




Représentation de connaissances



Interrogation de bases de connaissances

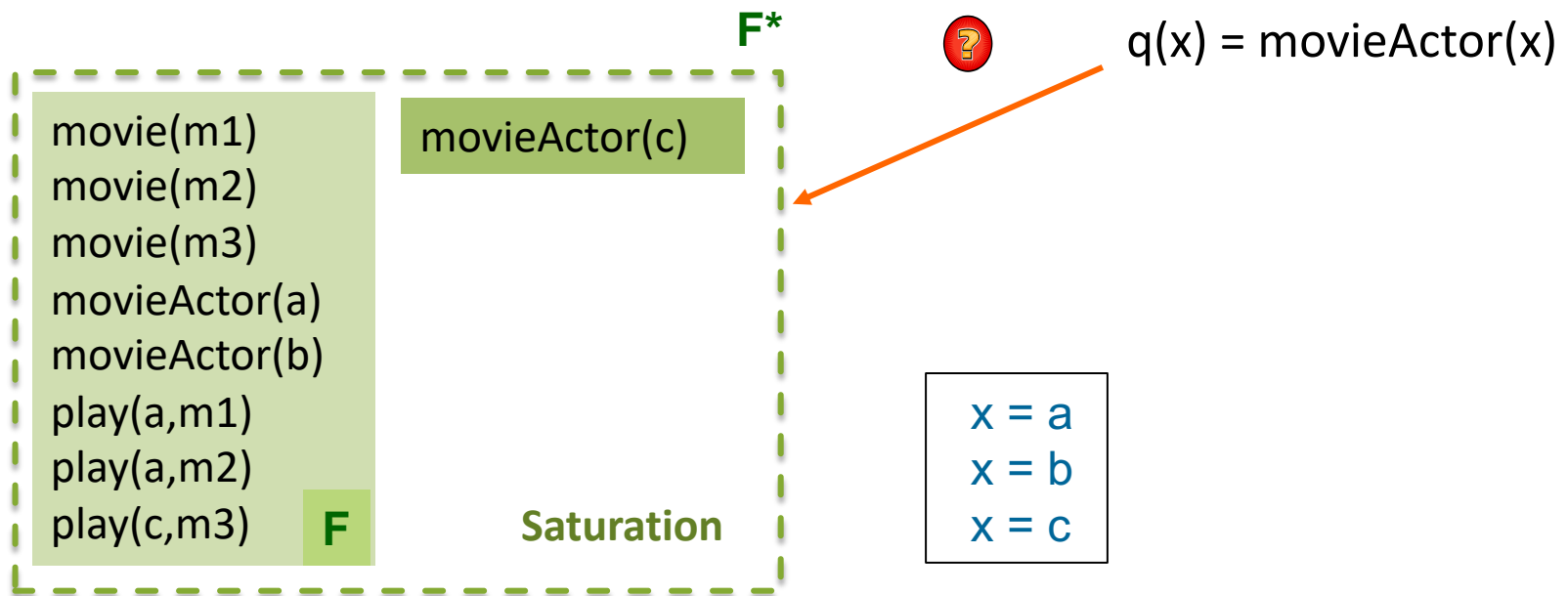
avec règles Datalog

- 
- Chaînage avant (Forward Chaining , « matérialisation »)
 - Chaînage arrière classique (Backward Chaining)
 - Variante de BC par réécriture de requête

EXAMPLE (FORWARD CHAINING)

$\mathcal{R} = \{ R: \text{play}(x,y), \text{movie}(y) \rightarrow \text{movieActor}(x) \}$

« find all movie actors »



Query Answering on (F, \mathcal{R})

1. **Forward chaining:** compute F^* (can be done offline)
2. Evaluate queries on F^* (based on homomorphisms to F^*)

EXAMPLE (BACKWARD CHAINING - CLASSICAL)

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

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

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

Rewrite queries with the **rules** and the **facts**
until obtaining empty sets of atoms
Compose substitutions to « trace » answers

- an atom of q unified with a fact

unifier $x \rightarrow a$

$\Rightarrow \emptyset$

$x = a$
 $x = b$

unifier $x \rightarrow b$

$\Rightarrow \emptyset$

- an atom of q unified
with head(R):

Renaming of R's variables (to avoid confusion)
R: $\text{play}(x',y'), \text{movie}(y') \rightarrow \text{movieActor}(x')$

unifier $x \rightarrow x'$

$\Rightarrow q_1: \text{play}(x',y'), \text{movie}(y')$

- an atom of q_1 unified with fact $\text{play}(c,m3)$

$x' \rightarrow c, y' \rightarrow m3$

$\Rightarrow q_2: \text{movie}(m3)$

- an atom of q_2 unified with fact $\text{movie}(m3)$

unifier (empty)

$\Rightarrow \emptyset$

$x = c$

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

and similarly with q_1 and the other facts

UNIFIERS

R1: $p(x1, y1), p(y1, z1) \rightarrow gp(x1, z1)$
R2: $mo(x2, y2) \rightarrow p(x2, y2)$
R3: $fa(x3, y3) \rightarrow p(x3, y3)$

$A = gp(x, a)$

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

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

$A = gp(x, a)$

$B = gp(x1, z1)$

$u1:$

$x \rightarrow x1$
 $z1 \rightarrow a$

$u2:$

$x1 \rightarrow x$
 $z1 \rightarrow a$

mgu

Il peut y avoir plusieurs mgu de A et B ,
mais ils sont « équivalents » :
on passe de l'un à l'autre par une substitution s qui
ne fait que renommer bijectivement les variables

$u2 = \{ x1 \rightarrow x \} \circ u1$

$u1 = \{ x \rightarrow x1 \} \circ u2$

$u3:$

$x \rightarrow a$
 $z1 \rightarrow a$
 $x1 \rightarrow a$

$u3 = \{ x1 \rightarrow a \} \circ u1$

DIRECT REWRITING

R1: $p(x1,y1), p(y1,z1) \rightarrow gp(x1,z1)$
R2: $mo(x2,y2) \rightarrow p(x2,y2)$
R3: $fa(x3,y3) \rightarrow p(x3,y3)$

Note:

Facts are seen as rules with an empty body

ex: F1 seen as $\rightarrow F1$

Basic step: computation of a **direct rewriting** of a query q

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

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

$q = gp(x,a)$

$head(R1) = gp(x1,z1)$

$u: x \rightarrow x1, z1 \rightarrow a$

$$rew(q,R1,u) = p(x1,y1), p(y1,a)$$

$q' = gp(x,a), p(x,b)$

$head(R1) = gp(x1,z1)$

$u: x \rightarrow x1, z1 \rightarrow a$

$$rew(q',R1,u) = p(x1,y1), p(y1,a), p(x1,b)$$

BACKWARD CHAINING POUR CQ **BOOLÉENNE**

Rappel : les faits sont vus ici comme des règles

Algorithme BC(K,Q)

// Donnée : $K = (F,R)$ et Q une requête booléenne (liste d'atomes)

// Remarque : on voit un fait $p(a_1, \dots, a_k)$ comme une règle $\rightarrow p(a_1, \dots, a_k)$

// Résultat : vrai si K répond positivement à Q , sinon faux

Début

Si $Q = \emptyset$, retourner vrai

Pour toute règle $R = B \rightarrow H$ de **F** et **R**

telle que premier(Q) s'unifie avec H par un unificateur u

// penser à renommer les variables de R si besoin

pour que les ensembles de variables de Q et R soient disjoints

$Q' \leftarrow u(B) \cup u(\text{reste}(Q))$

Si BC(K,Q') = vrai, retourner vrai

FinPour

Retourner faux

Fin

BACKWARD CHAINING POUR CQ NON BOOLÉENNE

Rappel : les faits sont vus ici comme des règles

// $Q(x_1, \dots, x_k)$ est une CQ vue comme un ensemble d'atomes

// $S = \{(x_1, x_1), \dots, (x_k, x_k)\}$ initialement

Algorithme BC(K, Q, S)

// Donnée : $K = (F, R)$ et Q une CQ, S substitution de $\{x_1, \dots, x_k\}$

// Résultat : Ensemble des $\{(x_1, a_1), \dots, (x_k, a_k) \mid (a_1, \dots, a_k) \text{ est une réponse à } Q \text{ dans } K\}$

Début

$\text{Rep} \leftarrow \emptyset$ // sera le résultat

Pour toute règle $R = B \rightarrow H$ de F et R telle que $\text{premier}(Q)$ s'unifie avec H par u

 // si besoin, renommer les variables de R pour que $\text{variables}(Q) \cap \text{variables}(R) = \emptyset$

$Q' \leftarrow u(B) \cup u(\text{reste}(Q))$

$S \leftarrow u \circ S$ // mise à jour de (x_i, y) par (x_i, z) si $(y, z) \in u$

Si $Q' = \emptyset$, ajouter S à Rep

Sinon $\text{Rep} \leftarrow \text{Rep} \cup \text{BC}(K, Q', S)$

FinPour

Retourner Rep

Fin

EXEMPLE BC

K :

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

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

« Qui sont tous les grand-parents de a ? »

F1: $fa(b, a)$

F2 : $fa(c, b)$

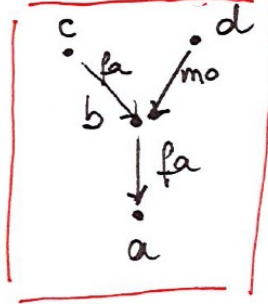
F3 : $mo(d, b)$

En chainage avant

$F^* = \{fa(b, a), fa(c, b), mo(d, b), p(b, a), p(c, b), p(d, b), gp(c, a), gp(d, a)\}$

$Q(K) = Q(F^*) = \{ (c), (d) \}$

$F_1 = fa(b, a)$
 $F_2 = fa(c, b)$
 $F_3 = mo(d, b)$



$$\varphi(x) = \underline{g p(x, a)} \quad S = \{(x, x)\}$$

$$R_1 \mid \mu: x \rightarrow x_1$$

$$z_1 \rightarrow a$$

$$(p(x_1, y_1), p(y_1, a)) \quad S = \{(x, x_1)\}$$

$$R_2 \mid \mu: x_1 \rightarrow x_2$$

$$y_1 \rightarrow y_2$$

$$(mo(x_2, y_2), p(y_2, a)) \quad S = \{(x, x_2)\}$$

$$F_3 \mid \mu: x_2 \rightarrow d$$

$$y_2 \rightarrow b$$

$$(p(b, a)) \quad S = \{(x, d)\}$$

$$R_2 \mid \mu: x_2 \rightarrow b$$

$$y_2 \rightarrow a$$

$$mo(b, a)$$

\times échec

$$R_3 \mid \mu: x_3 \rightarrow b$$

$$y_3 \rightarrow a$$

$$fa(b, a) \quad S = \{(x, d)\}$$

$$F_1$$

$$\emptyset \text{ succès } S = \{(x, d)\}$$

$$Rep = \{(x, d), (x, c)\}$$

$$R_3 \mid \mu: x_1 \rightarrow x_3$$

$$y_1 \rightarrow y_3$$

$$(fa(x_3, y_3), p(y_3, a)) \quad S = \{(x, x_3)\}$$

$$F_1 \mid \mu: x_3 \rightarrow b$$

$$y_3 \rightarrow a$$

$$(p(a, a)) \quad S = \{(x, b)\}$$

$$R_2 \mid \mu: x_2 \rightarrow a$$

$$y_2 \rightarrow a$$

$$mo(a, a)$$

\times

$$R_3 \mid \mu: x_3 \rightarrow a$$

$$y_3 \rightarrow a$$

$$fa(a, a)$$

\times

$$F_2 \mid \mu: x_3 \rightarrow c$$

$$y_3 \rightarrow b$$

$$(p(b, a)) \quad S = \{(x, c)\}$$

$$R_2 \mid \mu: x_2 \rightarrow b$$

$$y_2 \rightarrow a$$

$$mo(b, a)$$

\times

$$R_3 \mid \mu: x_3 \rightarrow a$$

$$y_3 \rightarrow a$$

$$fa(b, a)$$

$$F_1$$

$$\emptyset \quad S = \{(x, d)\}$$

$$M = \emptyset$$

APPROCHE ALTERNATIVE : RÉÉCRITURE DE REQUÊTE

- Le **chaînage arrière classique** nécessite de faire des unifications avec la base de faits « atome par atome »
- Si la base de faits est stockée dans une base de données, cela nécessite d'accéder fréquemment à la base de données pour faire des requêtes élémentaires
=> impraticable pour de grandes bases de faits
- Idée : décomposer le chaînage arrière en **deux phases**
 1. **Réécriture de requête** :
réécrire la requête avec la base de règles (les « vraies » règles)
=> on obtient un ensemble de requêtes conjonctives
que l'on voit comme une union de requêtes conjonctives (UCQ)
 2. **Evaluation de requête** : évaluer cette UCQ sur la base de données

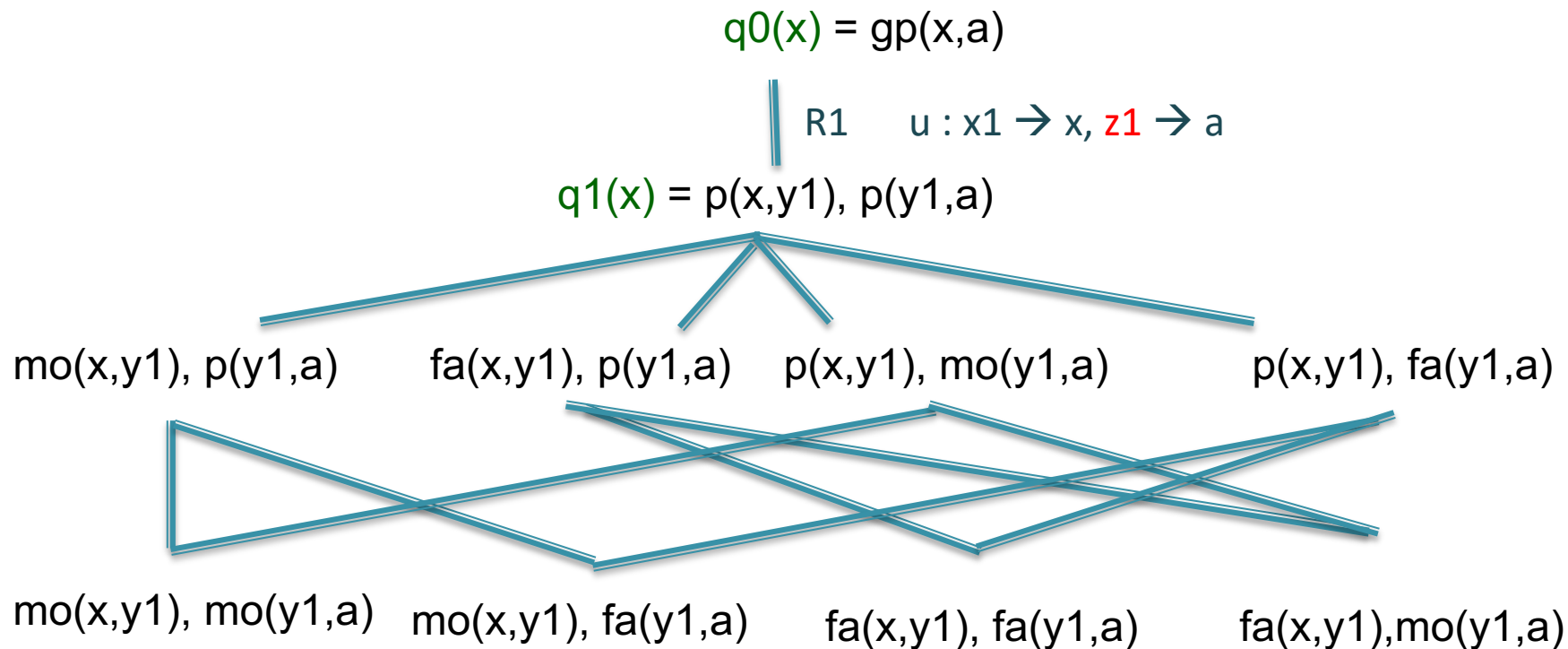
QUERY REWRITING WITH DATALOG RULES

\mathcal{R}

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

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

\mathcal{Q} set of rewritings of q with \mathcal{R} .



$\mathcal{Q} = \{ q_0, q_1, \dots, q_9 \}$

$UCQ(\mathcal{Q}) = q_0 \vee q_1 \vee \dots \vee q_9$

QUERY REWRITING WITH DATALOG RULES

\mathcal{R} 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} , set of rewritings of q with \mathcal{R}

$q0(x) = gp(x,a)$

$q1(x) = p(x,y1), p(y1,a)$

$q2(x) = mo(x,y1), p(y1,a)$

$q3(x) = fa(x,y1), p(y1,a)$

$q4(x) = p(x,y1), mo(y1,a)$

$q5(x) = p(x,y1), fa(y1,a)$

$q6(x) = mo(x,y1), mo(y1,a)$

$q7(x) = mo(x,y1), fa(y1,a)$

$q8(x) = fa(x,y1), mo(y1,a)$

$q9(x) = fa(x,y1), fa(y1,a)$

F

F1: $fa(b,a)$

F2 : $fa(c,b)$

F3 : $mo(d,b)$

$q((F, R)) = \mathcal{Q}(F)$

$= \{ (c), (d) \}$

$UCQ(\mathcal{Q}) = q0 \vee q1 \vee \dots \vee q9$

Let q be a Boolean CQ and \mathcal{Q} be its set of rewritings with \mathcal{R}

For any factbase F ,

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

Let $q(x_1, \dots, x_k)$ be a CQ and \mathcal{Q} be its set of rewritings with \mathcal{R}

For any factbase F , a tuple of constants (a_1, \dots, a_k) is an answer to q on (F, \mathcal{R})

iff (a_1, \dots, a_k) is an answer to the UCQ \mathcal{Q} on F

iff there is q_i in \mathcal{Q} such that (a_1, \dots, a_k) is an answer to q_i on F

PROBLÈME 1 : LA RÉÉCRITURE PEUT ÊTRE INFINIE

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

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

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

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

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

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

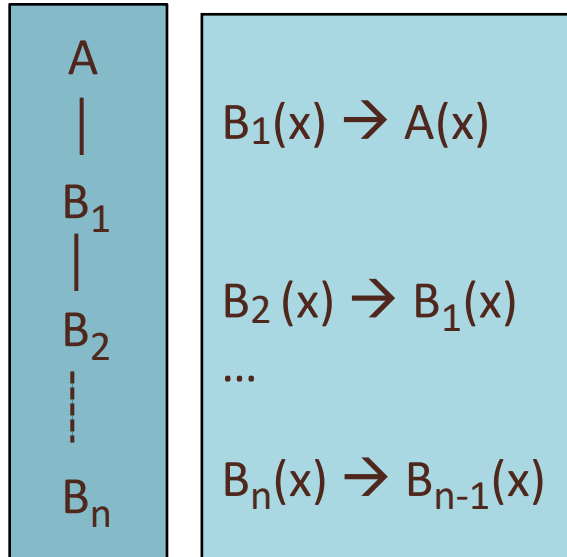
q_2 and $q_{2'}$
are equivalent

Etc.

Si on connaît la taille de la base de faits, on peut borner la taille des réécritures, mais cela donnera quand même de très grosses réécritures !

PROBLÈME 2 : EFFICACITÉ DE L'APPROCHE « RÉÉCRITURE » EN PRATIQUE ?

- 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

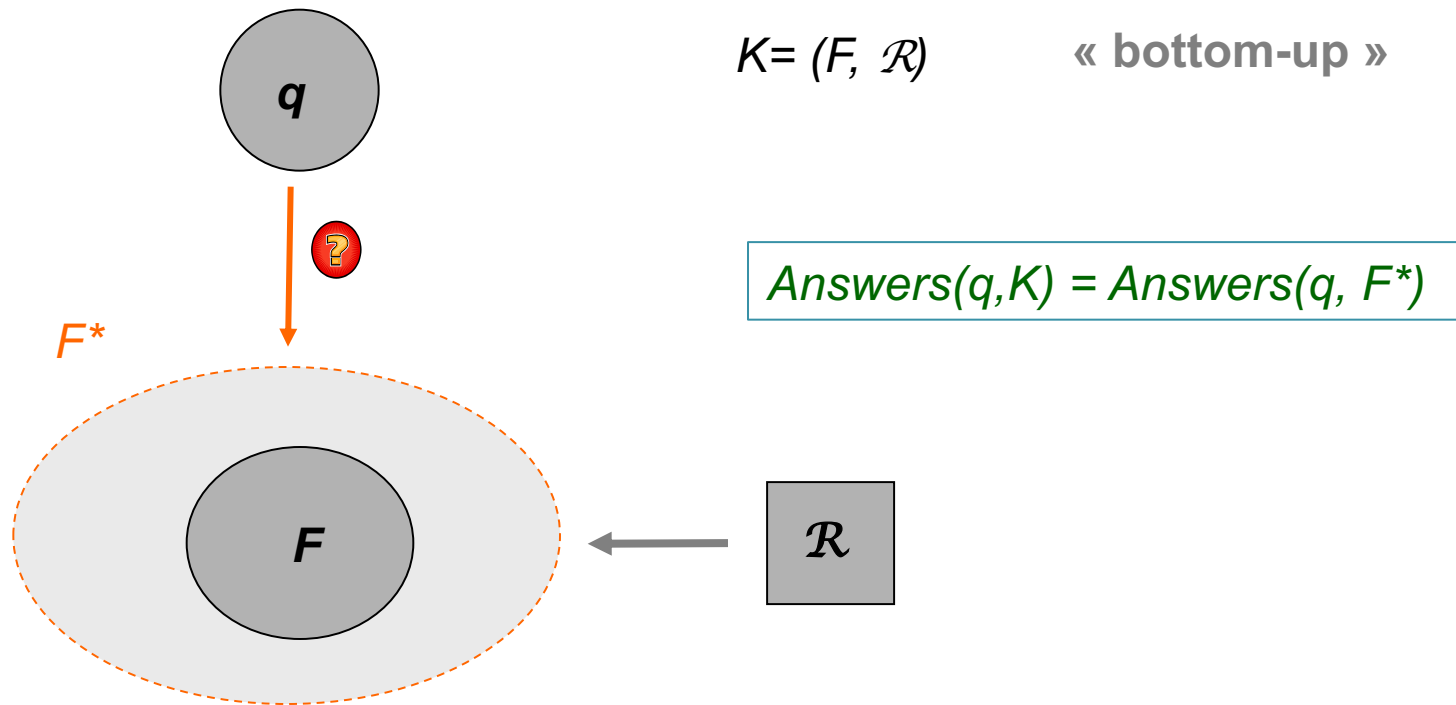
Ce n'est pas un « pire des cas » théorique :
se produit souvent en pratique

→ Réécriture en des formes de requêtes **plus compactes**

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

$(n+1) \times k$ atomes au lieu de $(n+1)^k$

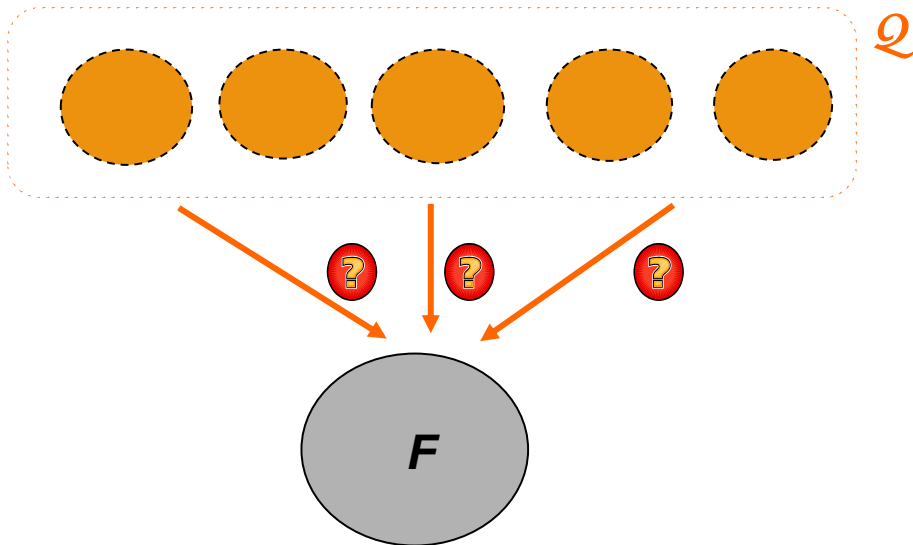
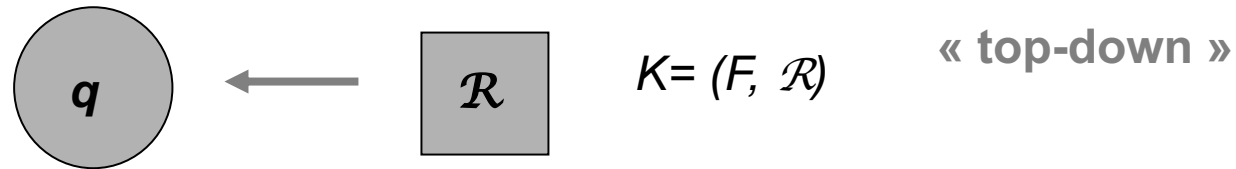
APPROACH 1 TO RULES : FORWARD CHAINING / MATERIALISATION



Pros: materialisation offline, then online query answering is fast

Cons: volume of the materialisation
not feasible if data is distributed among several databases
not adapted if data change frequently

APPROACH 2 TO RULES : QUERY REWRITING



Rewriting into a set of CQs, seen as a **union of conjunctive queries (UCQ)**

or to a more general first-order formula that still corresponds to an SQL query

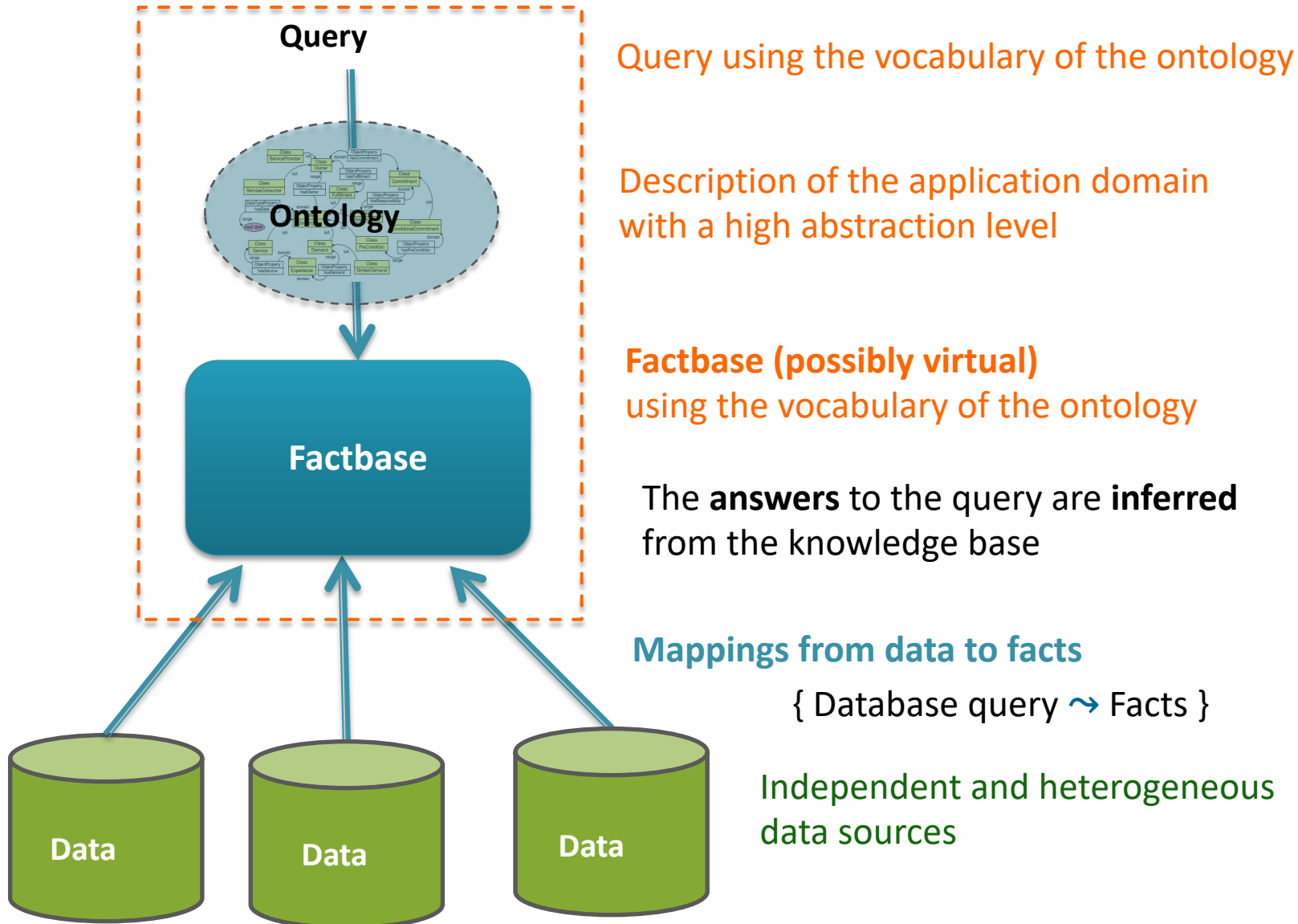
Query rewriting is independent from any factbase

For **any** F , $Answers(q, (F, \mathcal{R})) = Answers(\mathcal{Q}, F)$

Pros: independent from the data

Cons: rewriting done at query time, easily leads to huge and unusual queries

Conceptual level



MAPPINGS

Patient_T [ID_PATIENT, NAME,SSN]

Diagnosis_T[ID_PATIENT, DISORDER]

Patient /1

Diagnosis / 2

Influenza /1

Mapping = database query(X) \leadsto conjunction with free variables X

$q(x): \exists n \exists s \text{ Patient_T}(x,n,s) \leadsto \text{Patient}(x)$ **Datalog**

$q'(x): \exists n \exists s \text{ Patient_T}(x,n,s) \wedge \text{Diagnostic_T}(x,y) \wedge y = \text{« influenza »}$ **pas**
 $\leadsto \exists z (\text{diagnosis}(x,z) \wedge \text{Influenza}(z))$ **Datalog**

Patient_T			Diagnosis_T	
id	name	ssn	id	dis
P	P	influenza
..
..

\leadsto

Patient(P)
Diagnosis(P,M)
Influenza(M)

MAPPINGS CAN BE SEEN AS RULES

Patient_T [ID_PATIENT, NAME,SSN]

Diagnosis_T[ID_PATIENT, DISORDER]

$q(x): \exists n \exists s \text{ Patient_T}(x, n, s) \rightsquigarrow \text{Patient}(x)$ **GAV (Global-As-View)**
= Datalog

$q'(x): \exists n \exists s \text{ Patient_T}(x, n, s) \wedge \text{Diagnosis_T}(x, y) \wedge y = \text{« influenza »}$
 $\rightsquigarrow \exists z \text{ diagnosis}(x, z) \wedge \text{Influenza}(z)$

$\text{Patient_T}(x, n, s) \rightarrow \text{Patient}(x)$

Datalog Rule

$\text{Patient_T}(x, n, s), \text{Diagnosis_T}(x, \text{« influenza »}) \rightarrow \exists z \text{ diagnosis}(x, z), \text{Influenza}(z)$

Existential Rule

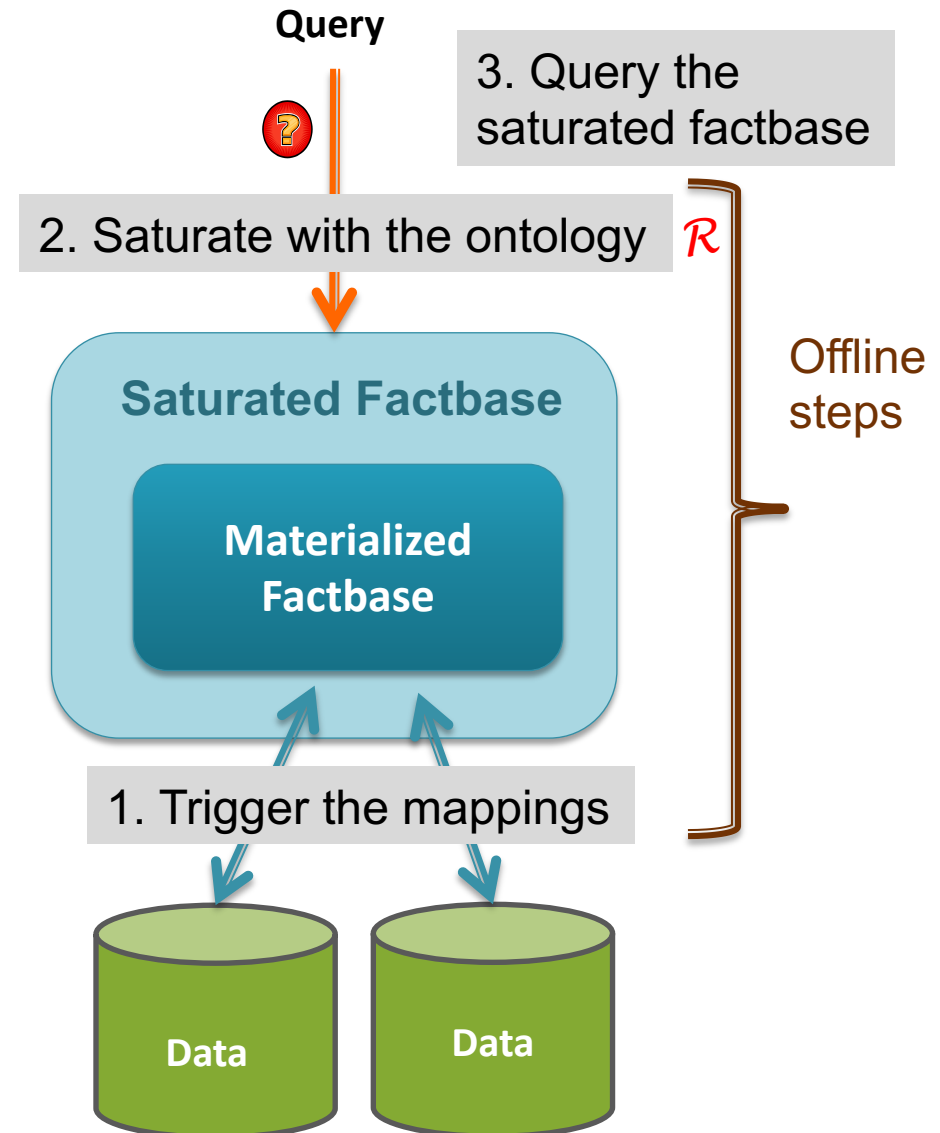
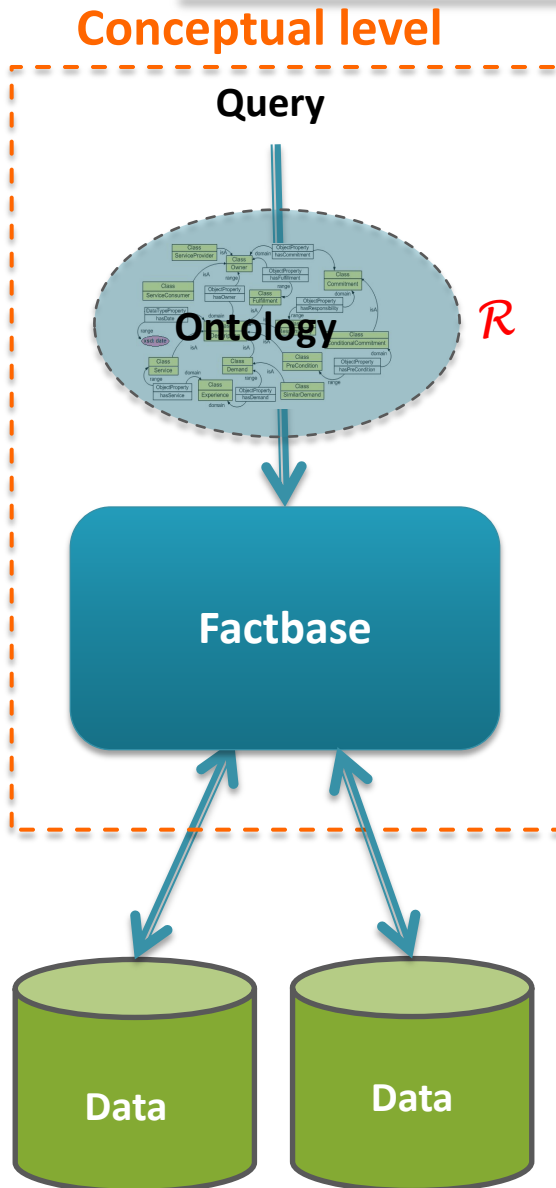
More generally: $q_1(X) \rightsquigarrow q_2(X)$ where q_1 is expressed in a native query language

Decomposition of a mapping into 2 mappings

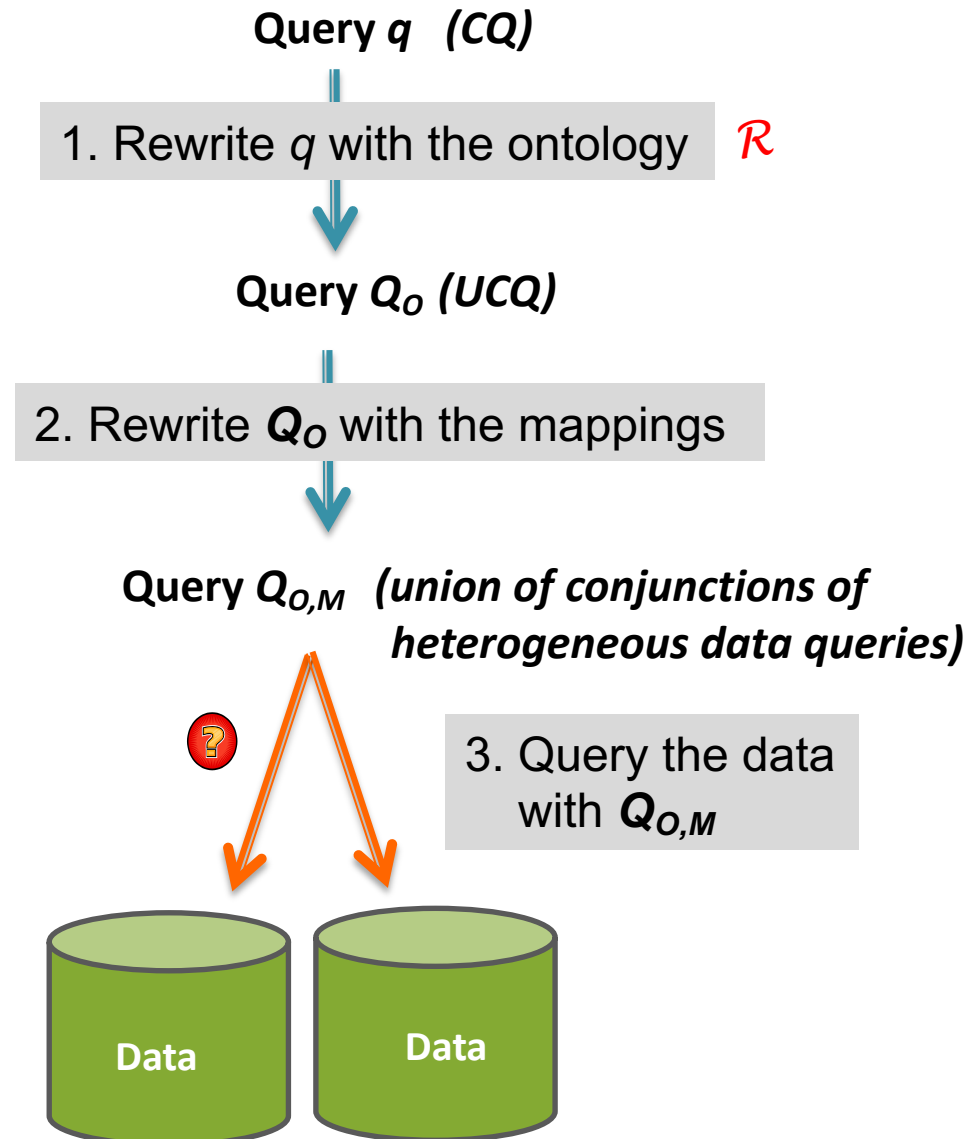
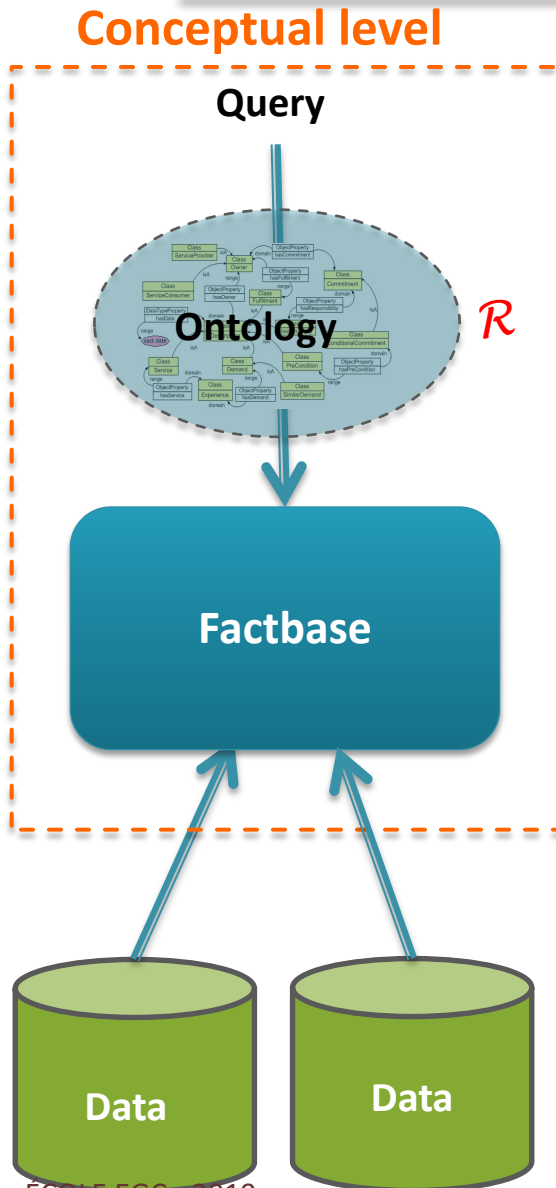
low level : $q_1(X) \rightsquigarrow \text{view}(X)$ **Result of the query stored in a view**

high level : $\text{view}(X) \rightarrow q_2(X)$ **Logical rule**

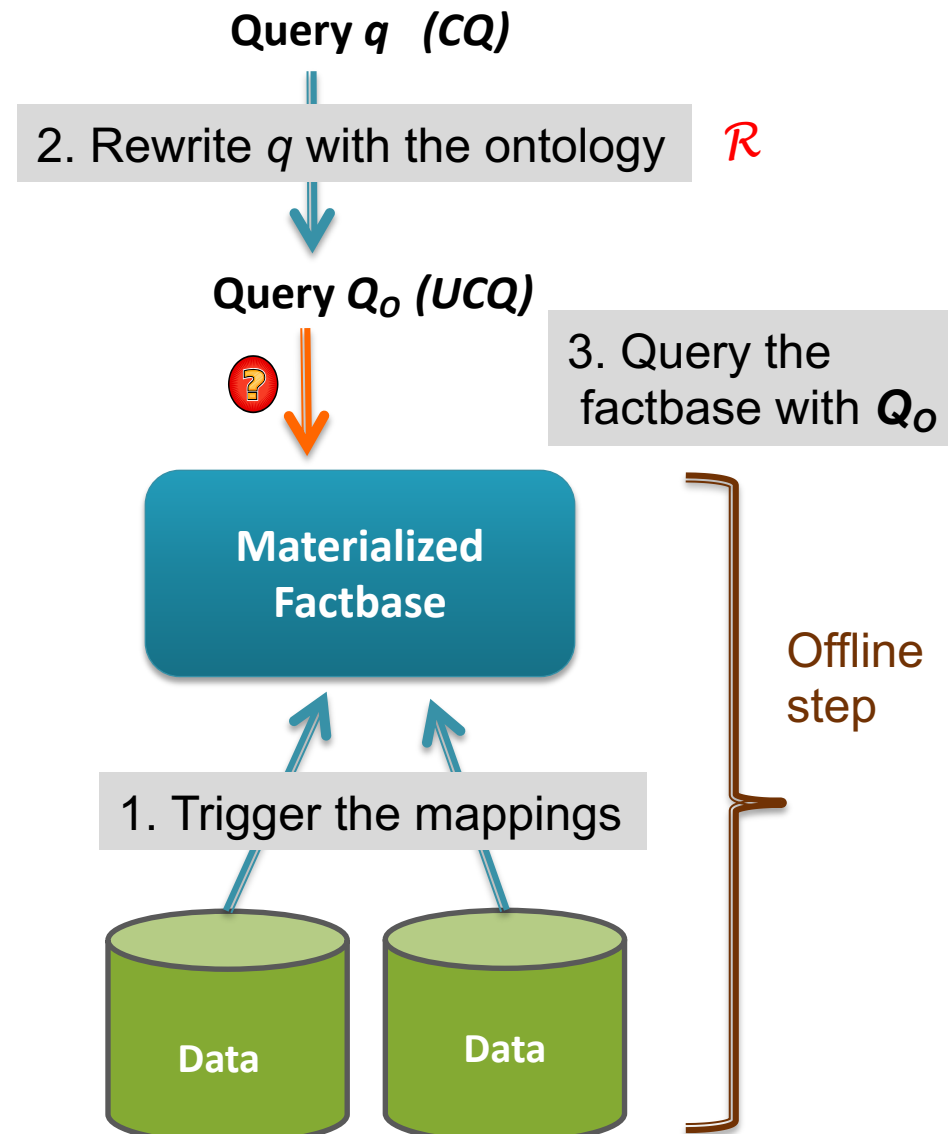
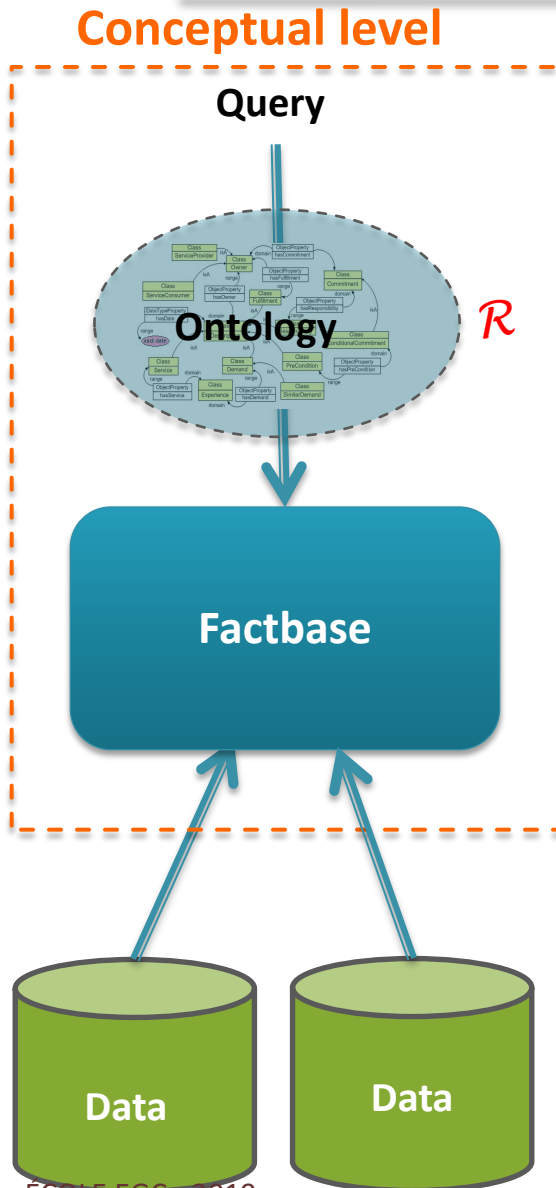
OBDA : TOTAL MATERIALIZATION (FORWARD CHAINING)



OBDA : TOTAL REWRITING



OBDA : EXAMPLE OF MIXED APPROACH



SYNTHÈSE

- **Bases de connaissances avec faits et règles datalog**
- Une **base de données** relationnelle peut être vue comme une base de faits et réciproquement
- Une **requête conjonctive (CQ)** correspond à une requête SQL de base
- Les réponses à une requête conjonctive sur une base de faits se calculent par **homomorphisme**
- Si on ajoute des **contraintes négatives**, la base de connaissances peut devenir insatisfiable
- Une base de connaissances (satisfiable) K a un **unique plus petit modèle**, qui suffit à calculer les réponses à une CQ sur K
- **Méthodes d'interrogation** de K :
 - par chaînage avant : la base de faits saturée correspond à un plus petit modèle de K
 - par chaînage arrière classique (« programmation logique », cf. Prolog)
 - par réécriture de requête puis évaluation de la requête réécrite
- Cadre plus général avec des **mappings**