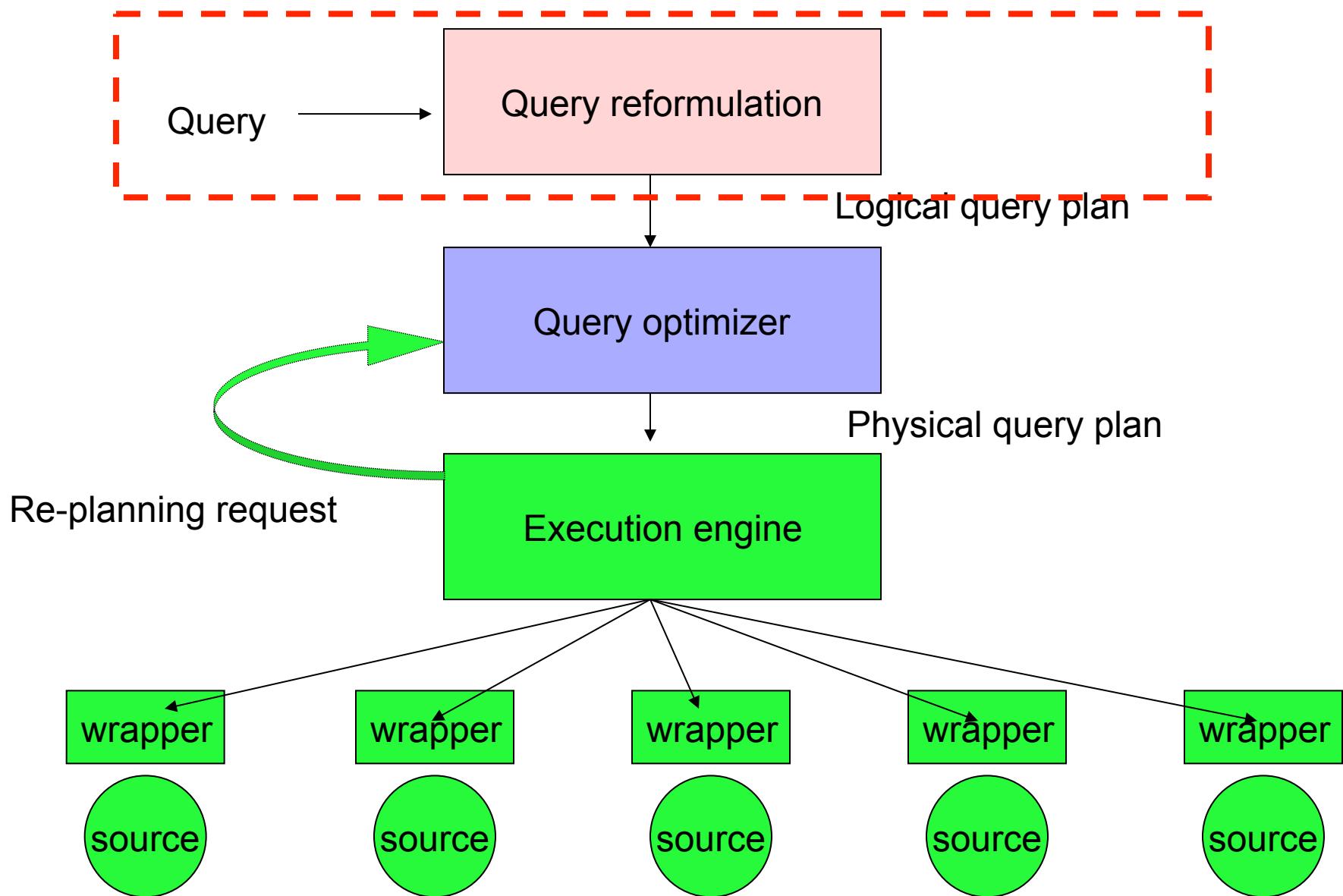


# Querying Heterogeneous Data Sources

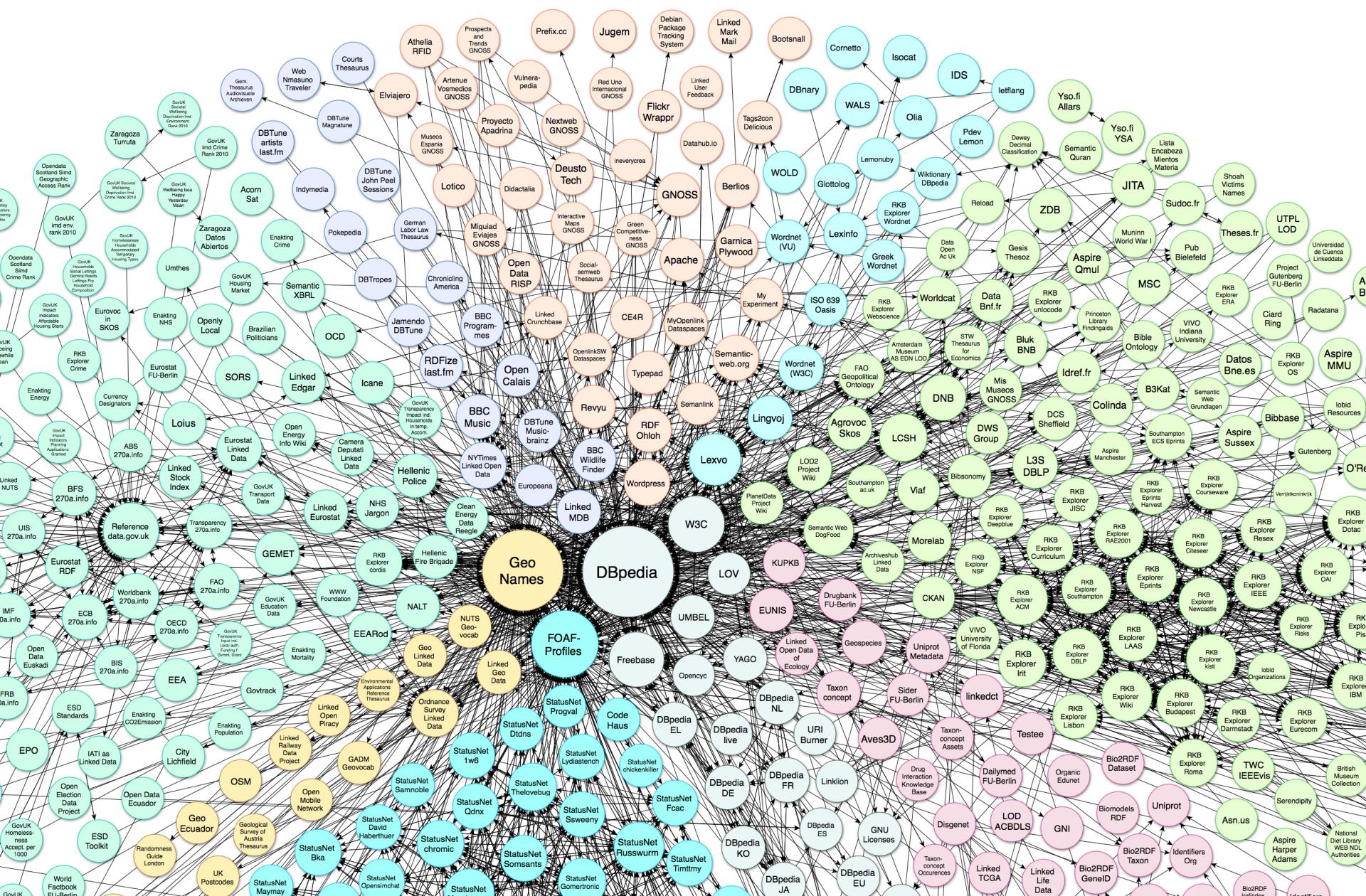
# Outline

- Federated Databases
- Source Descriptions
- Answering Queries Using Views

# Federated Databases



# A Possible Application: Linked Data



# Issues: Semantic Heterogeneity

- Differences in:
  - Naming of schema elements
  - Organization of tables
  - Coverage and detail of schema
  - Languages (XML vs. Relational)
  - Data-level representation (sculpture description in YAGO vs. Musée Rodin DB)

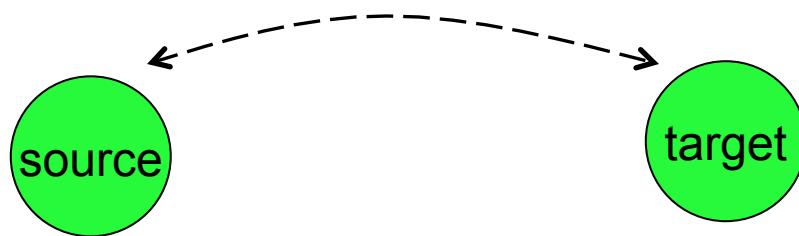
approach



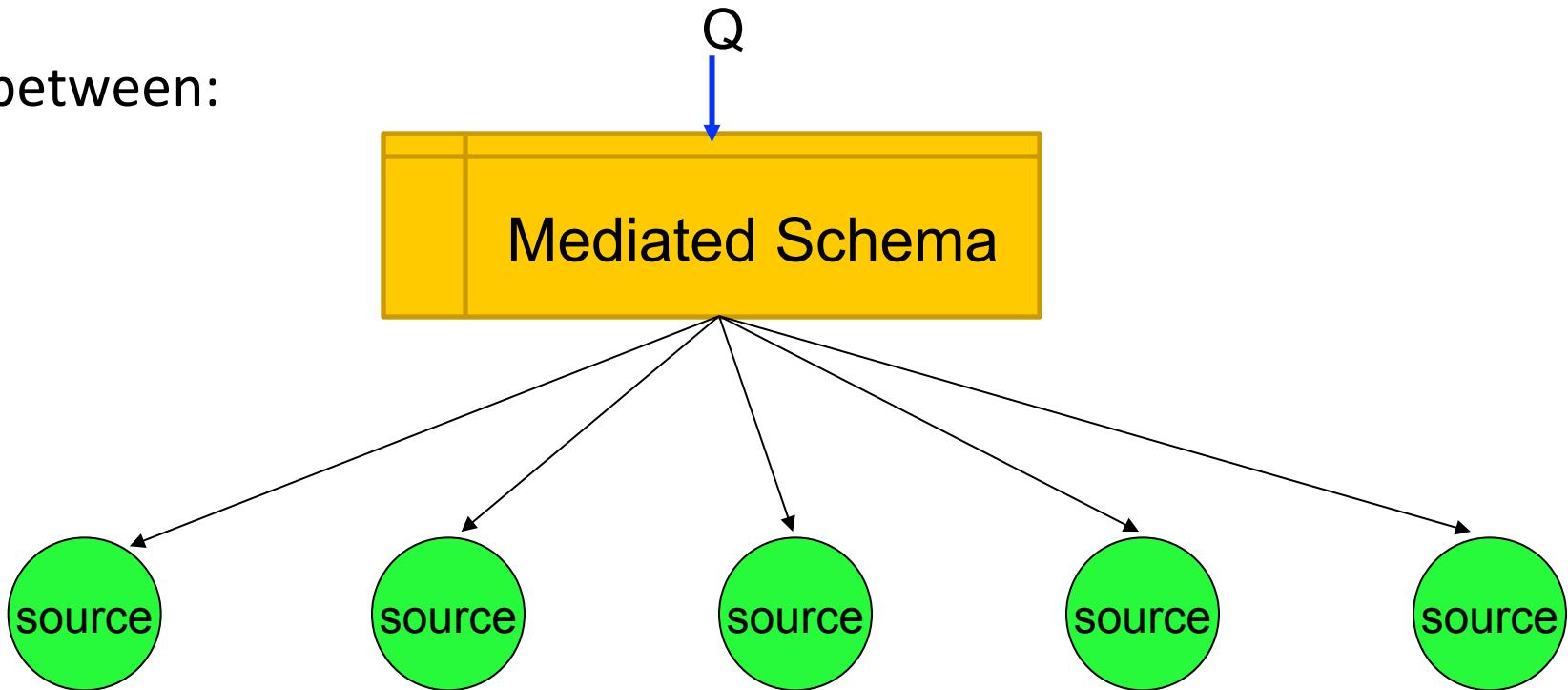
Schema Mappings

# Principles of Schema Mappings

Schema mappings describe the relationship between:



Or between:



# Question

How does the system know:

- Which source is relevant to the query?
- How to access the source appropriately?

# Source Description

Descriptions of data sources should enable the system to:

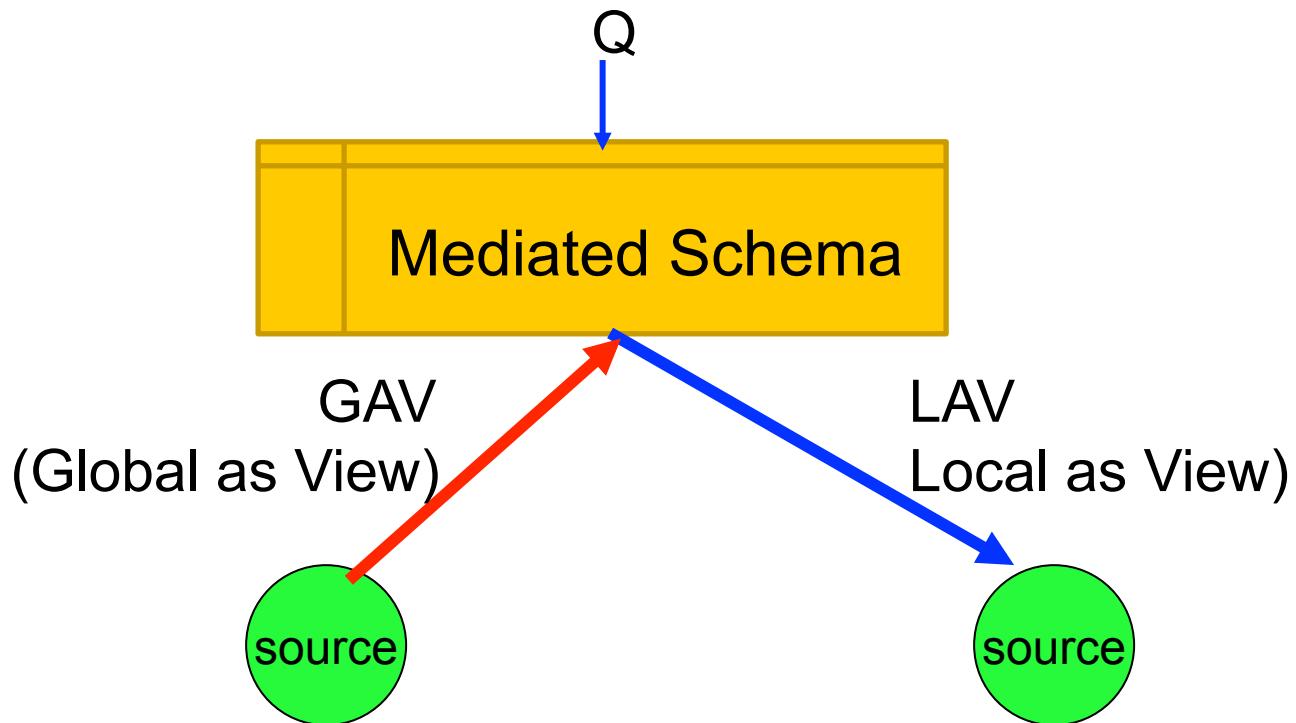
- Determine which sources are relevant to a query
- Access the sources appropriately
- Combine data from multiple sources
- Overcome limitations that specific sources may have
- Identify opportunities for more efficient query evaluation
- Possible access patterns to the data source:
  - Are there required inputs? (e.g., web forms, web services)
  - Can the source process database queries?
- Source completeness:
  - Is the source complete, or partially complete?
- Load restrictions, Quality Of the Service, Quality Of the Data ..

# Source Description Languages

Desiderata:

- **Effective:** it should be able to express relationships between real schemata
- **Enable efficient reformulation:** computational complexity of reformulation and finding answers
- **Easy update:** it should be easy to add and delete sources

# Languages for Schema Mapping



# Global-as-View (GAV)

Mediated schema defined as a set of views over the data sources

Film(title, director, year, genre)

$Film(title, director, year, genre) \supseteq$   
 $S1.Movie(MID, title),$   
 $S1.MovieDetail(MID, director, genre, year)$

S1

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

# Local-as-View (LAV)

Data sources defined as views over mediated schema!

Movie(title, director, year, genre)

Actors(title, name)

Plays(movie, location, startTime)

Reviews(title, rating, description)

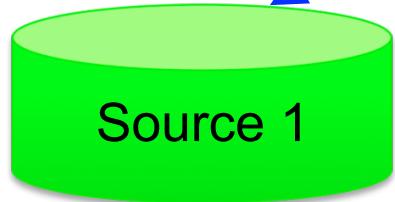
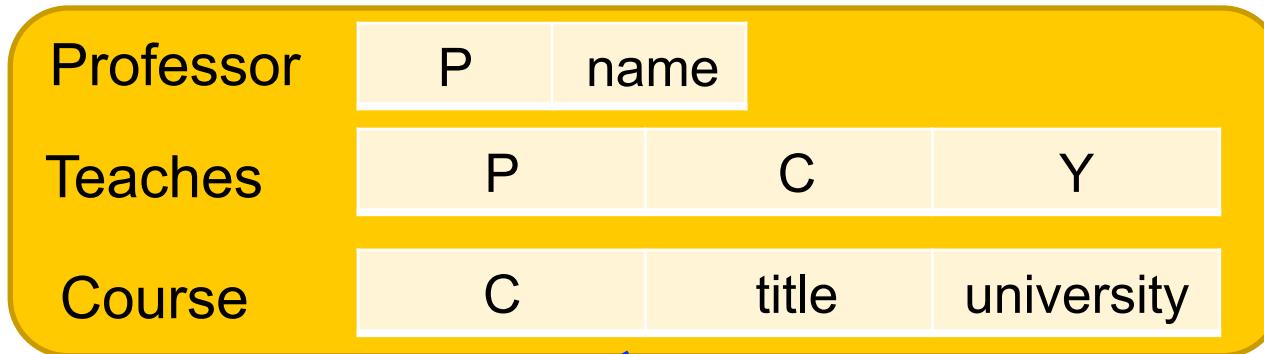
S2:

ActorDirectors(actor, dir)

$S2.\text{ActorDirectors}(\text{actor}, \text{dir}) \subseteq$   
 $\text{Movie}(\text{title}, \text{dir}, \text{year}, \text{genre}), \text{Actor}(\text{title}, \text{actor})$   
 $\text{year} \geq 1980$

# Local-As-View (LAV)

Data sources defined as views over mediated schema!



UVSQProfs

Create View  
Select  
From  
Where

UVSQProfs As  
Course.title, Professor.name  
Course, Teaches, Professor  
Course.C = Teaches.C and  
Teaches.P = Professor.P and  
Course.university = "UVSQ"

# Specification of Source Description

## GAV: Global as view

- Each mediator relation defined as a view over the sources
- Explains how to compute the extensions of mediator relations using the data of sources

## LAV: Local as view

- Each source defined as a view over the mediated schema
- Answer a query  $Q$  over global DB is answer  $Q$  using the views

# Answering Queries

Answering Queries over Possible Instances of Mediated Schema:

- **Certain Answer:** An answer is certain if it is true in every instance of the mediated schema that is consistent with: (1) the instances of the sources and (2) the mapping  $M$ .

# Examples: Certain Answers (1)

$S.\text{ActorDirectors}(actor, dir) \subseteq$

$\text{Movie}(title, dir, year, genre), \text{Actor}(title, actor), year \geq 1980$

S: Keaton Allen

$Q(actor, dir) :-$

$\text{Movie}(title, dir, year, genre), \text{Actor}(title, actor)$

Only one certain answer: (Keaton, Allen)

# Examples: Certain Answers (2)

$V_1(dir) = DirectorActor(ID, dir, actor)$  Keaton

$V_2(actor) = DirectorActor(ID, dir, actor)$  Allen

$Q(dir, actor) :- Director(ID, dir), Actor(ID, actor)$

Under closed-world assumption:  
single DB possible  $\Rightarrow$  (Allen,Keaton)

Under open-world assumption:  
no certain answers.

# Comparison

- **Query Reformulation:**
  - GAV: Query/View Unfolding (easy)
  - LAV: the problem of answering queries using views
- **Incompleteness:**
  - GAV: does not express incomplete information. Only a single instance of the mediated schema is consistent with sources.
  - LAV: expresses incomplete information

# Answering Queries using Views

after the survey Alon Y. Levy. [\*Answering Queries Using Views: A Survey.\*](#)

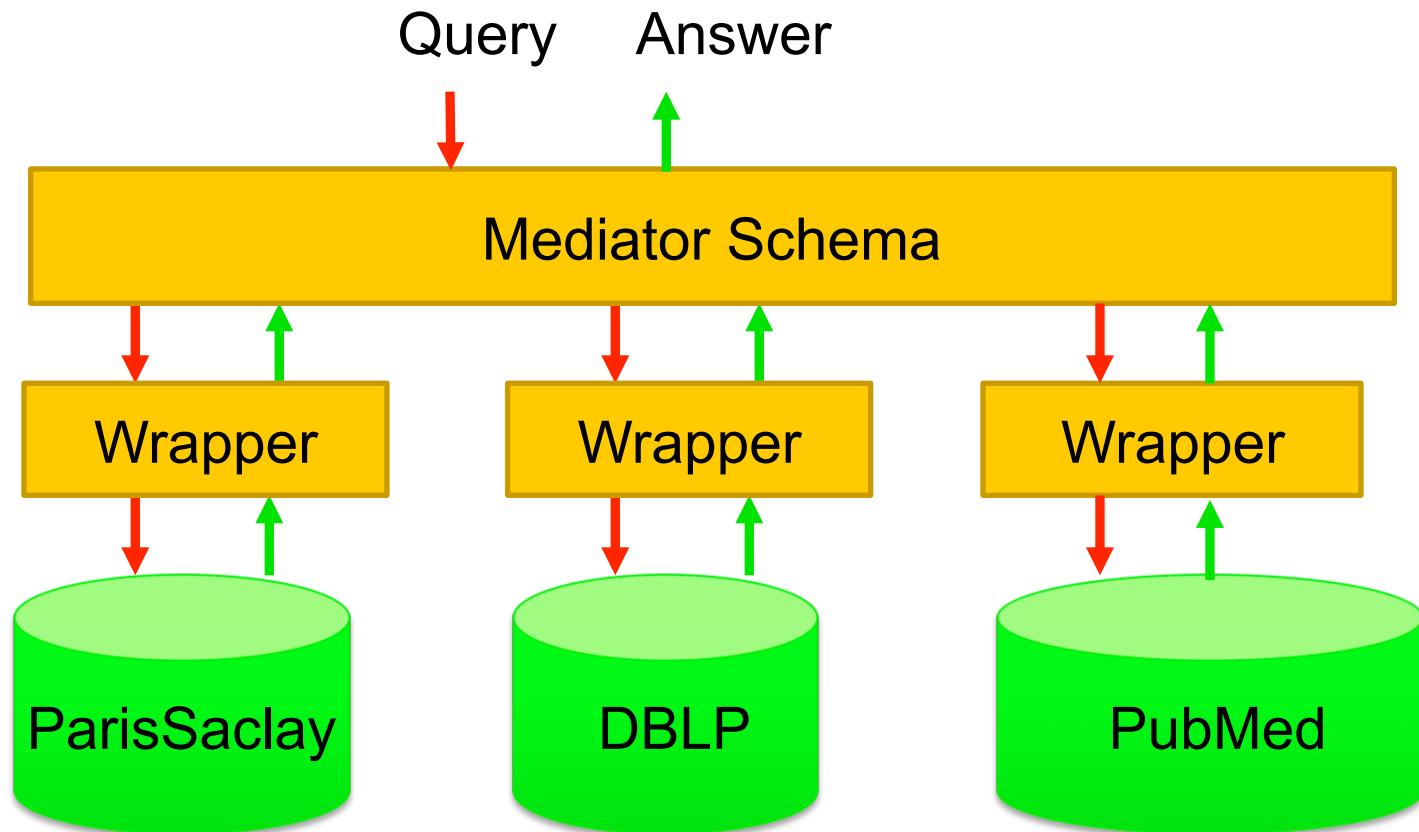
# Answering Queries Using Views



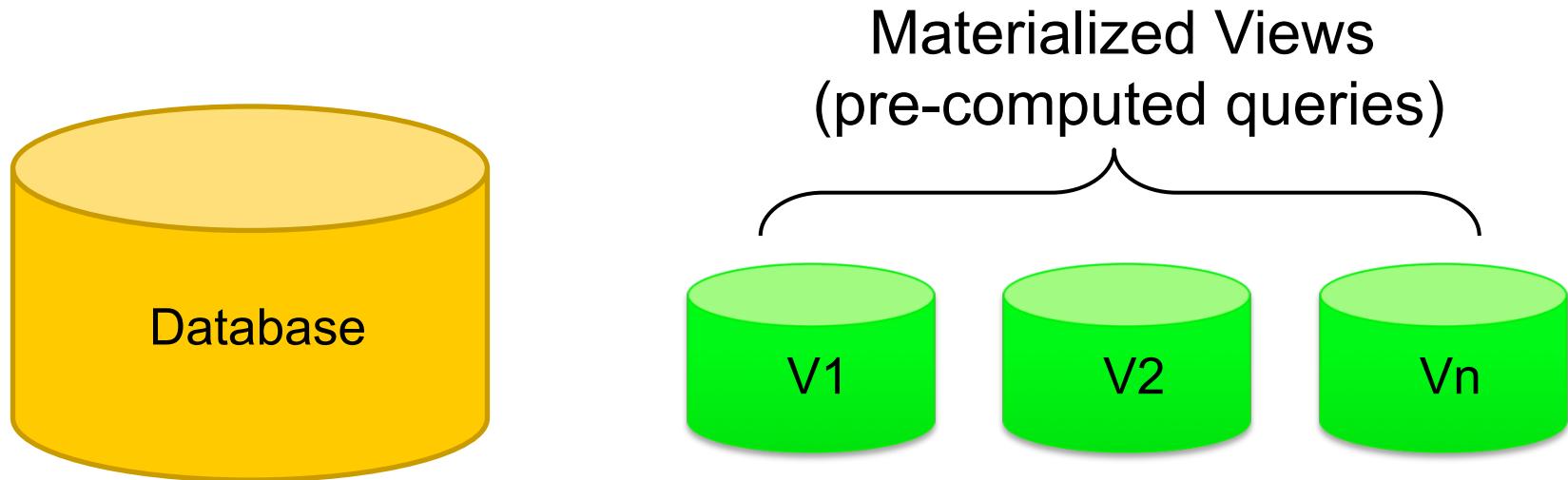
- Motivating Examples
- Problem Definition
- Two different settings:
  - Data Integration
  - Query Optimization & Physical Data Independence

# Scenario 1: Data Integration

Goal: Uniform access to data sources e.g. Amazon, IMDB



# Scenario 2: Query Optimization



Given a query  $Q$  against the Database use the pre-computed results of  $V_1, V_2 \dots V_n$  to reduce the evaluation time!

# Answering Queries Using Views

- 
- Motivating Examples
  - Problem Definition
  - Two different settings:
    - Data Integration
    - Query Optimization & Physical Data Independence

# Query Equivalence and Containment

- **Query Containment:**

$Q_1$  is contained in  $Q_2$  ( $Q_1 \subseteq Q_2$ ), if for all database instances  $D$ ,  
 $Q_1(D) \subseteq Q_2(D)$

- **Query Equivalence:**

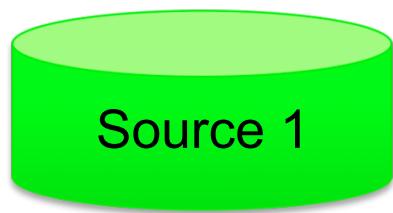
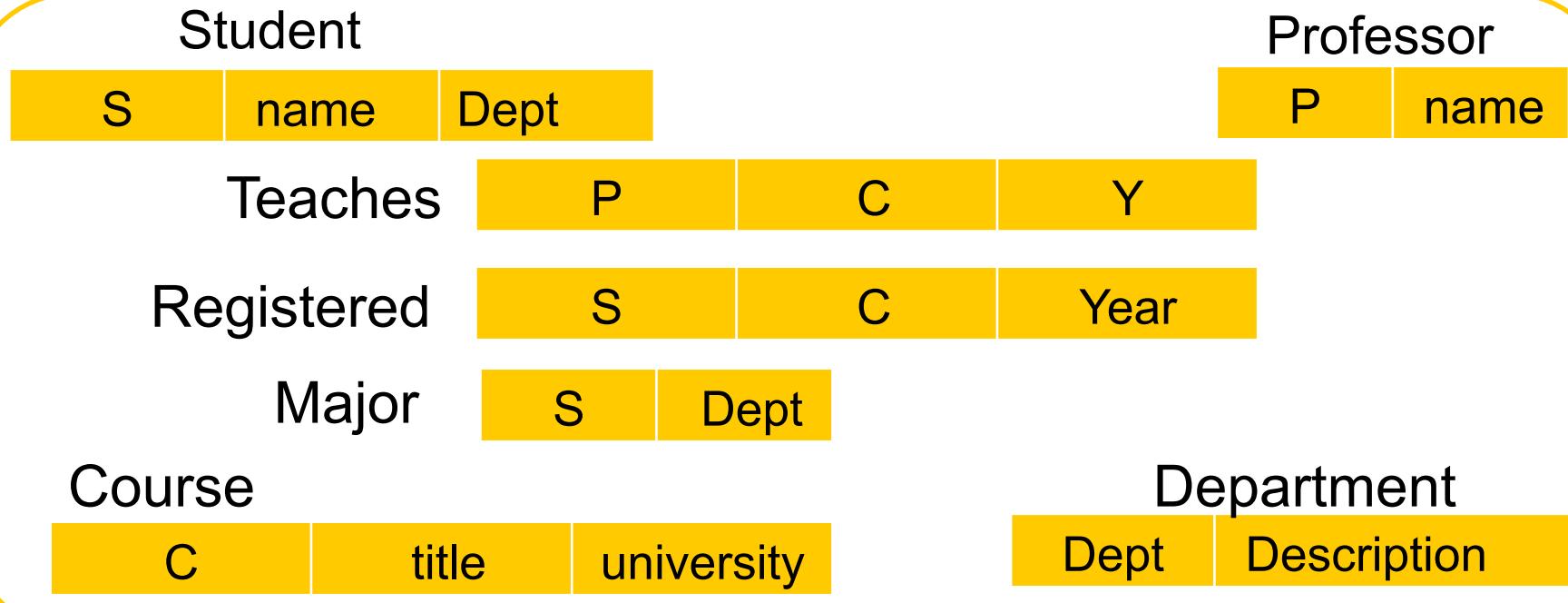
$Q_1$  is equivalent to  $Q_2$  if  $Q_1 \subseteq Q_2$  and  $Q_2 \subseteq Q_1$

# Equivalent and Containment Rewritings

Let  $A$  be a schema and  $V=\{V_1, V_2, \dots, V_m\}$  a set of views over  $A$  and  $Q$  be a query over  $A$

- $Q'$  is a **maximally contained rewriting** of  $Q$ , if
  - $Q'$  is a query expression **in a given language  $L$**
  - $Q'$  refers only to views in  $V$
  - $Q' \subseteq Q$
  - There is no other rewriting  $Q_1$ , such as  $Q' \subseteq Q_1 \subseteq Q$
- $Q'$  is a **equivalent rewriting** of  $Q$ , if
  - $Q'$  refers only to views in  $V$
  - $Q'$  is equivalent to  $Q$

# Example



DBProfs



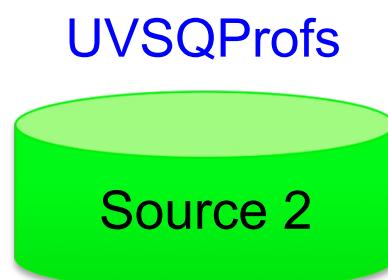
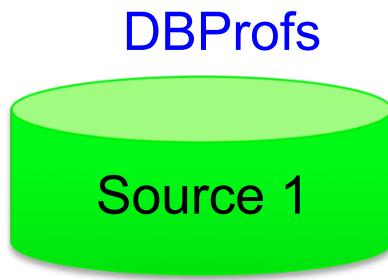
UVSQProfs

Create View  
Select  
From  
Where

UVSQProfs As  
Course.title, Professor.name  
Course, Teaches, Professor  
Course.C = Teaches.C and  
Teaches.P = Professor.P and  
Course.university = "UVSQ"

# Example: Query Containment

- Consider the query Q:  
“Select name from Professor”
- Let Q' be the query:  
“Select name from DBProfs union UVSQProfs”
- Q' is **maximally-contained** in Q
  - $Q' \subseteq Q$  but **Not**( $Q \subseteq Q'$ )
  - Let D be a instance of the global schema consisting (also) of “Professors of Algorithms from Paris XI”...

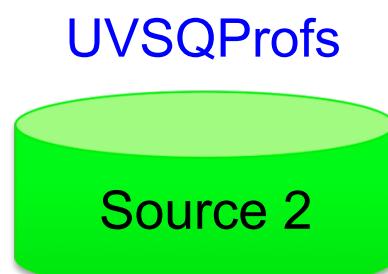
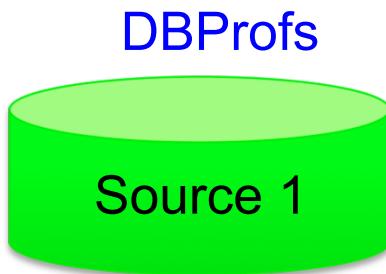


# Example: Query Containment (2)

- Consider the query Q:

```
Select Profesor.name  
From Course, Teaches, Professor  
Where Course.C = Teaches.C and  
      Teaches.P = Professor.P and  
      Course.university = "UVSQ" and  
      Course.title = "Databases"
```
- Let Q' be the query:

```
Select name  
from DBProfs, UVSQProfs  
where DBProfs.name= UVSQProfs.name
```
- Q' is **a equivalent** to Q (...if the sources are indeed complete)



# The variations of the problem

	Data integration	Query optimization	Database design
Equivalent rewriting	Only if the sources are complete	✓	✓
Maximally-contained rewriting	✓		
Views are materialized	✓	✓	✓
The db schema is virtual	✓		✓

# Usual Settings & Problems

## Data Integration:

Problem: Compute the **maximally contained rewriting**

Approach: Logical Rewriting (compute a query expression)

- Rewriting algorithms
- Query answering algorithms

## Query Optimization & Physical data independence:

Problem: Compute the **equivalent rewriting**

Approach: Cost Based Rewriting (compute an evaluation plan)

- System-R Style
- Transformational approaches

# Answering Queries Using Views

- Motivating Examples
- Problem Definition
- Two different settings:
  - ➔ Data Integration
  - Query Optimization & Physical Data Independence

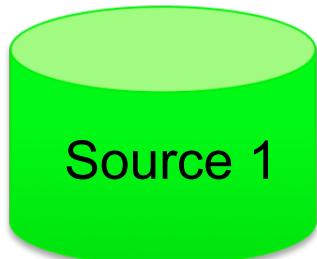
# Answering Queries Using Views In **Data Integration**

# Example: A Mediated Schema

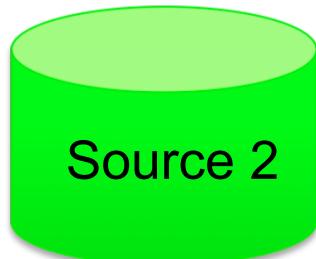
Student	P	C	Y
S	name	Dept	
Teaches	P	C	Y
Registered	S	C	Y
Major	S	Dept	

Course	Dept	Description
C	title	

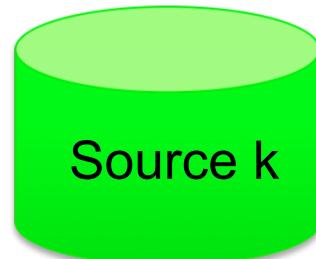
$V_1, V_2, \dots, V_i;$



$V_{i+1}, \dots, V_j$



$V_p, \dots, V_n$



# Conjunctive Queries

A conjunctive query has the form:

$$q(X, \dots) :- r_1(Y, \dots), \dots, r_n(Z, \dots)$$

where:

- $q, r_1, r_n$  are predicate names (database relations names or arithmetic comparison e.g.  $\leq, =, \geq$ )
- $X, Y, Z$  are variables or constants
- $q(X, \dots)$  is the head of the query
- $r_1(Y, \dots), r_n(Z, \dots)$  are subgoals
- the variables in the head should appear also in the body

- ✓ Conjunctive queries can express select-project-join queries
- ✓ Union is expressed as sets of conjunctive queries with the same head predicate

# Example

A Project-Select-Join SQL query as a conjunctive query:

Select      Advises.P, Advises.S  
From      Registered, Teaches, Advises  
Where     Registered.C = Teaches.C and  
              Registered.year=Teaches.year and  
              Advises.P = Teaches.P and  
              Advises.S = Registered.S  
              Registered.year > “98”

q(P, S)    :-    Registered(S, C, year),  
              Teaches(P, C, year),  
              Advises(P, S),  
              Year > “98”

# The Problem

Answering queries using views expressed as **conjunctive queries**

Given:

- A set of Database Relations (from the Mediated Schema)
- A set of views  $V_1, \dots, V_n$  declared as  
**conjunctive queries** over the Database Relations
- A conjunctive query  $q$  over the Database Relations

Compute **a new query expression** that is

- A disjunction of **conjunctive queries** over the view relations
- Is maximally contained rewriting of  $q$

# Algorithms

- Algorithms for **conjunctive queries**
  - The Bucket Algorithm
  - The Inverse-Rules Algorithm
  - The MiniCon algorithm

# Example

- ✓ The query:

$q(S, C, P) \quad :- \quad \text{Teaches}(P, C, \text{year}), \text{Registered}(S, C, \text{year}),$   
 $\quad \quad \quad \text{Course}(C, T), C \geq 300, \text{year} \geq 1995$

- ✓ The mediated schema (Database Relations):

$\text{Registered}(S, C, Y)$     $\text{Course}(C, \text{title})$     $\text{Teaches}(P, C, Y)$

- ✓ The views exported by the data sources:

- $V1(S, C, Y, \text{title}) \quad :- \quad \text{Registered}(S, C, Y), \text{Course}(C, \text{title}),$   
 $\quad \quad \quad C \geq 500, Y \geq 2000$
- $V2(S, P, C, Y) \quad :- \quad \text{Registered}(S, C, Y), \text{Teaches}(P, C, Y)$
- $V3(S, C) \quad :- \quad \text{Registered}(S, C, Y), \text{Year} \leq 1990$
- $V4(P, C, \text{title}, Y) \quad :- \quad \text{Registered}(S, C, Y), \text{Course}(C, \text{title}),$   
 $\quad \quad \quad \text{Teaches}(P, C, Y), Y \leq 1998$

# Complexity

- The number of possible rewritings is exponential in the size of the query
- Bucket Algorithm:
  - It is sufficient to consider for each subgoal the views that might be relevant

# The Bucket Algorithm

- $V1(S, C, Y, \text{title}) :- \text{Registered}(S, C, Y), \text{Course}(C, \text{title}), C \geq 500, Y \geq 2000$
- $V2(S, P, C, Y) :- \text{Registered}(S, C, Y), \text{Teaches}(P, C, Y)$
- $V3(S, C) :- \text{Registered}(S, C, Y), Y \leq 1990$
- $V4(P, C, \text{title}, Y) :- \text{Registered}(S, C, Y), \text{Course}(C, \text{title}), \text{Teaches}(P, C, Y), Y \leq 1998$

$q(S, C, P) :- \text{Teaches}(P, C, \text{year}), \text{Registered}(S, C, \text{year}), \text{Course}(C, T), C \geq 300, \text{year} \geq 1995$

$V4(P, C, T', \text{year})$   
 $V2(S', P, C, \text{year})$

$V2(S, P', C, \text{year})$   
 $V1(S, C, \text{year}, T')$

$V4(S', C, T, \text{year}')$   
 $V1(S', C, \text{year}', T)$

Teaches

Registered

Course

# The Bucket Algorithm

- I. For each subgoal of q, which is not a comparison predicate, create a 'bucket'
- II. Fill the bucket, with the **unifiers of relevant views**
- III. A solution must include a view from each bucket

Consider a query subgoal  $p(X_1, \dots, X_n)$ .

A subgoal  $p(Y_1, \dots, Y_n)$  of  $V$  is relevant to  $p(X_1, \dots, X_n)$ . :

- There is function  $\theta$  that unifies the variables & constants  
 $\theta(p(X_1, X_2, \dots)) = \theta(p(Y_1, Y_2, \dots))$
- The comparison predicates in  $V$  and  $q$  are compatible
- If  $X_i$  appears in the head of  $q$  then  $Y_i$  must be in the head of  $V$

# The Bucket Algorithm

- There are 8 combinations that have to be checked
- For each combination:
  - Make equal appropriate variables to enforce joins
  - If conflicting conditions, abandon
  - Eliminate duplicated goals

V4(P,C, $T'$ , year)

V2( $S'$ ,P,C, year)

V2(S, $P'$ ,C,year)

V1(S,C,year, $T'$ )

V4( $S'$ , C, T,  $year'$ )

V1( $S'$ , C,  $year'$ , T)

Teaches

Registered

Course

$q'(S,C,P) :- V2(S,P,C,year), V1(S,C,year,T'), V1(S',C,year',T)$



$q'(S,C,P) :- V2(S',P,C, year), V1(S, C, year, T')$  **The first solution!**

# Algorithms

- Algorithms for **conjunctive queries**
  - The Bucket Algorithm
  - The Inverse-Rules Algorithm
  - The MiniCon Algorithm



# The Inverse-Rules

**Idea:** Express the predicates of the mediated schema in terms of views: Invert the view definitions

**View Example:**  $V3(D, C) :- \text{Major}(S, D), \text{Registered}(S, C)$

**Intuition:** A tuple  $V3(D, C)$  in the view  $V3$  **witnesses** the existence of a value  $Z=f_S(D,C)$  such that  $\text{Major}(Z, D)$  and,  $\text{Registered}(Z, C)$

**Skolem functions:**

- $f_S$  is called Skolem function
- $f_S$  was introduced for the existential variable  $S$  of  $V3$

**Create Inverse Rules:**

- $\text{Major}(f_S(D,C), D) :- V3(D, C)$
- $\text{Registered}(f_S(D,C), C) :- V3(D, C)$

# Generate the Inverse-Rules

## Algorithm:

For each view V:

- I. For each existential variable in V, create a Skolem function
- II. For each subgoal of a view create a rule:
  - The head is the subgoal
    - ! Use the Skolem functions for existential variables in V
  - The body is the head of the view

# The Inverse Rule Algorithm

**Algorithm:** The rewriting of the query Q using a set of views V<sub>1</sub>, V<sub>2</sub>, .. V<sub>n</sub> is the Datalog program consisting of

- The inverse rules for V<sub>1</sub>, V<sub>2</sub>, .. V<sub>n</sub>
- The query Q

Observation: Datalog is a language similar to Prolog

# Inverse Rule Algorithm (2)

## Example:

Consider the query:  $q(D) :- \text{Major}(S,D), \text{Registered}(S, 444)$   
and the view:  $V3(D, C) :- \text{Major}(S, D), \text{Registered}(S, C)$

## Inputs:

The rules:

- $\text{Major}(f_S(D,C), D) :- V3(D, C)$
- $\text{Registered}(f_S(D,C), C) :- V3(D, C)$
- $q(D) :- \text{Major}(S,D), \text{Registered}(S, 444)$

Suppose that  $V3$  has the following tuples:

D	C
CS	444
EE	444
CS	333

# Inverse Rule Algorithm (3)<sub>v3</sub>

## Execute

- I. The inverse rules produce the tuples:

Major	
$f_S(D,C)$	C
$f_S(CS,444)$	444
$f_S(EE,444)$	444
$f_S(CS,333)$	333

Registered	
D	$f_S(D,C)$
CS	$f_S(CS,444)$
EE	$f_S(EE,444)$
CS	$f_S(EE,444)$

D	C
CS	444
EE	444
CS	333

- II. Apply the query q to the extensions Major & Registered

# Bucket versus Inverse-rules

- Both algorithms produce the **maximal contained rewritings**
- The same intuition:
  - Compute the views that are relevant to DB relations
- The inverse-rules are computed once and used anytime
- The bucket considers only the views relevant to the query

# The MiniCon Algorithm

## An improvement of Bucket

- Excludes much more views from participating in buckets
- There are fewer combination to check at the end

## New constraint:

- If a view  $V$  contains a ‘target’ subgoal of  $Q$  **and**
- The subgoal participates in a join in  $Q$

then  $V$  is relevant to  $Q$  if

- **Either** the ‘join’ variables appear in the head of  $V$
- **Or** the join takes place also within  $V$

# Bucket vs Minico Algorithm

- V1(S, C, title) :- **Registered(S, C, Y), Course(C, title),  
C ≥ 500, Y≥2000**
- V2(S, P, C) :- **Registered(S, C, Y), Teaches(P, C,Y)**
- V3 (S, C) :- **Registered(S, C, Y), Y≤1990**
- V4(P, C, title, Y) :- **Registered(S, C, Y), Course(C, title),  
Teaches(P, C, Y), Y≤1998**

q(S, C, P) :- **Teaches(P, C, year), Registered(S, C, year),  
Course(C, T), C ≥ 300, year≥1995**

V4(P,C,T', year)  
V2(S',P,C,year)

V2(S,P',C,year)  
V1(S,C,yar, T)

V4(S', C, T, year')  
V1(S', C, year', T)

Teaches

Registered

Course

# Further Reading

- Bucket Algorithm [Levy & Rajaraman & Ullman 1996]
- Inverse Rules Algorithm [Duscka & Genesereth 1997]
- MiniCom Algorithm [Pottinger & Halevy 2001]
- MDCSAT [Arvelo & Bonet & Vidal 2006]
- SSSAT [Izquierdo & Vidal & Bonet 2011]
- GQR [Konstantidis & Ambite, 2011]

# AQUV for other Languages

Rewriting language plays a key role:

- The maximal contained expression depends of it

Works considering other rewriting languages

- Object query languages (OQL) (Florescu 96)
- ▪ Views with binding patterns (Rajaraman 95)
  - Describe sources such as Web services, Web Forms

# Access-Pattern Limitations

- Often we can only access the data sources in specific ways:
  - Web forms: require certain inputs
  - Web services: interface definitions
  - Controlling load on systems: allow only limited types of queries
- Model limitations using adornments:
  - $V^{bf}(X, Y)$ : first argument must be bound

# Example Binding Patterns

Cites(X,Y)  
DBPapers(X)  
AwardPaper(X, P)

$S1 : CitationDB^{bf}(X, Y) \subseteq Cites(X, Y)$

$S2 : CitingPapers^f(X) \subseteq Cites(X, Y)$

$S3 : DBSource^f(X) \subseteq DBpapers(X)$

$S4 : AwardDB^{bf}(X, P) \subseteq AwardPaper(X, P)$

# Executable Plans

Conjunctive plan  $q_1(\overline{X}_1), \dots, q_n(\overline{X}_n)$

And let  $BF_i$  be the adornments for  $q_i$

We say that the plan is *executable* if there are  $bf_i \in BF_i$  s.t.

If a variable  $X$  is in position  $k$  of  $q_i$ , and  $bf_i$  has a 'b' in  $k$ 'th position, then  $X$  occurs earlier in the plan i.e., every variable that must be bound has a value.

# Example (1)

$S1 : CitationDB^{bf}(X, Y) \subseteq Cites(X, Y)$

$S2 : CitingPapers^f(X) \subseteq Cites(X, Y)$

$S3 : DBSource^f(X) \subseteq DBpapers(X)$

$S4 : AwardDB^{bf}(X, P) \subseteq AwardPaper(X, P)$

Cites(X, Y)  
DBPapers(X)  
AwardPaper(X, P)

$Q(X) :- Cites(X, 001)$

✗  $Q'(X) :- CitationDB(X, 001)$

$Q'(X) :- CitingPapers(X), CitationDB(X, 001)$

# Rewritings of unbound length!

$S1 : CitationDB^{bf}(X, Y) \subseteq Cites(X, Y)$

$S3 : DBSource^f(X) \subseteq DBpapers(X)$

$S4 : AwardDB^{bf}(X, P) \subseteq AwardPaper(X, P)$

Cites(X, Y)  
DBPapers(X)  
AwardPaper(X, P)

$Q(X) :- AwardPaper(X, "VLDB Best Paper")$

$Q'(X) :- DBSource(X), AwardDB(X, "VLDB Best Paper")$

$Q'(X) :- DBSource(Y), CitationDB(Y, X),$   
 $AwardDB(X, "VLDB Best Paper")$

$Q'(X) :- DBSource(Y), CitationDB(Y, X_1), \dots,$   
 $CitationDB(X_n, X), AwardDB(X, "VLDB Best Paper")$

# Recursive Plans to the Rescue

*papers(X) :- DBSource(X)*

*papers(X) :- papers(Y), CitationDB(Y, X)*

*Q'(X) :- papers(X), AwardDB(X, "VLDBBestPaper")*

Theorem: this can be done in general

[Duschka, Genesereth & Levy, 1997]

# Recursive Query Plans

[Duschka & al 2000]

- For every view with binding pattern  
k times

$$\underbrace{V^b, \dots b, f, \dots f}_{k \text{ times}}(x_1, \dots x_k \dots x_n) :- r_1(\dots), \dots r_m(\dots)$$

- Add **an inverse rule** for each of the terms  $r_1, \dots r_m$   
 $r_i(\dots) :- \text{dom}(x_1), \text{dom}(x_2), \dots \text{dom}(x_k), V^b, \dots b, f, \dots f(x_1, \dots x_k \dots x_n)$
- Add **a dom rule** for every free variable  $x_{k+1}, \dots x_n$   
 $\text{dom}(x_j) :- \text{dom}(x_1), \text{dom}(x_2), \dots \text{dom}(x_k), V^b, \dots b, f, \dots f(x_1, \dots x_k \dots x_n)$

- For every constant c in the query add **a rule**  
 $\text{dom}(c) :-$

The query is evaluated on the resulted Datalog program

# A Graph Perspective

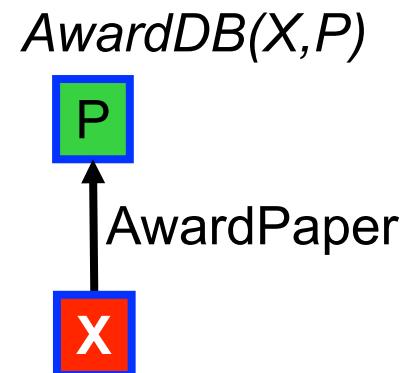
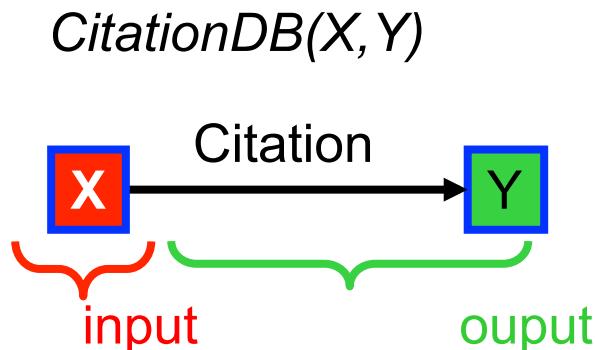
- If the mediated schema consists of binary relations (as in RDF), we can represent queries and plans as graphs

$$S1 : CitationDB^{bf}(X, Y) \subseteq Cites(X, Y)$$

Cites(X, Y)  
DBPapers(X, T)  
AwardPaper(X, P)

$$S2 : DBSource^f(X) \subseteq DBpapers(X, "Databases")$$

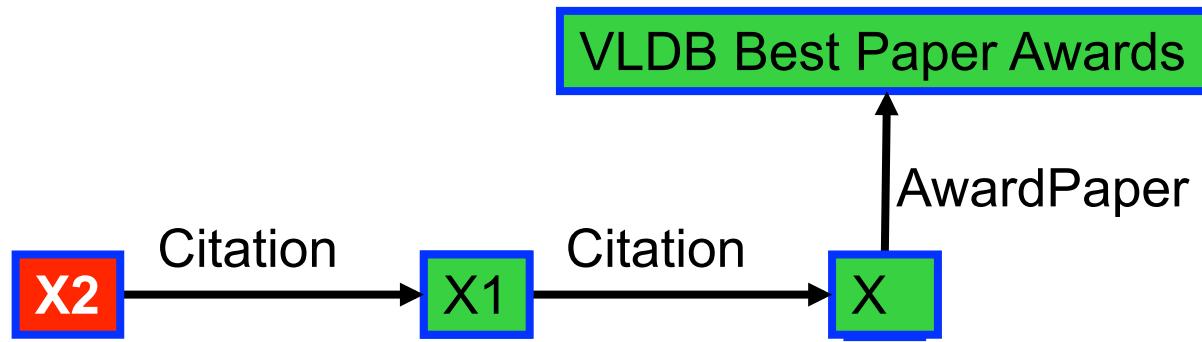
$$S4 : AwardDB^{bf}(X, P) \subseteq AwardPaper(X, P)$$



# Query Evaluation Example

- A query is a graph
- A query plan is also a graph

$Q(X) :- \text{AwardPaper}(X, "VLDB Best Paper")$



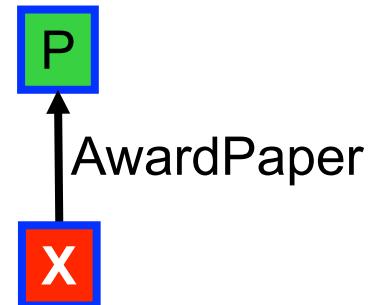
$DBPapers(X)$



$CitationDB(X, Y)$



$\text{AwardDB}(X, P)$



# Question

- Are the all the enumerated rewritings relevant (can they return new results) ?
- Answer: Unfortunately, yes as the sources are incomplete.

# Query Answering Problem

Given a query Q against

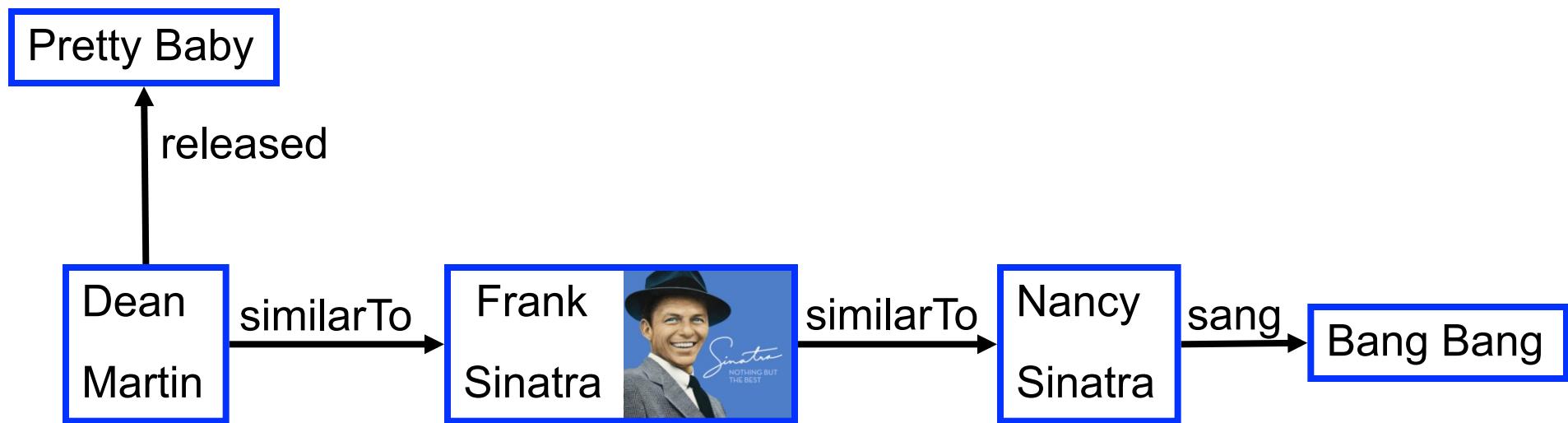
- a RDF base G
- a set of functions F

★ A limit **Max** for the number of Web calls

compute the larger number of answers using **Max** calls

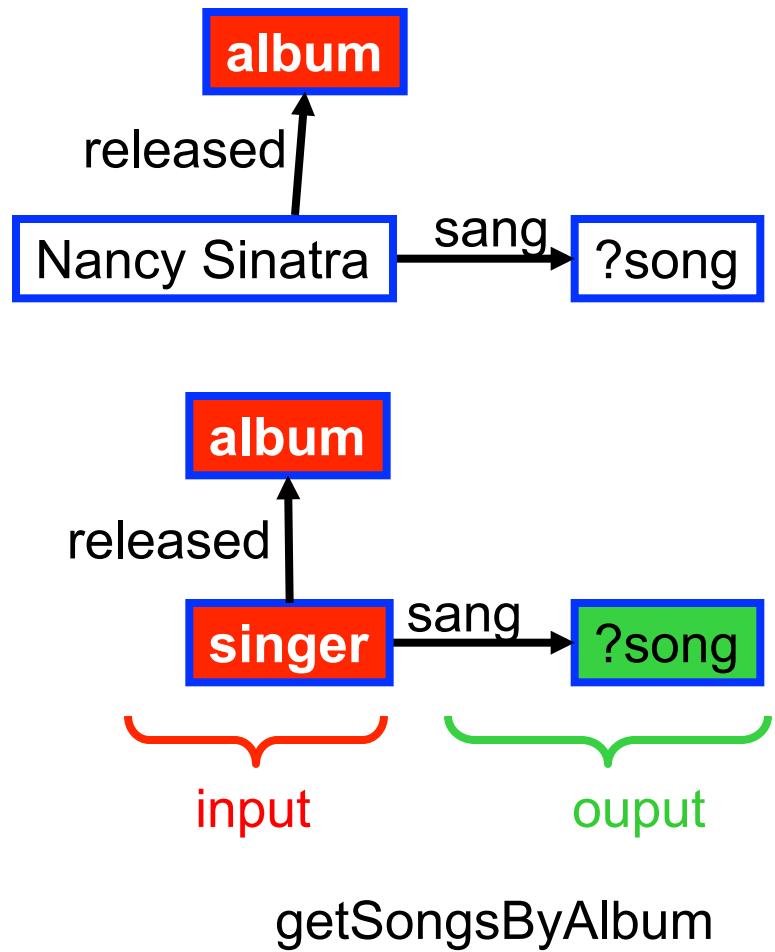
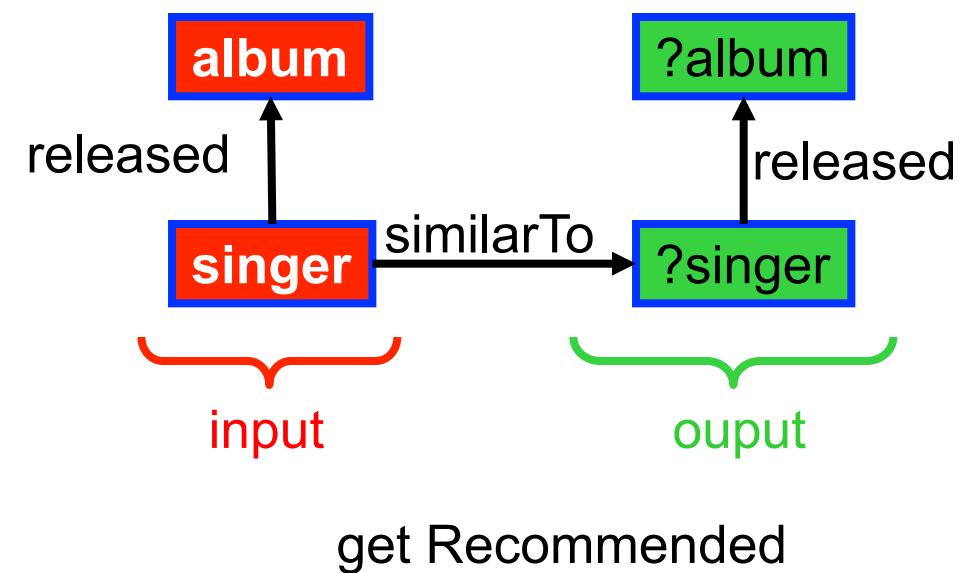
→ We need to prioritize the calls

# RDF Fragment

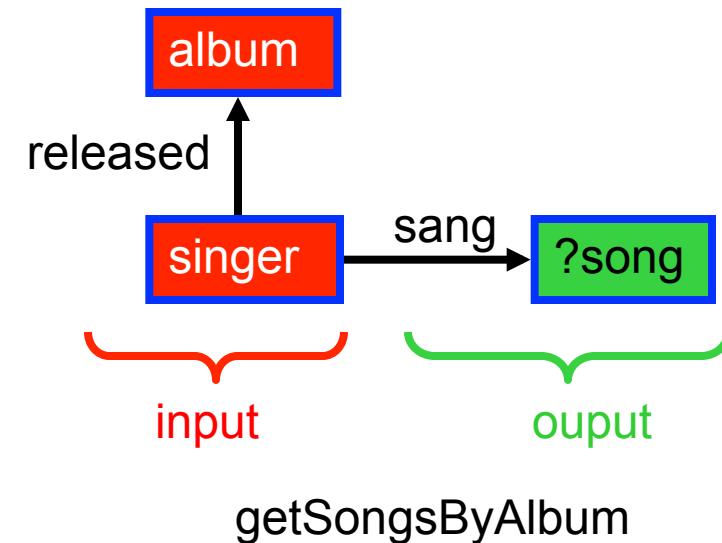
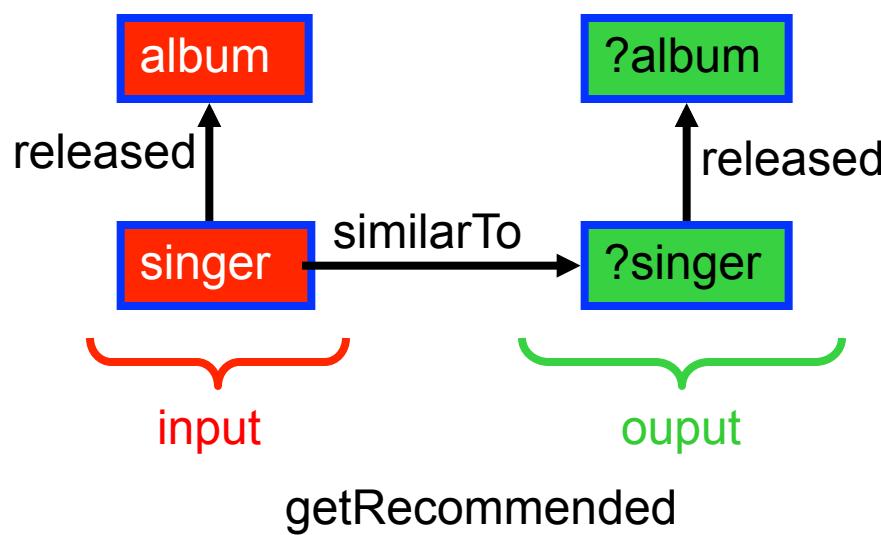
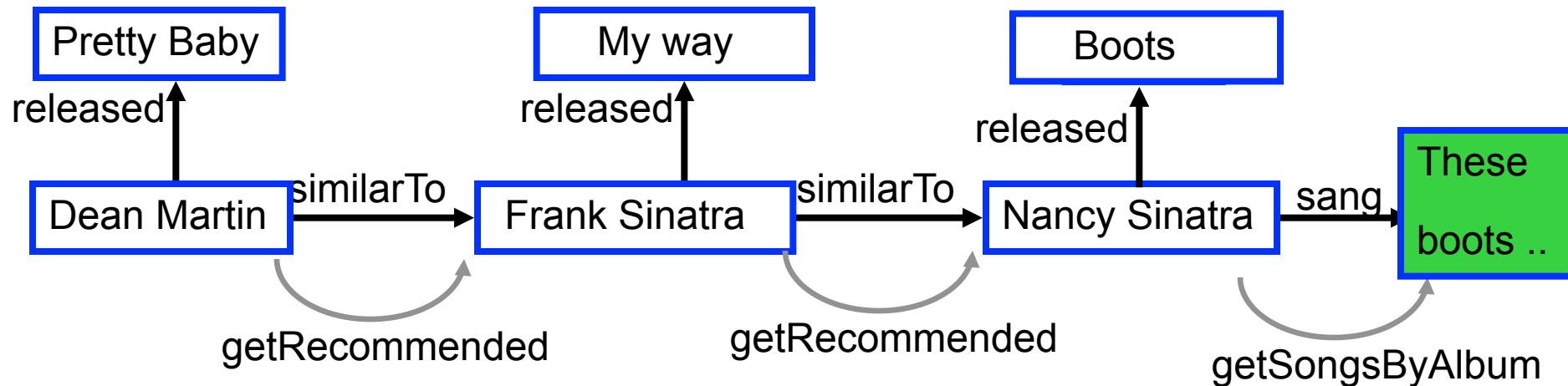


The RDF fragment itself can be described as LAV views!

# Query example

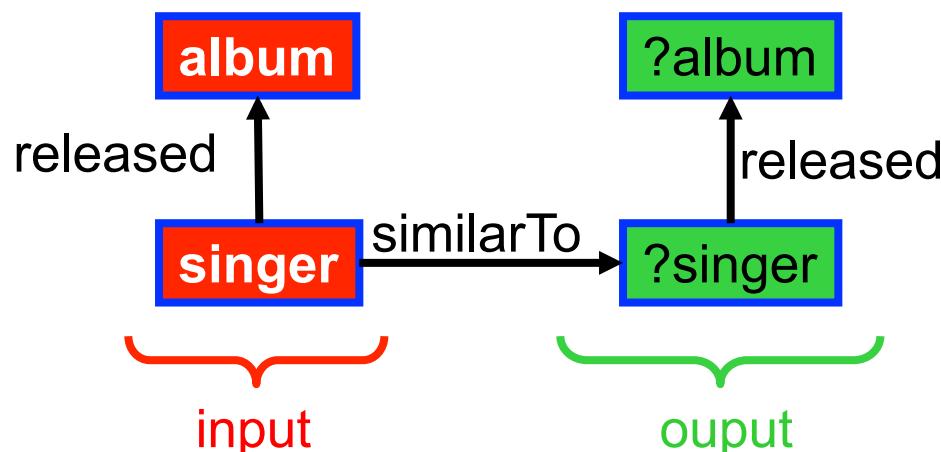
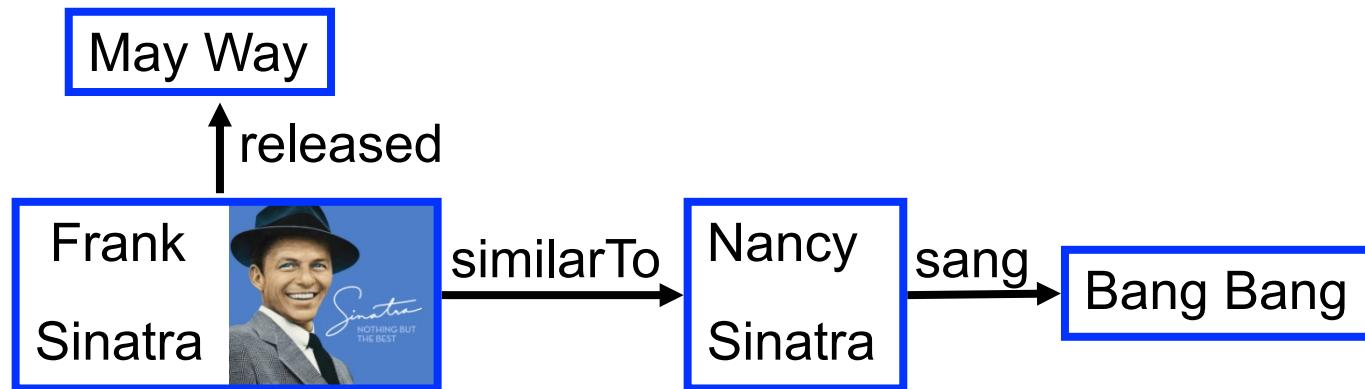


# ANGIE Algorithm



# Exercise:

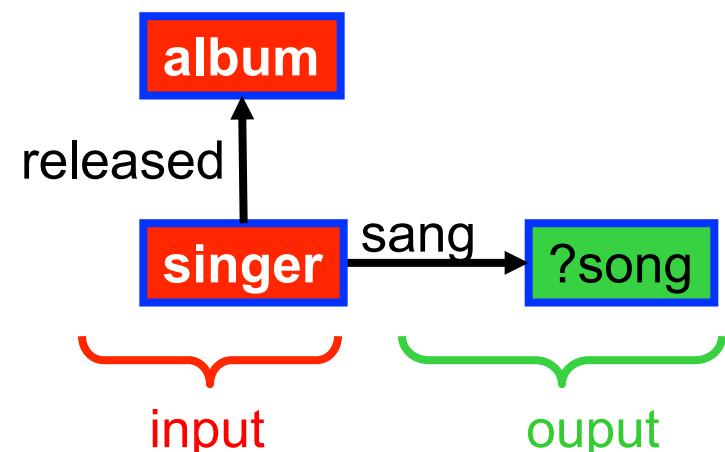
- Input:



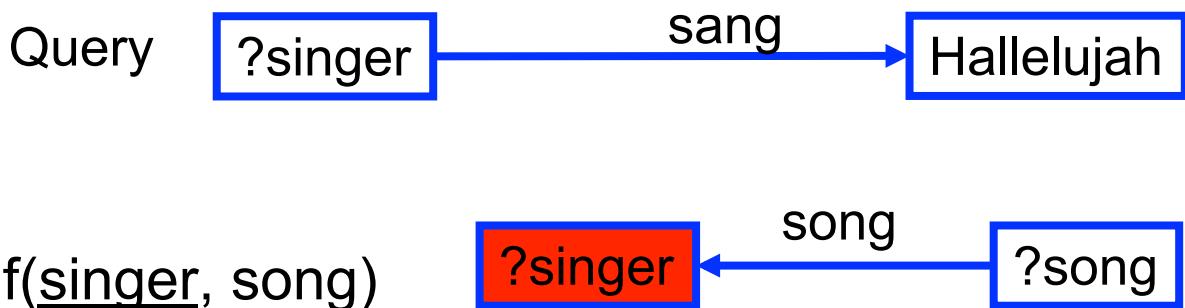
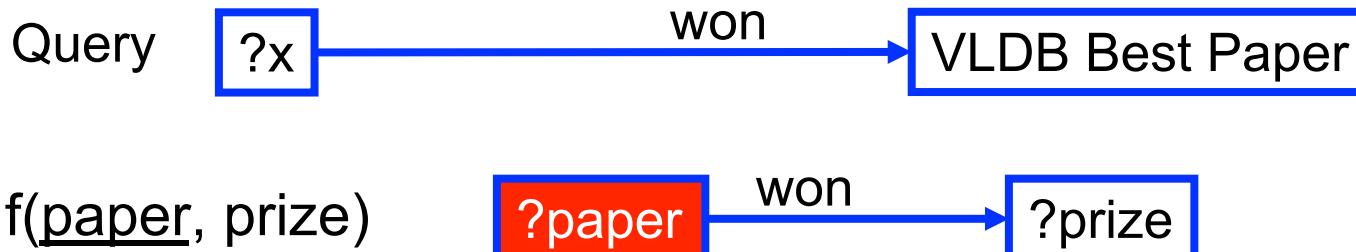
get Recommended



getSongsByAlbum

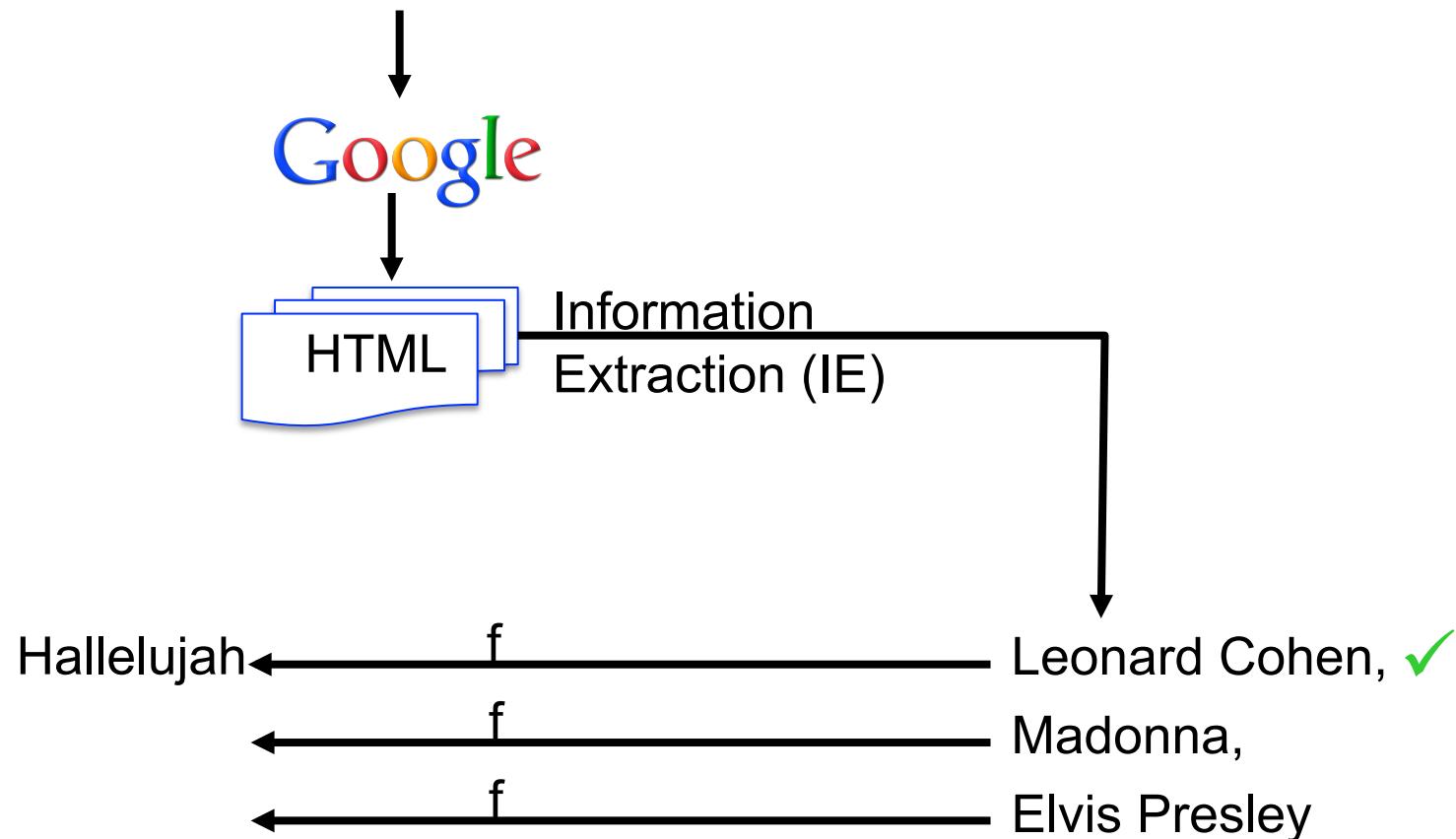


# Asymmetric Accesses

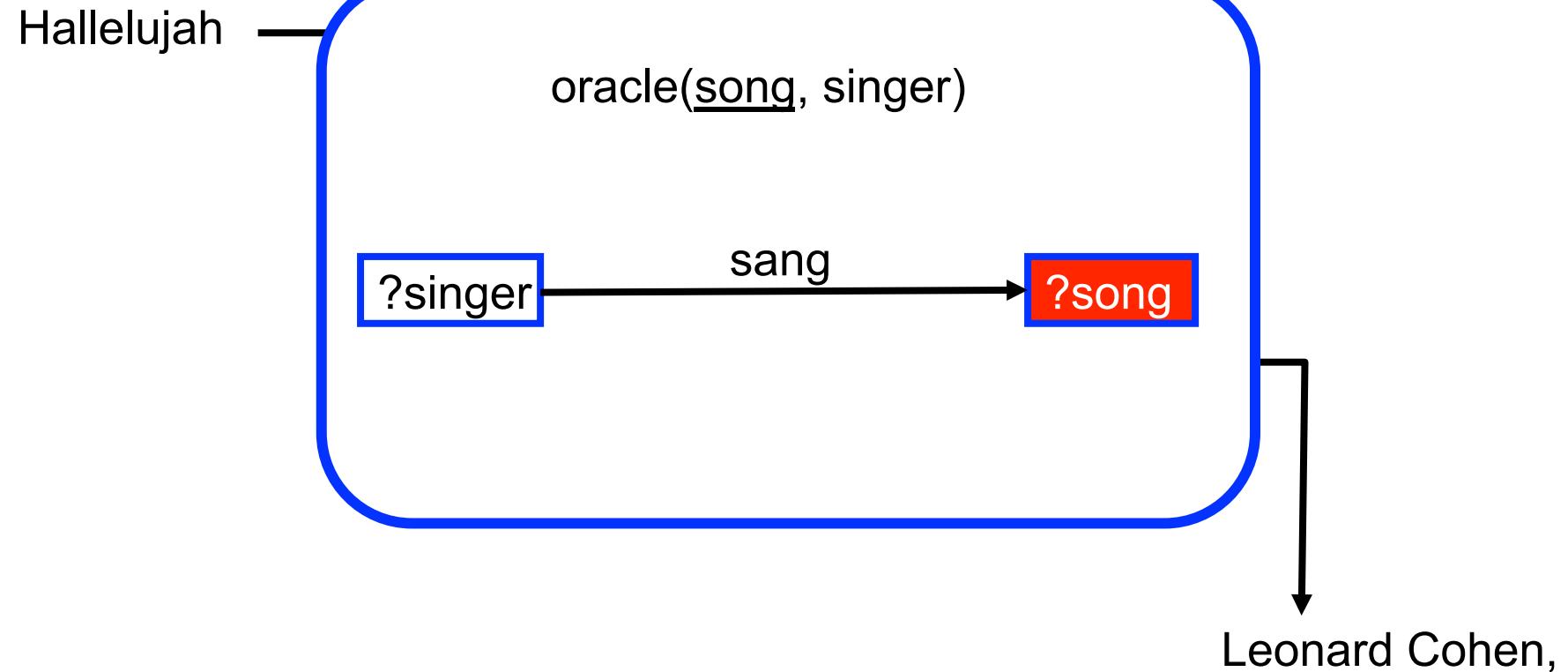


# Idea: Use the Web as an Oracle

Hallelujah → “Singers of Hallelujah”



# Model Oracles as Web services



# AQUV for other Languages

Rewriting language plays a key role:

- The maximal contained expression depends of it

Works considering other rewriting languages

- Object query languages (OQL) (Florescu 96)
- Views with binding patters (Rajaraman 95)
  - Describe sources such as Web services, Web Forms
- ▪ Conjunctive Queries with integrity constraints

# Integrity Constraints with LAV

Schedule(airline, flightNum, date, pilot, aircraft)

Functional dependencies:

$Pilot \rightarrow Airline$  and  $Aircraft \rightarrow Airline$

$S(date, pilot, aircraft) \sqsubseteq schedule(a, fN, date, pilot, aircraft)$

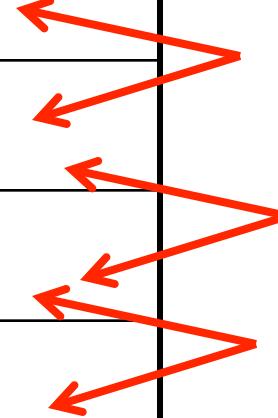
Query: (pilots in same airline as Mike)

$q(p) :- schedule(a, fN, date, "mike", aircraft)$   
 $schedule(a, f, d, p, a)$

*Without functional dependencies, no answers!*

# But We Do Have Answers...

Date	Pilot	Aircraft
1/1	Mike	111
5/2	Ann	111
1/3	Ann	222
4/3	John	222



Mike and Ann work for the same airline  
111 and 222 belong to the same airline  
John and Ann work for the same airline, ...

# No Limit to Rewriting Size

```
q'_n(p) :- S(D1, "mike", C1), S(D2, p2, C1),  
          S(D3, p2, C2), S(D4, p3, C2), ...,  
          S(D2n-2, pn, Cn-1), S(D2n-1, pn, Cn)  
          S(D2n, p, Cn)
```

Does this sound familiar?

Recursive plans to the rescue!

# Recursive Plan

- Start with inverse rules and get:

$schedule(f_1(1/1, Mike, \#111), f_2(1/1, Mike, \#111), 1/1, Mike, \#111)$

$schedule(f_1(5/2, Ann, \#111), f_2(5/2, Ann, \#111), 5/2, Ann, \#111)$

$schedule(f_1(1/3, Ann, \#222), f_2(1/3, Ann, \#222), 1/3, Ann, \#222)$

$schedule(f_1(4/3, John, \#222), f_2(4/3, John, \#222), 4/3, John, \#222)$

- But we know that:

$$f_1(1/1, Mike, \#111) = f_1(5/2, Ann, \#111)$$

$$f_1(5/2, Ann, \#111) = f_1(1/3, Ann, \#222)$$

$$f_1(1/3, Ann, \#222) = f_1(4/3, John, \#222)$$

# The e relation

- Let's define a relation  $e(X, Y)$  that stores all pairs that are actually equal.
- $e(x, x)$  – holds for every  $x$ .
- The other rules for inferring  $e(X, Y)$  are derived from the functional dependencies.
- Add:  $e(X, Y) :- e(X, Z), e(Z, Y)$

# The Rules

*Aircraft* → *Airline*: So...

```
e(X,Y) :- schedule(X,F,P,D,A),  
          schedule(Y,F',P',D',A'),  
          e(A,A')
```

*Pilot* → *Airline*: So...

```
e(X,Y) :- schedule(X,F,P,D,A),  
          schedule(Y,F',P',D',A'),  
          e(P,P')
```

# One More Step

- Reformulate the query to use  $e(X, Y)$ :

$q(P) :- \text{schedule}(AL, F', D', Mike, A),$

$\quad \text{schedule}(AL, F, D, P, A')$



$q(P) :- \text{schedule}(AL, F', D', P', A),$

$\quad \text{schedule}(AL', F, D, P, A'),$

$\quad e(AL, AL'), e(P', Mike)$

# Answering Queries Using Views

- Motivating Examples
  - Problem Definition
  - Two different settings:
    - Data Integration
    - Query Optimization & Physical Data Independence
- 

# Query answering: Optimization Versus Data Integration

- In the query optimization context
    - A variation of the Selinger Algorithm
    - Cost-based equivalent rewritings
  - In the data integration context
    - The bucket algorithm
    - The inverse-rules algorithm
    - The MiniCon algorithm
- 
- Equivalent rewritings
- Maximally-contained rewritings

# Materialized Views

Suppose that we have the following views:

- **create view** *JoinStudent&Registered* as  
select Name, C  
from Student, Registered  
where Student.S = Registered.S
- **create view** *JoinProfessor&Teaches* as  
select Name, P  
from Professor, Teaches  
where Professor.P = Teaches.P

## Observation:

- The views are evaluated and the result is materialized (stored)
- The materialized result can be used in the evaluation of queries

# Optimization Goal

Reduce the computation by using pre-computed joins  
stored as materialized views

# Query Optimization Example

- Consider the query:

Select Student.Name, Professor.Name

From Student, Registered, Teaches, Professor

Where Student.S=Registered.C and Teaches.P=Professor.P  
and Registered.C=Teaches.C

- A possible rewriting to reduces the query execution time:

Select Student.Name, Professor.Name

From JoinStudent&Registered, JoinProfessor&Teaches

Where JoinStudent&Registered.C= JoinProfessor&Teaches.C

2 Joins eliminated by using materialized views instead

# Query Optimization Problem

Given

- A query  $Q$
- A set of views  $V_1, V_2, \dots, V_n$
- A cost model  $C$

Compute **the equivalent query plan  $Q'$  ( $Q' = Q$ )**

that has **the cheapest cost**

# Difficulties

- Large number of many possible combinations
  - Large search space to find the solution
- In some cases, views may not help due to loss of indexes

# Relevant Views

A view **V** is useful in the evaluation of a query **Q** if:

- There is a mapping of relations from **V** to **Q** **and**
- **V** must project at least all attributes projected in **Q** **and**
- Each relation **R1** that appears in both **V** and **Q** satisfies:  
the join **and** the selection conditions

# Relevant Views: Join Conditions

Consider a relation R1 appearing in both Q and V and let R2 be the other relation in the join in Q.

R1 satisfies the join conditions if:

- In V, R1 and R2 are joined in the same attributes and conditions
- **OR** In V, R1 and R2 are joined in the same attributes and for weaker conditions and the conditions are projected in V
- **OR** In V, R1 and R2 are not joined, but the attributes participating in the join in Q are projected in V

# Relevant Views: Selection Conditions

Consider a relation R1 appearing in both Q and V and let A be an attribute of R1 for which Q applies a selection predicate:

R1 satisfies the selection conditions if :

- Either V applies the same selection predicate for the attribute A
- OR V applies a weaker selection predicate and A is projected in V
- OR V applies no selection predicate for A and A is projected in V

# Examples

Q: Select Student.Name, Professor.Name  
From Student, Registered, Teaches, Professor  
Where Student.S=Registered.S and Registered.C = Teaches.C  
and Teaches.P = Professor.P and Student.Year> 2000

V1: Select name  
from Student  
where Student.year>2000 

V5: Select Student.name, Course.name  
from Student, Registered, Course  
where Student.S = Registered.S  
and Registered.C = Course.C 

V2: Select name, S  
from Student  
where Year> 2000 

V6: Select Student.name, Student.year,  
Course.name,

Course.C  
from Student, Registered, Course  
where Student.S= Registered.S  
and Registered.C = Course.C 

V3: Select name, S  
from Student  
where Year >1996 

V4: Select name, S, Year  
from Student  
where Year >1996 

# Query Optimization – the Selinger Algorithm

- The Selinger Algorithm chooses the best join order
- A bottom-up dynamic programming algorithm
- Each iteration explores all join combinations of certain length
- For each combination of joins, it keeps only the cheapest plan, taking into account costs estimated in the previous iteration

# A variation of the Selinger Algorithm

- The Selinger Algorithm can be modified in order to choose the cheapest equivalent rewriting
- Start with views that seem useful for a certain query
- Continue in a bottom-up fashion exploring all join combinations
- For each combination, prune plans estimated to be more expensive or less contributing than others

# A variation of the Selinger Algorithm

- **First Iteration:**
  - Find all views that are relevant (useful) to the query
  - Choose the best access path for each view and evaluate its cost
- **Second Iteration:**
  - Find all pair combinations
  - For each pair find all possible Joins in all possible attributes that lead to useful plans
  - Plans that are complete (i.e. are sufficient to answer the query) are distinguished as final candidates
  - For the rest plans do the following:
    - Compare pairs of plans
    - For each pair  $(p, p')$ , if  $p'$  is cheaper and has greater or equal contribution (covers more relations), then  $p$  is pruned
- **Third Iteration:**
  - Find all possible pairs of plans derived from the previous iteration and the plans derived from the first iteration
  - ...

# A variation of the Selinger Algorithm

Query:  $R1 \bowtie R2 \bowtie R3 \bowtie R4$

Views:  $V1 = R1 \bowtie R2$ ,  $V2 = R2 \bowtie R3 \bowtie R4$

$V3 = R2 \bowtie R3$ ,  $V4 = R3 \bowtie R4$ ,

