

Logical Information Integration

Jose Luis Ambite

University of Southern California

Information Integration

- Information Integration: provide semantically *unified view* to collection of data stored in *multiple, autonomous, and heterogeneous sources*
- Unified view achieved by defining global schema
- Two approaches to realize global schema:
 - Data Exchange: materialized database
 - Data Integration: answer user queries over global schema at runtime, by rewriting and executing queries over sources

Information Integration Challenges

- Syntactic (Access/Format) heterogeneity:
 - Structured Sources: DBMS → ODBC, JDBC
 - Semistructured Sources: HTML, text, pdf, XML → Wrappers, Information Extraction
 - Web services → XML, SOAP, WSDL, REST
 - Geospatial Data → Maps, Images, Vector Data
- Semantic heterogeneity
 - Schema → Describe sources in common domain schema
 - Data → Record Linkage
- Scalability:
 - Mediation
 - Record Linkage
 - Efficient Query Execution

Schema Heterogeneity

- Schema heterogeneity is a fact of life
 - Whenever schemas are designed by different people/organizations, they will be different, even if they model the *same* domain!
- The goal of schema mappings is to reconcile schema heterogeneity
 - Mostly between the mediated schema and the schema of the data sources.

Schema Heterogeneity by Example

Mediated Schema

Movie: title, director, year, genre
Actors: title, name
Plays: movie, location, startTime
Reviews: title, rating, description

Sources

S1

Movie(title, director, year, genre)
Actor(AID, firstName, lastName,
nationality, yearofBirth)
ActorPlays(AID, MID)
MovieDetails(MID, director, genre, year)

S2

Cinemas(place, movie, start)

S3

NYCCinemas(name, title, startTime)

S5

MovieGenres(title, genre)

S6

MovieDirectors(title, dir)

S7

MovieYears(title, year)

S4

Reviews(title, date, grade, review)

Table and Attribute Names

Mediated Schema

Movie: title, director, year, genre

Actors: title, name

Plays: movie, location, startTime

Reviews: title, rating, description

Sources

Table and attribute names

S1

Movie(title, director, year, genre)

Actor(AID, firstName, lastName,
nationality, yearofBirth)

ActorPlays(AID, MID)

MovieDetails(MID, director, genre, year)

S2

Cinemas(place, movie, start)

S3

NYCCinemas(name, title, startTime)

S5

MovieGenres(title, genre)

S6

MovieDirectors(title, dir)

S7

MovieYears(title, year)

S4

Reviews(title, date, grade, review)

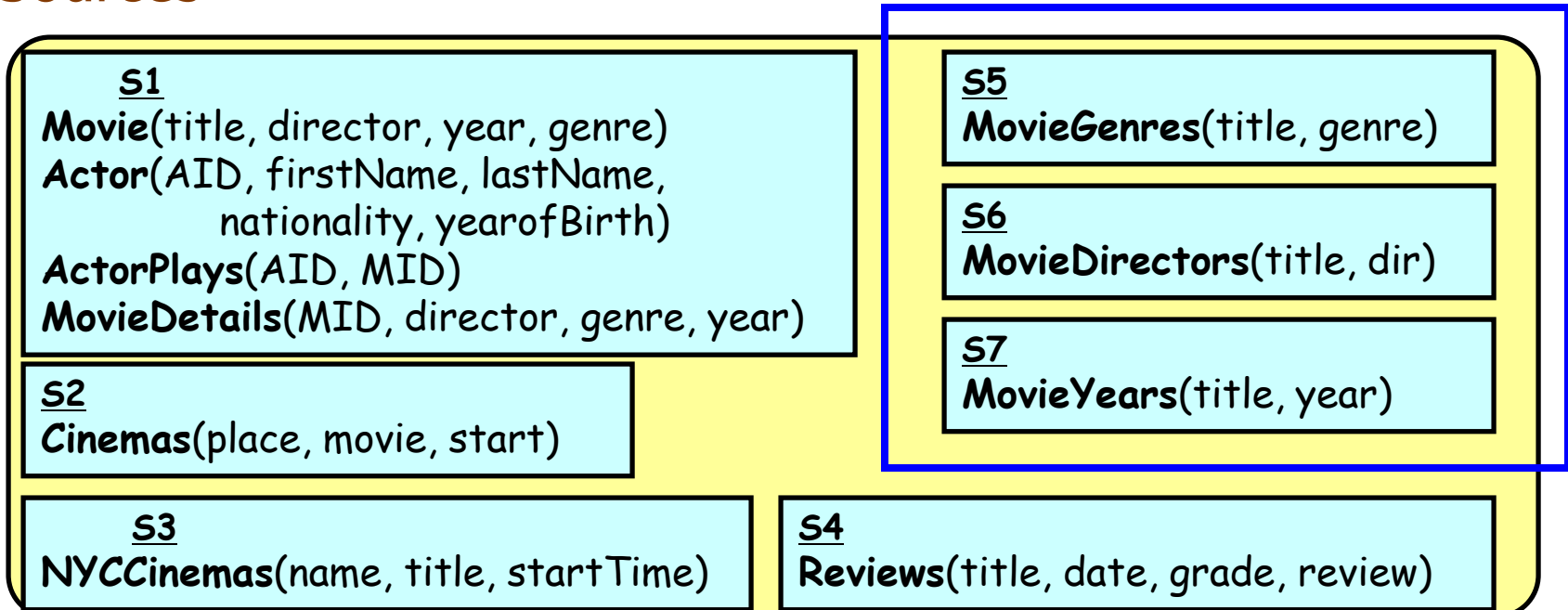
Organization of Schema

Mediated Schema

Movie: title, director, year, genre
Actors: title, name
Plays: movie, location, startTime
Reviews: title, rating, description

Different tabular organization

Sources



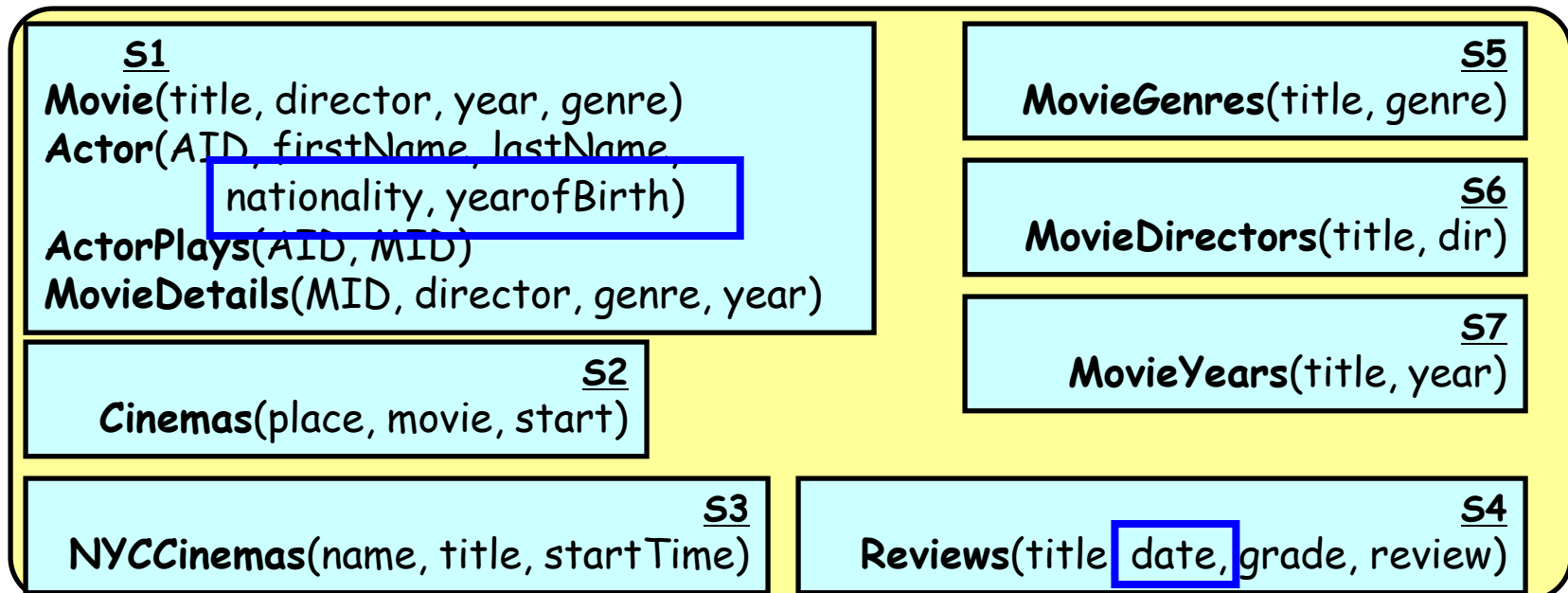
Schema Coverage

Mediated Schema

Movie: title, director, year, genre
Actors: title, name
Plays: movie, location, startTime
Reviews: title, rating, description

Sources

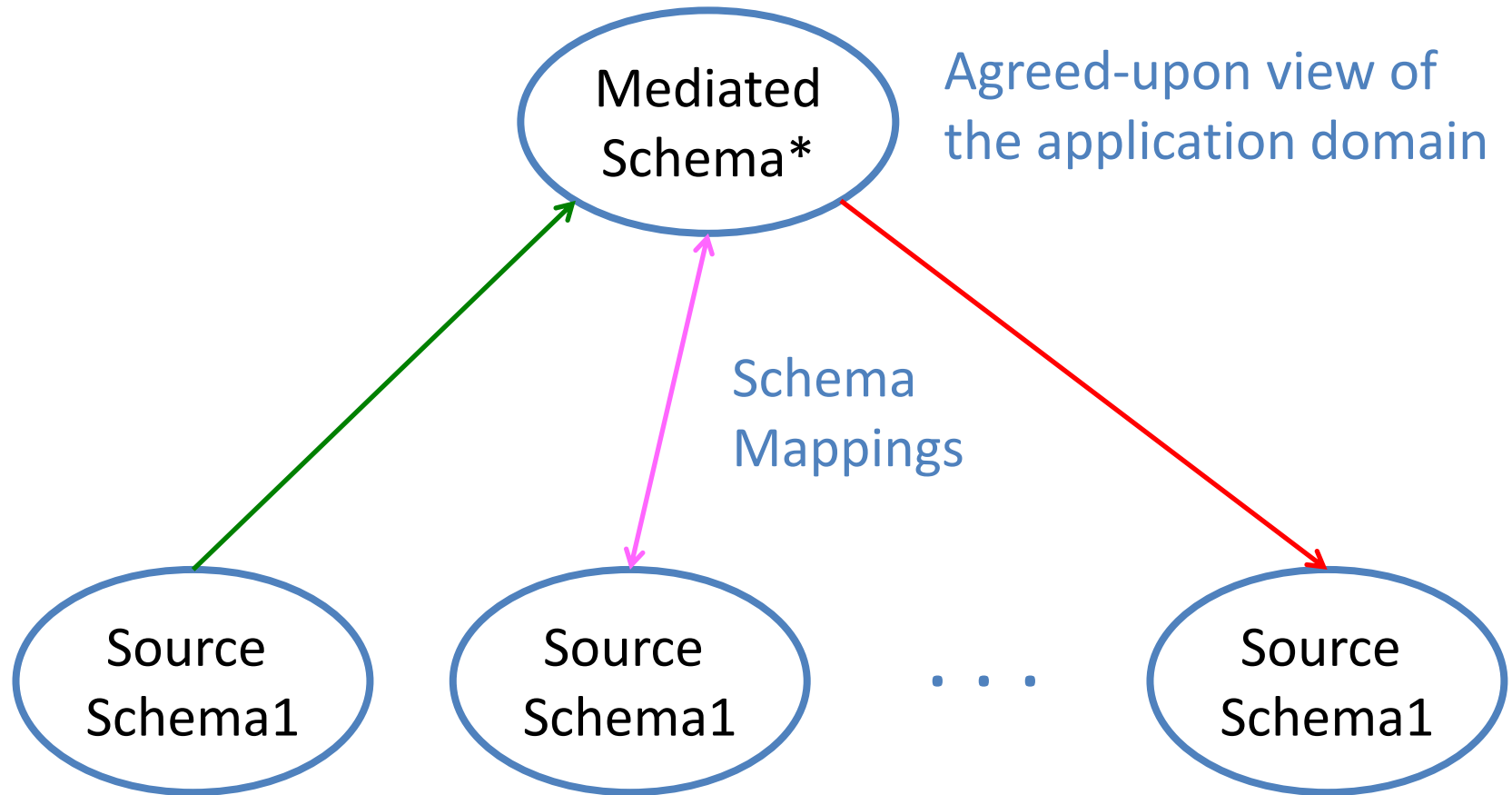
Different coverage and detail



Semantic Heterogeneity Summary

- Differences in:
 - Naming of schema elements
 - Organization of predicates (tables)
 - Coverage and detail of schema
 - Objects/Data values (IBM vs. International Business Machines) ... [Record Linkage]
- Reason: different perspectives
 - Schemas probably designed for different applications/contexts

Integration: Schema Mappings



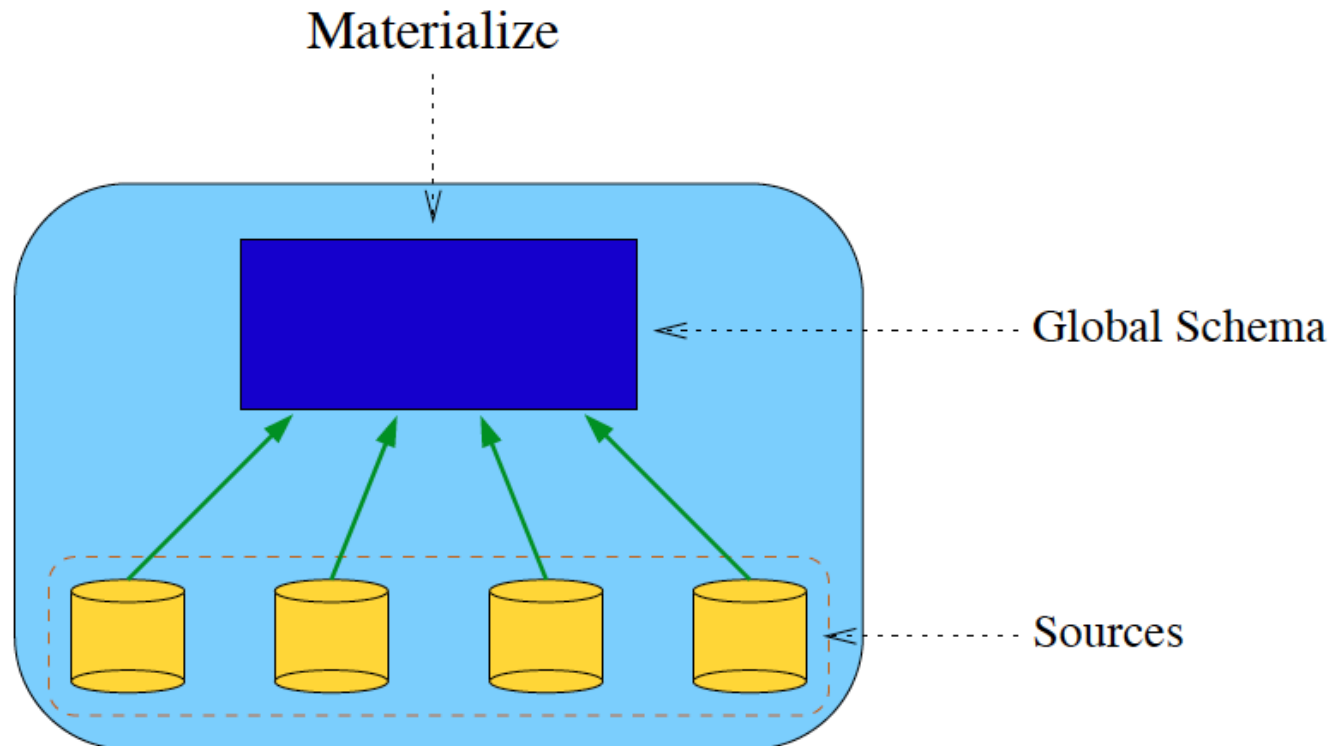
*Mediated/global/domain/target schema/ontology

Dimensions of Data Integration

- Query / Update
- Architecture
 - Warehouse/ETL
 - Virtual Data Integration
- Schema Mappings
 - Global-as-View (GAV)
 - Local-as-View (LAV)
 - GLAV (st-tgds)
- Language for mappings and queries
 - conjunctive queries (CQ)
 - union of CQs (UCQ)
 - first-order logic (\wedge, \vee, \neg)
 - Datalog (recursion)
 - XML, RDF, OWL (description logics) ...
- Source data models
 - Relational, XML, RDF, ...
- Global Schema Constraints
 - Inclusion, functional, ontologies, ...
- Source Capabilities:
 - Input/Output restrictions
- Source type
 - RDBMS, triplestore, XML DB, web services, formatted files, ...

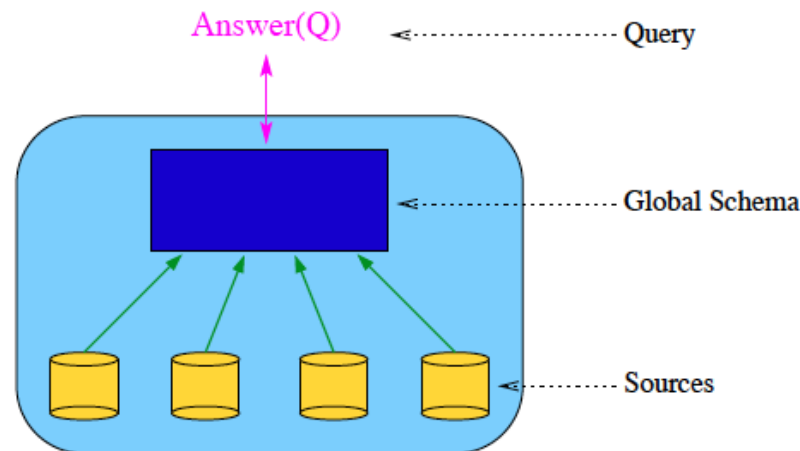
Data Exchange

- materialization of the global view
 - allows for query answering without accessing the sources
- data warehousing



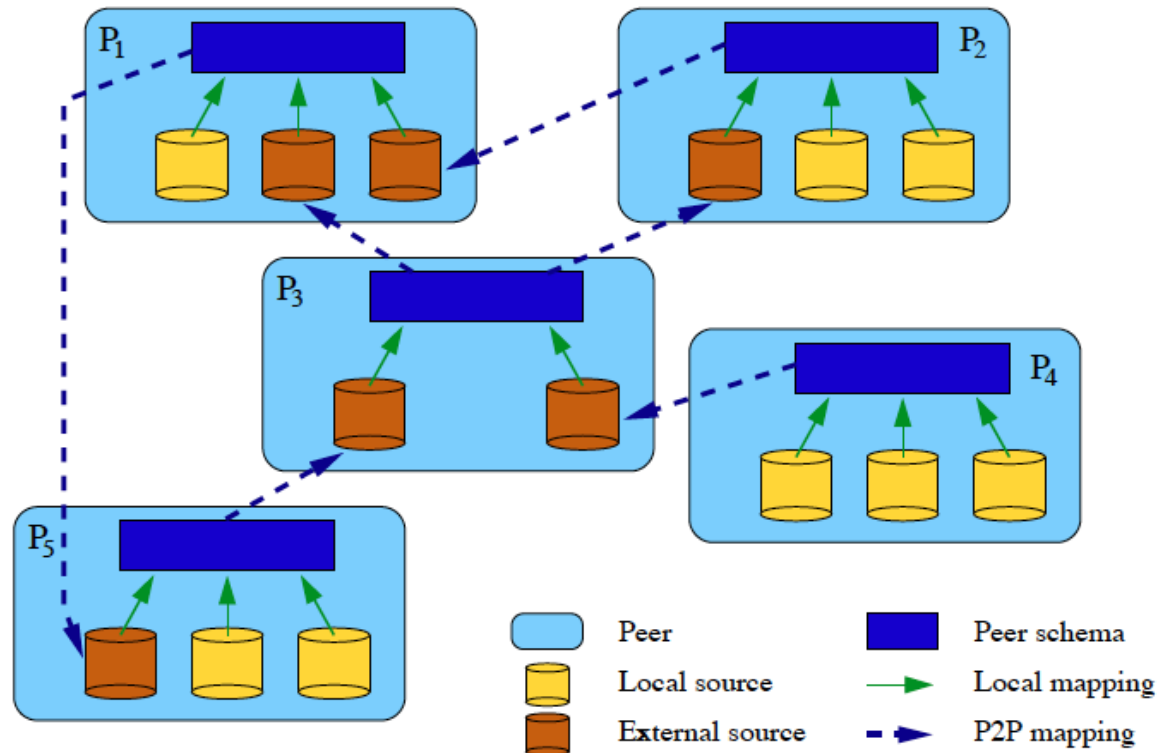
Virtual Data Integration

- Queries are expressed over a **global schema** (a.k.a. mediated schema, enterprise model, domain model, target schema/ontology/model)
- Data are stored in a set of sources
- **Wrappers** access the sources (provide a view in a uniform data model of the data stored in the sources)
- **Mediators** combine answers coming from wrappers and/or other mediators



Peer-to-Peer Data Integration

- several peers
- each peer with local and external sources
- queries over one peer



Virtual Data Integration: BIRN mediator

User quer
find all ma
patients ove
with t1

Results

HID XNAT Integrated Portal - Mozilla Firefox

http://hukadone.ai.edu:8080/birn-uf/

HID XNAT Integrated Portal | BIRN Data Integration Results | CENTRAL | Sessions: a1_MR1

BIRN The Center for Neurological Research

USC INFORMATION SCIENCE INSTITUTE UCI IRVINE

Subjects Experiments Subject Assessments Investigators

Experiment Characteristics

Scan Type: T1, T2, Structural MRI, Functional MRI, Any

Scanner: GE, Philips, Siemens, Any

Time Interval

Subject Characteristics

Race: White, Black, Asian, Native Hawaiian, Native American, Any

Ethnicity: Hispanic, Not Hispanic, Any

Age: TO

Submit Query

000934691111	Handesness_left_writing	0
000934691111	Handesness_right_writing	2
000934691111	Mother's Education	(null)
000934691111	Number of Children	(null)

Human Imaging Database(s)
Oracle DB

Results

Source	Subject ID	Age	Gender	Scan Type
XNAT	OAS1_0266_MR1	51	M	t1
XNAT	OAS1_0034_MR1	51	M	t1
XNAT	OAS1_0207_MR1	51	M	t1
HID	000998262706	51	M	t1
HID	000998262706	52	M	t1
HID	000998041611	53	M	t1
HID	000913186207	53	M	t1
HID	000960528669	53	M	t1
HID	000913186207	54	M	t1
HID	000947193547	55	M	t1
XNAT	OAS1_0389_MR1	55	M	t1

result
(XML)

```
<cell>76</cell>
<cell>active</cell>
</row>
- <row>
  <cell>OAS1_0016_MR1</cell>
  <cell><CENTRAL_OASIS_CS></cell>
  <cell>82</cell>
  <cell>active</cell>
</row>
```

XNAT

EXtensible Neuroimaging Archive Toolkit
Web service API

Query Rewriting (aka query reformulation)

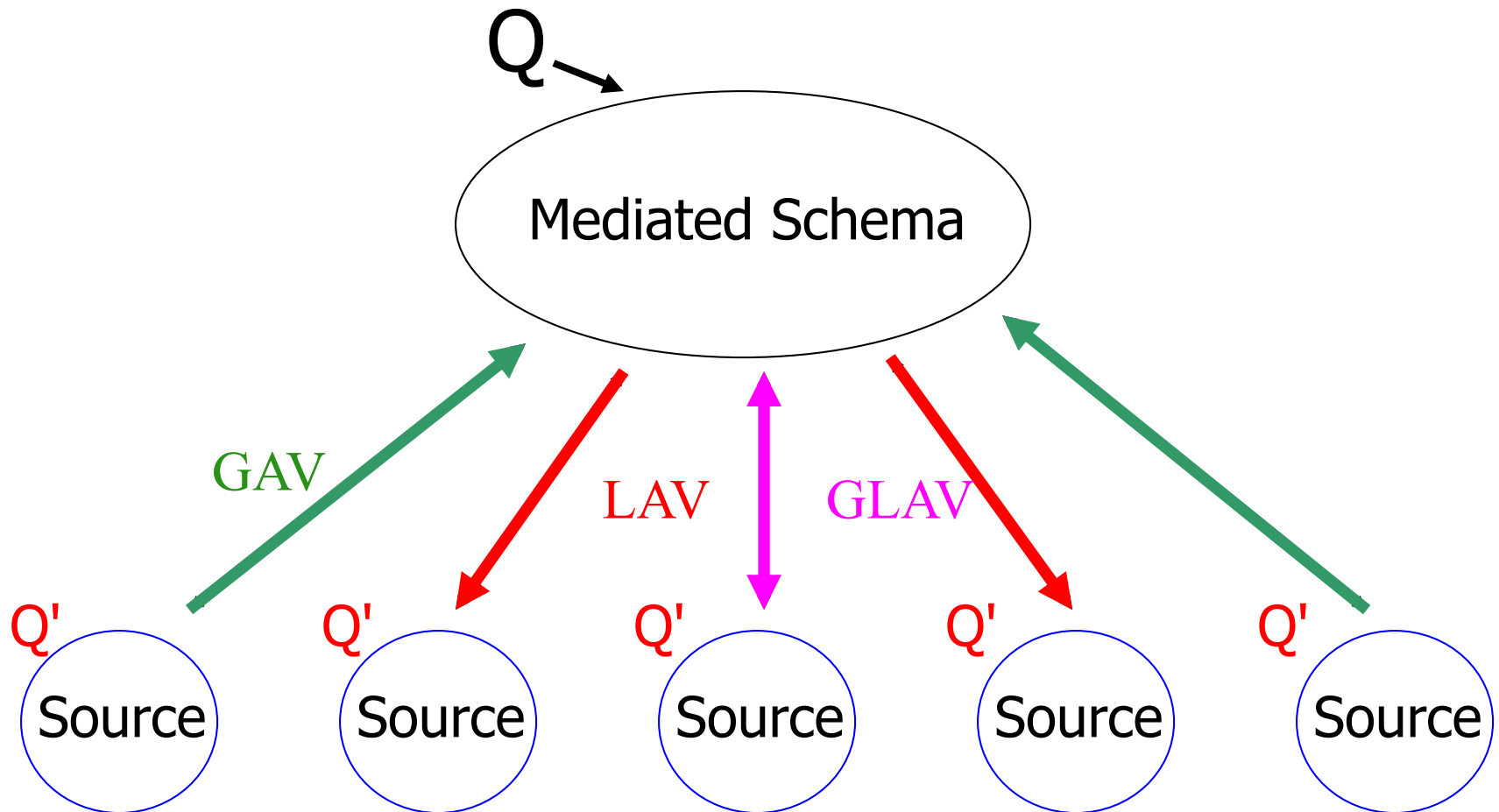
- Problem: rewrite user query expressed in mediated schema into query expressed in source schemas
- Given
 - Query Q in terms of the mediated-schema relations, and
 - Schema mappings (source descriptions)
 - (Mediated schema constraints)
- Find
 - Query Q' that uses only the source relations, such that
 - Q' provides only correct answers to Q , (i.e, $Q' \models Q$, $Q' \subseteq Q$), and
 - Q' provides all possible answers to Q given the sources

Note: I'll use the terms "schema mappings" and "source descriptions" interchangeably

Desiderata for Source Descriptions

- **Expressive power:** distinguish between sources with closely related data. Hence, be able to prune access to irrelevant sources.
- **Easy addition:** easy to add new data sources.
- **Reformulation:** be able to reformulate a user query into a query over the sources efficiently and effectively.
- **Nonlossy:** be able to handle all queries that can be answered by directly accessing the sources

Languages for Schema Mapping



Schema Mappings

- Global-as-View (GAV): $\Phi_s(X,Y) \rightarrow g(X)$
 - Mediator relation defined as a view over source relations
 - Ex: TSIMMIS (Stanford), HERMES (Maryland), SIMS, BIRN (USC)
- Local-as-View (LAV): $s(X) \rightarrow \exists Y \Phi_g(X,Y)$
 - Source relation defined as view over mediator relations
 - Ex: Information Manifold (AT&T), Tukwila (UW), InfoMaster (Stanford), Prometheus, BIRN Mediator/GQR (USC)
- General (GLAV): $\Phi_s(X,Y) \rightarrow \exists Z \Psi_g(X,Z)$
 - BIRN Mediator (USC)

View ~ named query ~ logical formula

Global-as-View (GAV) Schema Mapping

Each mediator relation is defined as a view over source relations: $\Phi_s(X,Y) \rightarrow g(X)$

$\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$

```
CREATE VIEW review AS  
SELECT Title, Year, Review  
FROM s1, s3  
WHERE s1.ID = s3.ID
```

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Domain Query: *Find reviews after 1997*

$q(T,R) \leftarrow \text{review}(T,Y,R), Y \geq 1997$

rewritten to

Source Query:

$q'(T,R) \leftarrow s1(I,T,D,Y), s3(I,R), Y \geq 1997$

using GAV schema mapping:

$$\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$$

Global-as-View (GAV)

Each mediator relation is defined as a view over source relations: $\Phi_s(X,Y) \rightarrow g(X)$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$

$\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$

Assume $s1$ and $s3$ are in the same database $db1$, and $s2$ is in a different source

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow s1(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $s1(I',T,D,Y), s3(I',R)$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow s1(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $s1(I,T,D,Y), s3(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then $I=I'$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$
 $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$
 $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow s1(I,T,\text{'Scott'},Y), Y \geq 1997,$
 ~~$s1(I,T,D,Y)$~~ , $s3(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then I=I' and 2nd s1
is redundant

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow s1(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $s3(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then I=I' and 2nd s1
is redundant

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow \text{s1}(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $\text{s3}(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then I=I' and 2nd s1
is redundant

2. $q'(T,R) \leftarrow \text{s2}(T,D,S,Y), Y \geq 1997,$
 $\text{s1}(I',T,D,Y), \text{s3}(I',R)$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s2}(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge \text{s3}(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow \text{s1}(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $\text{s3}(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then I=I' and 2nd s1
is redundant

2. $q'(T,R) \leftarrow \text{s2}(T,D,S,Y), Y \geq 1997,$
 $\text{s1}(I,T,D,Y), \text{s3}(I,R)$

Assume (T,Y) and (I)
are keys of s1, then I=I'

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s2}(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge \text{s3}(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

1. $q'(T,R) \leftarrow \text{s1}(I,T,\text{'Scott'},Y), Y \geq 1997,$
 $\text{s3}(I,R)$

Assume (T,Y) and (I)
are keys of s1,
then $I=I'$ and 2nd s1
is redundant

~~2. $q'(T,R) \leftarrow \text{s2}(T,D,S,Y), Y \geq 1997,$
 $\text{s1}(I,T,D,Y), \text{s3}(I,R)$~~

Assume (T,Y) and (I)
are keys of s1, then $I=I'$

If s1 and s3 have same movies,
then q2 is redundant wrt q1

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year})$ $\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow \text{s2}(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$ $\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow \text{s1}(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge \text{s3}(\text{ID}, \text{Review})$
--

Query Reformulation in GAV

Query reformulation = rule unfolding + simplification

Query: *Find reviews for movies directed by Scott after 1997*

$q(T,R) \leftarrow \text{movie}(T,Y, \text{'Scott'}), \text{review}(T,Y,R), Y \geq 1997$

$q'(T,R) \leftarrow s1(I,T,\text{'Scott'},Y), s3(I,R), Y \geq 1997,$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year})$

$\text{movie}(\text{Title}, \text{Year}, \text{Director}) \leftarrow s2(\text{Title}, \text{Director}, \text{Studio}, \text{Year})$

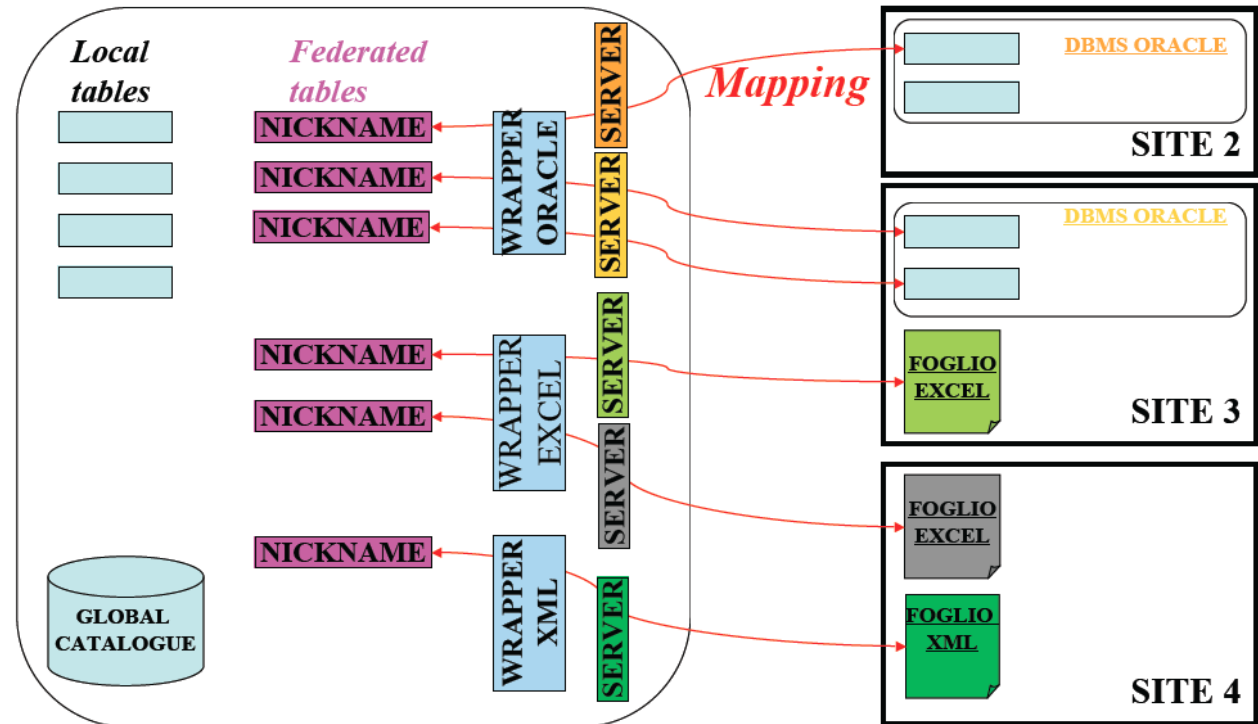
$\text{review}(\text{Title}, \text{Year}, \text{Review}) \leftarrow s1(\text{ID}, \text{Title}, \text{Director}, \text{Year}) \wedge s3(\text{ID}, \text{Review})$

Distributed/Federated Queries vs. Integration

- Some systems allow to query remote databases

- Ex:

IBM
Information
Integrator*



- However, If we just expose the source schemas, the user would quickly get overloaded
- For integration, we need to create global schema

(*) IBM marketing changes the name of the system quite often, today it is called InfoSphere Federation Server:
<http://www-01.ibm.com/software/data/infosphere/federation-server/>

Local-as-View (LAV) Schema Mappings

- Each source relation is defined as a view over mediator relations: $s(X) \rightarrow \exists Y \Phi_g(X, Y)$

$V1(\text{title}, \text{year}, \text{director}) \xrightarrow{\subseteq} \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$V2(\text{title}, \text{review}) \xrightarrow{\subseteq} \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{year} \geq 1990 \wedge \text{hasReview}(\text{title}, \text{review})$

Open world! Tuples in sources guaranteed to satisfy the definitions, but sources don't have all possible tuples

Existentials! Invent values! Skolem functions, aka labelled nulls

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Reformulated query:

$q'(\text{title}, \text{review}) \text{ :-}$

$s2(\text{title}, \text{review})$

Is s2 enough?

Only s2 provides hasReview

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Reformulated query:

$q'(\text{title}, \text{review}) \text{ :- } s1(\text{title}, \text{year}, \text{director}), s2(\text{title}, \text{review})$

Is s2 enough? No!
Need s1 to ensure Genre='Comedy'

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Reformulated query:

Are we done?

What about $\text{year} \geq 1950$ and $\text{Genre} = \text{'Comedy'}$?

$q'(\text{title}, \text{review}) \text{ :- } s1(\text{title}, \text{year}, \text{director}), s2(\text{title}, \text{review})$

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Reformulated query:

$q'(\text{title}, \text{review}) \text{ :- } s1(\text{title}, \text{year}, \text{director}), s2(\text{title}, \text{review})$

Are we done? *Yes!*

What about $\text{year} \geq 1950$ and $\text{Genre} = \text{'Comedy'}$?

Constraints are implied

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Does q' answers what q asks?

Reformulated query:

$q'(\text{title}, \text{review}) \text{ :- } s1(\text{title}, \text{year}, \text{director}), s2(\text{title}, \text{review})$

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}),$
 $\text{year} \geq 1950, \text{hasReview}(\text{title}, \text{review})$

Reformulated query:

Does q' answers what q asks?
Not quite, it gives a subset ...

$q' \subseteq q$

$q'(\text{title}, \text{review}) \text{ :- } s1(\text{title}, \text{year}, \text{director}), s2(\text{title}, \text{review})$

... but it's the best we can do given the available sources

$s1(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$s2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{hasReview}(\text{title}, \text{review})$

Answering queries using views

Given query q and view definitions $V = \{V_1 \dots V_n\}$

- q' is a Maximally-Contained Rewriting of q using V if:
 - q' refers only to views in V , and
 - $q' \subseteq q$, and
 - there is no rewriting q_1 , such that $q' \subseteq q_1 \subseteq q$ and $q_1 \neq q$
- q' is an Equivalent Rewriting of q using V if:
 - q' refers only to views in V , and
 - $q' = q$

GAV

vs.

LAV

$$\Phi s(X,Y) \rightarrow g(X)$$

- Not modular
 - adding new sources may require changes to all integration rules
- Query reformulation is easy
 - reduces to view unfolding
 - (though query simplification not so easy)
- Best when
 - Few, stable, data sources
 - well-known to the mediator (e.g. corporate integration)

$$s(X) \rightarrow \exists Y \Phi g(X,Y)$$

- Modular
 - adding new sources is easier, does not affect previous rules
- Expressive
 - power of entire query language available to describe sources
- Query reformulation is hard
 - answering queries using views
- Best when
 - Many, relatively unknown data sources
 - possibility of addition/deletion of sources

Query Reformulation in LAV: The Bucket Algorithm

- Given: user query q , source descriptions $\{V_i\}$
- Find relevant sources (fill buckets)
 - For each predicate g in query q
 - Find V_j that contains predicate g
 - Check that constraints in V_j are compatible with q
- Combine source relations $\{V_j\}$ from each bucket into a conjunctive query q' and check for containment ($q' \subseteq q$)

The Bucket Algorithm: Example

- Global Schema

Movie(ID, title, year, genre) Director(ID, director)

Actor(ID, actor)

Revenues(ID, Amount)

- LAV Schema Mappings

$V1(I, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A)$

$V3(I, A) \rightarrow \text{Revenues}(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Director}(I, D), I < 3000$

- User Query

$q(ID, Dir) \leftarrow \text{Movie}(ID, Title, Year, Genre), \text{Revenues}(ID, Amount),$
 $\text{Director}(ID, Dir), Amount > \$100M$

1. Filling the Buckets

A goal g_v in a view v is relevant to a goal g_q in a query q if:

1. The goals are on the same predicate
2. The interpreted predicates on variables of g_v , after appropriate variable substitutions, are satisfiable
3. If g_q includes a head variable of q , then the corresponding variable in g_v is also in the head of v

1. Filling the Buckets

$V1(I, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A)$

$V3(I, A) \rightarrow \text{Revenues}(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Director}(I, D), I < 3000$

$q(ID, Dir) \leftarrow$

$\text{Movie}(ID, \text{Title}, \text{Year}, \text{Genre}), \text{Revenues}(ID, \text{Amt}), \text{Director}(ID, Dir), \text{Amt} > \$100M$

1. Filling the Buckets

V1(I, Y) → Movie(I,T,Y,G), Revenues(I,A), I > 5000, A > \$200M

V2(I, A) → Movie(I,T,Y,G), Revenues(I,A)

V3(I, A) → Revenues(I,A), A < \$50M

V4(I, D, Y) → Movie(I,T,Y,G), Director(I,D), I < 3000

q(ID,Dir) ←

Movie(ID,Title,Year,Genre), Revenues(ID, Amt), Director(ID,Dir), Amt > \$100M

V1(ID,Year)

V2(ID,A')

V4(ID,D',Year)

1. Filling the Buckets

V1(I, Y) → Movie(I,T,Y,G), **Revenues**(I,A), I > 5000, A > \$200M

V2(I, A) → Movie(I,T,Y,G), **Revenues**(I,A)

V3(I, A) → **Revenues**(I,A), A < \$50M

V4(I, D, Y) → Movie(I,T,Y,G), Director(I,D), I < 3000

q(ID,Dir) ←

Movie(ID,Title,Year,Genre), **Revenues**(ID, Amt), Director(ID,Dir), Amt > \$100M

V1(ID,Year) **V1**(ID,Y')

V2(ID,A') **V2**(ID,Amt)

V4(ID,D',Year)

1. Filling the Buckets

$V1(I, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A)$

$V3(I, A) \rightarrow \text{Revenues}(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Director}(I, D), I < 3000$

$q(ID, Dir) \leftarrow$

$\text{Movie}(ID, \text{Title}, \text{Year}, \text{Genre}), \text{Revenues}(ID, \text{Amt}), \text{Director}(ID, Dir), \text{Amt} > \$100M$

$V1(ID, \text{Year}) \qquad V1(ID, Y') \qquad V4(ID, Dir, Y'')$

$V2(ID, A') \qquad V2(ID, \text{Amt})$

$V4(ID, D', \text{Year})$

1. Filling the Buckets

$V1(I, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow \text{Movie}(I, T, Y, G), \text{Revenues}(I, A)$

$V3(I, A) \rightarrow \text{Revenues}(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow \text{Movie}(I, T, Y, G), \text{Director}(I, D), I < 3000$

$q(ID, Dir) \leftarrow$

$\text{Movie}(ID, \text{Title}, \text{Year}, \text{Genre}), \text{Revenues}(ID, \text{Amt}), \text{Director}(ID, Dir), \text{Amt} > \$100M$

$V1(ID, \text{Year})$

$V1(ID, Y')$

$V4(ID, Dir, Y'')$

$V2(ID, A')$

$V2(ID, \text{Amt})$

$V4(ID, D', \text{Year})$

2. Checking Containment

- Construct a rewritten query q' by choosing one view from each bucket
- Check for containment: $q' \subseteq q$
- Collect all such rewritten queries $\{q'\}$, then maximally-contained rewriting is the union

2. Checking Containment

$q(ID, Dir) \leftarrow$

Movie(ID, Title, Year, Genre), *Revenues*(ID, Amt), *Director*(ID, Dir), Amt > \$100M

V1(ID, Year)

V1(ID, Y')

V4(ID, Dir, Y'')

V2(ID, A')

V2(ID, Amt)

V4(ID, D', Year)

V1(I, Y) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A), I > 5000, A > \$200M

V2(I, A) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A)

V3(I, A) \rightarrow *Revenues*(I, A), A < \$50M

V4(I, D, Y) \rightarrow *Movie*(I, T, Y, G), *Director*(I, D), I < 3000

2. Checking Containment

$q(ID, Dir) \leftarrow$

Movie(ID, Title, Year, Genre), *Revenues*(ID, Amt), *Director*(ID, Dir), Amt > \$100M

V1(ID, Year)

V1(ID, Y')

V4(ID, Dir, Y'')

V2(ID, A')

V2(ID, Amt)

V4(ID, D', Year)

V1(ID, Year), V1(ID, Y'), V4(ID, Dir, Y'') \rightarrow

Movie(ID, T, Year, G), *Revenues*(ID, A), ID > 5000, A > \$200M,

Movie(ID, T', Y', G'), *Revenues*(ID, A'), ID > 5000, A' > \$200M,

Movie(ID, T'', Y'', G''), *Director*(ID, Dir), ID < 3000,

V1(I, Y) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A), I > 5000, A > \$200M

V2(I, A) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A)

V3(I, A) \rightarrow *Revenues*(I, A), A < \$50M

V4(I, D, Y) \rightarrow *Movie*(I, T, Y, G), *Director*(I, D), I < 3000

2. Checking Containment

$q(ID, Dir) \leftarrow$

Movie(ID, Title, Year, Genre), *Revenues*(ID, Amt), *Director*(ID, Dir), Amt > \$100M

V1(ID, Year)

V1(ID, Y')

V4(ID, Dir, Y'')

V2(ID, A')

V2(ID, Amt)

V4(ID, D', Year)

V1(ID, Year), V1(ID, Y'), V4(ID, Dir, Y'') \rightarrow

Movie(ID, T, Year, G), *Revenues*(ID, A), ID > 5000, A > \$200M,

Movie(ID, T', Y', G'), *Revenues*(ID, A'), ID > 5000, A' > \$200M,

Movie(ID, T'', Y'', G''), *Director*(ID, Dir), ID < 3000,

Inconsistent!

V1(I, Y) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A), I > 5000, A > \$200M

V2(I, A) \rightarrow *Movie*(I, T, Y, G), *Revenues*(I, A)

V3(I, A) \rightarrow *Revenues*(I, A), A < \$50M

V4(I, D, Y) \rightarrow *Movie*(I, T, Y, G), *Director*(I, D), I < 3000

2. Checking Containment

$q(ID, Dir) \leftarrow$

$Movie(ID, Title, Year, Genre), Revenues(ID, Amt), Director(ID, Dir), Amt > \$100M$

$V1(ID, Year)$

$V1(ID, Y')$

$V4(ID, Dir, Y'')$

$V2(ID, A')$

$V2(ID, Amt)$

$V4(ID, D', Year)$

$V2(ID, A'), V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt),$

$Movie(ID, T'', Y'', G''), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$q(ID, Dir)$

$V1(I, Y) \rightarrow Movie(I, T, Y, G), Revenues(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow Movie(I, T, Y, G), Revenues(I, A)$

$V3(I, A) \rightarrow Revenues(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow Movie(I, T, Y, G), Director(I, D), I < 3000$

2. Checking Containment

$q(ID, Dir) \leftarrow$

$Movie(ID, Title, Year, Genre), Revenues(ID, Amt), Director(ID, Dir), Amt > \$100M$

$V1(ID, Year)$

$V1(ID, Y')$

$V4(ID, Dir, Y'')$

$V2(ID, A')$

$V2(ID, Amt)$

$V4(ID, D', Year)$

$V2(ID, A'), V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt),$

$Movie(ID, T'', Y'', G''), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$q(ID, Dir)$

Only movies
with $ID < 3000$

$V1(I, Y) \rightarrow Movie(I, T, Y, G), Revenues(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow Movie(I, T, Y, G), Revenues(I, A)$

$V3(I, A) \rightarrow Revenues(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow Movie(I, T, Y, G), Director(I, D), I < 3000$

2. Checking Containment

$q(ID, Dir) \leftarrow$

$Movie(ID, Title, Year, Genre), Revenues(ID, Amt), Director(ID, Dir), Amt > \$100M$

$V1(ID, Year)$

$V1(ID, Y')$

$V4(ID, Dir, Y'')$

$V2(ID, A')$

$V2(ID, Amt)$

$V4(ID, D', Year)$

$V2(ID, A'), V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$q'(ID, Amt) \leftarrow V2(ID, Amt), V4(ID, Dir, Y''), Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt),$

$Movie(ID, T', Y'', G''), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$Movie(ID, T', Y', G'), Revenues(ID, Amt), Director(ID, Dir), ID < 3000, Amt > \$100M \rightarrow$

$q(ID, Dir)$

$V1(I, Y) \rightarrow Movie(I, T, Y, G), Revenues(I, A), I > 5000, A > \$200M$

$V2(I, A) \rightarrow Movie(I, T, Y, G), Revenues(I, A)$

$V3(I, A) \rightarrow Revenues(I, A), A < \$50M$

$V4(I, D, Y) \rightarrow Movie(I, T, Y, G), Director(I, D), I < 3000$

Query Reformulation in LAV: Inverse-Rules Algorithm

- Idea: Construct an equivalent logic program which evaluation yields the answer to user query
- “Invert the rules”: simply use standard logical manipulations (conversion to clausal form)

The Inverse-Rules Algorithm

Conversion to Clausal Form

$V1(\text{dept}, \text{course}) \rightarrow \text{Enrolled}(\text{student}, \text{dept}) \wedge \text{Registered}(\text{student}, \text{course})$

$$a \rightarrow b \equiv \neg a \vee b$$

$$\forall D, C [v1(D, C) \rightarrow \exists S [e(S, D) \wedge r(S, C)]]$$

$$\equiv \neg v1(D, C) \vee [e(f(D, C), D) \wedge r(f(D, C), C)]$$

$$\equiv [\neg v1(D, C) \vee e(f(D, C), D)] \wedge [\neg v1(D, C) \vee r(f(D, C), C)]$$

$$\equiv [v1(D, C) \rightarrow e(f(D, C), D)] \wedge [v1(D, C) \rightarrow r(f(D, C), C)]$$

\equiv

$$e(f(D, C), D) \leftarrow v1(D, C)$$

$$r(f(D, C), C) \leftarrow v1(D, C)$$

This is essentially converting the rules to clausal form. See for example:
<http://logic.stanford.edu/classes/cs157/2010/lectures/lecture09.pdf>

The Inverse-Rules Algorithm: Example

$q(D) \leftarrow \text{Enrolled}(S,D) \wedge \text{Registered}(S,'DB')$
 $v1(D,C) \rightarrow \text{Enrolled}(S,D) \wedge \text{Registered}(S,C)$

$q(D) \leftarrow \text{Enrolled}(S,D) \wedge \text{Registered}(S,'DB')$
 $\text{Enrolled}(f(D,C),D) \leftarrow v1(D,C)$
 $\text{Registered}(f(D,C),C) \leftarrow v1(D,C)$

$q(D) \leftarrow v1(D,'DB')$

$\text{Ext}(v1) = \{('CS', 'DB'), ('EE', 'DB'), ('CS', 'AI')\}$
 $\text{Ext}(q) = \{('CS'), ('EE')\}$

GLAV mappings

$$\exists s(X,Y) \rightarrow \exists Z \exists g(X,Z)$$

Example:

$s1(id, title, year, director) \wedge s2(id, review) \rightarrow$
 $Movie(title, year, director, genre) \wedge American(director) \wedge$
 $year \geq 1960 \wedge genre = 'Comedy' \wedge year \geq 1990 \wedge$
 $hasReview(title, review)$

Note we did not model the `id` attribute (i.e., it's not in the consequent of the rule), perhaps it's an internal id to the sources that we don't care about

GLAV = LAV + GAV

A GLAV mapping $\Phi_s(X,Y) \rightarrow \exists Z \Psi_g(X,Z)$
can be converted to a LAV and a GAV mapping

- Just invent intermediate predicate v :
 - $\Phi_s(X,Y) \rightarrow v(X)$ [GAV]
 - $v(X) \rightarrow \exists Z \Psi_g(X,Z)$ [LAV]

So we can process GLAV rules by a combination
of LAV and GAV query rewriting algorithms

Modeling Source Capabilities

Negative capabilities:

- A web-site (HTML form) or a web service may require certain inputs in order to provide outputs.
 - Ex: to output current temperature, Yahoo weather needs ZIP as input
- Need to consider only valid query execution plans
 - Consistent orderings of source calls

Positive capabilities:

- A source may be database (understands SQL)
- Need to decide the placement of operations according to capabilities

Problem: how to describe and exploit source capabilities

Negative Capabilities: Binding Patterns

Sources:

$\text{AAAIdb}^{\mathbf{f}}(X) \rightarrow \text{AAAI Papers}(X)$

$\text{CitationDB}^{\mathbf{bf}}(X, Y) \rightarrow \text{Cites}(X, Y)$

$\text{AwardDB}^{\mathbf{b}}(X) \rightarrow \text{AwardPaper}(X)$

Query: find all the award winning papers:

$q(X) \leftarrow \text{AwardPaper}(X)$

Recursive Rewritings

$q(X) \leftarrow \text{AwardPaper}(X)$

- Problem: *Unbounded* union of conjunctive queries

$q_1(X) \leftarrow \text{AAAIdb}(X), \text{AwardDB}(X)$

$q_1(X) \leftarrow \text{AAAIdb}(X_1), \text{CitationDB}(X_1, X), \text{AwardDB}(X)$

...

$q_1(X) \leftarrow \text{AAAIdb}(X_1), \text{CitationDB}(X_1, X_2), \dots, \text{CitationDB}(X_n, X),$
 $\text{AwardDB}(X)$

- Solution: Recursive Rewriting

$\text{papers}(X) \leftarrow \text{AAAIdb}(X)$

$\text{papers}(X) \leftarrow \text{papers}(Y), \text{CitationDB}(Y, X)$

$q'(X) \leftarrow \text{papers}(X), \text{AwardDB}(X)$

$\text{AAAIdb}^{\mathbf{f}}(X) \rightarrow \text{AAAI Papers}(X)$
$\text{CitationDB}^{\mathbf{bf}}(X, Y) \rightarrow \text{Cites}(X, Y)$
$\text{AwardDB}^{\mathbf{b}}(X) \rightarrow \text{AwardPaper}(X)$

LAV query rewriting with sources with binding restrictions: Inverse-Rules Algorithm

Sources:

$\text{AAAIdb}^f(X) \rightarrow \text{AAAI Papers}(X)$

$\text{CitationDB}^{bf}(X,Y) \rightarrow \text{Cites}(X,Y)$

$\text{AwardDB}^b(X) \rightarrow \text{AwardPaper}(X)$

Query: find all the award winning papers:

$q(X) \leftarrow \text{AwardPaper}(X)$

Inverse-Rules Algorithm under binding restrictions: Domain predicate and rules

Invent a new predicate, $dom(X)$, that represents all objects in our domain
Use dom predicate to satisfy binding patterns

Given source descriptions (SD):

$AAAIdb^f(X) \rightarrow AAAIPapers(X)$

$CitationDB^{bf}(X,Y) \rightarrow Cites(X,Y)$

$AwardDB^b(X) \rightarrow AwardPaper(X)$

... and add *dom rules*, which
define the behavior of dom

$AAAIdb(X) \rightarrow dom(X)$

[we can get objects from AAAIdb, since it
has no restrictions]

Modify SD by introducing dom

$AAAIdb(X) \rightarrow AAAIPapers(X)$

[no restrictions on AAAIdb, just remove binding annotation]

$dom(X) \wedge CitationDB(X,Y) \rightarrow Cites(X,Y)$

[we need to input objects as 1st argument of CitationDB]

$dom(X) \wedge AwardDB(X) \rightarrow AwardPaper(X)$

[we need to input objects to AwardDB to test them]

$dom(X) \wedge CitationDB(X,Y) \rightarrow$
 $dom(Y)$

[if we input an object to CitationDB, we get
another object out]

Inverse-Rules Algorithm under binding restrictions

Evaluating this datalog program generates the answers to query q

$q(X) \leftarrow \text{AwardPaper}(X)$

} query

$\text{AAAI Papers}(X) \leftarrow \text{AAAI db}(X)$

$\text{Cites}(X, Y) \leftarrow \text{dom}(X) \wedge \text{CitationDB}(X, Y)$

$\text{AwardPaper}(X) \leftarrow \text{dom}(X) \wedge \text{AwardDB}(X)$

} inverse
rules
with *dom*

$\text{dom}(Y) \leftarrow \text{dom}(X) \wedge \text{CitationDB}(X, Y)$

$\text{dom}(X) \leftarrow \text{AAAI db}(X)$

} *dom* rules

Inverse-Rules Algorithm under binding restrictions

Simplify the program:

$$q(X) \leftarrow \text{dom}(X) \wedge \text{AwardDB}(X)$$
$$\text{dom}(Y) \leftarrow \text{dom}(X) \wedge \text{CitationDB}(X, Y)$$
$$\text{dom}(X) \leftarrow \text{AAAIdb}(X)$$