

Logic-based Techniques for Information Integration

Marie-Christine ROUSSET
University of Grenoble



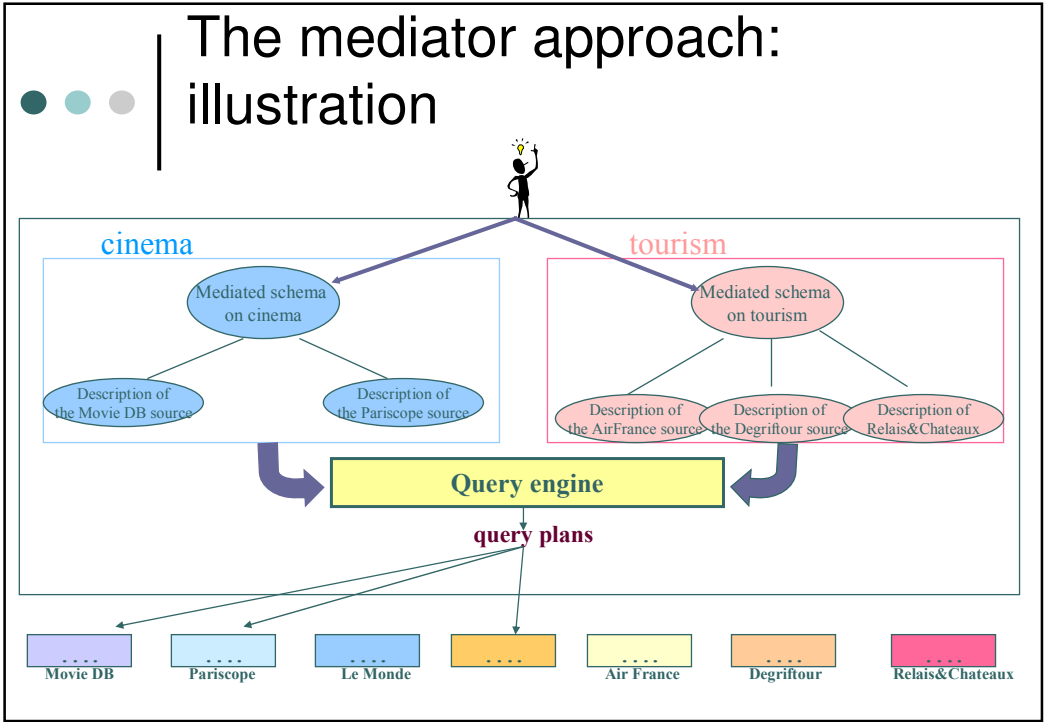
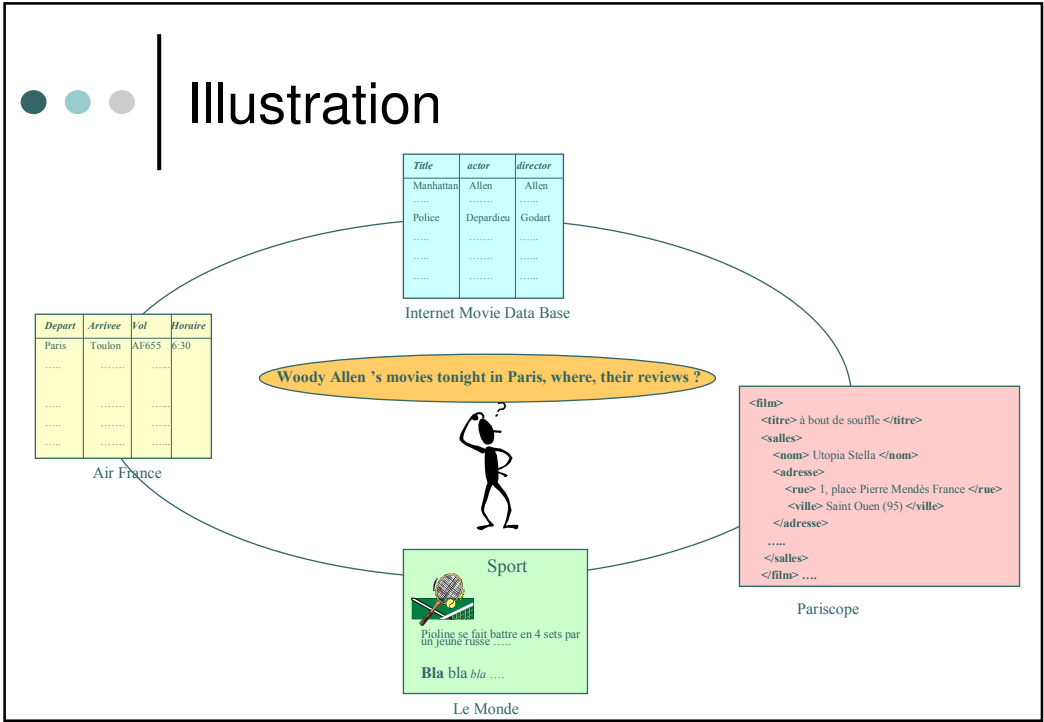
Introduction

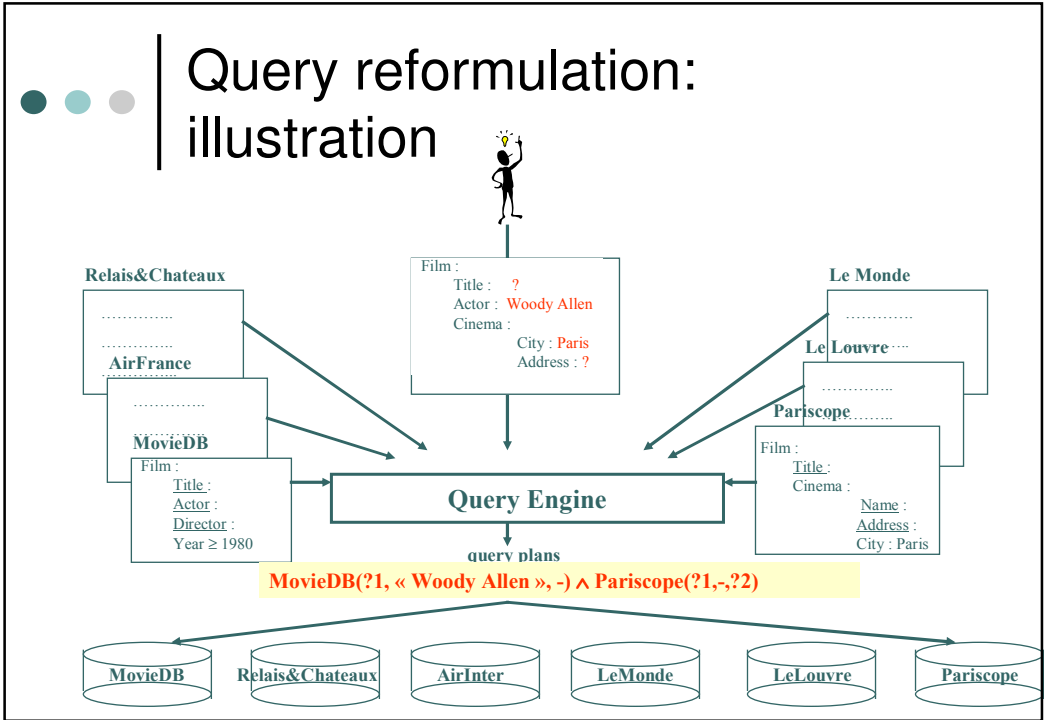
Information integration

- Core issue for the construction of modern information servers
 - **traditional DBMSs** : efficient management and querying of structured, centralized and reliable data
 - **data available through the Web** : distributed over autonomous and heterogeneous sources
- Requires innovative solutions
 - **mediation systems**

Key problems

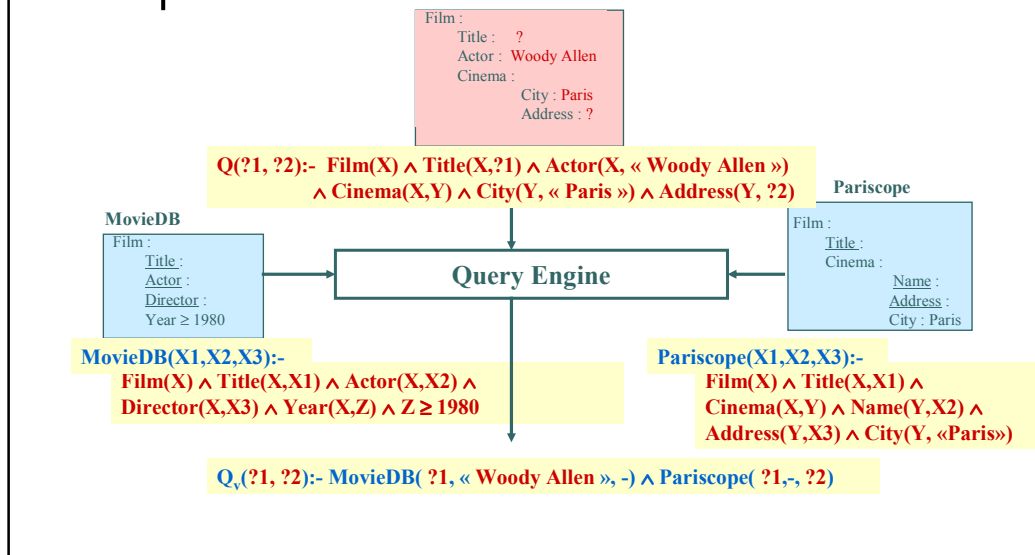
- Define a mediated schema
 - A structured vocabulary of the domain of interest serving as query interface for users
- Describe the content of pre-existing data sources in terms of a mediated schema
- Reformulate and decompose user's queries into queries executable against the relevant data sources
- Combine the answers to the local queries in order to get (certain) answers to the original user's query





- ## Principles of a logical approach
- Logical description of schema, queries and views
 - the mediated schema :
 - a domain specific knowledge base (domain ontology)
 - the views :
 - an abstract description of the content of the data sources to integrate
 - Reasoning for query reformulation
 - Query rewriting in terms of views
 - decidability and complexity depending on the logical formalism

Formal background: a logical semantics



Overview of the course

- Queries and views: logical semantics
- Query answering and rewriting: two reasoning problems
- Query rewriting using views in the simple relational model
 - Main algorithms
- Answering queries over ontologies by rewriting
 - First results on a complex problem

Queries

- Open formulae in first-order logic (FOL)

$$q(\underline{x}): \exists \underline{y} \Phi(\underline{x}, \underline{y})$$

- \underline{x} is a tuple of free (a.k.a distinguished) variables
- $\Phi(\underline{x}, \underline{y})$ is a FOL formula with as free variables those in \underline{x} and \underline{y}

- Example

$$q(X) : \exists A, C \text{ Flight}(X) \wedge \text{ArrivalAirport}(X, A) \wedge \text{Located}(A, C) \wedge \text{Capital}(C)$$

Datalog notation:

$$q(X) \leftarrow \text{Flight}(X) \wedge \text{ArrivalAirport}(X, A) \wedge \text{Located}(A, C) \wedge \text{Capital}(C)$$

Answers

- Answer set of a query q over a knowledge base K

$$\text{Ans}(q, K) = \{\underline{a} \mid K \models q(\underline{a})\}$$

- K is a set of closed FOL formulae
 - a database (a set of ground atoms of the form $R(a_1, \dots, a_n)$ where R is a relation of the schema) + possibly constraints on those relations
 - $A_{\text{box}} \cup T_{\text{box}}$ (in Description Logic knowledge bases)
- \underline{a} is a tuple of constants appearing in K
- $q(\underline{a})$: the formula $\exists \underline{y} \Phi(\underline{x}, \underline{y})$ where the variables of \underline{x} are replaced by the constants of \underline{a}

Views

o Named queries

- for representing by a formula a set of data
 - cached answers of a given query to be re-used for optimization of answering new queries
 - the content of a distant data source

o Example

$v1(X,Y1,Y2) : \text{Flight}(X) \wedge \text{DepartureAirport}(X,Y1) \wedge \text{ArrivalAirport}(X,Y2)$

$v2(X,Y) : \text{Place}(X) \wedge \text{Located}(X,Y) \wedge \text{Capital}(Y)$

o can be materialized or virtual

- Their extension can be stored (cached) or computed (by querying an external source)

Different semantics

of the mapping $v(\underline{x}) : \text{def}(\underline{x},\underline{y})$ between the view v and the query defining it

• exact semantics

$\text{Ext}(v) = \text{Ans}(\text{def},K)$

new axiom : $\forall \underline{x} [v(\underline{x}) \Leftrightarrow \exists y \text{def}(\underline{x},y)]$ added to K

• sound semantics

$\text{Ext}(v) \subseteq \text{Ans}(\text{def},K)$

new axiom : $\forall \underline{x} [v(\underline{x}) \Rightarrow \exists y \text{def}(\underline{x},y)]$ added to K

• complete semantics

$\text{Ans}(\text{def},K) \subseteq \text{Ext}(v)$

new axiom : $\forall \underline{x} [\exists y \text{def}(\underline{x},y) \Rightarrow v(\underline{x})]$ added to K

● ● ● | Answering queries using views

- Let q a query over a knowledge base K , V a set of views, K_v the set of axioms defining the semantics of the views w.r.t K , $\text{ext}(V)$ the union of the view extensions,
- the set of **certain** answers of q over a K w.r.t V is

$$\text{CertainAns}(q,K,V) = \{\underline{a} \mid K \cup K_v \cup \text{ext}(V) \models q(\underline{a})\}$$
 where \underline{a} is a tuple of constants appearing in $K \cup \text{ext}(V)$

● ● ● | a reasoning problem on data

- More complex than answering queries in DBMS
 - The views extensions provide **incomplete** data on the extensions of the query predicates.

Views :

$v1(X) : \text{Reservation}(X,X) \wedge \text{Reservation}(X,Y)$

$v2(X) : \text{Reservation}(X,Y)$

$v3(X) : \text{Reservation}(Y,X).$

Views extensions :

$v1 \quad v2 \quad v3$

$a \quad a \quad a$

$\quad b \quad c$

Query :

$q(X1,X2) \leftarrow \text{Reservation}(X1,X2).$

The only **certain** answer to q using views: **(a,a)**

Reservation

(a,a)

(a,?)

(a,?)

(b,?)

(?,a)

(?,c)

● ● ● | Rewriting queries using views

○ Given

- a query $q(\underline{x})$ defined over K using a language $L1$
- a set of views V defined on D in $L2$

$$V = \{v_1(\underline{x}_1) : Q_1(\underline{x}_1), \dots, v_k(\underline{x}_k) : Q_k(\underline{x}_k)\}$$

○ q_v is a (maximal) rewriting of q using V if it is:

- defined on $\{v_1, \dots, v_k\}$ in $L3$

- (maximally) contained in q

for every \underline{x} : $q_v(\underline{x}), K_v \cup K \models q(\underline{x})$

(if q'_v is a rewriting of q s.t for every \underline{x} $q_v(\underline{x}) \models q'_v(\underline{x})$ then for every \underline{x} $q_v(\underline{x}) \equiv q'_v(\underline{x})$)

● ● ● | Example

○ Query :

$q(X) \leftarrow \text{Flight}(X) \wedge \text{ArrivalAirport}(X,A) \wedge \text{Located}(A,C) \wedge \text{Capital}(C)$

○ Views (sound semantics) :

$v1(X,Y1,Y2) : \text{Flight}(X) \wedge \text{DepartureAirport}(X,Y1) \wedge \text{ArrivalAirport}(X,Y2)$

$v2(X,Y) : \text{Place}(X) \wedge \text{Located}(X,Y) \wedge \text{Capital}(Y)$

○ Rewriting :

$q_v(X) \leftarrow v1(X,Y1,Y2) \wedge v2(Y2,Y)$

\downarrow
 $\text{Flight}(X) \wedge \text{DepartureAirport}(X,Y1) \wedge \text{ArrivalAirport}(X,Y2)$

\searrow
 $\wedge \text{Place}(Y2) \wedge \text{Located}(Y2,Y) \wedge \text{Capital}(Y)$



Example with existential variables in the views definitions

Query :

$q(X) \leftarrow \text{Manager}(X,Y) \wedge \text{Head}(Y,Z)$

Views (exact semantics) :

$v1(X) : \text{Manager}(X,Y)$

$v2(Y,Z) : \text{Head}(Y,Z)$

$q_v(X) \leftarrow v1(X) \wedge v2(Y,Z)$ rewriting of q ?

1) $[\exists Y,Z v1(X) \wedge v2(Y,Z)], K_v \equiv [\exists Y,Z,Y' \text{Manager}(X,Y') \wedge \text{Head}(Y,Z)]$

2) $[\exists Y, Z,Y' \text{Manager}(X,Y') \wedge \text{Head}(Y,Z)] \not\equiv \exists Y, Z \text{Manager}(X,Y) \wedge \text{Head}(Y,Z)$

q_v is not a rewriting of q



Answering queries by rewriting

For any q_v rewriting of q using views V :

$\text{Ans}(q_v, \text{Ext}(V)) = \{\underline{a} \mid \text{Ext}(V) \models q_v(\underline{a})\}$

$q_v(\underline{a}), K_v \cup K \models q(\underline{a})$

$\text{CertainAns}(q,K,V) = \{\underline{a} \mid K \cup K_v \cup \text{ext}(V) \models q(\underline{a})\}$

Therefore: $\text{Ans}(q_v, \text{Ext}(V)) \subseteq \text{CertainAns}(q,K,V)$


Rewriting the query using views:

• A reasoning problem

- on the query, the view definitions and the axioms
- independent of the data


• Evaluating each rewriting against the views extensions

- A DBMS problem, known to be polynomial in data complexity (the size of the data)



Key question

- Can we compute all the certain answers by computing and evaluating the rewritings ?
 - depends on the query, view and rewriting languages and semantics
 - necessary condition :
 - The instance checking problem is polynomial in data complexity
 - a bunch of negative results coming from the complexity of reasoning in Description Logics
 - Sufficient condition :
 - There exists a finite number of maximal conjunctive rewritings
 - Positive results in the simple setting of relational databases



Answering queries by rewriting
using views in relational setting

● ● ● | Query rewriting using views in relational databases

- Extensively studied in DB theory
 - central for query optimization
 - closely related to the problem of query containment
- Results of decidability and complexity
 - depending on the (relational) languages used for the queries, rewritings and views
 - NP-complete when queries, rewritings and views are conjunctive queries

● ● ● | The «Bucket » algorithm

- Input
 - a conjunctive query (possibly with comparison predicates) over a relational schema
 - a set of sound views defined by conjunctive queries (possibly with comparison predicates)
- output : a set of conjunctive rewritings
 - maximally contained if the query is without comparison predicates (complete algorithm)
- Implemented in Information Manifold

A.Levy, A. Rajaraman, J.Ordille. Querying heterogeneous information sources using source descriptions. Proceedings of the Int.Conference on Very Large Data bases (VLDB 96)

● ● ● | Principles : two steps

- Create a « bucket » for each atom g in the query
 - store there all the view-atoms having a conjunct unifiable with g in its definition (without violating the comparison constraints)
- Build the set of candidate rewritings
 - pick one view-atom per bucket and take their conjunction
 - For each of them check whether their expansion is contained in the query
 - If yes : return it in the output
 - If no : add comparison atoms if it makes containment possible

● ● ● | Creation of bucket(g)

- let V be a view-atom such that one atom g' in its definition is unifiable with g
 - let ψ be the substitution of the variables of g' onto the variables of g
- let $C(V)$ (resp. $C(Q)$) be the conjunction of comparison atoms in the definition of V (resp. Q)
 - If $\psi(C(V))$ and $C(Q)$ are compatible
add the atom $\psi'(V)$ to bucket(g)
where ψ' extends ψ to the variables in V and not in g' by substituting them with new variables

Example

views

- V1(Student,Number,Year) : registered(Student,Course,Year) \wedge course(Course,Number)
 \wedge Number \geq 500 \wedge Year \geq 1992
- V2(Student,Dept,Course) : registered(Student,Course,Year) \wedge enrolled(Student,Dept)
- V3(Student,Course) : registered(Student,Course,Year) \wedge Year \leq 1990
- V4(Student,Course,Number) : registered(Student,Course,Year) \wedge course(Course,Number)
 \wedge enrolled(Student,Dept) \wedge Number \leq 100

query : $q(S,D) \leftarrow$ enrolled(S,D) \wedge registered(S,C,Y) \wedge course(C,N) \wedge N \geq 300 \wedge Y \geq 1995

Bucket(enrolled(S,D)) = {V2(S,D,C'), V4(S,C'',N')}

V3(S,C) \notin Bucket(registered(S,C,Y))
because Y \leq 1990 and Y \geq 1995 are not compatible

V4(S,C,N'') \in Bucket(registered(S,C,Y))
because N'' \leq 100 and N \geq 300 are compatible

Example(ctd)

Bucket(enrolled(S,D)) = {V2(S,D,C'), V4(S,C'',N')}

Bucket(registered(S,C,Y)) = {V1(S,N'',Y), V2(S,D',C), V4(S,C,N'')}

Bucket(course(C,N)) = {V1(S',N,Y')}

Candidate rewriting number 1 :

Q1(S,D) \leftarrow V2(S,D,C') \wedge V1(S,N'',Y) \wedge V1(S',N,Y')

● ● ● | Example (ctd)

Verification of that candidate as a rewriting

$Q1(S,D) \leftarrow V2(S,D,C') \wedge V1(S,N'',Y) \wedge V1(S',N,Y')$

Containment checking between the query resulting from the candidate expansion :

$\text{Exp}(Q1)(S,D) \leftarrow \text{registered}(S,C',Y1) \wedge \text{enrolled}(S,D) \wedge \text{registered}(S,C1,Y) \wedge \text{course}(C1,N'')$
 $\wedge N' \geq 500 \wedge Y \geq 1992 \wedge \text{registered}(S',C2,Y') \wedge \text{course}(C2,N) \wedge N \geq 500 \wedge Y' \geq 1992$

and the initial query :

$q(S,D) \leftarrow \text{enrolled}(S,D) \wedge \text{registered}(S,C,Y) \wedge \text{course}(C,N) \wedge N \geq 300 \wedge Y \geq 1995$

- The variables of exp(Q1) are « frozen » :

$S \rightarrow a1, C' \rightarrow a2, Y1 \rightarrow a3, D \rightarrow a4, C1 \rightarrow a5, Y \rightarrow a6, N'' \rightarrow a7, S' \rightarrow a8, C2 \rightarrow a9, Y' \rightarrow a10,$
 $N \rightarrow a11$

- q(S,D) is evaluated against the resulting database :

FAILURE

SUCCESS if we add the comparison atoms : $N=N''$ and $Y \geq 1995$

● ● ● | Minicon : optimization of the bucket algorithm

- Query containment checking is avoided by a more elaborate verification of the atoms to add to the buckets.
 - When the definition of a view V contains an atom g' such that : $\sigma(g') = g$
 - If an existential variable Y of g also appears in other atoms $g1, g2, \dots, gk$ of the query
 - if $Y' = \sigma(Y)$ is also existential in the view definition
 - $\sigma(V)$ is added to $\text{Bucket}(g)$ only if $g1, g2, \dots, gk$ are also covered by the definition of $\sigma(V)$



Illustration by example

$V4(X) :- cite(X,Y) \wedge cite(Y,X)$

$V5(X,Y) :- sameTopic(X,Y)$

$V6(X,Y) :- cite(X,Z) \wedge cite(Z,Y) \wedge sameTopic(X,Z)$

Query : $Q(U) :- cite(U,V) \wedge cite(V,U) \wedge sameTopic(U,V)$

Bucket (cite(U,V)) ?

- **V4(U) is not added** because $sameTopic(U,V)$ is not covered by the body of $V4(U)$

- **V6 ?**

$\sigma(X)=U$ and $\sigma(Z)=V$

covering of $cite(V,U)$ by the body of $V6(U,Y) \Rightarrow \sigma(Y)=U$

covering of $sameTopic(U,V)$ by the body of $\sigma(V6(X,Y))$? yes

Bucket(cite(U,V)) = {V6(U,U)}

cover(V6(U,U)) = {cite(U,V), cite(V,U), sameTopic(U,V)}



Advantages of Minicon

- The rewritings are directly obtained by conjuncting the view-atoms of each bucket which have pairwise pairwise disjoint covering.
- Results
 - theoretical :
 - Same worst-case complexity as the Bucket algorithm (exponential in the size of the query)
 - experimental :
 - Scales up much better than the Bucket algorithm in presence of many views.

● ● ● | Summary

- When queries and views are (union of) conjunctive queries over a simple relational schema, the number of maximal conjunctive rewritings is finite and several algorithms exist for computing them.
- No guarantee that this property holds when some knowledge is added to the global schema or to the views
 - to express some binding access patterns to the sources
 - when queries are expressed w.r.t an ontology (schema with axioms or constraints)

● ● ● | Binding access patterns

- Some sources have some constraints to access to their content
 - to get movie titles from Internet Movie Database, we have to provide the name of an actor or of the director
 - for each view $v(\text{Att}_1, \dots, \text{Att}_n)$ describing a source, for each attribute Att_i we have to indicate if it must be instantiated for the source to provide answers
 - a simple way of modeling this knowledge: associate to each view a word of size n of b or f
 - a b in position i indicates that the attribute Att_i must be instantiated for the source to return answers



Example

Source 1: returns papers published in AAAI conferences

PaperSource^f(X) : AAAIpapers(X)

Source 2: returns papers cited by a given paper that must be provided as input

CitationSource^{bf}(X,Y) : Cites(X,Y)

Source3: for a given paper provided as input, answers if it has been awarded or not

AwardSource^b(X) : AwardPaper(X)

Query : find all the awarded papers

Q(X) ← AwardPaper(X)

- $Q_v(X) \leftarrow \text{AwardSource}(X)$ is not an executable rewriting
- $Q'_v(X) \leftarrow \text{PaperSource}(X) \wedge \text{AwardSource}(X)$ is an executable rewriting
- $Q''_v(X) \leftarrow \text{PaperSource}(V) \wedge \text{CitationSource}(V,X_1) \wedge \dots \wedge \text{CitationSource}(X_n,X) \wedge \text{AwardSource}(X)$

is also an executable rewriting




The problem

- Infinite number of maximal conjunctive rewritings
- A maximal rewriting in Datalog :


Papers(X) ← PaperSource(X)

Papers(X) ← Papers(Y) ∧ CitationSource(Y,X)

Q_v(X) ← Papers(X) ∧ AwardSource(X)



Answering queries over ontologies by rewriting



Ontologies

- vocabulaires structurés
 - noms de concepts/classes
 - noms de propriétés
- définis à l'aide de langages formels
 - pour la définition et typage des concepts et des propriétés
 - permettant de faire des inférences fondées sur une sémantique logique
- OWL :Ontology Web Language
 - fondé sur les Logiques de Description

● ● ● | Logiques de Description

- Fragments décidables de la logique du premier ordre avec égalité
- Algorithmes de raisonnement avec une étude approfondie de leur complexité
- Des systèmes implémentés (RACER, PELLET)

● ● ● | Les constructeurs de classes

Constructor	DL Syntax	Example	Modal Syntax
intersectionOf	$C_1 \sqcap \dots \sqcap C_n$	Human \sqcap Male	$C_1 \wedge \dots \wedge C_n$
unionOf	$C_1 \sqcup \dots \sqcup C_n$	Doctor \sqcup Lawyer	$C_1 \vee \dots \vee C_n$
complementOf	$\neg C$	\neg Male	$\neg C$
oneOf	$\{x_1\} \sqcup \dots \sqcup \{x_n\}$	{john} \sqcup {mary}	$x_1 \vee \dots \vee x_n$
allValuesFrom	$\forall P.C$	\forall hasChild.Doctor	$[P]C$
someValuesFrom	$\exists P.C$	\exists hasChild.Lawyer	$\langle P \rangle C$
maxCardinality	$\leq nP$	≤ 1 hasChild	$[P]_{n+1}$
minCardinality	$\geq nP$	≥ 2 hasChild	$\langle P \rangle_n$

Les constructeurs de classes peuvent être imbriqués

- $\text{Person} \sqcap \exists \text{haschild} (\forall \text{haschild}.\text{Doctor})$



Les contraintes (axiomes)

Axiom	DL Syntax	Example
subClassOf	$C_1 \sqsubseteq C_2$	Human \sqsubseteq Animal \sqcap Biped
equivalentClass	$C_1 \equiv C_2$	Man \equiv Human \sqcap Male
disjointWith	$C_1 \sqsubseteq \neg C_2$	Male $\sqsubseteq \neg$ Female
sameIndividualAs	$\{x_1\} \equiv \{x_2\}$	{President_Bush} \equiv {G_W_Bush}
differentFrom	$\{x_1\} \sqsubseteq \neg\{x_2\}$	{john} $\sqsubseteq \neg$ {peter}
subPropertyOf	$P_1 \sqsubseteq P_2$	hasDaughter \sqsubseteq hasChild
equivalentProperty	$P_1 \equiv P_2$	cost \equiv price
inverseOf	$P_1 \equiv P_2^-$	hasChild \equiv hasParent ⁻
transitiveProperty	$P^+ \sqsubseteq P$	ancestor ⁺ \sqsubseteq ancestor
functionalProperty	$\top \sqsubseteq \leq 1P$	$\top \sqsubseteq \leq 1$ hasMother
inverseFunctionalProperty	$\top \sqsubseteq \leq 1P^-$	$\top \sqsubseteq \leq 1$ hasSSN ⁻



FOL axiomatization

$$(C_1 \cap C_2)(X) \equiv C_1(X) \wedge C_2(X)$$

$$(C_1 \cup C_2)(X) \equiv C_1(X) \vee C_2(X)$$

$$(\forall R.C)(X) \equiv \forall Y (R(X, Y) \Rightarrow C(Y))$$

$$(\exists R.C)(X) \equiv \exists Y (R(X, Y) \wedge C(Y))$$

$$(\geq n R)(X) \equiv \exists Y_1 \dots Y_n (R(X, Y_1) \wedge \dots \wedge R(X, Y_n) \wedge \bigwedge_{\{i,j/i \neq j\}} Y_i \neq Y_j)$$

$$(\leq n R)(X) \equiv \exists Y_1 \dots Y_{n+1} ((R(X, Y_1) \wedge \dots \wedge R(X, Y_{n+1})) \Rightarrow \bigvee_{\{i,j/i \neq j\}} Y_i = Y_j)$$

● ● ● | Problèmes d'inférence étudiés

- Test de subsomption
 - Entre 2 expressions de concepts
 - Entre 2 expressions de concepts modulo un ensemble de contraintes définies dans une Tbox T:

$$T \models C1 \subseteq C2 ?$$
- Test d'appartenance d'une instance à un concept
 - Étant donné un concept C, une constante a, une Tbox T, une Abox A (un ensemble de faits de la forme A(b) et P(b,c))

$$T \cup A \models C(a) ?$$
- Nombreux résultats de décidabilité et de complexité en fonction des constructeurs et des axiomes autorisés dans T

● ● ● | Quelques résultats de complexité

Constructeurs	Complexité de la subsomption
ALN (\cap , forall, atleast, atmost)	P
ALE (\cap , forall, exists)	NP-complet
ALNE (\cap , forall, atleast, atmost, exists)	NP-complet
ALN + conjonction de rôles	co-NP-hard
ALC (\cap , forall, not)	Pspace-complet
ALCNR	Pspace-complet

● ● ● | La famille DL-Lite

- Restriction sur les constructeurs et les contraintes autorisées pour une complexité polynômiale du calcul des réponses à une requête conjonctive
- Approche suivie
 - Trouver des logiques de description (ou des langages logiques) pour lesquels le calcul de réponses est réductible à l'évaluation de requêtes sur une BD relationnelle simple
 - Condition nécessaire: la reconnaissance qu'un tuple est une réponse est LOGSPACE (et donc P)

● ● ● | DL-Lite : constructeurs et axiomes

- DL-Lite_{core} :
un ensemble d'axiomes d'inclusion $B \subseteq C$ où B et C sont des expressions de concept sur des concepts atomiques A et des rôles atomiques P de la forme

$$B \rightarrow A \mid \exists R \quad R \rightarrow P \mid P^- \quad C \rightarrow B \mid \neg B$$
- DL-Lite_R :
DL-Lite_{core} + axiomes d'inclusion de rôles $R \subseteq E$ où E est de la forme R ou $\neg R$
- DL-Lite_F :
DL-Lite_{core} + axiomes de fonctionnalité : (funct R)

Expressivité de DL-Lite

- Capture les principales contraintes utilisées en bases de données et génie logiciel
 - Relation ISA : $A1 \subseteq A2$
 - Disjonction : $A1 \subseteq \neg A2$
 - typage :
 - $\exists P \subseteq A1$ (le premier attribut de P est typé par A1)
 - $\exists P^- \subseteq A2$ (le second attribut de P est typé par A2)
 - Propriétés obligatoires ou interdites :
 - $A \subseteq \exists P$ $A \subseteq \exists P^-$ $A \subseteq \neg \exists P$ $A \subseteq \neg \exists P^-$
 - Restriction de fonctionnalité pour les relations binaires (les rôles)

Exemple

$\text{Professor} \subseteq \exists \text{TeachesTo}$
 $\text{Student} \subseteq \exists \text{HasTutor}$
 $\exists \text{TeachesTo}^- \subseteq \text{Student}$
 $\exists \text{HasTutor}^- \subseteq \text{Professor}$
 $\text{Professor} \subseteq \neg \text{Student}$

$\text{HasTutor}^- \subseteq \text{TeachesTo}$
 (func HasTutor)

DL-Lite_R
 DL-Lite_F

Requêtes

- (Union de) requêtes conjonctives sur des concepts et rôles atomiques

$$q_0(x) \leftarrow \text{TeachesTo}(x,y) \wedge \text{HasTutor}(y,z)$$

- Algorithme de réécriture correct et complet

$$q_1(x) \leftarrow \text{TeachesTo}(x,y) \wedge \text{Student}(y)$$

$$q_2(x) \leftarrow \text{TeachesTo}(x,y) \wedge \text{TeachesTo}(z',y)$$

$$q_3(x) \leftarrow \text{TeachesTo}(x,y')$$


$$q_4(x) \leftarrow \text{Professor}(x)$$

$$q_5(x) \leftarrow \text{HasTutor}(u,x)$$

$$\text{Ans}(q_0, T \cup A) = \bigcup_i \text{Ans}(q_i, A)$$


Résumé

- Répondre (de façon complète) à des requêtes par réécriture n'est « tractable » que pour très peu de Logiques de Description
- La reconnaissance d'instance (et donc le calcul des réponses à une requête) est coNP-hard en complexité des données pour des extensions mineures of $\text{DL-Lite}_{\text{core}}$ telles que :
 - $\text{DL-Lite}_{\text{FR}}$
 - $\forall R.A$ ou $\neg A$ apparaissent en partie gauche d'un axiome d'inclusion
 - $\forall R.A$ ou $\neg A$ apparaissent dans la requête
 - Union de concepts en partie gauche d'un axiome d'inclusion



References

- Querying heterogeneous information sources using source descriptions. A. Levy, A. Rajaraman, J. Ordille. Proceedings of the Int. Conference on Very Large Data bases (VLDB 96)
- Information Integration Using Logical Views, J. Ullman, Proceedings ICDT '97
- Logic-based techniques in data integration. A. Halevy, in Logic Based Artificial Intelligence, Ed. Jack Minker, Kluwer Publishers, 2000.
- Query Containment for Data Integration Systems, T. Millstein, A. Levy, M. Friedman, Proceedings PODS 2000
- Complexity of Answering Queries Using Materialized Views, S. Abiteboul, O. Duschka, Proceedings PODS 1998
- Theory of Answering Queries Using Views, A. Halevy, Proceedings of SIGMOD 2000
- The Use of CARIN Language and Algorithms for Information Integration: the PICSEL System, F. Goasdoué, V. Lattès, -C. Rousset. International Journal of Cooperative Systems, Vol 9, Number 4 (2000)
- Minicon: a scalable algorithm for answering queries using views. R. Pottinger, A. Halevy, VLDB Journal, Volume 10 (2-3), 2001.
- Query Rewriting and Answering under Constraints in Data Integration Systems, A. Cali, D. Lembo, and R. Rosati, Proceedings of IJCAI 2003
- D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, R. Rosati, DL-Lite: Tractable Description Logics for Ontologies, Proceedings of AAAI 2005.
- D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, R. Rosati, Data Complexity of Query Answering in Description Logics, Proceedings of KR 2006



Logique des prédicats (rappels)

- Vocabulaire :
 - prédicats (ou relations), fonctions, constantes, variables
- Connecteurs logiques :
 - usuels : $\wedge, \vee, \neg, \Rightarrow, \dots$
 - quantificateurs : \forall, \exists
- Règles de bonne formation de formules :
 - Termes
 - de base (variables, constantes)
 - fonctionnels (application de fonctions à des termes)
 - Formules atomiques : $R(t_1, t_2, \dots, t_n)$
 - Quantification de variables :
 - Variables libres versus liées
 - Formules :
 - de base : formules atomiques
 - application des connecteurs logiques à des formules



Exemples

$\forall x R(x,a)$ R: prédicat, a: constante

$\exists y \forall x R(x,y)$ x,y: variables liées

$\forall x R(f(x),x)$ f: fonction

$\forall x \forall y \forall z (R(x,y) \wedge R(y,z) \Rightarrow R(x,z))$

$\forall x \exists y \neg R(x,y)$

$\exists y \neg R(y,x)$ x: variable libre



Interprétations et modèles de formules

- Une interprétation I de φ :
 - un domaine d'interprétation Δ^I
 - une fonction d'interprétation associant
 - aux constantes a de φ des éléments a^I de Δ^I
 - aux fonctions f (prédicats R) de φ des fonctions f^I de Δ^I dans Δ^I (relations R^I sur Δ^I).
 - des règles d'interprétations pour interpréter une formule à partir de l'interprétation de ses sous-formules
 - Interprétation d'une formule close (dont toutes les variables sont quantifiées) : valeur de vérité (vrai ou faux)
 - Interprétation d'une formule ouverte (n variables libres) : une relation n-aire sur Δ^I
- Un modèle d'une formule φ : une interprétation I
 - $\varphi^I = \text{vrai}$ (si φ est close)
 - $\varphi^I \neq \emptyset$ (si φ est ouverte)

● ● ● | Règles d'interprétation des quantificateurs

- Soit φ une formule close de la forme $\forall x \psi$
 $[\forall x \psi(x)]^I = \text{vrai}$ ssi pour tout $e \in \Delta^I$, $\psi^I(e)$ est vrai
- Soit φ une formule close de la forme $\exists x \psi$
 $[\exists x \psi(x)]^I = \text{vrai}$ ssi il existe $e \in \Delta^I$, $\psi^I(e)$ est vrai
- Soit $\varphi(x_1, \dots, x_n)$ une formule avec n variables libres
 $[\varphi(x_1, \dots, x_n)]^I = \{(e_1, \dots, e_n) \in \Delta^I \times \dots \times \Delta^I / \varphi^I(e_1, \dots, e_n) \text{ est vrai}\}$

● ● ● | Exemples (suite)

$\Delta^I : \mathbb{N}$ (ensemble des entiers naturels)

$a^I = 0$; $f^I = \text{succ}$; $R^I = \geq$

$[\forall x R(x, a)]^I = \text{vrai}$ (car pour tout $e \in \mathbb{N}$, $e \geq 0$)

$[\exists y \forall x R(x, y)]^I = \text{vrai}$ (car il existe $e \in \mathbb{N}$ ($e=0$) tel que pour tout $e' \in \mathbb{N}$ $e' \geq e$)

$[\forall x R(f(x), x)]^I = \text{vrai}$ (car pour tout $e \in \mathbb{N}$, $\text{succ}(e) \geq e$)

$[\forall x \forall y \forall z (R(x, y) \wedge R(y, z) \Rightarrow R(x, z))]^I = \text{vrai}$ car pour tout $e_1, e_2, e_3 \in \mathbb{N}$, si $e_1 \geq e_2$ et $e_2 \geq e_3$ alors $e_1 \geq e_3$)

$[\forall x \exists y \neg R(x, y)]^I = \text{faux}$ (car $[\exists y \forall x R(x, y)]^I = \text{vrai}$)

$[\exists y \neg R(y, x)]^I = \{e \in \mathbb{N} / \text{il existe } e' \in \mathbb{N} \text{ tel que } e' < e\}$
 $= \mathbb{N}^*$



Problèmes d'inférence

- Satisfiabilité :
 - existence d'un modèle
- Déterminer si une formule φ est conséquence logique d'un ensemble de formules $\varphi_1, \dots, \varphi_n$:
 - tout modèle de $\varphi_1, \dots, \varphi_n$ est modèle de φ
 - $\varphi_1, \dots, \varphi_n, \neg\varphi$
- Problèmes semi-décidables :
 - on ne connaît pas d'algorithmes qui décident si une formule quelconque est satisfiable ou non satisfiable (est conséquence logique ou non d'un ensemble de formules)
- L'ensemble des interprétations d'une formule quelconque est infini



Fragments décidables

Des sous-ensembles de formules pour lesquels les problèmes d'inférence sont décidables

- Les règles logiques sans fonction :
 - Fondements logiques des systèmes experts avec variables
 - Datalog /BD déductives
- Les logiques de description :
 - Fondements logiques des langages à base de classes ou de concepts