in Σ

Step 2: build the ontology by means of DLR-DB axioms K and generate view in R for each ontology element

Classification of relations

Classification of relations

- Base, if K and FK don't share attributes

Classification of relations

- **Base**, if K and FK don't share attributes
- **Relationship**, if K is entirely composed of FK and $|\text{FK}| > 1$

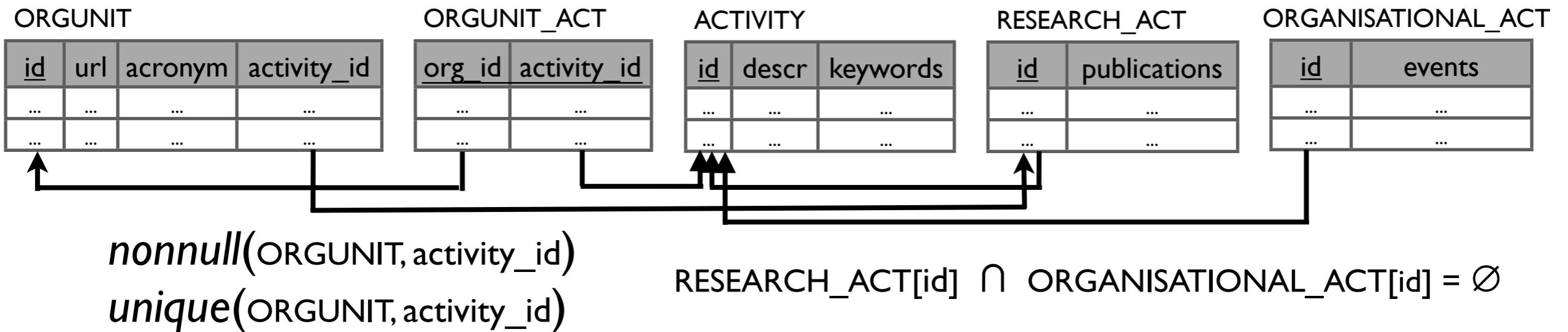
Classification of relations

- **Base**, if K and FK don't share attributes
- **Relationship**, if K is entirely composed of FK and $|\text{FK}| > 1$
- **Specific**, if K is among FK and either $|\text{FK}| = 1$, or K is referenced by some other relation

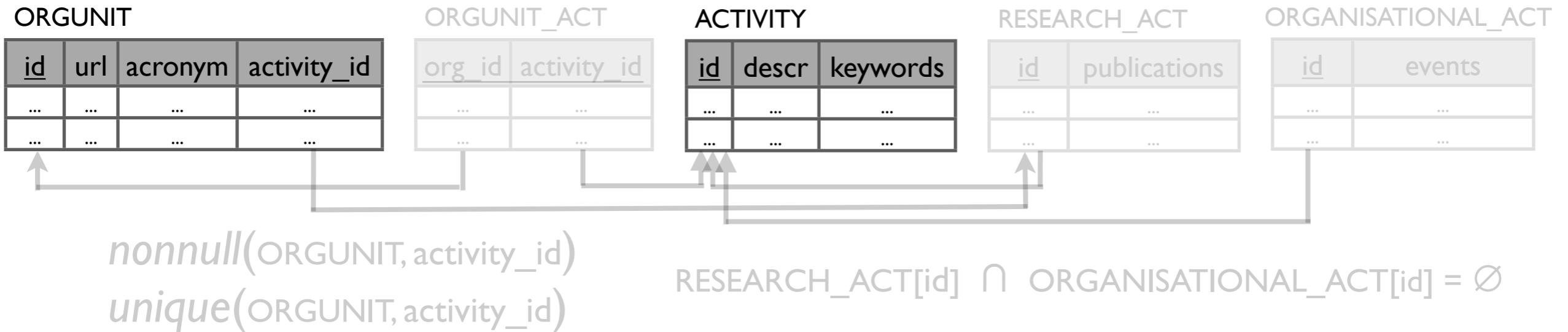
Classification of relations

- **Base**, if K and FK don't share attributes
- **Relationship**, if K is entirely composed of FK and $|\text{FK}| > 1$
- **Specific**, if K is among FK and either $|\text{FK}| = 1$, or K is referenced by some other relation
- **Ambiguous**, otherwise

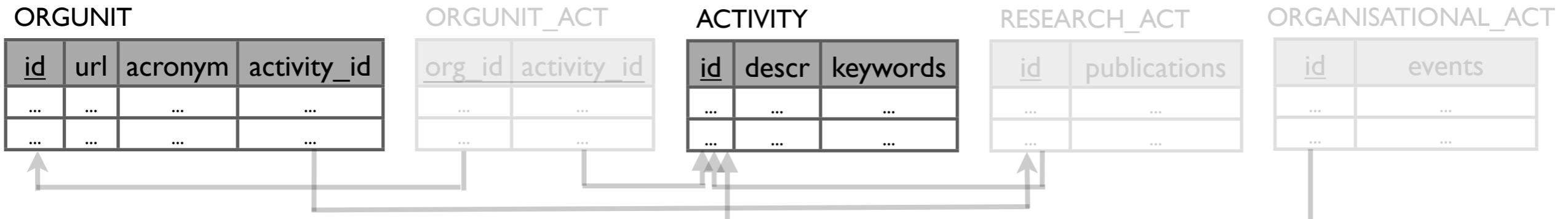
Ontology extraction by example



Ontology extraction by example



Ontology extraction by example



nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit

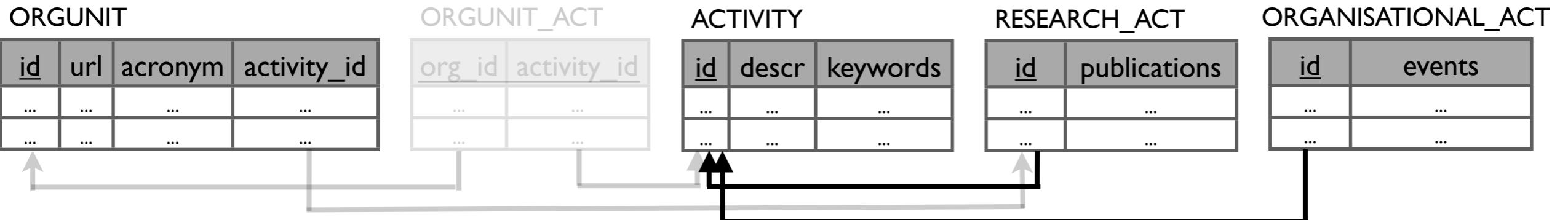
<u><u>id</u></u>	<u><u>url</u></u>	<u><u>acronym</u></u>
...
...

activity

<u><u>id</u></u>	<u><u>descr</u></u>	<u><u>keywords</u></u>
...
...

key(Orgunit[id])
 key(Activity[id])

Ontology extraction by example



nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit

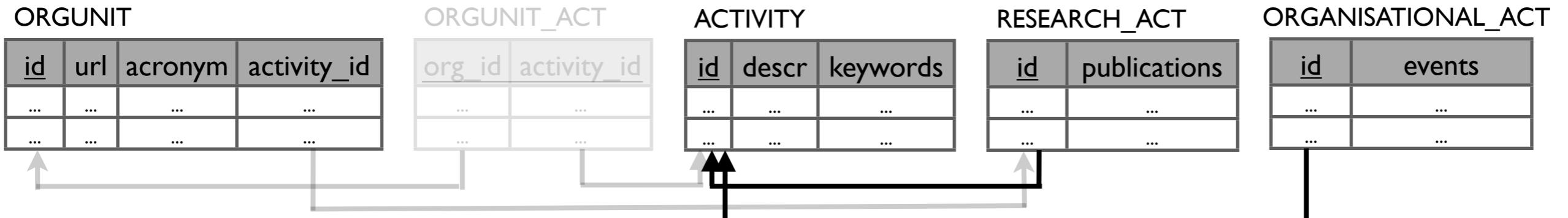
<u><u>id</u></u>	url	acronym
...
...

activity

<u><u>id</u></u>	descr	keywords
...
...

key(Orgunit[id])
key(Activity[id])

Ontology extraction by example



*nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)*

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit
<code>id</code>
<code>url</code>
<code>acronym</code>

activity
<code>id</code>
<code>descr</code>
<code>keywords</code>

research_act
<code>id</code>
<code>publications</code>

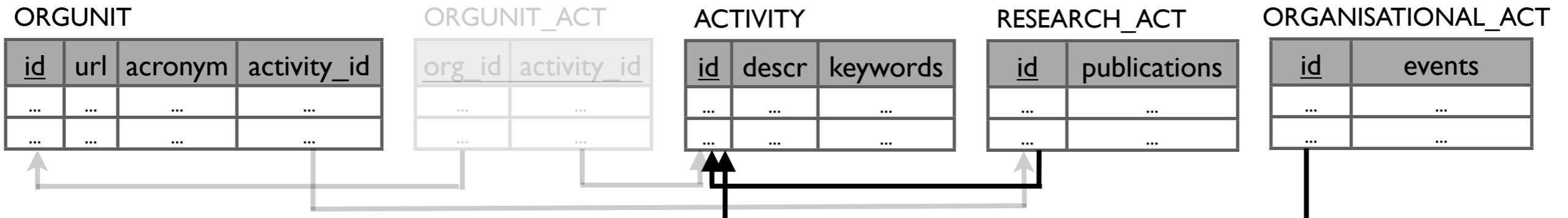
organisational_act
<code>id</code>
<code>events</code>

key(Orgunit[id]) **key(ResearchAct[id])**
key(Activity[id]) **key(OrganisationalAct[id])**

ResearchAct[id] ⊑ Activity[id]

OrganisationalAct[id] ⊑ Activity[id]

Ontology extraction by example



*nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)*

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit
<code>id</code>
<code>url</code>
<code>acronym</code>

activity
<code>id</code>
<code>descr</code>
<code>keywords</code>

research_act
<code>id</code>
<code>publications</code>

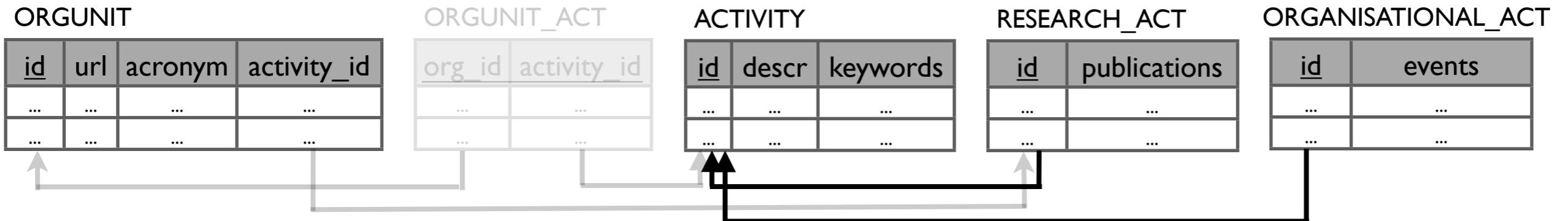
organisational_act
<code>id</code>
<code>events</code>

key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$

$\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$

Ontology extraction by example



*nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)*

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit
<code>id</code>
<code>url</code>
<code>acronym</code>

activity
<code>id</code>
<code>descr</code>
<code>keywords</code>

research_act
<code>id</code>
<code>publications</code>

organisational_act
<code>id</code>
<code>events</code>

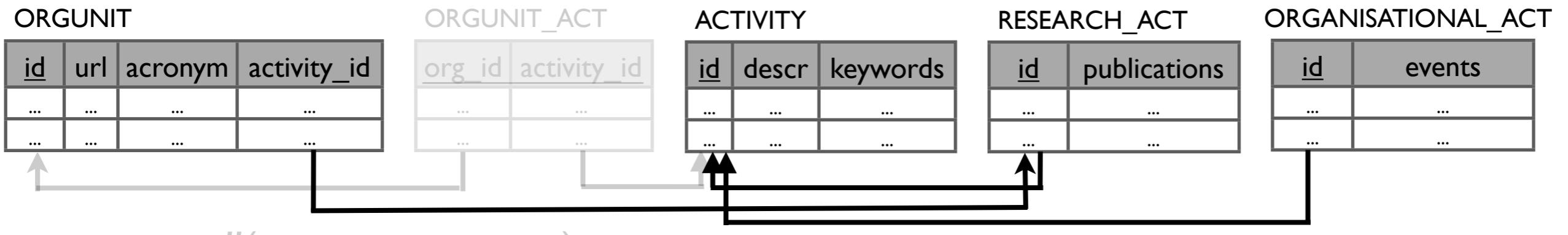
key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$

$\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$

$\text{ResearchAct[id]} \text{ disj } \text{OrganisationalAct[id]}$

Ontology extraction by example



*nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)*

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

orgunit
<code>id</code>
<code>url</code>
<code>acronym</code>

activity
<code>id</code>
<code>descr</code>
<code>keywords</code>

research_act
<code>id</code>
<code>publications</code>

organisational_act
<code>id</code>
<code>events</code>

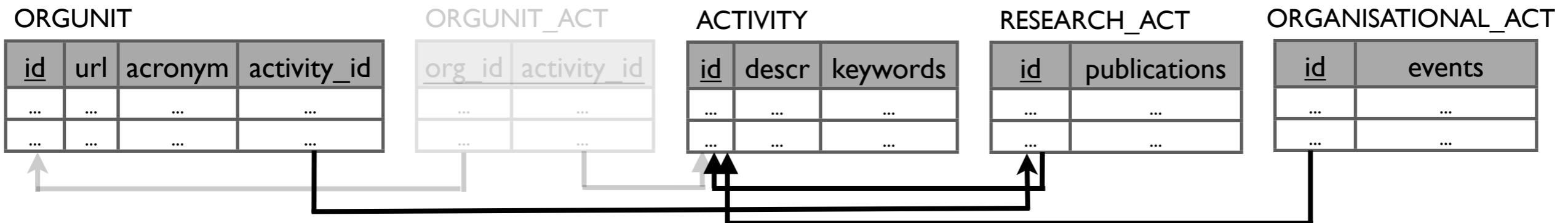
key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$

$\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$

$\text{ResearchAct[id]} \text{ disj } \text{OrganisationalAct[id]}$

Ontology extraction by example



$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

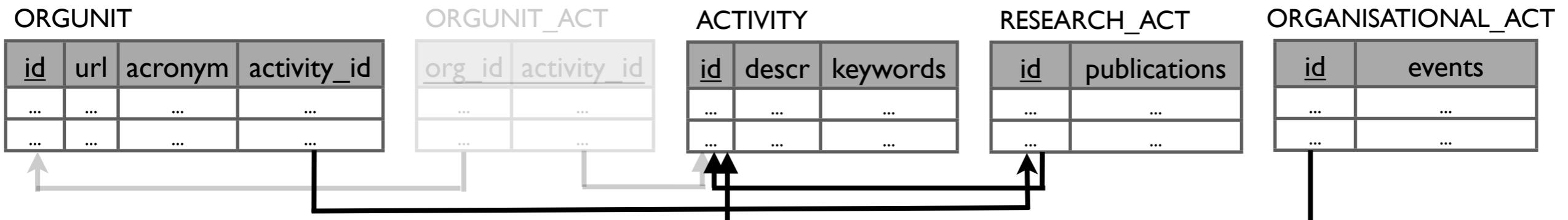


key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{orgunit_research[org_id]} \sqsubseteq \text{Orgunit[id]}$
 $\text{orgunit_research[res_act_id]} \sqsubseteq \text{ResearchAct[id]}$

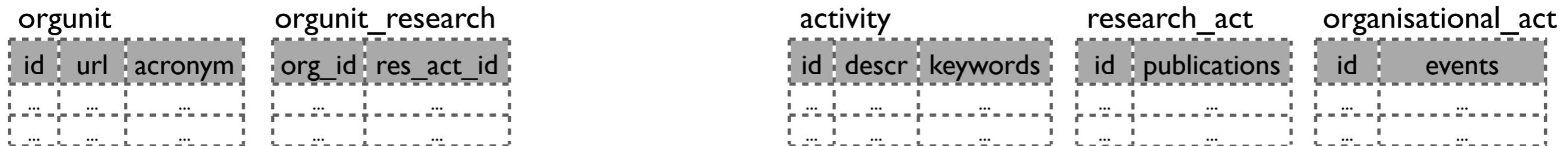
$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{ResearchAct[id]} \text{ disj } \text{OrganisationalAct[id]}$

Ontology extraction by example



$\text{nonnull}(\text{ORGUNIT}, \text{activity_id})$
 $\text{unique}(\text{ORGUNIT}, \text{activity_id})$

$\text{RESEARCH_ACT}[id] \cap \text{ORGANISATIONAL_ACT}[id] = \emptyset$



$\text{key}(\text{Orgunit}[id])$
 $\text{key}(\text{Activity}[id])$

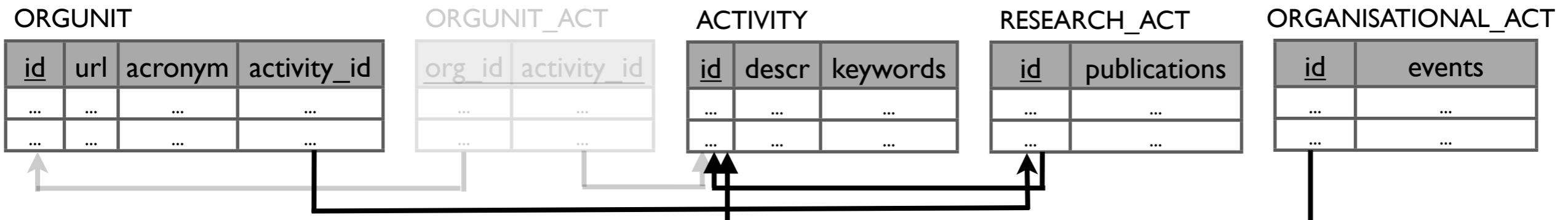
$\text{key}(\text{ResearchAct}[id])$
 $\text{key}(\text{OrganisationalAct}[id])$

$\text{ResearchAct}[id] \sqsubseteq \text{Activity}[id]$

$\text{OrganisationalAct}[id] \sqsubseteq \text{Activity}[id]$

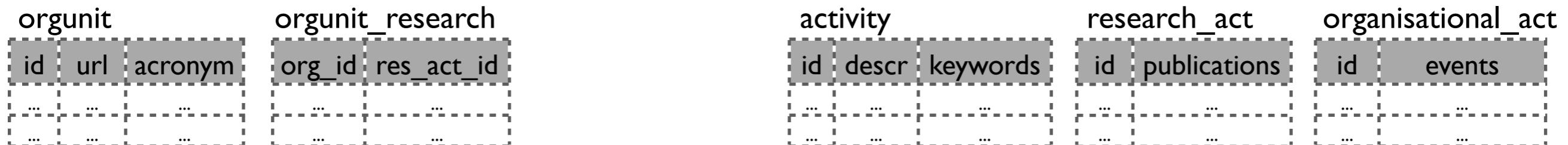
$\text{ResearchAct}[id] \text{ disj } \text{OrganisationalAct}[id]$

Ontology extraction by example



$\text{nonnull}(\text{ORGUNIT}, \text{activity_id})$
 $\text{unique}(\text{ORGUNIT}, \text{activity_id})$

$\text{RESEARCH_ACT}[id] \cap \text{ORGANISATIONAL_ACT}[id] = \emptyset$

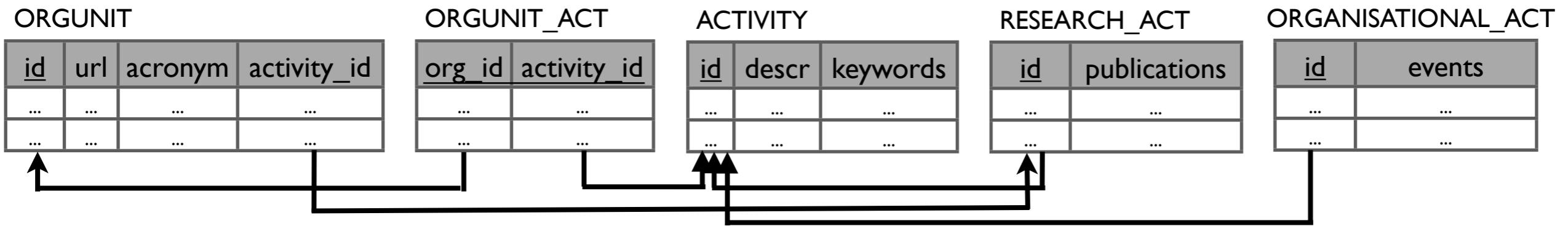


$\text{key}(\text{Orgunit}[id])$
 $\text{key}(\text{Activity}[id])$

$\text{ResearchAct}[id] \sqsubseteq \text{Activity}[id]$
 $\text{OrganisationalAct}[id] \sqsubseteq \text{Activity}[id]$
 $\text{ResearchAct}[id] \text{ disj } \text{OrganisationalAct}[id]$

$\text{orgunit_research}[org_id] \sqsubseteq \text{Orgunit}[id]$
 $\text{orgunit_research}[res_act_id] \sqsubseteq \text{ResearchAct}[id]$
 $\text{Orgunit}[id] \sqsubseteq \text{orgunit_research}[org_id]$
 $\text{key}(\text{orgunit_research}[org_id])$

Ontology extraction by example



nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$

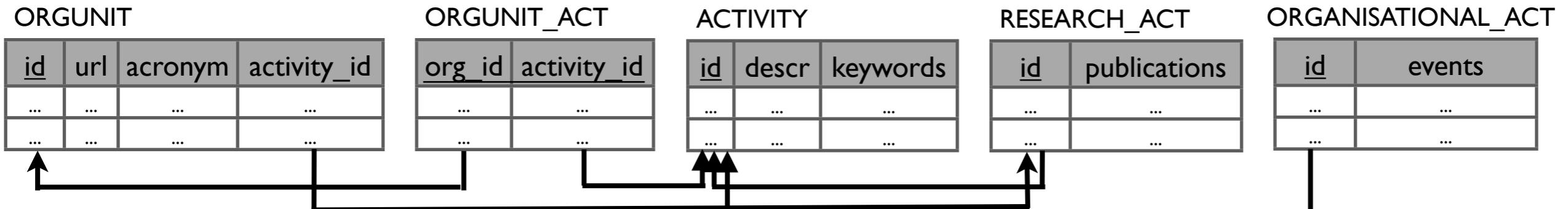


key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{ResearchAct[id]} \text{ disj OrganisationalAct[id]}$

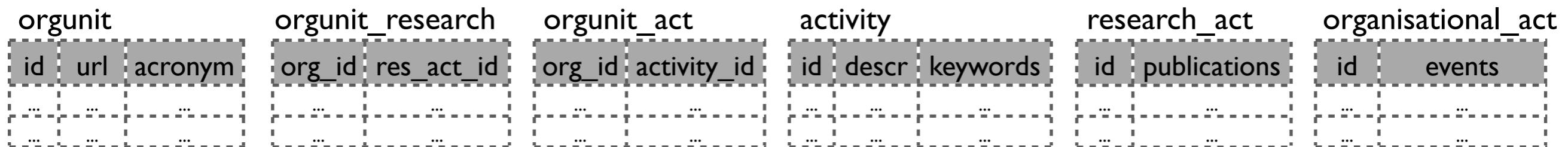
$\text{orgunit_research[org_id]} \sqsubseteq \text{Orgunit[id]}$
 $\text{orgunit_research[res_act_id]} \sqsubseteq \text{ResearchAct[id]}$
 $\text{Orgunit[id]} \sqsubseteq \text{orgunit_research[org_id]}$
 $\text{key(orgunit_research[org_id])}$

Ontology extraction by example



nonnull(ORGUNIT, activity_id)
unique(ORGUNIT, activity_id)

$\text{RESEARCH_ACT}[\text{id}] \cap \text{ORGANISATIONAL_ACT}[\text{id}] = \emptyset$



key(Orgunit[id]) $\text{key(ResearchAct[id])}$
 key(Activity[id]) $\text{key(OrganisationalAct[id])}$

$\text{ResearchAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{OrganisationalAct[id]} \sqsubseteq \text{Activity[id]}$
 $\text{ResearchAct[id]} \text{ disj OrganisationalAct[id]}$

$\text{orgunit_research[org_id]} \sqsubseteq \text{Orgunit[id]}$
 $\text{orgunit_research[res_act_id]} \sqsubseteq \text{ResearchAct[id]}$
 $\text{Orgunit[id]} \sqsubseteq \text{orgunit_research[org_id]}$
 $\text{key(orgunit_research[org_id])}$
 $\text{orgunit_act[org_id]} \sqsubseteq \text{Orgunit[id]}$
 $\text{orgunit_act[activity_id]} \sqsubseteq \text{Activity[id]}$

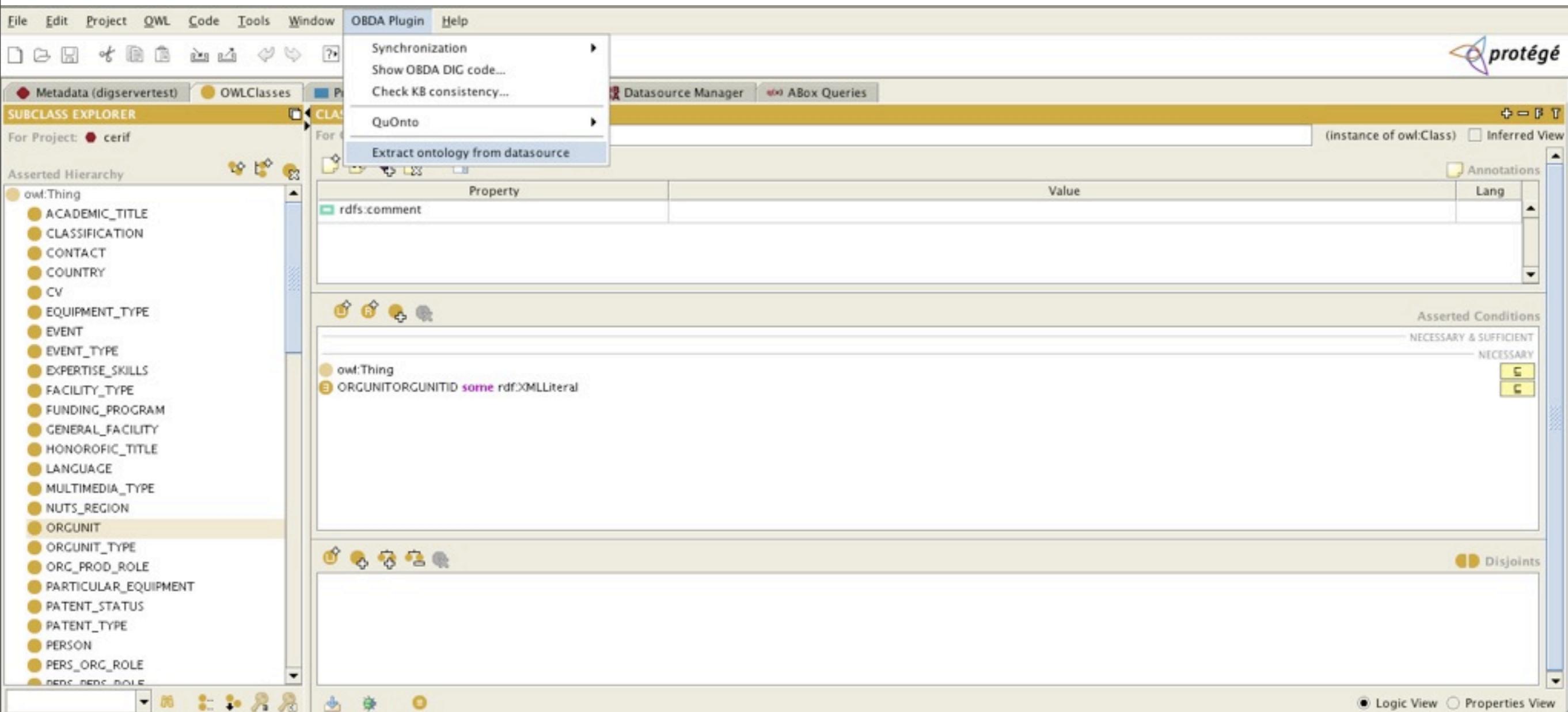
Correctness and completeness

- Correctness and completeness using the general concept of **information capacities** [Miller et al.;'93]
- **Theorem:** the ontology extraction procedure is an **equivalence preserving schema transformation**
 - no information loss
 - queries expressed over the ontology can be answered by simply expanding the generated views!

Ontology extraction in Protégé

[Rodríguez-Muro et al., IIMAS'08; Poggi et al., OWLED'08; Calvanese et al., OWLED'08],

Ontology extraction in Protégé



[Rodríguez-Muro et al., IIMAS'08; Poggi et al., OWLED'08; Calvanese et al., OWLED'08],

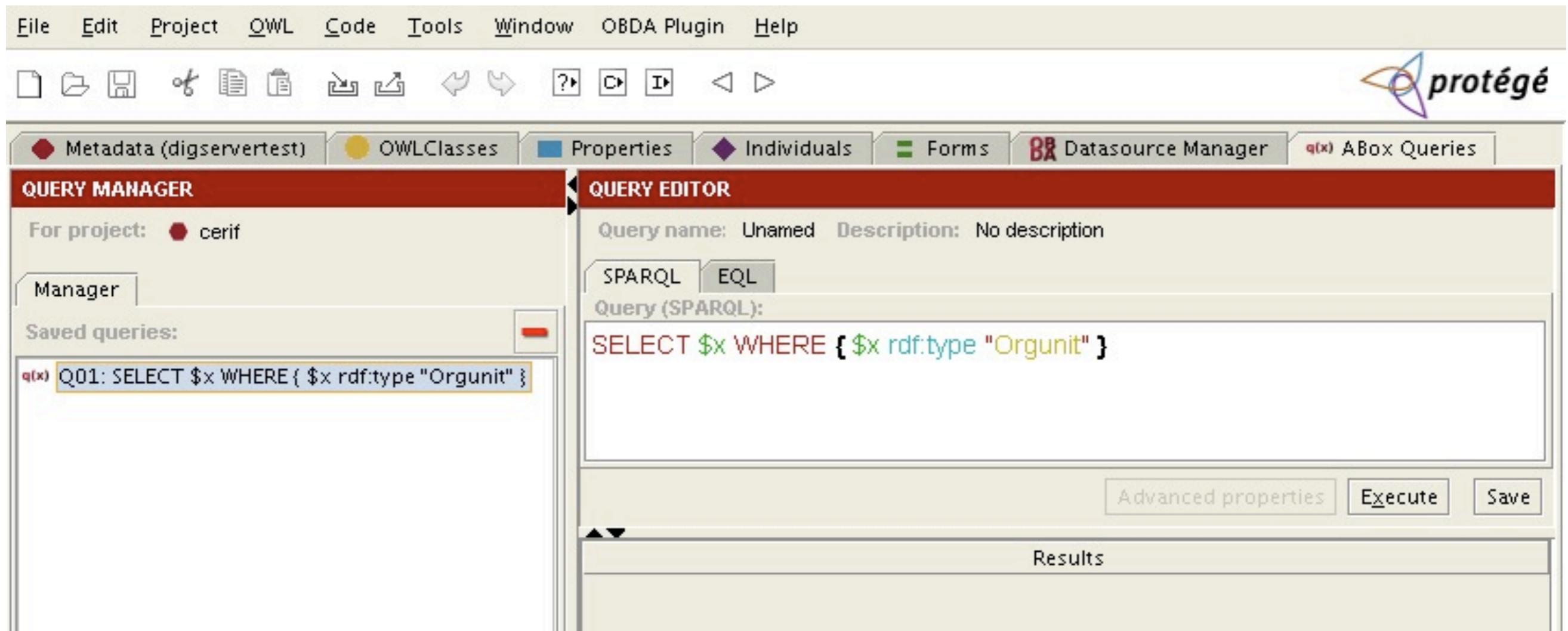
Ontology extraction in Protégé

The screenshot shows the Protégé interface with the 'Datasource Manager' tab selected. The left sidebar displays project metadata and database connection details for the 'cerif' project. The main pane lists 16 mappings (M:f through M:16) with their corresponding SQL queries. The queries involve various ontology classes like NUTS_REGION, ORGUNIT, and PATENT_STATUS, and their properties or relationships.

Mapping ID	Description
M:f	SELECT "NUTS_REGION"."REGION_CODE" AS REGION_CODE, "NUTS_REGION"."REGION_NAME" AS REGION_NAME, "NUTS_REGION"."NAME" AS NAME
M:10	SELECT "ORGUNIT"."ORGUNITID" AS ORGUNITID, "ORGUNIT"."ORG_ACRONYM" AS ORG_ACRONYM, "ORGUNIT"."ORG_TYPE" AS ORG_TYPE
M:11	SELECT "ORGUNIT_TYPE"."ORG_TYPE" AS ORG_TYPE, "ORGUNIT_TYPE"."ORG_TYPE_FULL" AS ORG_TYPE_FULL FROM
M:12	SELECT "ORG_PROD_ROLE"."ORG_PROD_ROLE" AS ORG_PROD_ROLE, "ORG_PROD_ROLE"."ORG_PROD_ROLE_FULL" AS ORG_PROD_ROLE_FULL
M:13	SELECT "PARTICULAR_EQUIPMENT"."EQUIP_ID" AS EQUIP_ID, "PARTICULAR_EQUIPMENT"."EQUIP_OW_INV_ID" AS EQUIP_OW_INV_ID
M:14	SELECT "PATENT_STATUS"."PATENT_STATUS" AS PATENT_STATUS, "PATENT_STATUS"."PATENT_STATUS_FULL" AS PATENT_STATUS_FULL
M:15	SELECT "PATENT_TYPE"."PATENT_TYPE" AS PATENT_TYPE, "PATENT_TYPE"."PATENT_TYPE_FULL" AS PATENT_TYPE_FULL
M:16	SELECT "PERSON"."PER_ID" AS PER_ID, "PERSON"."FAMILY_NAMES" AS FAMILY_NAMES

[Rodriguez-Muro et al., IIMAS'08; Poggi et al., OWLED'08; Calvanese et al., OWLED'08],

Ontology extraction in Protégé



[Rodriguez-Muro et al., IIMAS'08; Poggi et al., OWLED'08; Calvanese et al., OWLED'08],

CERIF case study

- CERIF (Common European Research Information Framework) database schema
 - <http://cordis.europa.eu/cerif/src/toolkit.htm>
 - 123 tables; rich with constraints
 - correct ontology constructs were generated for 77 tables

Relation type	#Classified
base	52
relationship	25
specific	0
ambiguous	46

Axiom	#Extracted
inclusion	68
key	63
disjointness	0
covering	0

CERIF: ambiguous relations

- All 46 ambiguous relations fall into one of the two types:

PERSON_RESEARCH_INTEREST

res_int_language	res_int_trans_type	per_id	keywords
	
...

PROJ_PERSON

proj_id	per_id	proj_per_role
...
...

PERSON

per_id
...
...

PROJECT

proj_id
...
...

PROJ_PERS_ROLE

proj_per_role	...
...	...
...	...

CERIF: ambiguous relations

- All 46 ambiguous relations fall into one of the two types:

PERSON_RESEARCH_INTEREST

res_int_language	res_int_trans_type	per_id	keywords
	
...
...

RESEARCH_INTEREST

language	trans_type
...	...
...	...
...	...

PERSON

per_id
...
...

PROJ_PERSON

proj_id	per_id	proj_per_role
...
...
...

PROJECT

proj_id
...
...

PROJ_PERS_ROLE

proj_per_role	...
...	...
...	...

CERIF: ambiguous relations

- All 46 ambiguous relations fall into one of the two types:

PERSON_RESEARCH_INTEREST

res_int_language	res_int_trans_type	per_id	keywords
	
...
...

RESEARCH_INTEREST

language	trans_type
...	...
...	...
...	...

PERSON

per_id
...
...

PROJ_PERSON

proj_id	per_id	proj_per_role
...
...
...

PROJECT

proj_id
...
...

PROJ_PERS_ROLE

proj_per_role	...
...	...
...	...

→ a slight extension of the algorithm correctly generates CERIF ontology

Related work

Related work

- Database reverse engineering (DBRE) in the 90s
 - [Chiang, Barron et al.;‘94], [Johannesson;‘94], etc.

Related work

- Database reverse engineering (DBRE) in the 90s
 - [Chiang, Barron et al.;'94], [Johannesson; '94], etc.
- RDBs and ontologies for the SW and Information Integration
 - [Astrova; '04], [Volz et al.; '04] use DBRE to migrate DB content into an ontology and to deeply annotate a DB
 - SWARD and Virtuoso automatically generate RDF views over RDBs

Related work

- Database reverse engineering (DBRE) in the 90s
 - [Chiang, Barron et al.;'94], [Johannesson; '94], etc.
- RDBs and ontologies for the SW and Information Integration
 - [Astrova; '04], [Volz et al.; '04] use DBRE to migrate DB content into an ontology and to deeply annotate a DB
 - SWARD and Virtuoso automatically generate RDF views over RDBs
- Linking RDBs and ontologies
 - D2R MAP, R2O, [Poggi et al.;'08] allow to declaratively state mappings
 - [An et al; '05] presents semi-automatic mapping discovery

Advantages of our approach

Advantages of our approach

- Extraction of a **formal ontology wrapping relational sources**
 - automatic generation of mappings
 - enables **automated reasoning** to support the designer

Advantages of our approach

- Extraction of a **formal ontology wrapping relational sources**
 - automatic generation of mappings
 - enables **automated reasoning** to support the designer
- Ensuring **faithfulness** of the extracted model
 - no data loss

Current work

- Tool for exploring the logical schema of a DB and for facilitating the annotation of the schema by means of constraints
- Experiments with other real-world DB schemas to identify other peculiarities that schemas exhibit