



---

# Ontology-Based Data Access

HAI824 – Traitement sémantique des données  
ML Mugnier

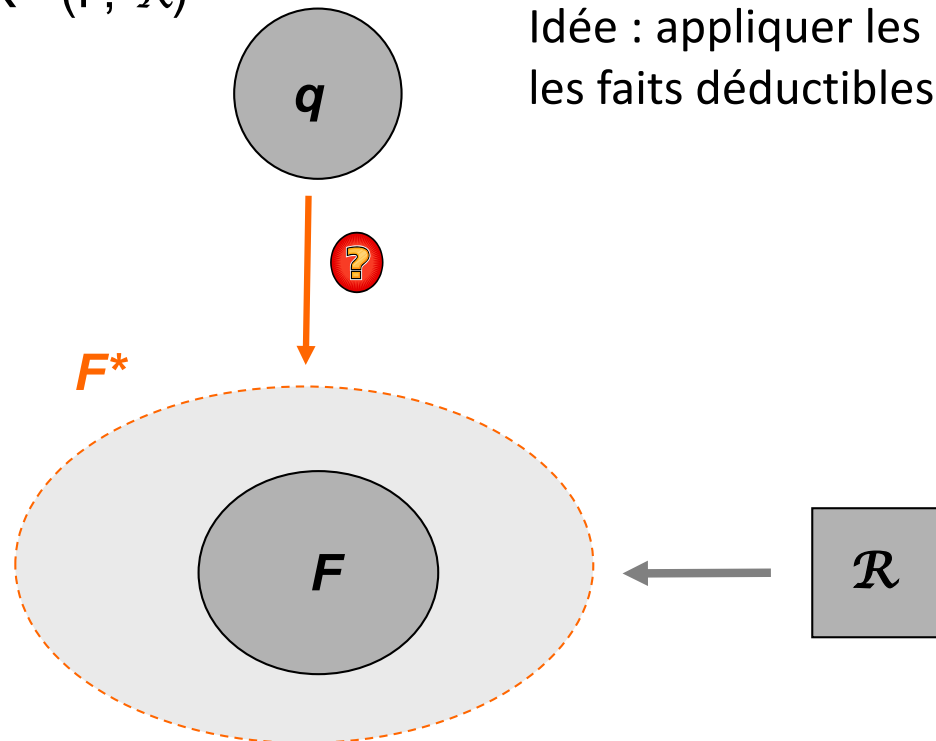
# PLAN

---

- Deux grandes approches pour l'interrogation d'une base de connaissances
  - **Saturation** de la base de faits (chaînage avant)
  - **Récriture** de requête (forme de chaînage arrière)
- Si la base de faits provient d'une base de données
  - **Matérialisation** ou **virtualisation** de la base de faits ?
- Intégration classique de données : les mappings à la rescousse
  - Matérialisation (**warehouse, Extract-Transform-Load**) ou virtualisation (**médiation**) ?
- **Ontology-Based Data Access**  
ou le mariage de l'intégration de données et de la représentation de connaissances

# SATURATION DE LA BASE DE FAITS (CHAÎNAGE AVANT) « bottom-up »

$K = (F, \mathcal{R})$



Idée : appliquer les règles sur les faits pour obtenir tous les faits déductibles de la base de connaissance

$F^*$  (la saturation de  $F$  par  $\mathcal{R}$ )  
est contenue dans tous les modèles de  $K$

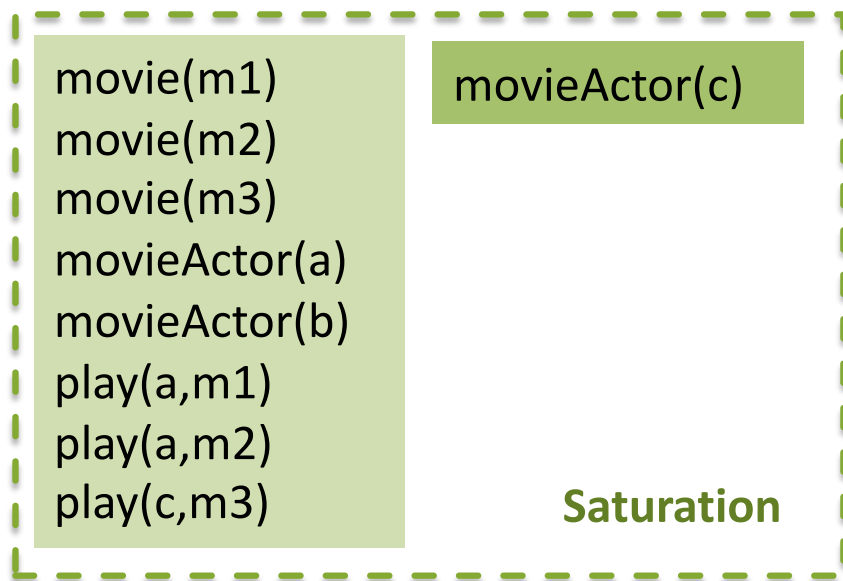
Si  $q$  est une (U)CQ :

$K \models q$  **ssi**  $F^* \models q$   
**ssi** il y a un homomorphisme de  $q$  dans  $F^*$

- + raisonnement effectué **offline**,  
l'interrogation elle-même ne considère qu'une base de faits (saturée)
- **volume** de la base de faits saturée  
inadapté si la base de faits **change** fréquemment  
(en particulier si on doit recalculer la saturation à chaque requête)

# EXAMPLE (SATURATION)

$\text{play}(x,y) \wedge \text{movie}(y) \rightarrow \text{movieActor}(x)$



$q(x) = \text{movieActor}(x)$

« *find all movie actors* »

$x = a$   
 $x = b$   
 $x = c$

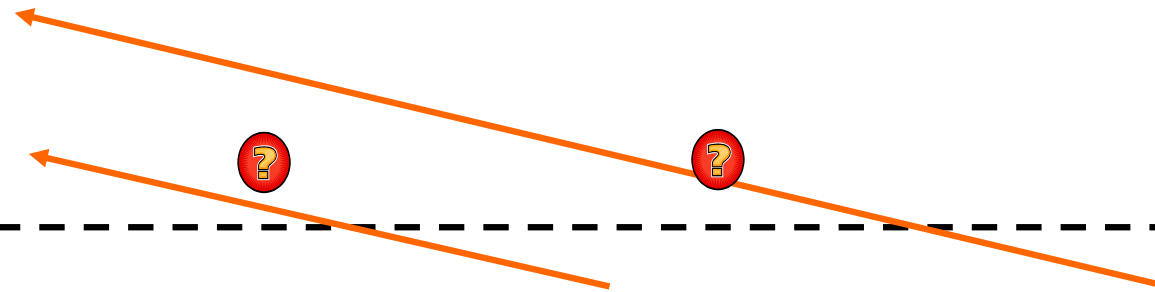
# EXEMPLE (RÉÉCRITURE DE REQUÊTE)

$\text{play}(x,y) \wedge \text{movie}(y) \rightarrow \text{movieActor}(x)$

movie(m1)  
movie(m2)  
movie(m3)  
movieActor(a)  
movieActor(b)  
play(a,m1)  
play(a,m2)  
play(c,m3)

$q(x) = \text{movieActor}(x)$

« find all movie actors »



Rewriting( $q$ ) =  $\exists y (\text{movie}(y) \wedge \text{play}(x, y)) \vee \text{movieActor}(x)$

{ ans(x) :- movie(y), play(x,y)  
ans(x) :- movieActor(x) }

**Réécriture de requête**

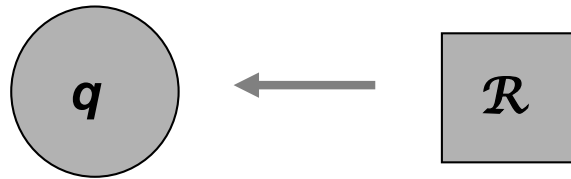
x = a    y = m1  
x = a    y = m2  
x = c    y = m3

x = a  
x = b

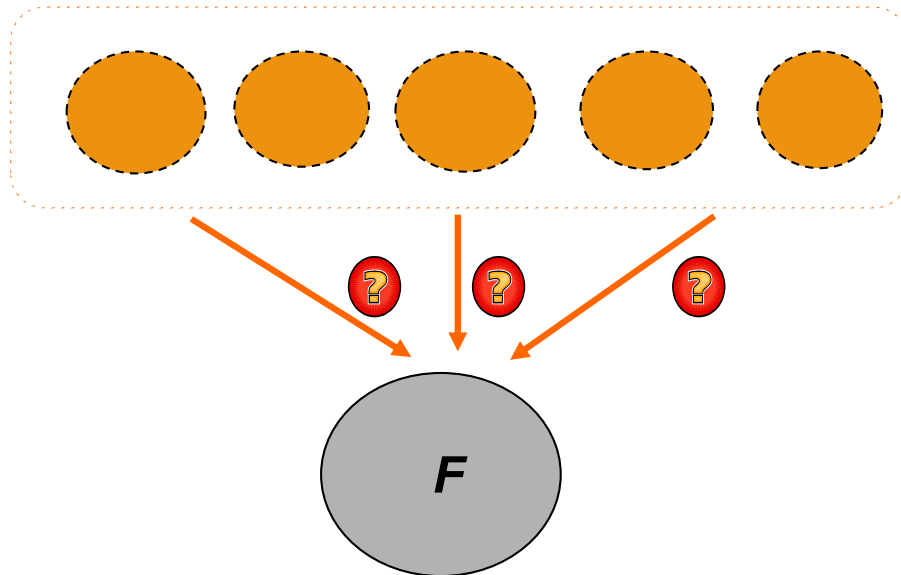
x = a  
x = b  
x = c

# RÉÉCRITURE DE REQUÊTE (SEMI CHAINAGE ARRIÈRE) « top-down »

$K = (F, \mathcal{R})$



Idée : chercher toutes les façons de répondre à la requête en « remontant » le long des règles



$rewriting(q)$

Technique essentiellement applicable aux CQ (et UCQ)

Réécriture en un ensemble de CQ vu comme une UCQ

Pour une CQ  $q$  :  
 $K \models q$  ssi  $F \models rewriting(q, \mathcal{R})$

La réécriture de la requête est indépendante de toute base de faits

- + processus **indépendant des changements** dans la base de faits
- raisonnement **pendant** l'interrogation ( $\rightarrow$  plus lente)  
la réécriture peut conduire à une requête **complexe**  
(très grande et d'une forme inhabituelle pour les SGBD)

# LIEN AVEC LE CHAÎNAGE ARRIÈRE ?

Rappels cours sur les règles positives  
en logique des propositions

## Chaînage avant :

BF = {A, B, C}  
BR = {  
     $R_1 : A \wedge B \rightarrow E$   
     $R_2 : C \wedge E \rightarrow D$   
}

$BF^* = BF \cup \{E, D\}$

## Chaînage arrière :

- principe : prouver un but (atome) en «remontant» le long des règles
- le but initial est prouvé lorsqu'on arrive à une liste de buts vide

But initial : <b>D</b> ?	liste de buts { <b>D</b> }
Avec $R_2$ :	{ <b>C</b> , <b>E</b> }
Avec le fait C (vu comme $\rightarrow C$ ) :	{ <b>E</b> }
Avec $R_1$ :	{ <b>A</b> , <b>B</b> }
Avec le fait A	{ <b>B</b> }
Avec le fait B	{}

# LIENS AVEC LE CHAÎNAGE ARRIÈRE ?

Qu'est-ce qui change ?

## 1. Décomposition du chaînage arrière en **deux étapes** :

- Réécriture avec les « vraies » règles
- puis recherche dans la base de faits

BF = {A, B, C}

BR = {  
     $R_1 : A \wedge B \rightarrow E$   
     $R_2 : C \wedge E \rightarrow D$   
}

Partant de la « requête »  $q_0 = D$  :

- $q_0$  se réécrit avec  $R_2$  : on obtient  $q_1 = C \wedge E$
- $q_1$  se réécrit avec  $R_1$  : on obtient  $q_2 = C \wedge A \wedge B$

D'où la requête réécrite :

$q_0 \vee q_1 \vee q_2 = D \vee (C \wedge E) \vee (C \wedge A \wedge B)$

## 2. En logique du premier ordre, on n'a pas une simple égalité entre un atome d'une requête et une tête de règle, mais une **unification** de ces atomes

R:  $p_2(x) \rightarrow p_1(x)$        $Q() = p_1(a) ?$       on réécrit :  $Q'() = p_2(a)$

R:  $p_2(x,y) \rightarrow p_1(x,x)$        $Q() = \exists u \exists v. p_1(a,u) \wedge p_3(u,v)$   
on réécrit :  $Q'() = p_2(a,y) \wedge p_3(a,v)$



# QUERY REWRITING WITH DATALOG RULES (1)

R1:  $p(x_1, y_1), p(y_1, z_1) \rightarrow gp(x_1, z_1)$

R2:  $mo(x_2, y_2) \rightarrow p(x_2, y_2)$

R3:  $fa(x_3, y_3) \rightarrow p(x_3, y_3)$

$q(x) = gp(x, a)$

A **unifier**  $u$  of atoms  $A$  and  $B$  is a substitution (of variables) such that  $u(A) = u(B)$

A **most general unifier** (mgu) of  $A$  and  $B$  is a unifier  $u$  of  $A$  and  $B$  such that every other unifier of  $A$  and  $B$  can be written as  $s \circ u$  where  $s$  is a substitution

Basic step: computation of a **direct rewriting** of a CQ  $q$  with a rule  $R$   
*[Important: we always assume that  $q$  and  $R$  have disjoint sets of variables; if it is not the case, rename some variables with fresh ones]*

1. look for a mgu  $u$  of an atom  $A$  in  $q$  and the head of  $R$
2. the direct rewriting of  $q$  with  $R$  according to  $u$  is

$$rew(q, R, u) = u(q \setminus A) \cup u(body(R))$$

The **naive rewriting** of  $q$  with  $\mathcal{R}$  is the set of all CQs obtained by a (possibly empty) sequence of direct rewritings starting from  $q$  and using the rules in  $\mathcal{R}$

It is logically seen as the **disjunction** of all the CQs it contains

## QUERY REWRITING WITH DATALOG RULES (2)

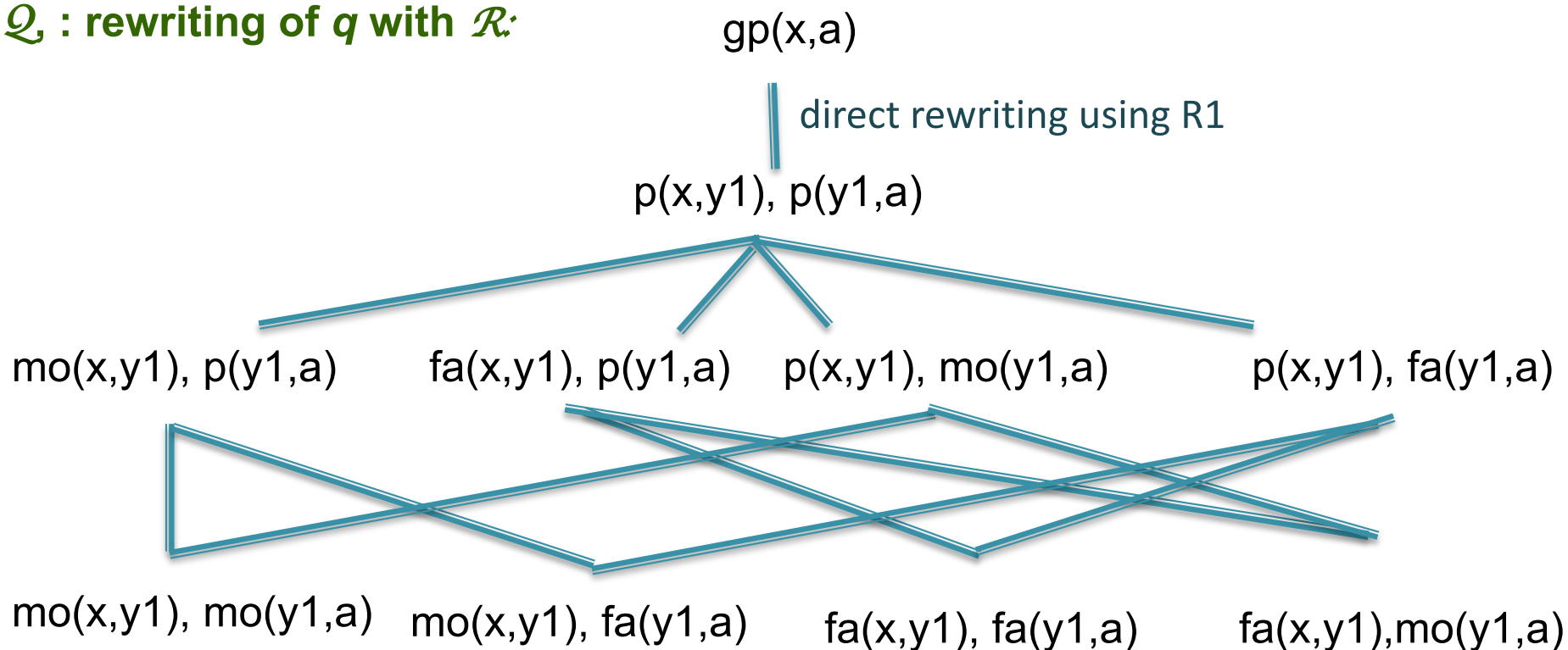
R1:  $p(x1,y1), p(y1,z1) \rightarrow gp(x1,z1)$

R2:  $mo(x2,y2) \rightarrow p(x2,y2)$

R3:  $fa(x3,y3) \rightarrow p(x3,y3)$

$q(x) = gp(x,a)$

$\mathcal{Q}$  : rewriting of  $q$  with  $\mathcal{R}$ :



Let  $q$  be a Boolean CQ and  $\mathcal{Q}$  be a rewriting of  $q$  with  $\mathcal{R}$

For any factbase  $F$ ,  $F, \mathcal{R} \models q$  iff  $F \models \mathcal{Q}$  (where  $\mathcal{Q}$  is seen as a disjunction of CQs)  
 iff there is  $q_i$  in  $\mathcal{Q}$  such that  $F \models q_i$

# ALGORITHME (BASIQUE) DE RÉÉCRITURE

**Réécriture**( $q, \mathcal{R}$ ) : retourne l'ensemble des réécritures de  $q$  avec  $\mathcal{R}$

**Début**

Résultat  $\leftarrow \emptyset$

AExplorer  $\leftarrow \{q\}$  // *requêtes pas encore réécrites*

**Tant que** AExplorer  $\neq \emptyset$

Retirer  $q_i$  de AExplorer

Ajouter  $q_i$  dans Résultat

**Pour** toute règle  $R_j \in \mathcal{R}$

et tout unificateur le plus général  $u$  de tête( $R_j$ ) avec un atome de  $q_i$

Calculer  $q_k = \text{rewriting}(q_i, R_j, u)$

Ajouter  $q_k$  dans AExplorer si  $q_k$  n'est pas déjà dans AExplorer  $\cup$  Résultat  
(à un isomorphisme près)

**Fin Pour**

**Fin Tant que**

Retourner Résultat

**Fin**

# HOMOMORPHISME / ISOMORPHISME

Soient  $S1$  et  $S2$  deux ensembles d'atomes

- Un **homomorphisme**  $h$  de  $S1$  dans  $S2$  est une **application** des variables de  $S1$  dans les **termes** de  $S2$  telle que  $h(S1) \subseteq S2$
- Un **isomorphisme**  $h$  de  $S1$  dans  $S2$  est une **bijection** des variables de  $S1$  dans les **variables** de  $S2$  telle que  $h(S1) = S2$ . On dit que  $S1$  et  $S2$  sont isomorphes

Un isomorphisme est donc un homomorphisme, mais l'inverse est faux

**Propriété** : un homomorphisme  $h$  de  $S1$  dans  $S2$  est un isomorphisme  
ssi  $h^{-1}$  est un homomorphisme de  $S2$  dans  $S1$

**Remarque** : on peut avoir un homomorphisme  $h1$  de  $S1$  dans  $S2$  et  
un homomorphisme  $h2$  de  $S2$  dans  $S1$  sans que  $S1$  et  $S2$  soient isomorphes

$S1 = \{ p(x,y), p(x,z) \}$   
 $S2 = \{ p(u,v) \}$

## REMOVING REDUNDANCES IN THE REWRITING (« MINIMIZATION »)

Let  $Q = q_1 \vee q_2$ , where  $q_1$  maps by « query-homomorphism » to  $q_2$

i.e.: there is a homomorphism  $h$  from  $q_1(x_1 \dots x_n)$  to  $q_2(y_1 \dots y_n)$  such that  $h(x_i) = y_i$  for all  $i$

Then  $q_2$  is useless because every answer to  $q_2$  is also an answer to  $q_1$  (notation :  $q_2 \sqsubseteq q_1$ )

Example:

$$q_1(x) = \exists y. \text{parent}(x, y)$$

$$q_2(x) = \exists y. \text{parent}(y, x)$$

$$q_3(x) = \exists y. \text{parent}(x, y) \wedge \text{girl}(y)$$

$$Q = q_1 \vee q_2 \vee q_3$$

$q_3$  is useless, but not  $q_2$

$Q$  is equivalent to  $q_1 \vee q_2$

# REMOVING REDUNDANCES IN THE REWRITING (« MINIMIZATION »)

We can keep only a ``cover'' of the rewriting  $Q$

A **cover** of  $Q$  is a subset  $Q_c$  of  $Q$  such that:

1. for any  $q$  in  $Q$ , there is  $q'$  in  $Q_c$  such that  $q \sqsubseteq q'$
2. elements of  $Q_c$  are pairwise incomparable with respect to  $\sqsubseteq$

We can also suppress redundancies **inside** each conjunctive query  $q$   
(computation of the **core** of  $q$ : the minimal subset of  $q$  that is equivalent to  $q$  by **query homomorphism**)

**Minimal UCQ** : no CQ can be made smaller or removed (without losing equivalence)

If the naive rewriting of  $q$  with  $\mathcal{R}$  is equivalent to a **finite** set of CQs  
we can find it by taking one of its **covers**  
(then minimizing each CQ if we want a minimal UCQ)

# QUERY REWRITING CAN BE INFINITE !

$R = \text{friend}(u,v) \wedge \text{friend}(v,w) \rightarrow \text{friend}(u,w)$

No finite minimal rewriting

$q = \text{friend}(\text{Giorgos}, \text{Maria})$

$q_1 = \text{friend}(\text{Giorgos}, v_0) \wedge \text{friend}(v_0, \text{Maria})$

$q_2 = \text{friend}(\text{Giorgos}, v_1) \wedge \text{friend}(v_1, v_0) \wedge \text{friend}(v_0, \text{Maria})$

$q_{2'} = \text{friend}(\text{Giorgos}, v_0) \wedge \text{friend}(v_0, v_1) \wedge \text{friend}(v_1, \text{Maria})$

$q_2$  and  $q_{2'}$   
are isomorphic

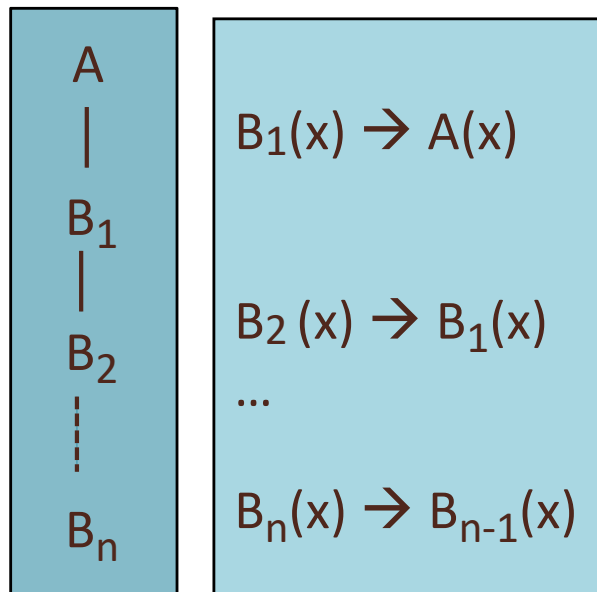
$q_3 = \text{friend}(\text{Giorgos}, v_2) \wedge \text{friend}(v_2, v_1) \wedge \text{friend}(v_1, v_0) \wedge \text{friend}(v_0, \text{Maria})$

*Etc.*

- If we know the number of facts we can bound the size of a maximal « path of friends », hence number of atoms in a rewriting, (however, this will result in very large rewritings!)

# MÊME QUAND LA RÉÉCRITURE EST FINIE, ELLE PEUT ÊTRE (TRÈS) GRANDE !

- Intérêt de l'approche : indépendance vis à vis de la base de faits
- Cependant, la taille de la réécriture peut être prohibitive en pratique



$$q = A(x_1) \wedge \dots \wedge A(x_k)$$

UCQ produite :  $(n+1)^k$  CQ,  $k \times (n+1)^k$  atomes

Ce n'est pas un « pire des cas » théorique :  
se produit souvent en pratique

→ Réécriture en des formes de requêtes **plus compactes** (mais pas toujours plus facile à évaluer !)

$$(A(x_1) \vee B_1(x_1) \dots \vee B_n(x_1)) \wedge \dots \wedge (A(x_k) \vee B_1(x_k) \dots \vee B_n(x_k))$$

$k \times (n+1)$  atomes

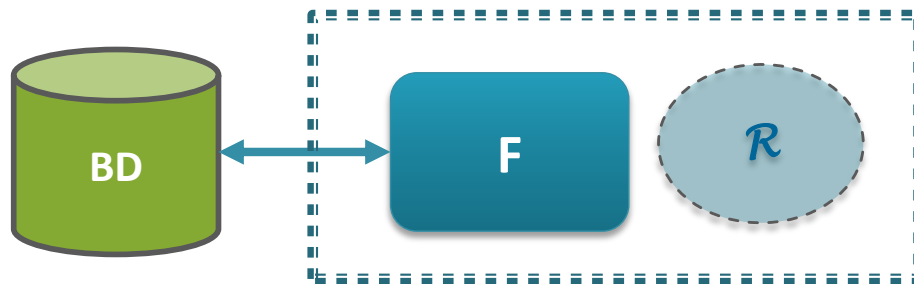
→ Utilisation **combinée** de saturation et de réécriture de requête



# QUAND LA BASE DE FAITS EST LA TRADUCTION D'UNE BASE DE DONNÉES

BD relationnelle,  
triplets RDFS

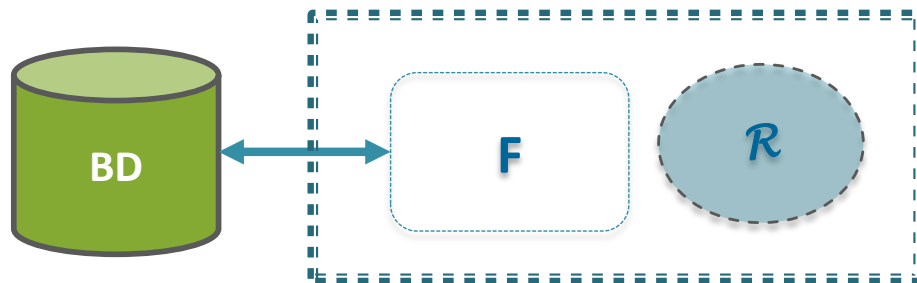
La base de faits peut être **matérialisée**



Interrogation avec une requête  $q$  :

- soit :
  - calcul de la saturation de  $F$  avec  $\mathcal{R}$
  - évaluation de  $q$  sur  $F^*$
- soit :

ou elle peut rester **virtuelle**



Interrogation avec une requête  $q$  :

- calcul de la réécriture de  $q$  avec  $\mathcal{R}$
- traduction de cette réécriture en une requête sur la BD source
- évaluation de cette requête sur la BD  
[Ceci permet de **tirer parti des optimisations du SGBD**]