# RÈGLES EXISTENTIELLES

THÉORIE DES BASES DE DONNÉES ET DE CONNAISSANCES
HAI933I
COURS DE ML MUGNIER

# Retour aux Bases de Connaissances

**Requête**



**Base de connaissances K**

Requête du premier ordre
(par ex une CQ ou une UCQ)

Ensemble de formules logiques, par ex :
    Règles Datalog (M1)
    Règles existentielles (extension vue ici)

Ensemble de faits :
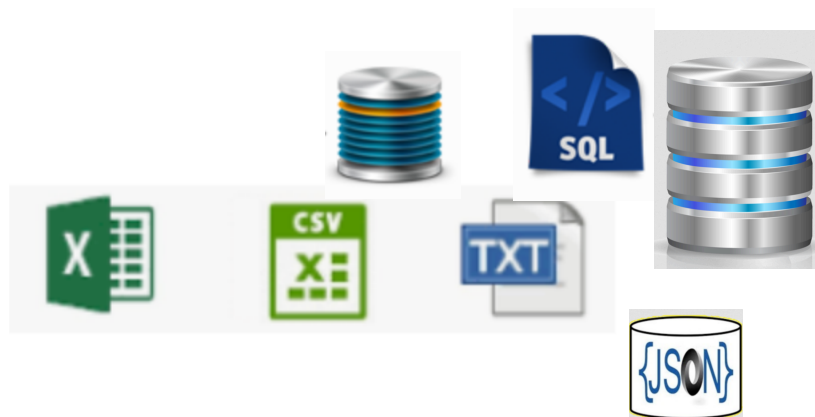    instanciés (M1)
    avec variables existentielles
    (extension vue au cours précédent)

La réponse à une requête booléenne *q* est oui
si q est **conséquence logique** de K
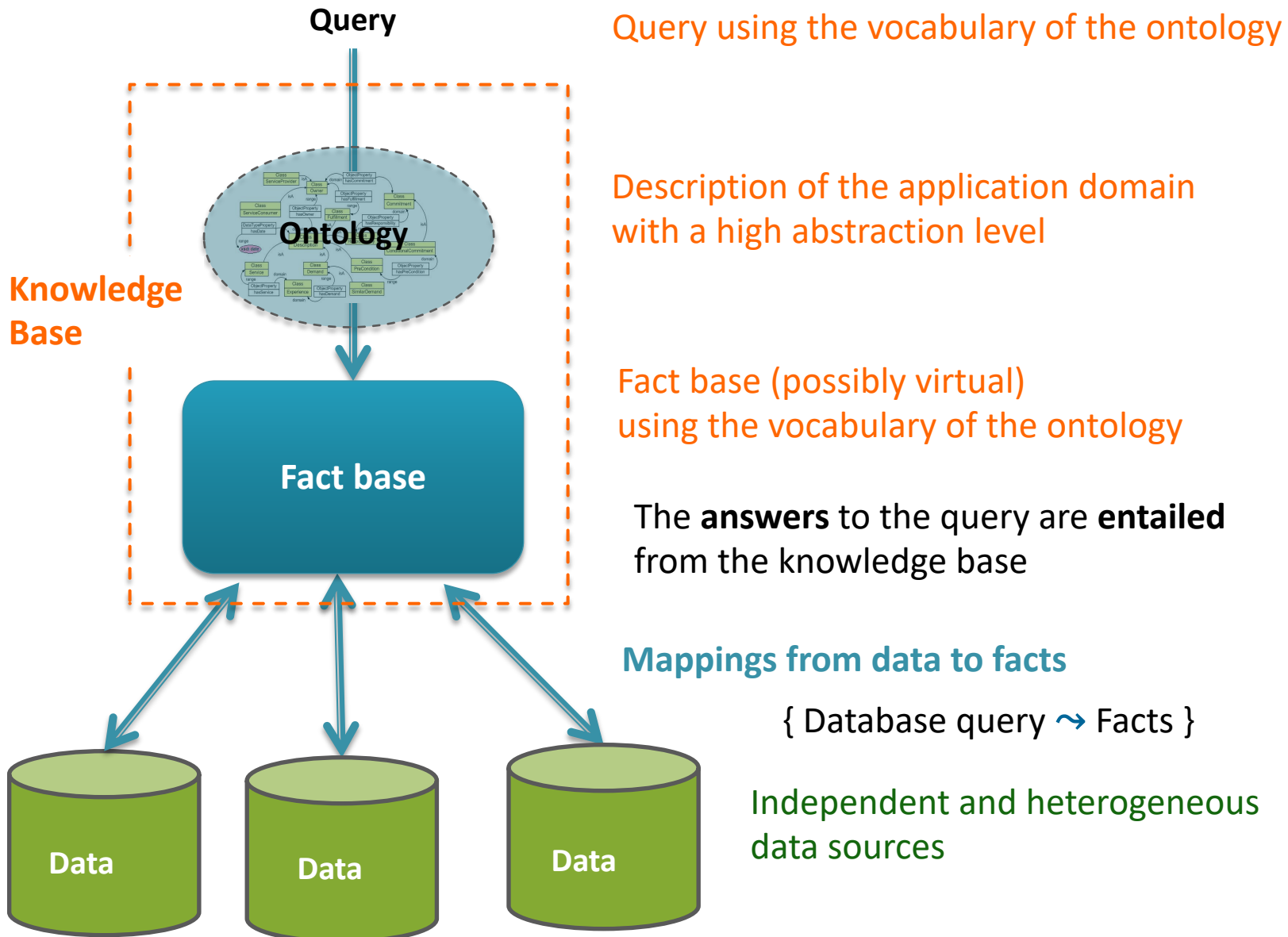
# Example: Detecting Conflicts of Interest

○ Problem: detect links between scientific experts involved in public interest studies and industry that are likely to constitute a conflict of interest



Data on the scientific commissions and experts
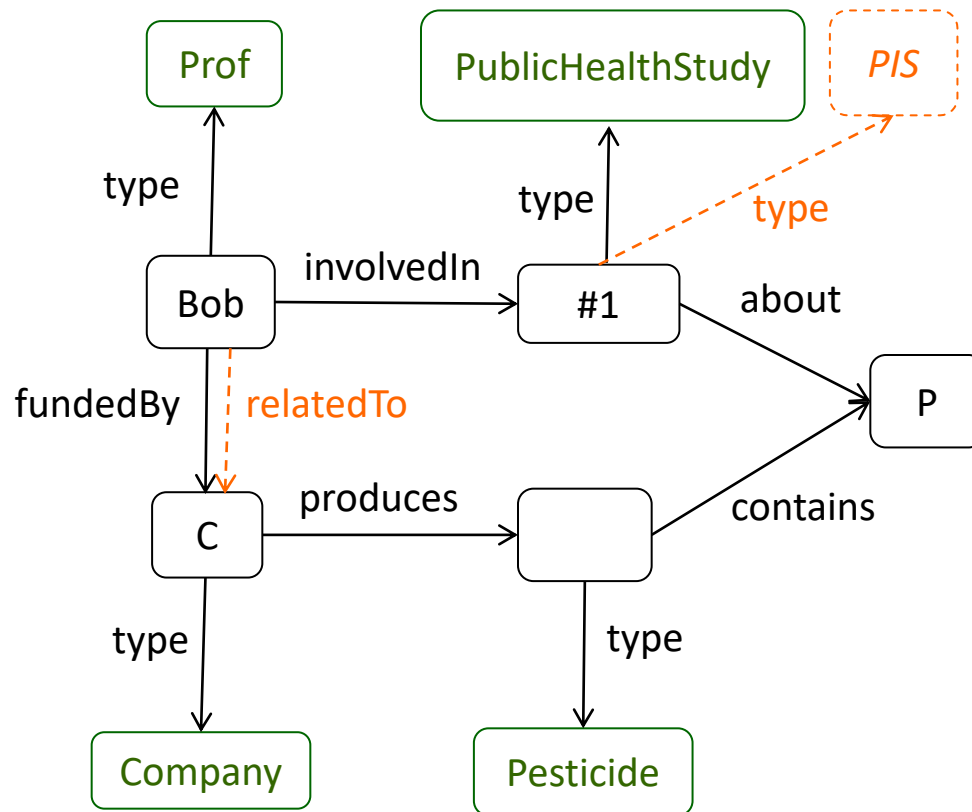Data on companies and their products
Data on research laboratories, scientific projects and their funding
Data on lobbies declarations: members, participants to meetings
Data on scientific publications
.....

**Dream query**: ''Find all persons **x**, studies **y** and companies **z** such that **x** has a conflict of interest for **y** because of its relationships with **z**''

# ONTOLOGY-BASED DATA ACCESS



**Query**

Query using the vocabulary of the ontology

**Ontology**

Description of the application domain with a high abstraction level

**Knowledge Base**

**Fact base**

Fact base (possibly virtual) using the vocabulary of the ontology

The **answers** to the query are **entailed** from the knowledge base

**Mappings from data to facts**

{ Database query ↝ Facts }

Independent and heterogeneous data sources

**Data**   **Data**   **Data**

# ASSUME WE HAVE BUILT A FACT BASE (« KNOWLEDGE GRAPH »)



**Facts**

∃x (
  Prof(Bob)         ∧
  PHS(#1)           ∧
  Comp(C)           ∧
  Pest(x)           ∧
  involvedIn(Bob,#1) ∧
  fundedBy(Bob,C)   ∧
  about(#1,P)       ∧
  produces(C,x)     ∧
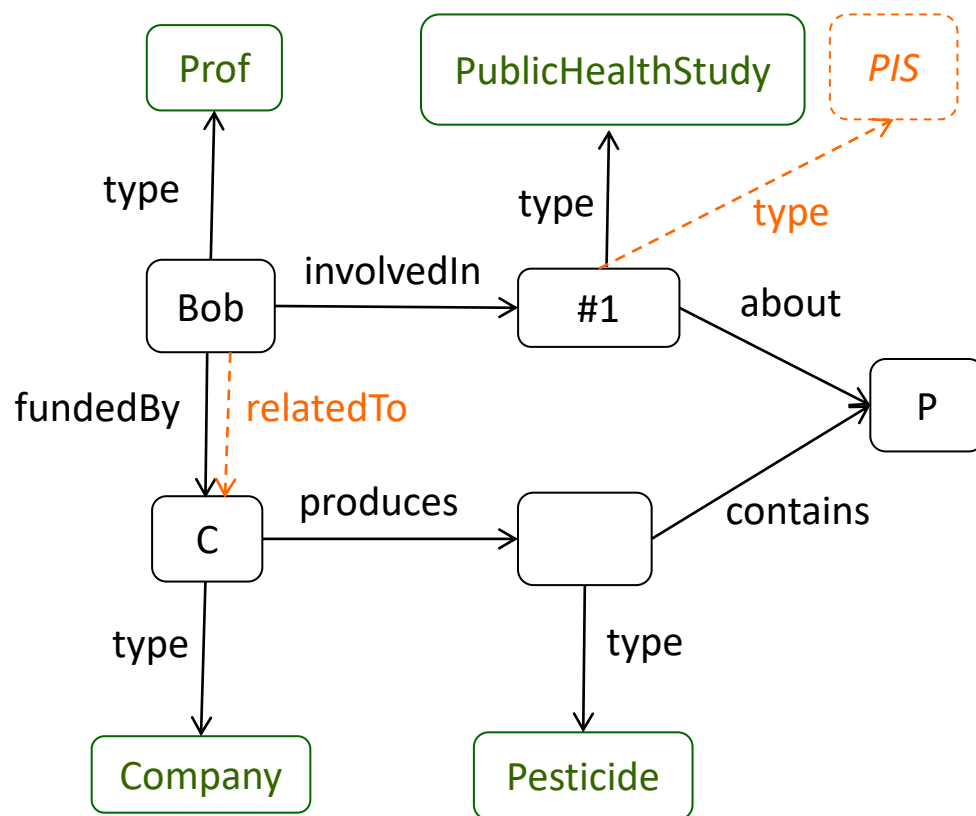  contains(x,P)     )

+ Basic ontological knowledge

PublicHealthStudy **subclass of** PublicInterestStudy
fundedBy **subproperty of** relatedTo

**Rules**

∀x (PHS(x) → PIS(x))
∀x∀y (fundedBy(x,y) → relatedTo(x,y))

**Allow to infer:**

PIS(#1), relatedTo(Bob,C)

# How to Infer Conflicts of Interest (C.o.I.) ?

Prof

PublicHealthStudy

*PIS*

type

type

type

Bob → involvedIn → #1 → about → P

fundedBy | relatedTo

C → produces → ◻ → contains

type

type

Company

Pesticide

What kind of **ontological knowledge** would allow to infer conflicts of interest?

PIS

involvedIn

*x* → *y*

relatedTo

type

about

*z* → hasInterest → *u*

type

Company

C.o.I. pattern

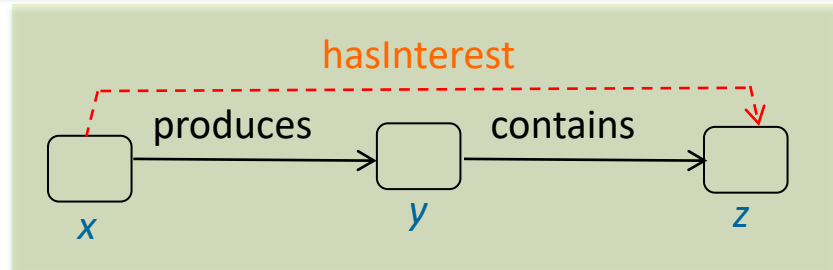**Query**: ''Find all **x, y, z** such that **x** has a conflict for study **y** because of its relationships with company **z**''

# DEFINING CONFLICTS OF INTEREST

$R_1$: $\forall x \forall y \forall z$ ( produces(x,y) $\land$ contains(y,z)
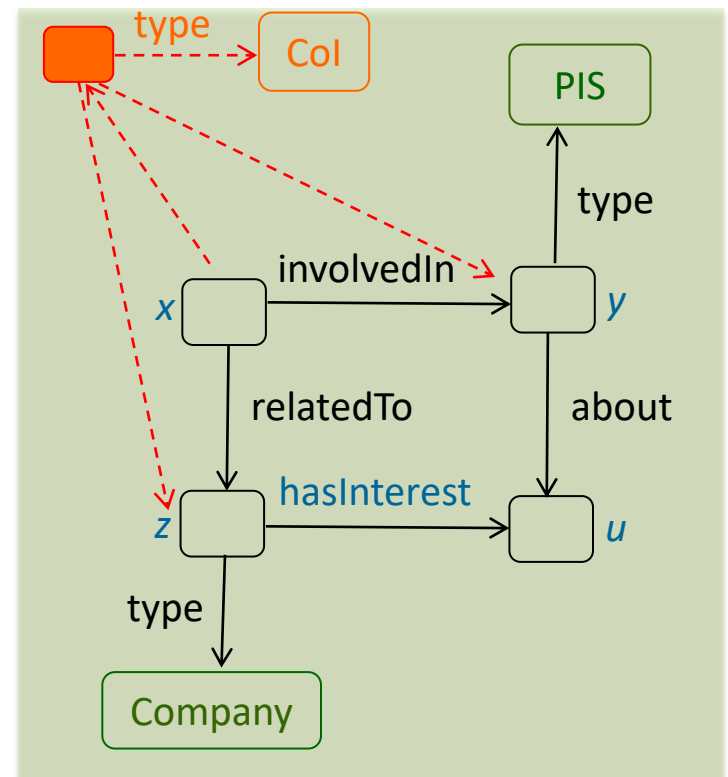$\rightarrow$ hasInterest(x,z) )



$R_2$: $\forall x \forall y \forall z \forall u$ ( involvedIn(x,y) $\land$ PIS(y) $\land$ about(y,u) $\land$ relatedTo(x,z) $\land$ Company(z) $\land$ hasInterest(z,u)
$\rightarrow$ CoI(x,y,z) )

What if we only have unary and binary predicates
i.e. graphs and not hypergraphs ?

Reification: new object of type CoI

$R_2$: $\forall x \forall y \forall z \forall u$ ( *body[x,y,z,u]* $\rightarrow \exists o$
(CoI(o) $\land$ in(x,o) $\land$ on(o,y) $\land$ with(o,z))
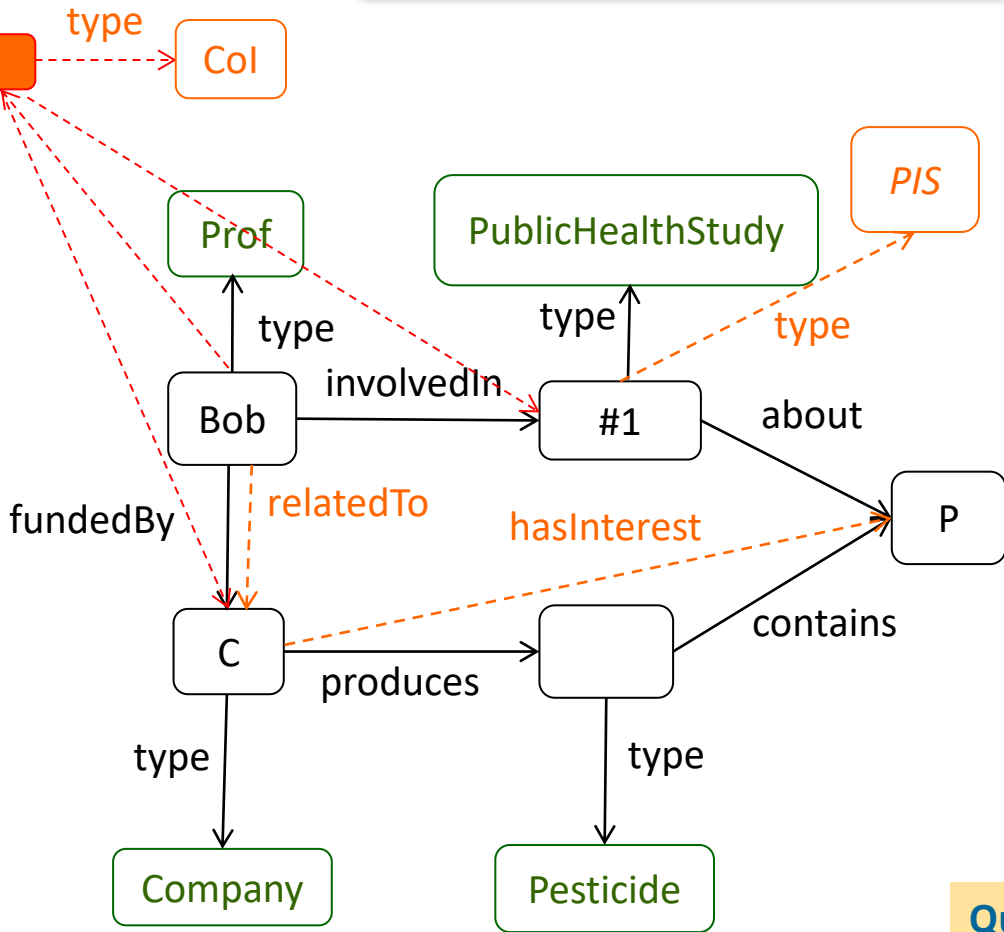)

# CREATING NEW OBJECTS

$R_2$: $\forall x \forall y \forall z \forall u$ ( *body[x,y,z,u]* $\rightarrow$ $\exists o$ (CoI(o) $\wedge$ in(x,o) $\wedge$ on(o,y) $\wedge$ with(o,z)) )

Interest of creating a new object:

- Flexible description of C.o.I. instead of a fixed arity predicate
  Not all C.o.I. need to be described by the same properties

- Ability to talk about C.o.I. because they become objects (reification)

E.g. $R_3$: $\forall x \forall z$ (CoI(x) $\wedge$ with(x,z) $\wedge$ ChemicalCompany(z) $\rightarrow$ toBeInvestigated(x) )

# INFERRING CONFLICTS OF INTEREST



**Rules** (universal quantifiers omitted)

$PHS(x) \rightarrow PIS(x)$
$fundedBy(x,y) \rightarrow relatedTo(x,y)$

$R_1$: $produces(x,y) \wedge contains(y,z)$
$\rightarrow hasInterest(x,z)$

$R_2$: $involvedIn(x,y) \wedge PIS(y) \wedge about(y,u) \wedge$
$relatedTo(x,z) \wedge Company(z) \wedge$
$hasInterest(z,u)$
$\rightarrow \exists o\ CoI(o) \wedge in(x,o) \wedge on(o,y) \wedge with(o,z)$

**Inferred facts**

$PIS(\#1)$, $relatedTo(Bob,C)$, $hasInterest(C,P)$
$CoI(o_1)$, $in(Bob,o_1)$, $on(o_1,\#1)$, $with(o_1,C)$

**Query**: find (x,y,z) such that
$\exists o\ CoI(o) \wedge in(x,o) \wedge on(o,y) \wedge with(o,z)$

**Answer**: (Bob,#1,C)

**Facts**

$Prof(Bob)$, $PHS(\#1)$, $Comp(C)$, $Pest(x)$
$involvedIn(Bob,\#1)$, $fundedBy(Bob,C)$
$about(\#1,P)$, $produces(C,x)$, $contains(x,P)$

# EXISTENTIAL RULES

$$\forall X \; \forall Y \; ( \; Body \; [X,Y] \rightarrow \exists Z \; Head \; [X,Z] \; )$$

X, Y, Z : sets *of variables*
(possibly empty)

any **positive conjunction** (without functional symbols)

$$\forall x \; ( \; actor(x) \rightarrow \exists z \; (movie(z) \wedge play(x,z) \; )$$

$$\forall x \; \forall y \; ( \; siblingOf(x,y) \rightarrow \exists z \; (parentOf(z,x) \wedge parentOf(z,y)) \; )$$

Key point: ability to assert **the existence of unknown entities**

Crucial for representing ontological knowledge in « **open domains** »

[Open domain = we do not assume that the only existing objects are those known in the factbase]
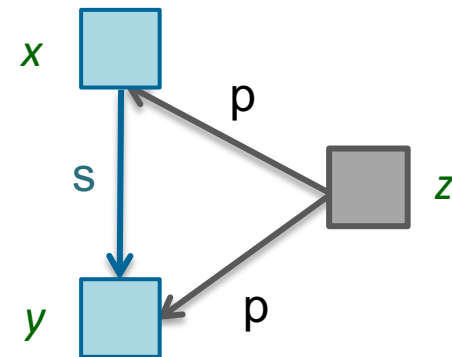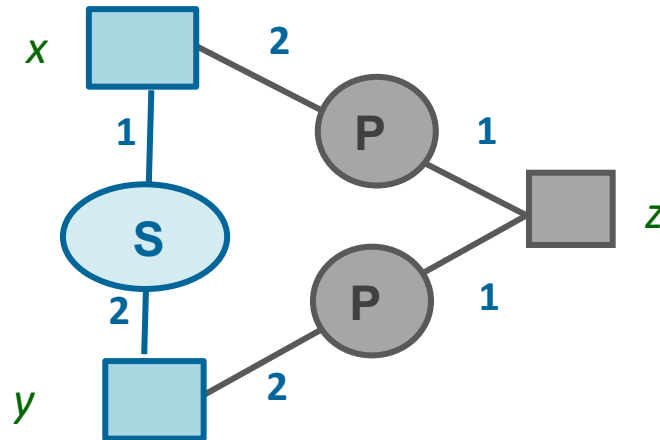
# GRAPH VIEW OF (EXISTENTIAL) RULES

$$\forall X \, \forall Y \, ( \, \text{Body} \, [X,Y] \rightarrow \exists Z \, \text{Head} \, [X,Z] \, )$$

**graph**          **graph**

$\forall x \, \forall y \, ( \, \text{siblingOf}(x,y) \rightarrow \exists z \, (\text{parentOf}(z,x) \wedge \text{parentOf}(z,y)) \, )$



The rule head has 2 kinds of variables (or unlabelled term nodes):
- **frontier**: shared with the body (**X**)          {x,y} on the example
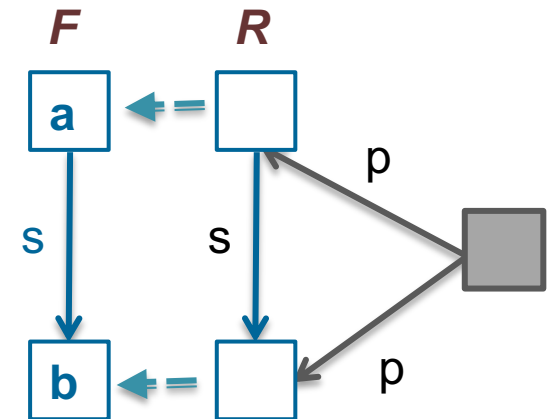- **existential**: (**Z**)          {z} on the example

# GENERATION OF FRESH (UNKNOWN) INDIVIDUALS

$R = \forall x \, \forall y \, (siblingOf(x,y) \rightarrow \exists z \, (parentOf(z,x) \land parentOf(z,y)))$

$F = siblingOf(a,b)$

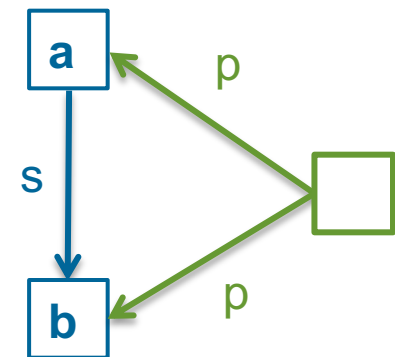*R* is **applicable** to *F* if there is a **homomorphism** *h* from *body(R)* to *F*

$$x \rightarrow a$$
$$y \rightarrow b$$

Applying *R* to *F* w.r.t. *h* produces *F* ∪ *h(head(R))*

*where a **fresh** variable (a «null») is created for each existential variable in R*

$F' = \exists z0 \, ( \, siblingOf(a,b) \land parentOf(z0,a) \land parentOf(z0,b) \, )$

Formal notation (when needed) : **F ∪ h^safe(head(R))**

where **h^safe** is a substitution of variables(head(R))
such that:      $h^{safe}(x) = h(x)$ if x is in frontier(R)
                     otherwise $h^{safe}(x)$ is a fresh variable (a null)

# Retour sur Datalog

- Les règles Datalog sont un cas particulier de règles existentielles

$$\forall X \; \forall Y \; ( \; Body \; [X,Y] \; \rightarrow \; \exists \; Z \; Head \; [X,Z] \; ) \; avec \; Z = \emptyset$$

- Soit une base de connaissances **K =(F,$\mathcal{R}$)** où F est une **base de faits sans variables** et $\mathcal{R}$ est un ensemble de règles **Datalog**.
  Alors:

  - K possède un **unique plus petit modèle** qui est l'**intersection** de tous ses modèles (« plus petit » au sens de l'inclusion et du cardinal)

  - Donc, étant donnée une CQ Booléenne q, pour déterminer si **K ⊨ q** il **suffit** de vérifier si le plus petit modèle de K est un modèle de q

  - Le plus petit modèle de K se calcule en **saturant** F avec $\mathcal{R}$ (« chainage avant »)

Qu'est-ce qui change quand on passe aux **règles existentielles** ?

RAPPEL

Vocabulaire V = (P, C)

Base de faits *F* (sans variables) sur V

Modèle **canonique** de *F* *(ou : modèle isomorphe à F)*

M:        $D_M = C$

pour tout *p* ∈ P d'arité *k*, $p^M = \{ (c_1, ..., c_k) \mid p(c_1, ...,c_k) \in F \}$

Le modèle canonique de *F* correspond à l'**intersection** de tous les modèles de *F*

V = ( {$r_{/3}$, $p_{/2}$, $q_{/1}$ }, {a, b, c, d, e} )

*F* = { p(a,b), p(b,c), q(c) }

M:        $D_M = \{a,b,c,d,e\}$

$p^M = \{ (a,b), (b,c) \}$

$q^M = \{ c \}$

$r^M = \varnothing$

Qu'est-ce qui change quand la base de faits peut avoir des variables ?

# MODEL "ISOMORPHIC" TO A CLOSED FOL($\exists$,$\wedge$) FORMULA

To a closed formula *F* in FOL($\exists$,$\wedge$), we assign its **isomorphic model**
(also called **canonical model**)*:*

$\mathcal{M}$:

- $D^\mathcal{M} = C \cup$ *variables(f)*      *We add a domain element for each variable*

- for all *p* in $\mathcal{P}$, $p^\mathcal{M} = \{(t_1 \ldots t_k) \mid p(t_1 \ldots t_k) \text{ in } f\}$,

$\mathcal{V} = (\{s_{/1}, p_{/2}, r_{/3}\}, \{a,b\})$

$F = \exists x \exists y \exists z\ (\ p(x,y) \wedge p(y,z) \wedge r(x,z,a)\ )$

*M*:       $D_M = \{a, b, x, y, z\}$
          $p^M = \{\ (x,y),\ (y,z)\ \}$
          $r^M = \{\ (x, z, a)\ \}$
          $s^M = \emptyset$

Reciprocally, any **interpretation *I*** can be seen as a **closed FOL($\exists$,$\wedge$) formula**

Each element from $D_I \setminus C$ is translated into a new variable

# « INTERSECTION DES MODÈLES » : ?

Le modèle canonique d'une base de faits F avec variables n'est plus un « plus petit modèle » au sens de l'inclusion ☹

$\mathcal{V}$ = ({p$_{/2}$}, {a,b})

$F$ = ∃x ∃y (p(a,x) ∧ p(a,b) ∧ p(b,y))

$M$:    D = {a, b, x, y}
        p$^M$ = { (a,x), (a,b), (b,y)}

Plus petit modèle au sens de l'inclusion ?

$D_{M'}$ = {a, b, y}
$p^{M'}$= { (a,b), (b,y) }

Ceci est dû au redondances
($F$ n'est pas un core)

Et ce n'est plus l'intersection des modèles de $F$

$\mathcal{V}$ = ({p$_{/2}$}, {a,b})

$F$ = ∃x p(a,x)

$M$:    D = {a, b, x}
        p$^M$ = { (a,x)}

Autres modèles de $F$ ?

$D_{M'}$ = {a, b}          $D_{M'}$ = {a, b}          Etc.
$p^{M'}$= { (a,b) }        $p^{M'}$= { (a,a) }

# MODÈLES UNIVERSELS

Le modèle canonique d'une base de faits F avec variables n'est plus un « plus petit modèle » au sens de l'inclusion, ce n'est plus non plus l'intersection des modèles de F

**Mais** c'est un modèle « le plus général » :

Le modèle canonique d'une formule close *F* de FOL($\exists$,$\wedge$) est un modèle **universel** de *F :*

il s'envoie par homomorphisme dans tous les modèles de F

**Conséquence** :

Si un CQ *Q* est satisfaite dans un modèle universel de *F*

alors Q est satisfaite dans tous les modèles de *F*

Donc *F* $\vDash$ *Q*

# HOMOMORPHISMS AGAIN AND AGAIN

- One can define **homomorphisms** between **interpretations**

  Homomorphism $h$ from $\mathcal{I}_1=(D_1, .^{\mathcal{I}1})$ to $\mathcal{I}_2 = (D_2, .^{\mathcal{I}2})$:
  
  mapping from $D_1$ to $D_2$ such that:

  for all $c$ in $\mathcal{C}$, h(c) = c

  for all $p$ in $\mathcal{P}$ and $(t_1 \ldots t_k)$ in $p^{\mathcal{I}1}$, $(h(t_1) \ldots h(t_k))$ in $p^{\mathcal{I}2}$

- Homomorphisms between interpretations correspond to homomorphisms between the associated fact bases

- If $\mathcal{I}_1$ maps by homomorphism to $\mathcal{I}_2$ then,
  
  for any $F$ in FOL($\exists,\wedge$), $\mathcal{I}_1$ model of $f \Rightarrow \mathcal{I}_2$ model of $f$

  Indeed: $f$ maps to $\mathcal{I}_1$ and $\mathcal{I}_1$ maps to $\mathcal{I}_2$, hence $f$ maps to $\mathcal{I}_2$

# Nice semantic properties of FOL($\exists,\wedge$)

- For any f in FOL($\exists,\wedge$), the canonical model of $f$ is **universal:**

    for all $M'$ model of $f$, $M_f$ maps by homomorphism to $M'$

- $g \vDash f$   (i.e., every model of $g$ is a model of $f$) iff

    $M_g$ *is a model of f* (the canonical model of $g$ is a model of $f$) iff

    *f maps to g*     (there is a homomorphism from $f$ to g)

Donc : pour déterminer si **F $\vDash$ Q** lorsque F a des variables, on peut toujours se reposer sur l'homomorphisme

Ajoutons un ensemble $\mathcal{R}$ de règles existentielles :

- peut-on saturer F avec $\mathcal{R}$ ?

- le résultat correspond-il à un modèle universel de (F,$\mathcal{R}$) ?

# KNOWLEDGE BASES WITH EXISTENTIAL RULES

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

$\mathcal{R}$ is a set of **existential rules**

*F* is a set of **facts** (rules with an empty body): existential conjunctions of atoms

Forward chaining called « **chase** » (we still denote by *F\** the result of the chase)

Main change with respect to Datalog rules: *F\** can be **infinite**

> $R$ = person(x) $\rightarrow$ $\exists$ y hasParent(x,y) $\wedge$ person(y)
>
> $F$ = person(a)
>
> $\wedge$ hasParent(a, y0) $\wedge$ person(y0)
>
> $\wedge$ hasParent(y0, y1) $\wedge$ person(y1)     Etc.
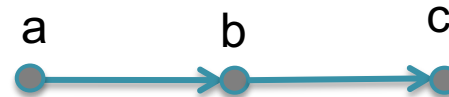
but it remains a **universal model**

Hence, for Boolean CQs :  **K ⊨ q  iff  q maps to F\***

Other changes: **F\* is not unique** (but all F\*we will see are logically **equivalent**)
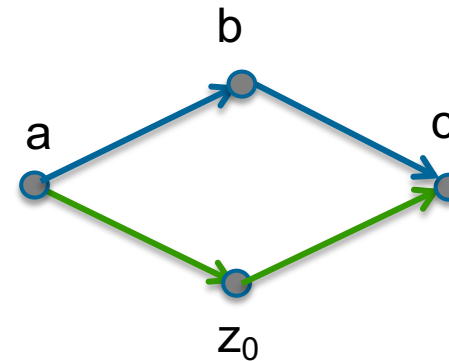
All chase variants we will see compute **universal models** of the KB

but they differ on how they handle **redundancies** possibly caused by nulls

p(a,b), p(b,c)

p(a,b), p(b,c),
$\exists z_0$ p(a,$z_0$) p($z_0$,c)

$z_0 \mapsto b$

**Core**: set of atoms without homomorphism to one of its strict subsets

# DERIVATION

- **Trigger** for a factbase F: **(R,h)** | h homomorphism from body(R) to F

- **Derivation**: **$(F_0 = F)$ $(R_1,h_1)$ $F_1$ $(R_2,h_2)$ $F_2$ , ...**
  
  where for all i,    $(h_i, R_i)$ trigger for $F_{i-1}$
  
  and $F_i = F_{i-1} \cup h_i^{safe}(head(R_i))$

  When the triggers are not needed, we note **$(F_0=F)$, $F_1$, $F_2$, ...**

- **Different chase variants** with their own rule application criteria

$\rightarrow$ different notions of *active* trigger $(R_i, h_i)$

A chase variant considers only derivations with active triggers

# OBLIVIOUS CHASE

**Oblivious** (or **naive**): « performs all rule applications according to all new triggers »

A trigger (R,h) to $F_i$ is *active on $F_i$* iff this trigger has *not already been used*
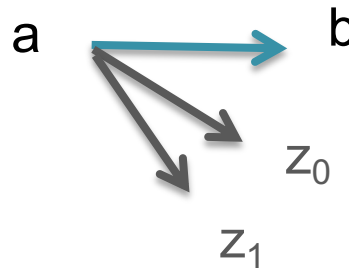in the derivation from $F_0$ to $F_{i-1}$

$R = p(x,y) \rightarrow \exists z\ p(x,z)$

$F = p(a,b)$

$p(a,z_0)$

$p(a,z_1)$    ...

a ⟶ b

$z_0$

$z_1$

*stupid rules to keep examples simple!*
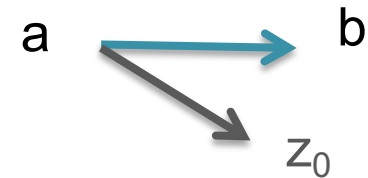
infinite derivation

# SEMI-OBLIVIOUS = SKOLEM CHASE

**Semi-oblivious**: consider only homomorphisms that differ on the rule frontier (x)

A trigger (R,h) to $F_i$ is *active on $F_i$* iff there is *no* trigger (R,h')
such that $h'(x) = h(x)$ for all $x$ in frontier(R)
in the derivation from $F_0$ to $F_{i-1}$

$F = p(a,b)$    $R = p(x,y) \rightarrow \exists z\ p(x,z)$

a ⟶ b

$z_0$

**Skolem chase**: similar behavior

(1) skolemize rules: in R, replace each existential
   variable z by a function $f_R^z(\text{frontier}(R))$

(2) perform the oblivious chase on skolemized rules

$R = p(x,y) \rightarrow p(x,f(x))$

p(a,b)    p(a,f(a))

Skolemization can be seen as a way of naming existential variables and
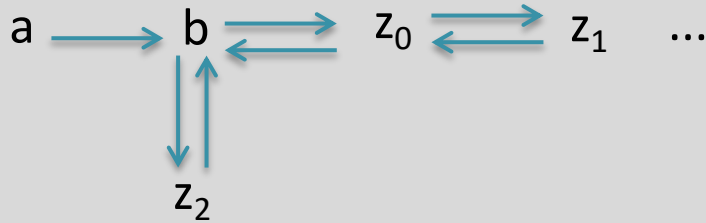« tracking » the nulls created during the semi-oblivious chase

# RESTRICTED (ALSO KNOWN AS STANDARD) CHASE

**Restricted**: do not perform a rule application that brings *only* redundant information

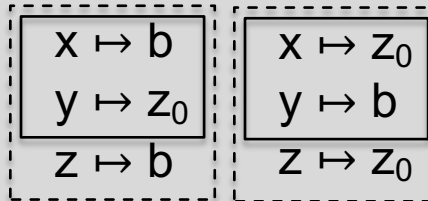A trigger (R,h) to $F_i$ is *active on $F_i$* iff
*h cannot* be extended to homomorphism *h'*: body ∪ head → $F_i$

F : p(a,b)     R : p(x,y) → ∃z p(y,z), p(z,y)

a ⟶ b ⇄ $z_0$ ⇄ $z_1$   ...

$z_2$

(semi-) oblivious chase:

 infinite

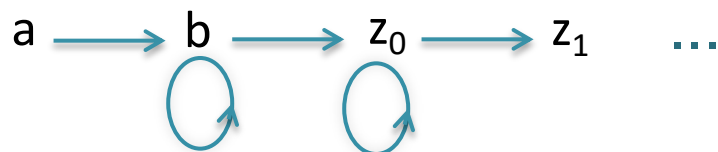| | |
|---|---|
| x ↦ b | x ↦ $z_0$ |
| y ↦ $z_0$ | y ↦ b |
| z ↦ b | z ↦ $z_0$ |

a ⟶ b ⇄ $z_0$

restricted chase:

 halts after one rule application

# RESTRICTED CHASE: NATURAL BUT TRICKY
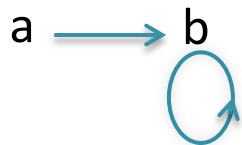
- For the same KB, some derivations may halt while others may not

$F$ : p(a,b)     $R_1$: p(x,y) $\rightarrow$ $\exists$z p(y,z)
$R_2$: p(x,y) $\rightarrow$ p(y,y)

If $R_1$ is always applied before $R_2$ for a given homomorphism of p(x,y):

$$a \longrightarrow b \longrightarrow z_0 \longrightarrow z_1 \quad \ldots$$

If $R_2$ is applied first:
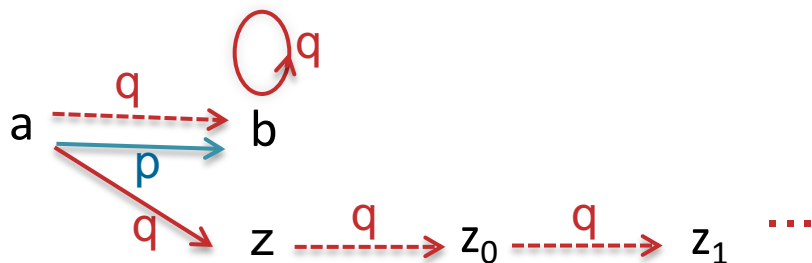
$$a \longrightarrow b$$

Iterate:

(1)    perform a finite number of rule applications as in the restricted chase

(2)    compute the core of the result

$F$ : p(a,b), q(b,b), q(a,z)                    where z is a variable

$R_1$: p(x,y) → q(x,y)
$R_2$: q(x,y) → ∃z q(y,z)



The restricted chase only checks redundancy of **newly** added atoms ⇒ infinite here

The core chase outputs { p(a,b), q(b,b), q(a,b) }

The core chase allows to detect **global** redundancies

- **Terminating** derivation:
  (1) finite and (2) there is no active trigger on the last factbase

- A chase derivation has to be **fair**: no active trigger is indefinitely delayed
  Formally:  if $(R,h)$ is an active trigger on $F_i$

  then there is $F_j$ with $j > i$ such that $F_j$ is obtained by applying $(R,h)$
  or $(R,h)$ is not active anymore on $F_j$

**Terminating = finite and fair**

$R_1$: $p(x,y) \rightarrow \exists z \, p(y,z)$
$R_2$: $p(x,y) \rightarrow p(y,y)$

$F= p(a,b)$

*unfair* infinite derivation: apply only $R_1$ ...

**(semi-) oblivious**: all fair derivations are infinite

**restricted**: some terminating derivations,
some infinite fair derivations

**core**: all fair derivations are terminating

For a chase variant C, C **halts** on a KB K if all **fair** derivations on K are finite

# IN SHORT

All previous chase variants compute **universal** models of a KB

They can be **strictly ordered wrt termination**:

<p align="center"><strong>oblivious &lt; semi-oblivious = skolem &lt; restricted &lt; core</strong></p>

[X &lt; Y means that: for any KB K, if X-chase halts on K then Y-chase halts on K
and there is a KB on which Y-chase halts but not X-chase ]

Only the **core** chase halts if and only if the KB admits a **finite** universal model
but it is costly (involves homomorphisms from the whole factbase)

The **O**, **S-O** and **core** chases yield a **unique** result  (up to the name of nulls):
all fair derivations for a given chase variant yield the same result on a given KB
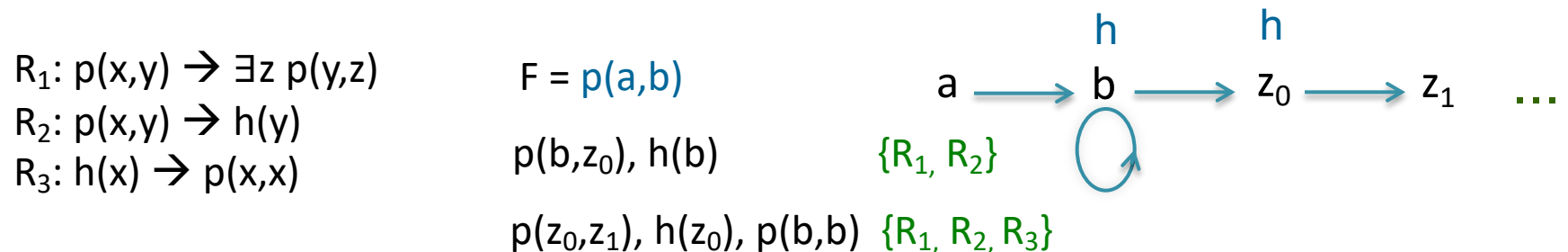but not the **R chase**: we can even have finite and infinite fair derivations

The **R chase** seems to achieve a good tradeoff
redundancy elimination / efficiency of computation (when it stops)
but its behavior is difficult to control

**Open question:**
**is there an ordering strategy that terminates more often than the others?**

- **Breadth-first ordering** is a natural candidate (iterate:
  (1) compute all rule body homomorphisms to the current factbase,
  (2) apply all active triggers according to these homomorphisms)

- however, it is **not optimal for restricted chase** termination

$R_1$: $p(x,y) \rightarrow \exists z\ p(y,z)$
$R_2$: $p(x,y) \rightarrow h(y)$
$R_3$: $h(x) \rightarrow p(x,x)$

F = $p(a,b)$

$p(b,z_0)$, $h(b)$  $\{R_1, R_2\}$

$p(z_0,z_1)$, $h(z_0)$, $p(b,b)$  $\{R_1, R_2, R_3\}$

$$a \xrightarrow{} b \xrightarrow{h} z_0 \xrightarrow{h} z_1 \quad \ldots$$

Optimal order: apply $R_2$ then $R_3$ (ie delay application of $R_1$)

$$a \xrightarrow{} b \quad (\text{with } h \text{ self-loop})$$

- Usual heuristic: at each step, first saturate with all datalog rules, then apply an existential rule

→ would be optimal on this example but it is not true in general