



# Reasoning on Data:

## The Ontology-Mediated Query Answering Problem

**Marie-Laure Mugnier**

University of Montpellier



UNILOG School – Vichy – 2018

# KNOWLEDGE REPRESENTATION AND REASONING (KR)

---



- A field historically at the heart of **Artificial Intelligence**
  - Study **formalisms (or languages)** to
    - **represent** various kinds of human knowledge
    - do **reasoning** on these representations
  - along the tradeoff **expressivity / tractability** of reasoning
- KR languages based on **computational logic**

In this talk: classical first-order logic (FOL)

**Major conferences:  
IJCAI, AAI, KR**

# OVERVIEW OF THE TUTORIAL

---

## **Part 1: Basics**

Knowledge bases, Ontologies

Logical view of Queries and Data

Main KR formalisms to represent and reason with ontologies

Ontology-Mediated Query Answering

## **Part 2: KR formalisms and algorithmic approaches**

## **Part 3: Decidability issues in the existential rule framework**

# KNOWLEDGE BASED SYSTEMS

Knowledge Base  
(KB)



Reasoning  
Services



- **General knowledge on the application domain**  
« *Cats are Mammals* »

## Ontology

- **Factual Knowledge**  
Description of specific individuals, situations, ...

*Félix is a Cat*

## Factbase, Database

## Fundamental tasks

- **Checking the consistency** of the KB
- **Computing answers** to a query over the KB

...

Reasoning algorithms associated with the KR language

Knowledge expressed in a KR language

# WHAT IS AN ONTOLOGY?

---

In computer science:

a **formal specification of the knowledge of a particular domain**

- which allows for machine processing
- that relies on the semantics of knowledge

> **automated reasoning**

Such a specification consists of

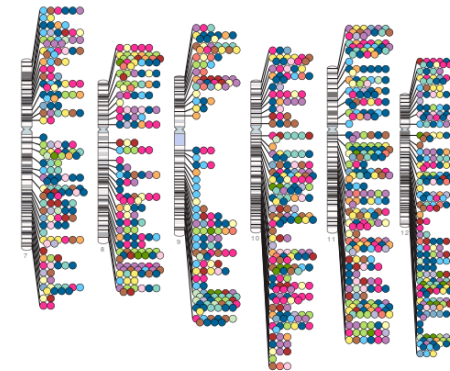
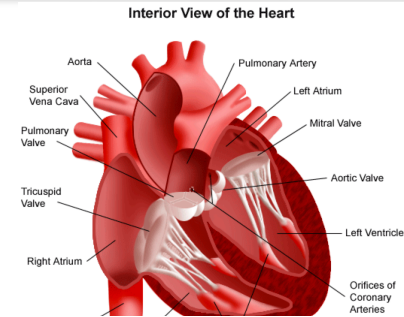
- a **vocabulary** in terms of concepts and relations
- **semantic relationships** between these elements

# EXAMPLES OF ONTOLOGIES

- **Medecine and life sciences :**  
hundreds of available ontologies

- general medical ontologies  
SNOMED CT (400 000 terms)  
GALEN (> 30 000 terms)
- specialized medical ontologies  
FMA (anatomy)  
NCI (cancer), ...
- biology
- agronomy

- **Information systems**  
**of large organizations and corporations**



- Banana
- Barley
- Cassava
- Chickpea
- Common bean
- Cowpea
- Groundnut
- Lentil
- Maize
- Oat (*Global Triticeae*)
- Pearl millet
- Pigeon Pea
- Potato
- Soybean (*USDA & IIT*)
- Sweet Potato
- Rice
- Sorghum
- Vitis (*INRA*)
- Wheat
- Yam

Crop Ontology  
[www.cropontology.org](http://www.cropontology.org)

The screenshot shows the 'Crop Ontology Curation Tool' interface. It features a navigation menu with 'Home', 'About', 'Users', and 'Feedback'. The main content area displays the 'Common Bean Ontology' with a list of terms and their associated identifiers. A sidebar on the right shows the 'Crop Ontology Curation Tool' logo and the 'Integrated Breeding Platform' logo. The interface is designed for managing and curating ontological data.

# AT THE CORE OF ONTOLOGIES: CONCEPTS / CLASSES

- **Concept / class** : a category of entities (objects) that share properties  
In FOL: unary predicate: **Cat**, **Mammal**
- **Instance** of class: a specific member of this class  
In FOL: a term (variable or constant)

- Fundamental relation on classes : **specialization (subsumption)**  
(«is a kind of », «subclass of »)

C2



C1

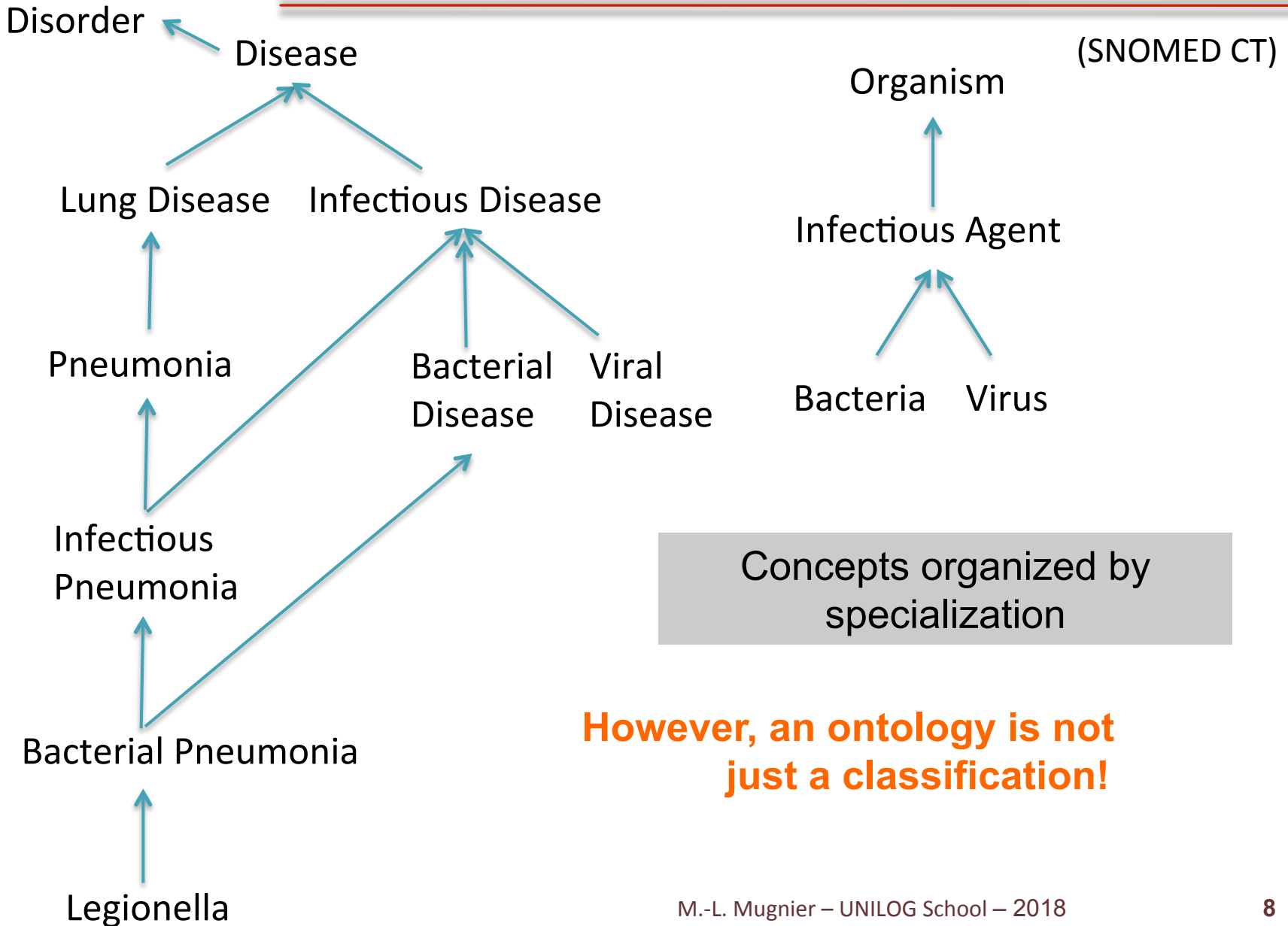
Semantics : every instance of C1 is also instance of C2

$$\forall x (C1(x) \rightarrow C2(x))$$

$$\forall x (Cat(x) \rightarrow Mammal(x))$$

C1 **subclass** of C2

# AT THE HEART OF ONTOLOGIES: CONCEPTS / CLASSES





# ONTOLOGIES ARE MUCH MORE THAN CLASSIFICATIONS

---

« An ontology specifies the **vocabulary** of an application domain and **semantic relationships** between the terms of the vocabulary»

## Vocabulary

1. **concepts / classes**

+ semantic relationships between concepts

2. relations (between instances)

+ semantic relationships between relations

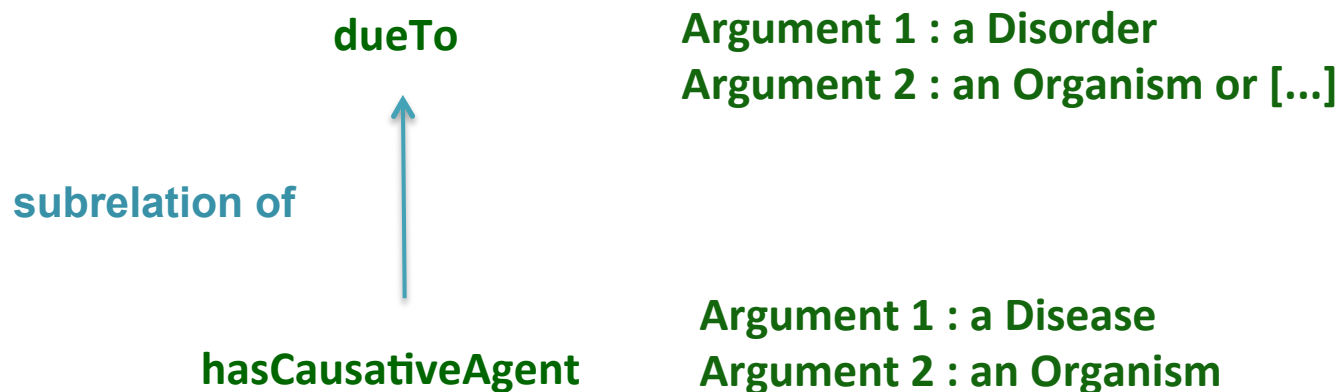
+ properties of concepts

+ properties of relations

+ other axioms that more generally express domain knowledge

# RELATIONS BETWEEN INSTANCES

Often these are binary relations (also called « roles » or « properties »)



$$\forall x \forall y (\text{hasCA}(x,y) \rightarrow \text{dueTo}(x,y))$$

**Signature** of a relation : assigns a maximum concept to each argument  
(« domain » and « range » in OWL)

$$\forall x \forall y (\text{hasCA}(x,y) \rightarrow \text{Disease}(x) \wedge \text{Organism}(y))$$

# EXAMPLES OF OTHER FREQUENT TYPES OF AXIOMS

---

- Negative constraints (disjointness between concepts, relations, ...)

$$\mathbf{Bacteria} \cap \mathbf{Virus} = \emptyset$$

$$\forall x (\mathbf{Bacteria}(x) \wedge \mathbf{Virus}(x) \rightarrow \perp)$$

$$\forall x (\mathbf{Bacteria}(x) \rightarrow \neg \mathbf{Virus}(x))$$

- Necessary and/or sufficient properties of concepts (ex: BacterialDisease)

*A bacterial disease is caused by a bacteria*

$$\forall x (\mathbf{BacterialDisease}(x) \rightarrow \exists y (\mathbf{Bacteria}(y) \wedge \mathbf{hasCausativeAgent}(x,y)))$$

- Properties of relations

**inverse relations:**  $\forall x \forall y (\mathbf{hasPart}(x,y) \leftrightarrow \mathbf{isPartOf}(y,x))$

**symmetry, transitivity, ...**

**functional relation:**  $\forall x \forall y \forall z (\mathbf{isPartOf}(x,y) \wedge \mathbf{isPartOf}(x,z) \rightarrow y = z)$

# WHAT KINDS OF LANGUAGES TO EXPRESS ONTOLOGIES?

---

## Very light languages

Hierarchies of classes

Hierarchies of binary relations (called « properties »)

Signatures of these relations (« domain » and « range »)

→ **OWL DL fragment of RDF Schema** (Semantic Web)

## More expressive fragments of first-order logics

**Description Logics**

**Rule-based languages**

Datalog, existential rules,

RDF deductive rules, Answer Set Programming ...

**From a logical viewpoint:** an ontology is composed of

a **finite set of predicates** (to express concepts and relations)

a **finite set of (closed) formulas** over these predicates

of the form  $\forall X (\text{condition}[X] \rightarrow \text{conclusion}[X])$

# WHAT ARE ONTOLOGIES GOOD FOR?

---

- **provide a common vocabulary**
  - it is easier to share information  
(typically between experts of several domains)
- **constrain the meaning of terms**
  - forces to explicit not-said things and to remove ambiguities  
hence less misunderstandings
- to do **automated reasoning**, basis of high-level services
  - find **implicit links** between pieces of knowledge
  - check the **consistency** of the KB, find errors in modeling
  - enrich **data query answering**

# ONTOLOGY-MEDIATED QUERY ANSWERING (EX: MEDICAL RECORDS)

---

**Query (SQL, SPARQL, MongoDB ...)**

« find all patients affected by a lung disease  
due to a bacteria »



??

Patient P : Diagnosis = « legionella »

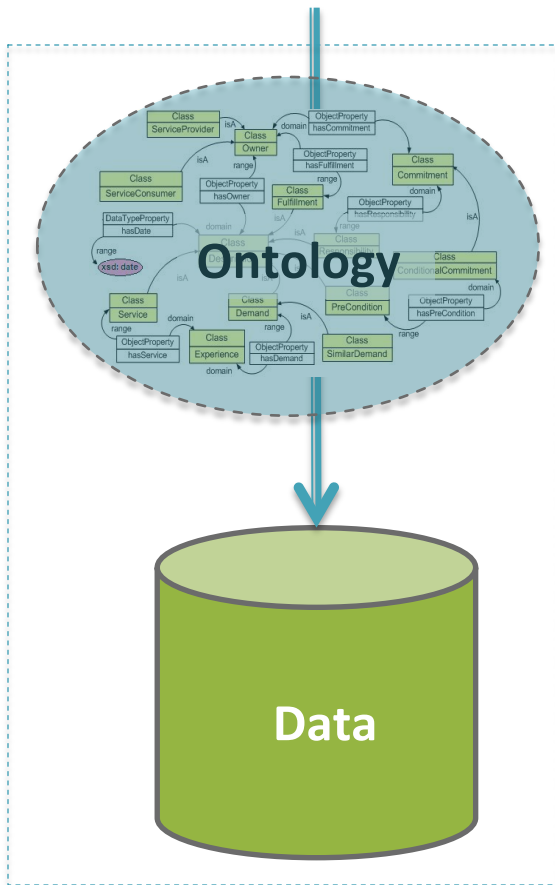
**Database (relational, RDF, NoSQL, ...)**



# ONTOLOGY-MEDIATED QUERY ANSWERING (OMQA)

## Adding an ontological layer on top of data

Query



Knowledge base

1- **Enrich** the vocabulary

allowing to **abstract** from a specific data storage

2 - **Infer** new facts, not explicitly stored,

allowing for **incomplete data** representation





# OMQA EXAMPLE: ONTOLOGICAL KNOWLEDGE

---

A legionella is bacterial pneumonia

$$\forall x (\text{Legionella}(x) \rightarrow \text{BacterialPneumonia}(x))$$

A bacterial pneumonia is a pneumonia

A pneumonia is a lung disease

A bacterial pneumonia is caused by a bacteria

$$\forall x (\text{BacterialPneumonia}(x) \rightarrow \exists y (\text{hasCausativeAgent}(x,y) \wedge \text{Bacteria}(y)))$$

If  $x$  is caused by  $y$  then  $x$  is due to  $y$

$$\forall x \forall y (\text{hasCausativeAgent}(x,y) \rightarrow \text{dueTo}(x,y))$$

If the diagnosis of a patient  $x$  contains a disease  $y$  then  $x$  is affected by  $y$

$$\forall x \forall y ((\text{Diagnosis}(x,y) \wedge \text{Disease}(y)) \rightarrow \text{isAffectedBy}(x,y))$$

# FACTBASE

**Factbase** : usually a set of **ground atoms** (on the ontological vocabulary)  
seen as the conjunction of these atoms

« *The diagnosis for the patient P is legionella* »

$F = \{ \text{Patient}(P), \text{Diagnosis}(P,M), \text{Legionella}(M) \}$

**A relational database** may naturally be viewed as a factbase

Relational schema :            finite set R of relations     $\rightarrow$  predicates  
   infinite domain of values    $\rightarrow$  constants

Instance of a relation r :    finite set of tuples on r     $\rightarrow$  atoms on r

r	
attr1	attr2
a1	a2
a2	a3
a1	a1

$\{ r(a1,a2), r(a2,a3), r(a1,a1) \}$

Database instance = { instance for each r in R }     $\rightarrow$  factbase

# CONJUNCTIVE QUERIES (CQ)

---

« *find all patients affected by a lung disease due to a bacteria* »

$q(x) = \exists y \exists z (\text{Patient}(x) \wedge \text{isAffBy}(x,y) \wedge \text{LungDisease}(y) \wedge \text{dueTo}(y,z) \wedge \text{Bacteria}(z))$

A **CQ** is an **existentially quantified conjunction of atoms**

The **free variables** are the **answer variables**

If closed formula: **Boolean CQ**

Datalog notation

$\text{ans}(x) \leftarrow \text{Patient}(x), \text{isAffBy}(x,y), \text{LungDisease}(y), \text{dueTo}(y,z), \text{Bacteria}(z)$

Select-Join-Project queries in relational algebra (SQL)

SELECT ... FROM ... WHERE *<join conditions>*

SPARQL (semantic web queries)

SELECT ... WHERE *<basic graph pattern>*

# ANSWERS TO A CONJUNCTIVE QUERY

---

- The **answer** to a Boolean CQ  $q$  in  $F$  is *yes* if  $F \models q$     *yes* = ()
- Let the CQ  $q(x_1, \dots, x_k)$ . A tuple  $(a_1, \dots, a_k)$  of *constants* is an **answer** to  $q$  with respect to a factbase  $F$  if  $F \models q[a_1, \dots, a_k]$ ,  
where  $q[a_1, \dots, a_k]$  is obtained from  $q(x_1, \dots, x_k)$  by replacing each  $x_i$  by  $a_i$
- Let  $F$  and  $q$  be seen as sets of atoms. A **homomorphism**  $h$  from  $q$  to  $F$  is a mapping from *variables*( $q$ ) to *terms*( $F$ ) such that  $h(q) \subseteq F$

$F \models q()$  iff  $q$  can be mapped by **homomorphism** to  $F$

$(a_1, \dots, a_k)$  is an answer to  $q(x_1, \dots, x_k)$  w.r.t.  $F$  iff  
there is a **homomorphism** from  $q$  to  $F$  that maps each  $x_i$  to  $a_i$

# KEY NOTION: HOMOMORPHISM

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

$$\begin{array}{l} \text{movie}(y) \\ \text{play}(x, y) \end{array}$$

$$F \quad \begin{array}{l} \text{movie}(m1) \\ \text{movie}(m2) \\ \text{movie}(m3) \\ \text{actor}(a) \\ \text{actor}(b) \\ \text{actor}(c) \\ \text{play}(a, m1) \\ \text{play}(a, m2) \\ \text{play}(c, m3) \end{array}$$

**Homomorphism**  $h$  from  $q$  to  $F$ :  
substitution of  $\text{var}(q)$  by  $\text{terms}(F)$   
such that  $h(q) \subseteq F$

$$\begin{array}{l} h1 : x \rightarrow a \\ \quad y \rightarrow m1 \end{array}$$

$$h1(q) = \text{movie}(m1) \wedge \text{play}(a, m1)$$

$$\begin{array}{l} h2 : x \rightarrow a \\ \quad y \rightarrow m2 \end{array}$$

$$h2(q) = \text{movie}(m2) \wedge \text{play}(a, m2)$$

$$\begin{array}{l} h3 : x \rightarrow c \\ \quad y \rightarrow m3 \end{array}$$

$$h3(q) = \text{movie}(m3) \wedge \text{play}(c, m3)$$

**Answers:** obtained by restricting the domains of homomorphisms  
to answer variables

$$\begin{array}{l} x = a \\ x = c \end{array}$$

# ON THE OMQA EXAMPLE

$$q(x) = \exists y \exists z (\text{Patient}(x) \wedge \text{isAffectedBy}(x,y) \wedge \text{LungDisease}(y) \wedge \text{dueTo}(y,z) \wedge \text{Bacteria}(z))$$

« find all patients affected by a lung disease due to a bacteria »

Factbase = { Patient(P), Diagnosis(P,M), Legionella(M) }

« *The diagnosis for the patient P is legionella* »

Legionella *specialisation of LungDisease and BacterialDisease (and Disease)*

**hence LungDisease(M)**

**hence BacterialDisease(M),  
Disease(M)**

$\forall x (\text{BacterianDisease}(x) \rightarrow \exists y (\text{hasCausativeAgent}(x,y) \wedge \text{Bacteria}(y)))$

**hence hasCausativeAgent(M,b) and Bacteria(b)**

$\forall x \forall y (\text{hasCausativeAgent}(x,y) \rightarrow \text{dueTo}(x,y))$

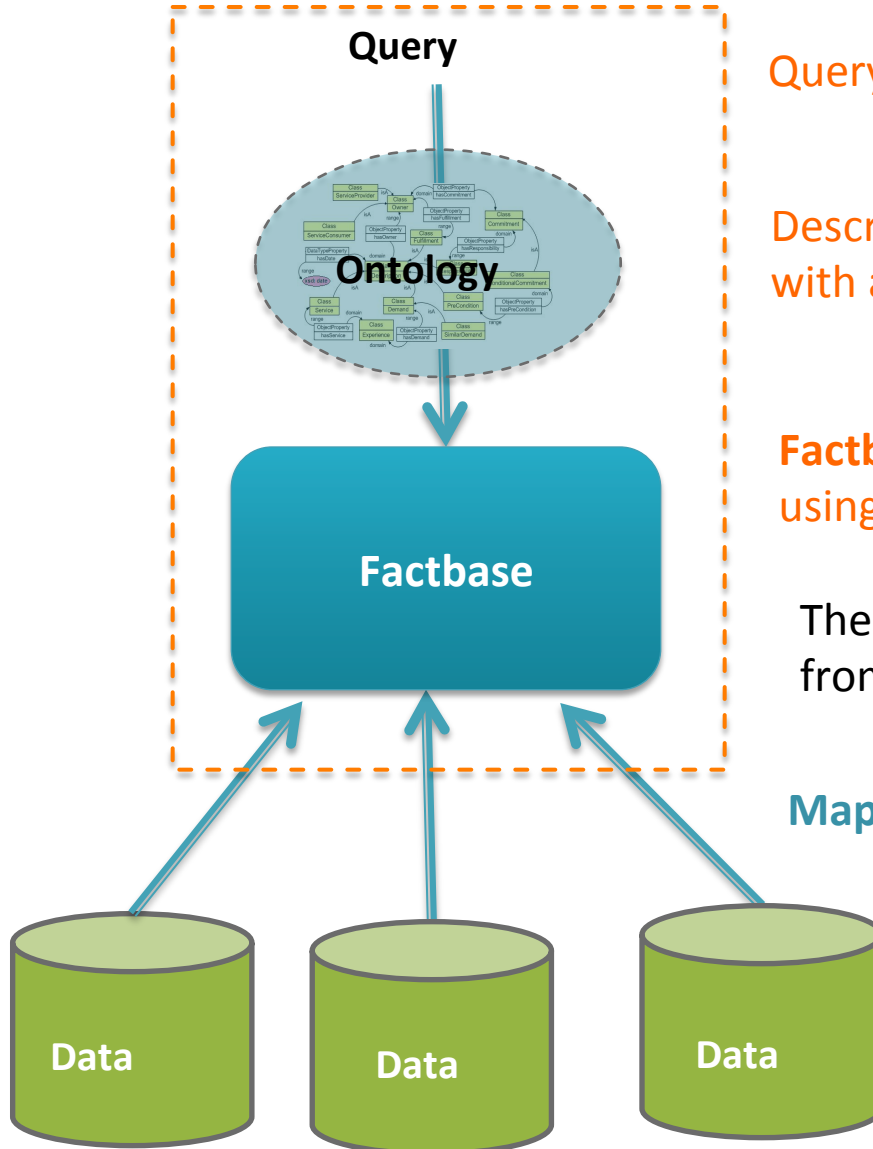
**hence dueTo(M,b)**

$\forall x \forall y ((\text{Diagnosis}(x,y) \wedge \text{Disease}(y)) \rightarrow \text{isAffectedBy}(x,y))$

**hence isAffectedBy(P,M)**

**Answer : x = P**

## Conceptual level



Query using the vocabulary of the ontology

Description of the application domain with a high abstraction level

**Factbase (possibly virtual)**  
using the vocabulary of the ontology

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

**Mappings from data to facts**

{ Database query  $\rightsquigarrow$  Facts }

Independent and heterogeneous data sources



# MAPPINGS

Patient\_T [ID\_PATIENT, NAME,SSN]

Diagnosis\_T[ID\_PATIENT, DISORDER]

Patient /1

Diagnosis / 2

Legionella /1

**Mapping: database query(X)  $\rightsquigarrow$  conjunction with free variables X**

$q(x): \exists n \exists s \text{ Patient\_T}(x,n,s) \rightsquigarrow \text{Patient}(x)$

$q'(x): \exists n \exists s \text{ Patient\_T}(x,n,s) \wedge \text{Diagnostic\_T}(x,y) \wedge y = \text{« Legionella »}$   
 $\rightsquigarrow \exists z (\text{diagnosis}(x,z) \wedge \text{legionella}(z))$

Patient_T			Diagnosis_T	
id	name	ssn	id	dis
P	..	..	P	« Leg. »
..	..	..	..	..
..	..	..	..	..

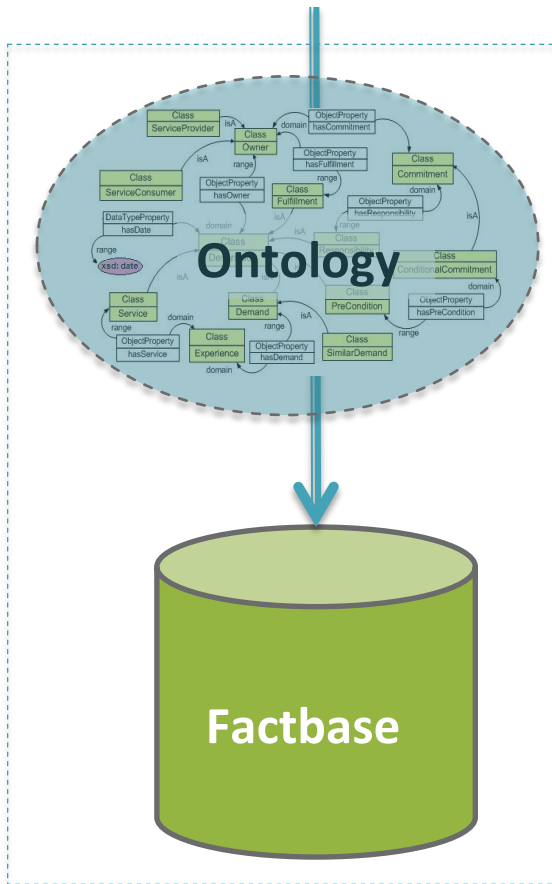
$\rightsquigarrow$

Patient(P)  
Diagnosis(P,M)  
Legionella(M)

# ONTOLOGY-MEDIATED QUERY ANSWERING (OMQA)

Query

(Boolean) conjunctive query  $q$



Theory  $O$  in a suitable FOL fragment

Set of ground atoms (or existentially closed formula)  $F$

Fundamental **decision problem**

$$O, F \models q \quad ?$$

Knowledge base

# OVERVIEW OF THE LECTURE

---

## **Part 1: Basics**

## **Part 2: KR formalisms and algorithmic approaches**

Outline of description logics – Horn DLs

Existential Rules

Materialization approach (forward chaining)

Query rewriting approach (related to backward chaining)

## **Part 3: Decidability issues in the existential rule framework**

# DESCRIPTION LOGICS

---

- A family of KR languages for representing and reasoning with ontologies
- Mostly correspond to **decidable fragments of FOL**  
(related to modal propositional logic, the guarded fragment of FOL, ...)
- **Variable-free syntax**
- Used to be called « **concept languages** »:
  - from concept and role names (unary and binary predicates)
  - and a **set of constructors**
  - define complex concepts (more recently: complex roles)
- An ontology is a set of axioms that state **inclusions between concepts**  
(and between roles)

# DESCRIPTION LOGICS: BUILDING BLOCKS (SYNTAX)

---

## Vocabulary

Atomic concepts:      Human, Parent, Student ...      (unary predicates)

Atomic roles:      parentOf, siblingOf, ...      (binary predicates)

**Complex concepts and roles** can be built using a **set of constructors**  
**(which depends on each particular DL)**

conjunction ( $\sqcap$ ), disjunction ( $\sqcup$ ), negation ( $\neg$ )

Human  $\sqcap$   $\neg$ Parent      Female  $\sqcup$  Male

restricted forms of existential and universal quantification ( $\exists$ ,  $\forall$ )

$\exists$ parentOf.(Female  $\sqcap$  Student)       $\forall$ parentOf.Male

inverse of a role ( $\bar{\phantom{x}}$ ), composition of roles ( $\circ$ )

$\exists$ parentOf $\bar{\phantom{x}}$       parentOf  $\circ$  parentOf

# DESCRIPTION LOGICS: BUILDING BLOCKS (SEMANTICS)

---

To each **concept** is assigned a **FOL sentence with free variable x**

Human

Human(x)

Human  $\sqcap$   $\neg$ Parent

Human(x)  $\wedge$   $\neg$ Parent(x)

$\exists$ parentOf.(Female  $\sqcap$  Student)

$\exists y$  (parentOf(x,y)  $\wedge$  Female(y)  
 $\wedge$  Student(y))

$\forall$  parentOf.Female

$\forall y$  (parentOf(x,y)  $\rightarrow$  Female(y))

To each **role** is assigned a **FOL sentence with 2 free variables x and y**

parentOf o parentOf

$\exists z$  (parentOf(x,z)  $\wedge$  parentOf(z,y))

# DESCRIPTION LOGICS: KNOWLEDGE BASE

---

Knowledge Base = TBox (ontology) + ABox (factbase)

**Tbox:** axioms of the form  $C1 \sqsubseteq C2$   $\forall x ( \text{fol}(C1) \rightarrow \text{fol}(C2) )$   
or  $r1 \sqsubseteq r2$   $\forall x \forall y ( \text{fol}(r1) \rightarrow \text{fol}(r2) )$

$\text{Human} \sqsubseteq \text{Male} \sqcup \text{Female}$   $\forall x ( \text{Human}(x) \rightarrow \text{Male}(x) \vee \text{Female}(x) )$

$\text{Adult} \sqsubseteq \neg \text{Child}$   $\forall x ( \text{Adult}(x) \wedge \text{Child}(x) \rightarrow \perp )$

$\text{Parent} \sqsubseteq \exists \text{parentOf}$   $\forall x ( \text{Parent}(x) \rightarrow \exists y \text{parentOf}(x,y) )$

$\text{HappyFather} \sqsubseteq \forall \text{parentOf.Female}$   $\forall x ( \text{HP}(x) \rightarrow ( \forall y ( \text{parentOf}(x,y) \rightarrow \text{Female}(y) ) )$

$\text{Human} \sqsubseteq \exists \text{parentOf}^{\neg} . \text{Human}$   $\forall x ( \text{Human}(x) \rightarrow \exists y ( \text{parentOf}(y,x) \wedge \text{Human}(y) ) )$

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

**Abox** : set of ground facts  $\text{parentOf}(A,B), \text{Female}(A), \dots$

# DESCRIPTION LOGICS: STANDARD REASONING TASKS

## Standard reasoning tasks on a KB $(\mathcal{T}, \mathcal{A})$

- Concept subsumption  $\mathcal{T} \models C \sqsubseteq D ?$
- Concept satisfiability is C satisfiable w.r.t.  $\mathcal{T} ?$
- KB satisfiability is  $(\mathcal{T}, \mathcal{A})$  satisfiable ?
- Instance checking  $(\mathcal{T}, \mathcal{A}) \models C(b)$ , where b is a constant?

All these tasks can be expressed in terms of KB (un)satisfiability provided that the constructors in the considered DL allow for it

Concept subsumption	$\mathcal{T} \models C \sqsubseteq D$ iff $(\mathcal{T}, \{C(a), \neg D(a)\})$ unsatisfiable
Concept satisfiability	C satisfiable w.r.t. $\mathcal{T}$ iff $(\mathcal{T}, \{C(a)\})$ satisfiable
Instance checking	$(\mathcal{T}, \mathcal{A}) \models C(b)$ iff $(\mathcal{T}, \mathcal{A} \cup \{\neg C(b)\})$ unsatisfiable

Query answering beyond instance checking?  
cannot be reduced to the standard reasoning tasks



# EVOLUTION OF DLs

---

## Standard expressive DL *ALC*

- Concepts:

$$C := \top \mid A \mid C_1 \sqcap C_2 \mid \exists R.C \mid \neg C \mid C_1 \sqcup C_2 \mid \forall R.C$$

- TBox axioms: only concept inclusions

Satisfiability and instance checking in *ALC* are:  
EXPTIME-complete in combined complexity  
coNP-complete in data complexity

Even worse if we add inverse roles: 2EXPTIME-complete in combined complexity

# TWO COMPLEXITY MEASURES FOR QUERY ANSWERING PROBLEMS

**Problem:** Given a KB =  $(O, F)$ , with  $O$  the ontology and  $F$  the factbase, and a query  $q$ , is  $q$  entailed by the KB?

**Combined complexity** (usual complexity measure)

The input is  $O, F$  and  $q$

*E.g.,*  
 $q$  Boolean CQ,  $F$  factbase  
Does  $F \models q$  ?

**Data complexity**

The input is  $F$   
( $O$  and  $q$  supposed to be fixed)

NP-complete (combined)  
PTime (data)

This distinction comes from database theory:  
the size of the query is negligible compared to the size of the data

# EVOLUTION OF DLs

---

## Standard expressive DL *ALC*

- Concepts:

$$C := \top \mid A \mid C_1 \sqcap C_2 \mid \exists R.C \mid \neg C \mid C_1 \sqcup C_2 \mid \forall R.C$$

- TBox axioms: only concept inclusions

Satisfiability and instance checking in *ALC* are:  
EXPTIME-complete in combined complexity  
coNP-complete in data complexity

Even worse if we add inverse roles: 2EXPTIME-complete in combined complexity

Two factors led to the **evolution of description logics**:

1. **practical use** (e.g. SNOMED CT): people mostly use conjunction and existential quantification
2. **complexity too high** for query answering problems

# NEW DLs WITH LOWER COMPLEXITY

DL-Lite<sub>R</sub>

$$B_1 \sqsubseteq B_2 \quad B_1 \sqsubseteq \neg B_2 \quad S_1 \sqsubseteq S_2 \quad S_1 \sqsubseteq \neg S_2$$

where

$$B := A \mid \exists S \quad S := R \mid R^-$$

Large ABoxes  
Query answering

EL

$$C_1 \sqsubseteq C_2$$

where

$$C := \top \mid A \mid C_1 \sqcap C_2 \mid \exists R.C$$

Large TBoxes  
Classification

Common feature: no disjunction (no « true » negation)

Then a satisfiable KB has a **unique canonical model  $M$** :

For any Boolean CQ  $q$ ,  $\text{KB} \models q$  iff  $M$  is a model of  $q$

Reasoning techniques for these lighter DLs  
are very similar to forward or backward chaining in rule-base systems

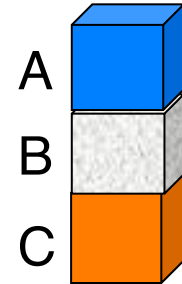
# COMPLEXITY INTRODUCED BY DISJUNCTION OR NEGATION

---

KB ( $\mathcal{T}, \mathcal{A}$ )

$\mathcal{T}$ :  $T \sqsubseteq \text{Blue} \sqcup \text{Other}$

$\mathcal{A}$ :  $\text{Blue}(A), \text{Other}(C), \text{on}(A,B), \text{on}(B,C)$



$q()$ :  $\exists x \exists y (\text{Blue}(x) \wedge \text{on}(x,y) \wedge \text{Other}(y))$

To answer  $q$ , we have to consider two cases:

in each model of the KB, either  $\text{Blue}(B)$  or  $\text{Other}(B)$  holds

Similarly if we replace  $\mathcal{T}$  by:  $\neg\text{Blue} \sqsubseteq \text{Other}$  (equivalent axiom)

Note that  $\text{Other} \sqsubseteq \neg\text{Blue}$  is harmless: it is just a disjointness constraint

## IN SUMMARY

---

**DL ontology (TBox)** has axioms of the form

$$\forall x (\text{fol}(C_1) \rightarrow \text{fol}(C_2))$$

$\forall x \forall y (\text{fol}(r_1) \rightarrow \text{fol}(r_2))$  where  $\text{fol}(r)$  is a path of atomic roles or their inverses

DLs essentially satisfy the **tree model property**:

if a KB is satisfiable then it has a « tree-shaped » model

**With the new DLs:** left and right parts of the implication are both **existentially quantified conjunctions of atoms**

called « **Horn description logics** »

# WHY « HORN DLS » ON AN EXAMPLE

---

$\mathcal{EL}$  Axiom

$$C \sqcap \exists R.T \sqsubseteq \exists S.\exists R.B$$

FOL translation  $\forall x((C(x) \wedge \exists y R(x, y)) \rightarrow \exists u(S(x, u) \wedge \exists v(R(u, v) \wedge B(v))))$

$$\forall x \forall y((C(x) \wedge R(x, y)) \rightarrow \exists u \exists v(S(x, u) \wedge R(u, v) \wedge B(v)))$$

prenex form

$$\forall x \exists u \exists v \forall y(\neg C(x) \vee \neg R(x, y) \vee (S(x, u) \wedge R(u, v) \wedge B(v)))$$

Let us skolemize ( $u$  and  $v$  resp. replaced by  $f_1(x)$  and  $f_2(x)$ ):

$$\forall x \forall y(\neg C(x) \vee \neg R(x, y) \vee (S(x, f_1(x)) \wedge R(f_1(x), f_2(x)) \wedge B(f_2(x))))$$

we obtain a set of 3 Horn clauses (with skolem terms)

$$(\neg C(x) \vee \neg R(x, y) \vee S(x, f_1(x))) \quad (\neg C(x) \vee \neg R(x, y) \vee R(f_1(x), f_2(x))) \quad (\neg C(x) \vee \neg R(x, y) \wedge B(f_2(x)))$$

Hence the name **Horn description logics**

# EXISTENTIAL RULES

---

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

$x, y, z :$   
*sets of variables*

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

$\forall x ( \text{actor}(x) \rightarrow \exists z \text{play}(x,z) )$

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

*we often simplify by omitting universal quantifiers*

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

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

See « **value invention** » in databases

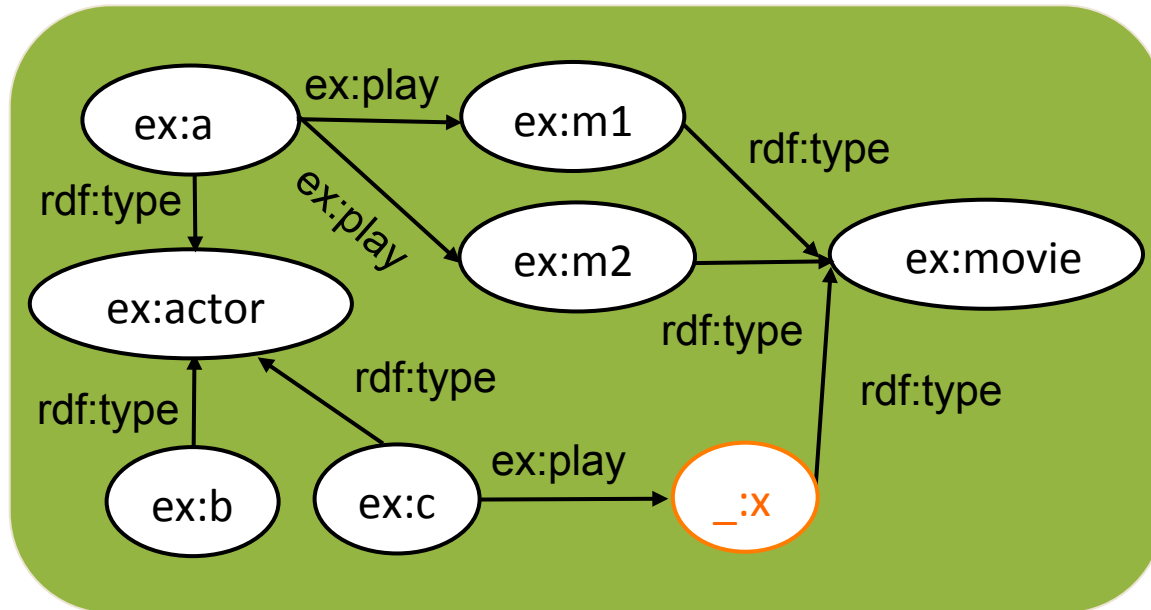


# DATA / FACTS

## Relational database

Movie		Actor		Play	
m_id		a_id		m_id	a_id
m1	...	a	...	a	m1
m2	...	b	...	a	m2
?x	...	c	...	c	?x

## RDF



## Etc.

## Abstraction in first-order logic (FOL)

$\exists x ( \text{movie}(m1) \wedge \text{movie}(m2) \wedge \text{movie}(x) \wedge$   
 $\text{actor}(a) \wedge \text{actor}(b) \wedge \text{actor}(c) \wedge$   
 $\text{play}(a,m1) \wedge \text{play}(a,m2) \wedge \text{play}(c,x) )$

We generalize here the classical notion of a fact by allowing existential variables

**fact / factbase = existentially closed conjunction of atoms**

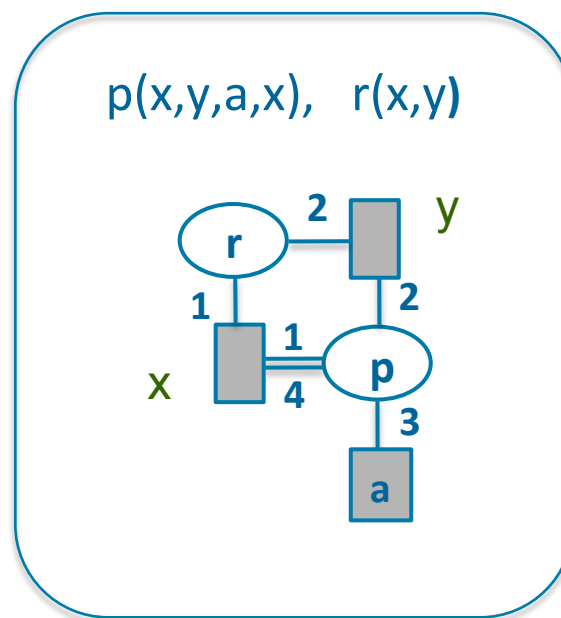
# LABELLED HYPERGRAPH / GRAPH REPRESENTATION

- A fact or a set of facts can be seen as a **set of atoms**

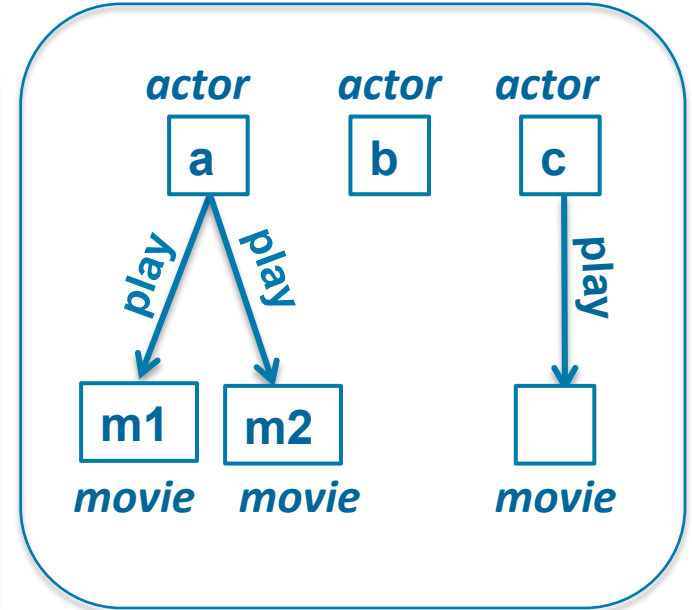
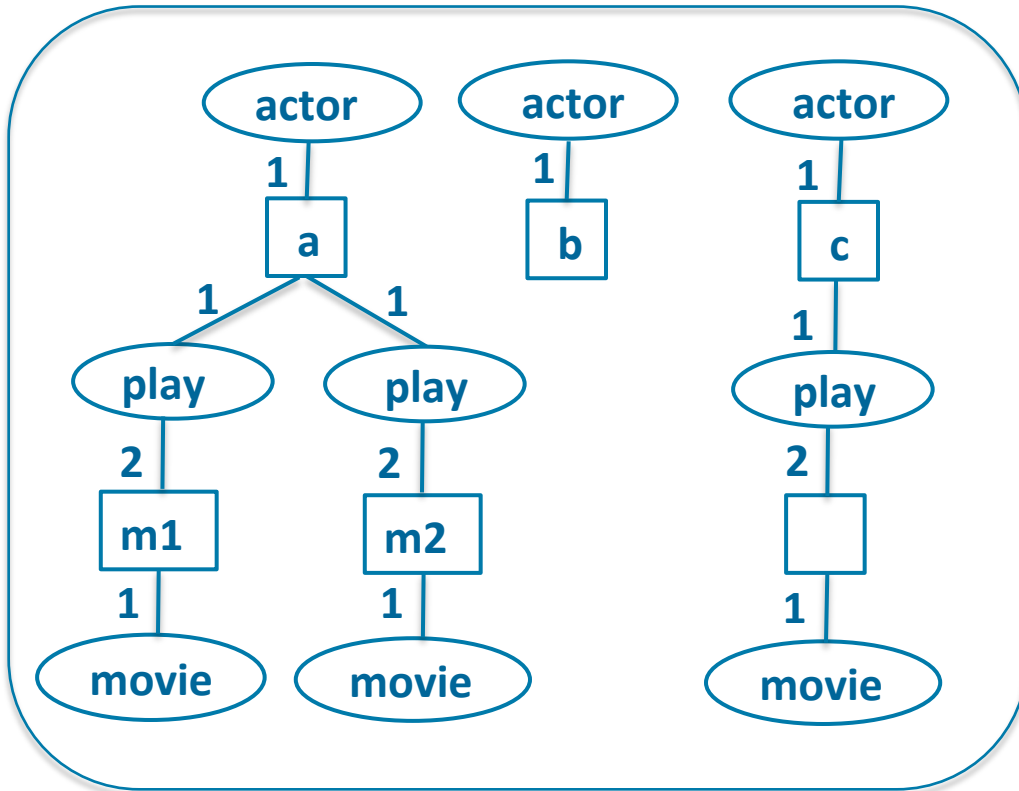
movie(m1), movie(m2), movie(x), actor(a), actor(b), actor(c),  
play(a,m1), play(a,m2), play(c,x)

→ hence a **hypergraph**  
or its associated **bipartite (multi-)graph**

- one (labelled) node per term
- one (labelled) node per atom (~ hyperedge)
- totally ordered edges



movie(m1), movie(m2), movie(x), actor(a), actor(b), actor(c),  
 play(a,m1), play(a,m2), play(c,x)

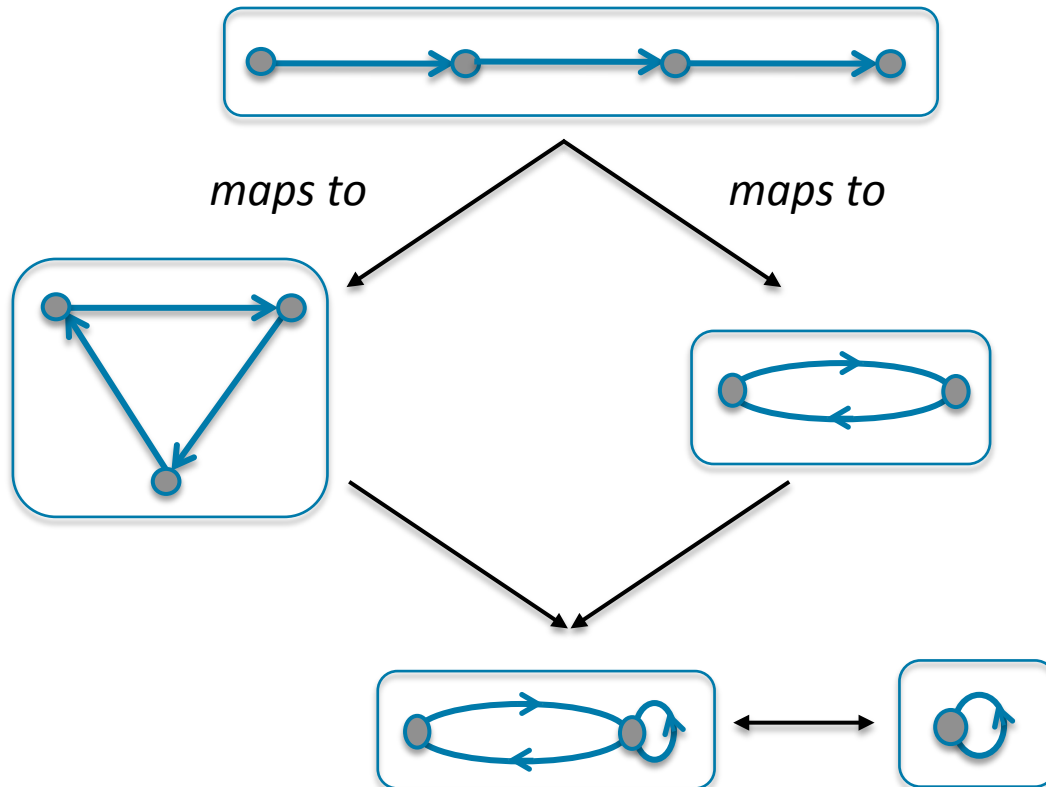


If predicates are at most binary:  
 atom nodes can be replaced by **labels** and **directed edges**

# GRAPH HOMOMORPHISMS (1)

- Let  $G_1=(V_1,E_1)$  to  $G_2=(V_2,E_2)$  be classical graphs.

**Homomorphism**  $h$  from  $G_1$  to  $G_2$ : mapping from  $V_1$  to  $V_2$  s. t.  
for every edge  $(u,v)$  in  $E_1$ ,  $(h(u),h(v))$  is in  $E_2$

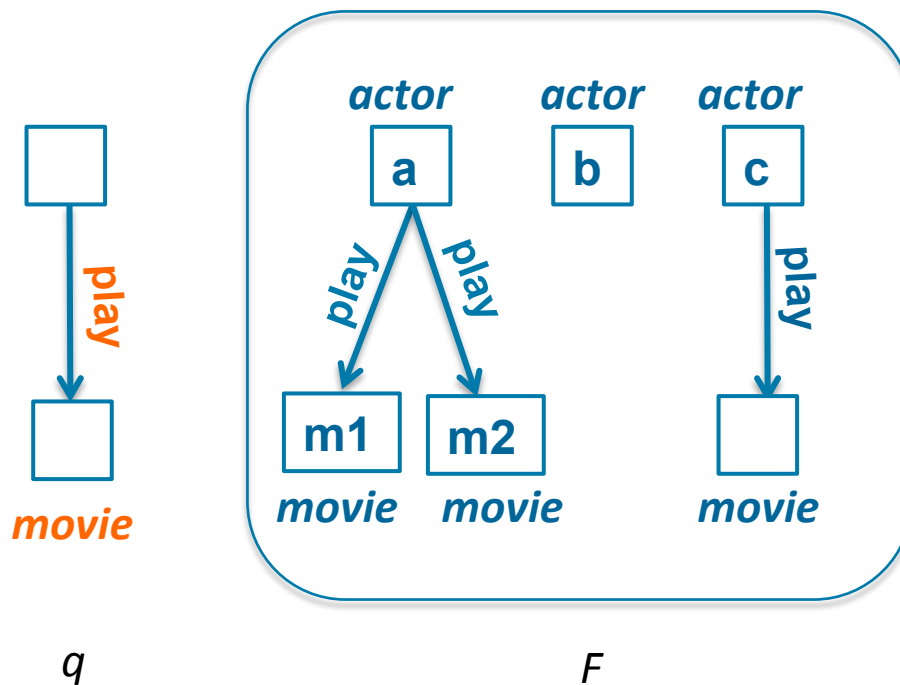


# GRAPH HOMOMORPHISMS (2)

- Let  $G_1=(V_1,E_1)$  to  $G_2=(V_2,E_2)$  be classical graphs.

**Homomorphism**  $h$  from  $G_1$  to  $G_2$ : mapping from  $V_1$  to  $V_2$  s. t.  
for every edge  $(u,v)$  in  $E_1$ ,  $(h(u),h(v))$  is in  $E_2$

- If there are labels: they have to be “kept” as well

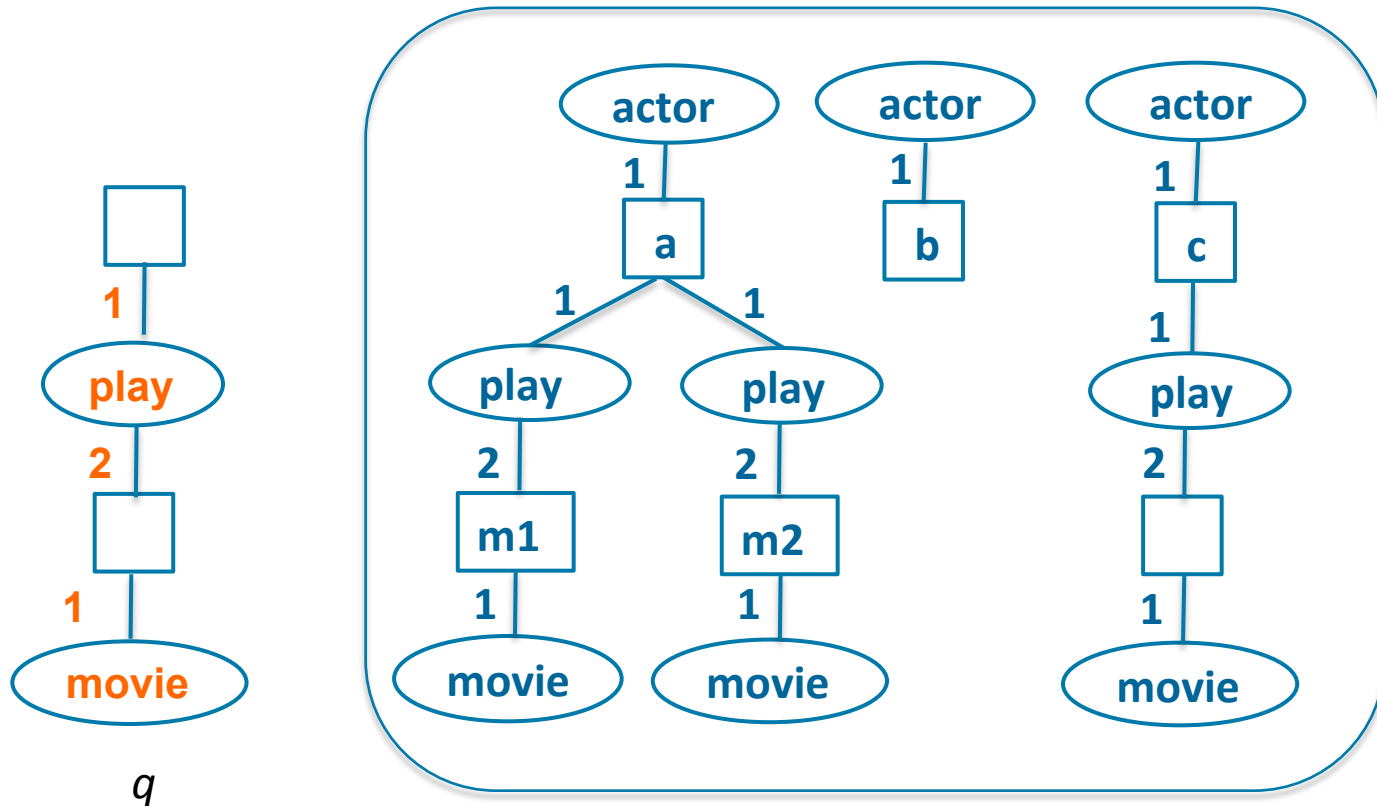


# GRAPH HOMOMORPHISMS (3)

- Let  $G_1=(V_1,E_1)$  to  $G_2=(V_2,E_2)$  be classical graphs.

**Homomorphism**  $h$  from  $G_1$  to  $G_2$ : mapping from  $V_1$  to  $V_2$  s. t.  
for every edge  $(u,v)$  in  $E_1$ ,  $(h(u),h(v))$  is in  $E_2$

- If there are labels: they have to be "kept" as well



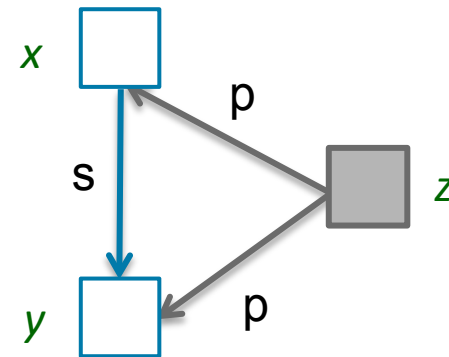
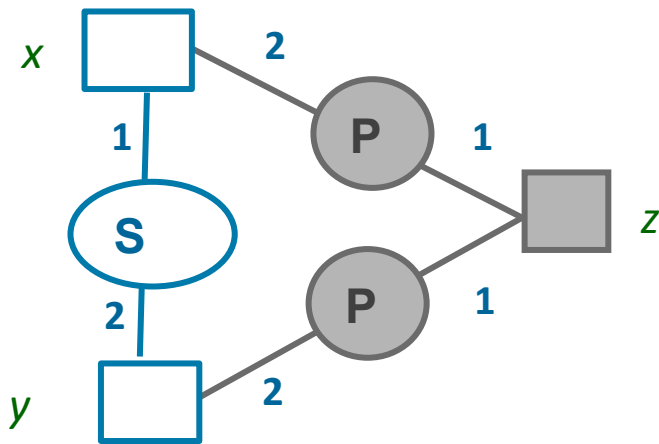
# GRAPH VIEW OF EXISTENTIAL RULES

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

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:

- **frontier**: shared with the body
- **existential** (new ``blank'' nodes)

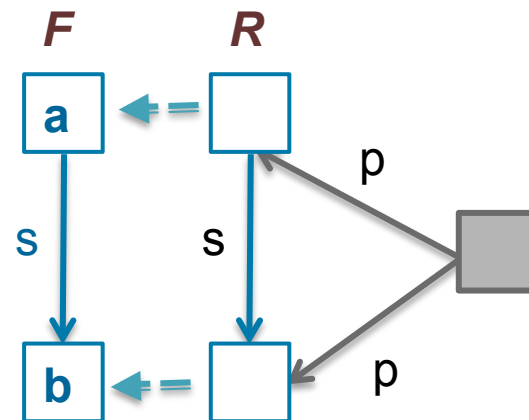
# GENERATION OF FRESH (UNKNOWN) INDIVIDUALS

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

$$F = \text{siblingOf}(a,b)$$

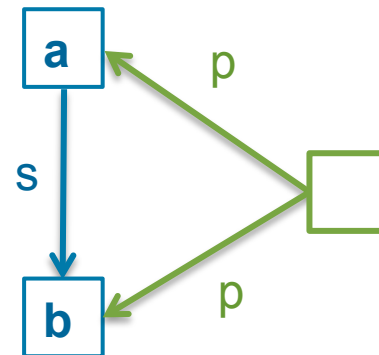
$R$  is **applicable** to  $F$  if there is a **homomorphism**  $h$  from  $\text{body}(R)$  to  $F$

$$\begin{array}{l} x \rightarrow a \\ y \rightarrow b \end{array}$$



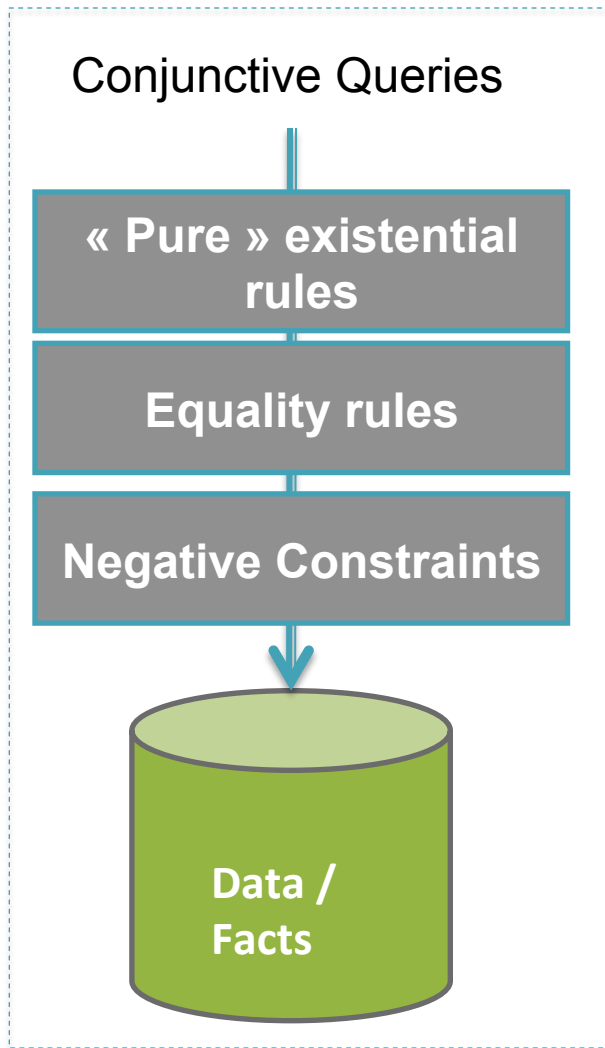
Applying  $R$  to  $F$  w.r.t.  $h$  produces  $F \cup h(\text{head}(R))$  where a new variable is created for each existential variable in  $R$

$$F' = \exists z_0 (\text{siblingOf}(a,b) \wedge \text{parentOf}(z_0,a) \wedge \text{parentOf}(z_0,b))$$





# EXISTENTIAL RULE FRAMEWORK (LOGICAL / GRAPHICAL)



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

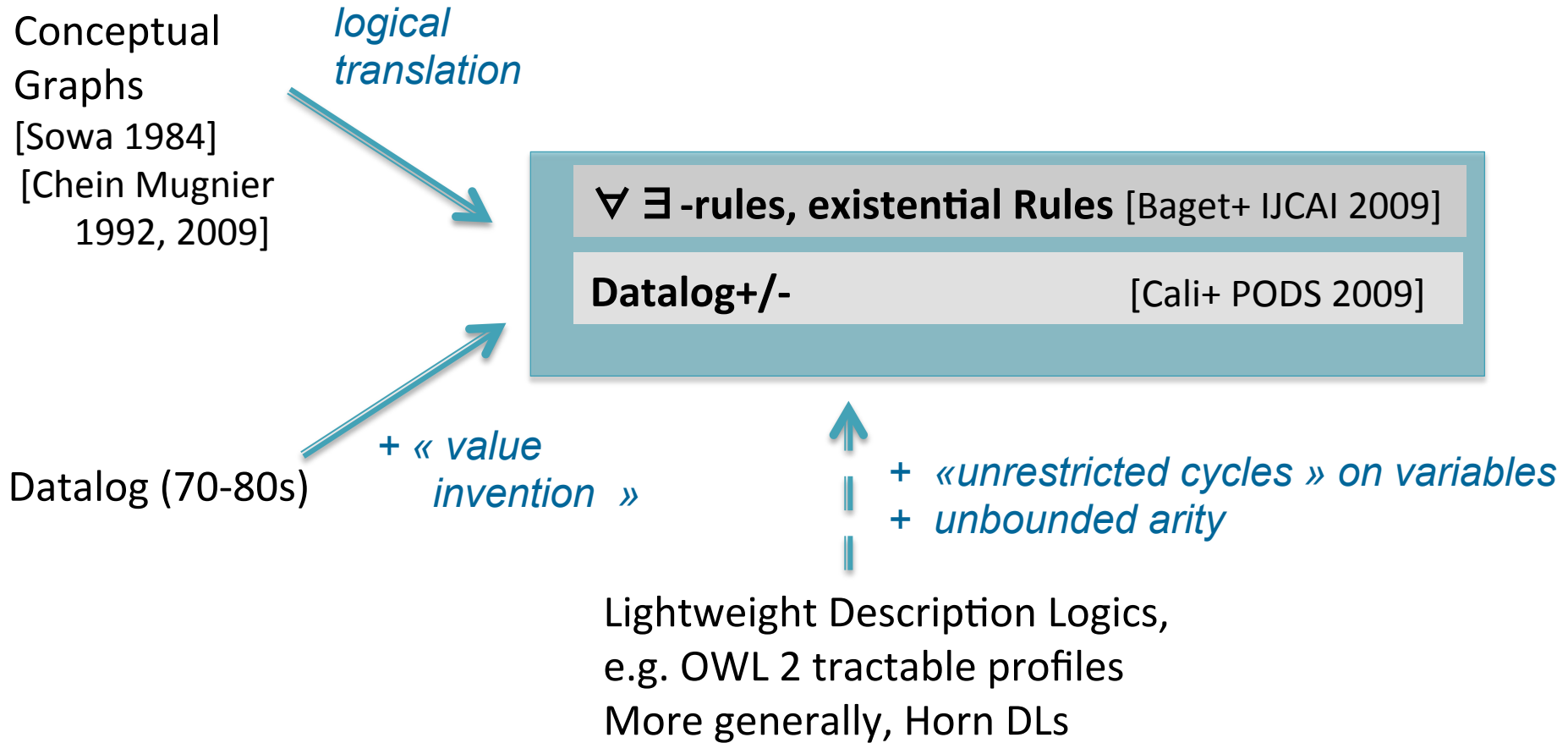
$$\forall x ( \text{actor}(x) \rightarrow \exists z ( \text{movie}(z) \wedge \text{play}(x,z) ) )$$

$$\forall x \forall y \forall z ( \text{movie}(y) \wedge \text{director}(x,y) \wedge \text{director}(z,y) \rightarrow x = z )$$

$$\forall x ( \text{movie}(x) \wedge \text{person}(x) \rightarrow \perp )$$

movie(m1)  
play(a,m1)  
play(c, x)  
...

# MULTIPLE THEORETICAL FOUNDATIONS



- Same logical form as « Tuple-Generating Dependencies » (TGDs) long studied in relational databases

# EXISTENTIAL RULES ARE MORE EXPRESSIVE THAN HORN-DLS

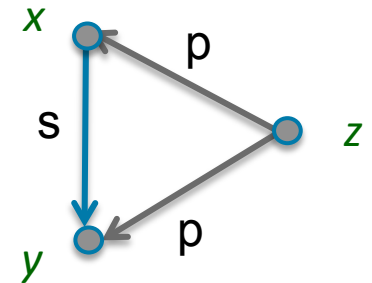
- ◆ The FOL translation of Horn DLs yields existential rules
- ◆ Existential rules are strictly **more expressive**:

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

*cannot be expressed in most DLs because of the « cycle on variables »*

*(needs role composition:  $s \sqsubseteq p \circ p$ )*

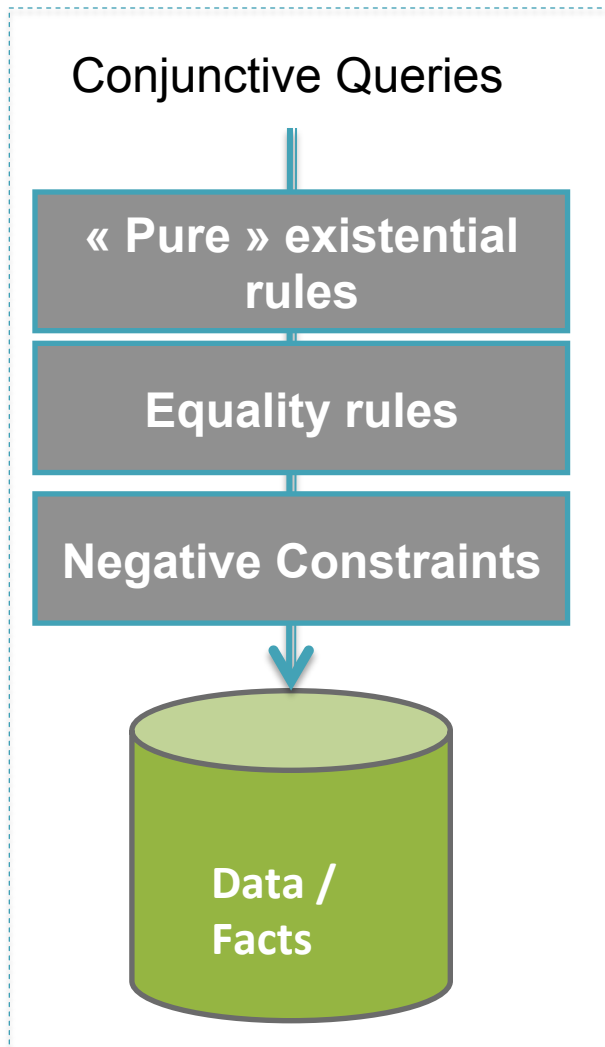
*More complex interactions between variables cannot be expressed at all in DLs*



- ◆ The **unbounded predicate arity** allows for more flexibility:
  - direct translation of database relations
  - adding contextual information is easy (provenance, trust, etc.)

Unsurprisingly, this added expressivity has a cost

# EXISTENTIAL RULE FRAMEWORK



- **Fundamental decision problem**

Input:  $\mathcal{K} = (F, \mathcal{R})$  knowledge base  
 $q$  Boolean conjunctive query

Question: is  $q$  entailed by  $\mathcal{K}$ ?

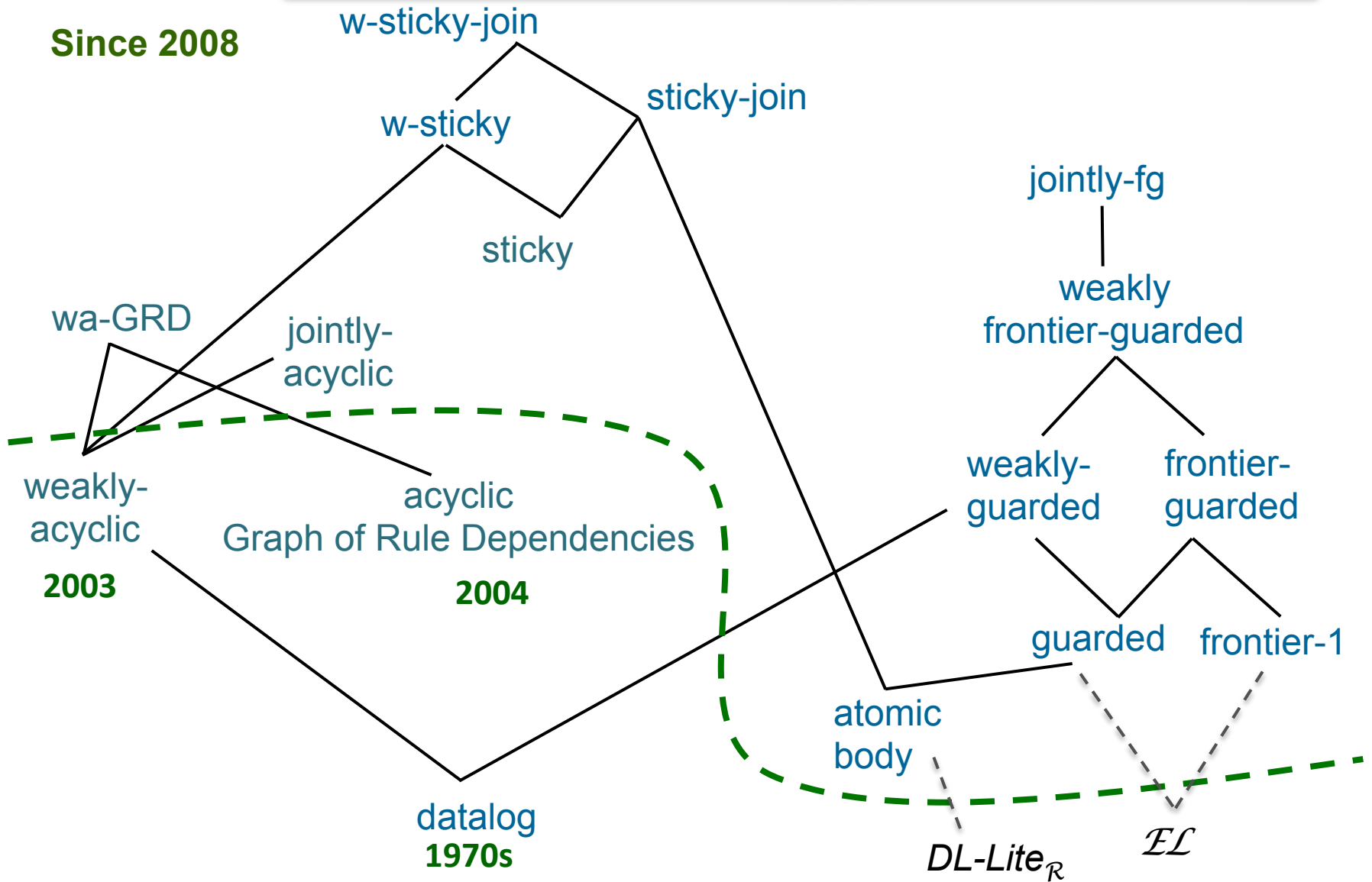
- This problem is **not decidable**

*f.i.* [Beeri Vardi ICALP 1981] on TGDs

even with a single rule [Baget & al. KR 2010]

→ find « decidable » classes of rules  
with good expressivity/tractability tradeoff

# (PARTIAL) MAP OF DECIDABLE CLASSES



# FUNDAMENTAL NOTIONS FOR REASONING IN $\text{FOL}(\exists, \wedge)$

---

- Back to the **positive conjunctive existential fragment** of FOL:  $\text{FOL}(\exists, \wedge)$
- Allows to express **facts** and (Boolean) **conjunctive queries**
- Such formulas can be seen as sets of atoms, labelled graphs, relational structures, ...
- **Homomorphism** is a fundamental notion in this fragment:
  - An interpretation  $\mathcal{I}$  is a model of a sentence  $f$  **iff** there is a **homomorphism** from  $f$  to  $\mathcal{I}$
  - One can define **homomorphisms** between **interpretations**. Then:  
If  $\mathcal{I}_1$  maps to  $\mathcal{I}_2$  then, for any  $f$ ,  $\mathcal{I}_1$  model of  $f \Rightarrow \mathcal{I}_2$  model of  $f$
  - To a formula  $f$ , we assign its **isomorphic model**  $M(f)$  (aka **canonical model**)

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

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

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

$$\begin{aligned} M(f): \quad D &= \{dx, dy, dz, a\} \\ p^{M(f)} &= \{ (dx,dy), (dy,dz) \} \\ r^{M(f)} &= \{ (dx, dz, da) \} \end{aligned}$$

The canonical model  $M(f)$  is **universal**: for all  $M'$  model of  $f$ ,  $M(f)$  maps to  $M'$

for any  $f$  and  $g$  in FOL( $\exists, \wedge$ ),  $g \models f$  iff  $M(g)$  is a model of  $f$  iff  $f$  maps to  $g$

# ADDING RANGE-RESTRICTED (= DATALOG) RULES TO FACTS

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

$\mathcal{R}$  is a set of **range-restricted rules** (i.e.,  $\text{var}(\text{head}) \subseteq \text{var}(\text{body})$ )

$F$  is a factbase (rules with an empty body): ground atoms

By applying rules from  $\mathcal{R}$  starting from  $F$ , a unique result is obtained:  
the **saturation** of  $F$  (denoted by  $F^*$ )

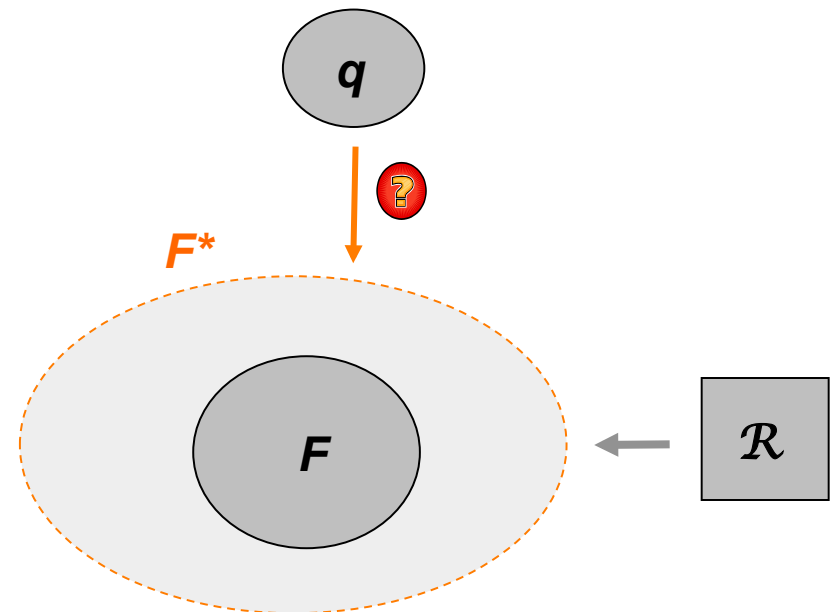
**$F^*$  is finite** since no new variable is created

**$F^*$  is a core** (no redundancies)

The nice properties of  $\text{FOL}(\exists, \wedge)$  are kept:

**$F^*$  is a universal model** of  $\mathcal{K}$

Hence: for any CQ  $q$ ,  $\mathcal{K} \models q$  iff  $q$  maps to  $F^*$





# KNOWLEDGE BASES WITH EXISTENTIAL RULES

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

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

$F$  is a factbase (rules with an empty body): existential conjunctions of atoms

Main change:  $F^*$  can be **infinite**

$R = \text{person}(x) \rightarrow \exists y \text{ hasParent}(x,y) \wedge \text{person}(y)$

$F = \text{person}(a)$

$\wedge \text{person}(y_0) \wedge \text{hasParent}(a, y_0)$

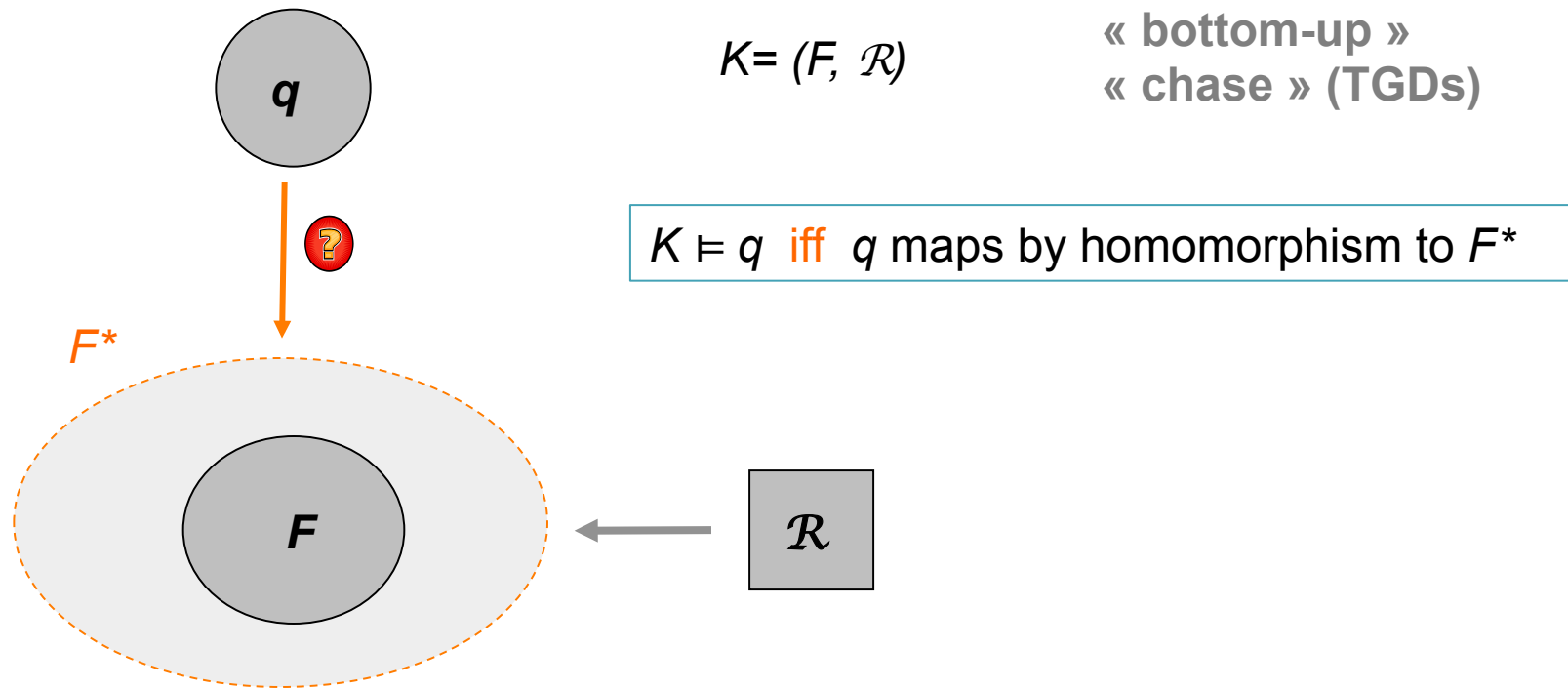
$\wedge \text{person}(y_1) \wedge \text{hasParent}(y_0, y_1)$

Etc.

but it remains a **universal model**

hence  $\mathcal{K} \models q$  iff  $q$  maps to  $F^*$

# APPROACH 1 TO RULES : FORWARD CHAINING / MATERIALISATION

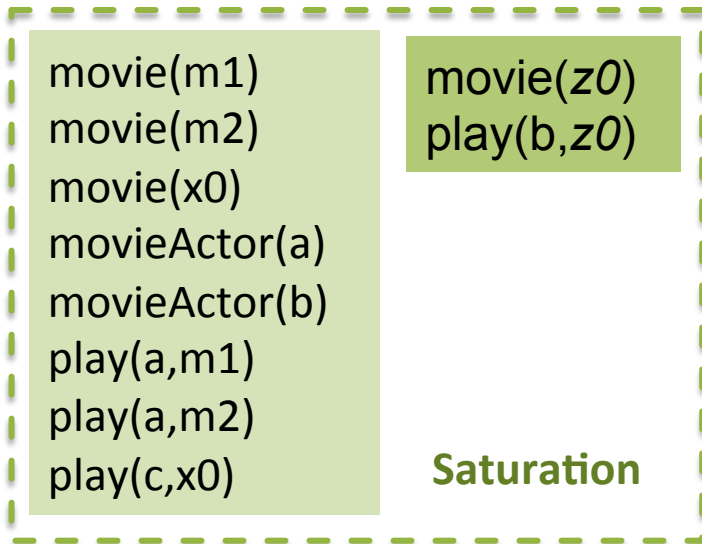


**Pros:** materialisation offline, then online query answering is fast

**Cons:** volume of the materialisation  
needs writing access rights to the data  
not feasible if data is distributed among several databases  
not adapted if data change frequently

# EXAMPLE (MATERIALIZATION)

$\forall x (\text{movieActor}(x) \rightarrow \exists z (\text{movie}(z) \wedge \text{play}(x,z)))$

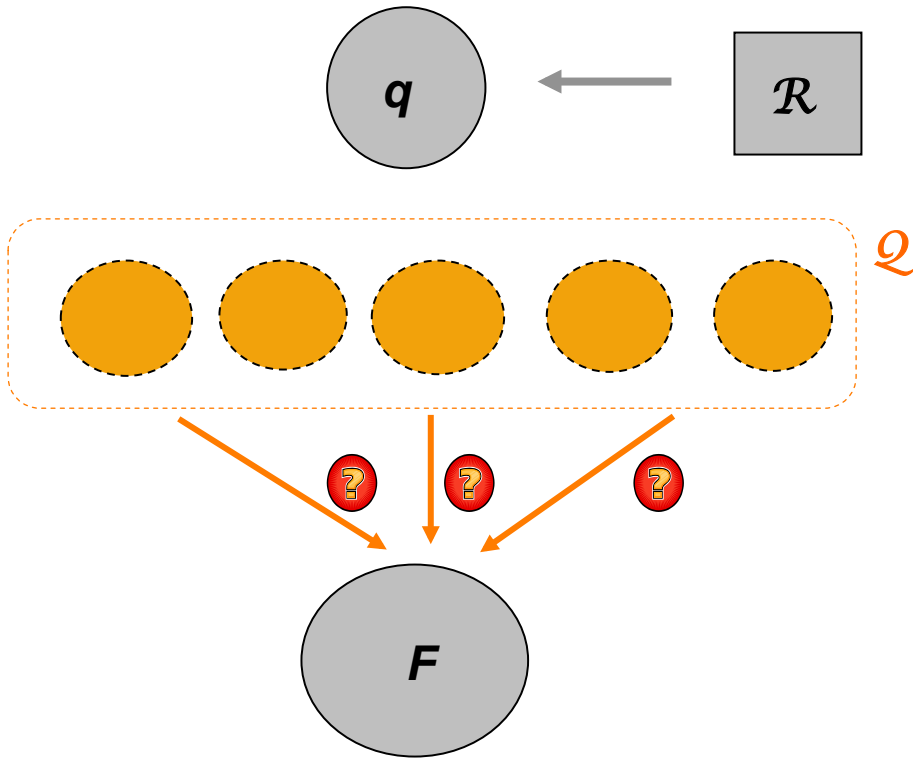


$q(x) = \exists y (\text{movie}(y) \wedge \text{play}(x, y))$

« find those who play in a movie »

x = a	y = m1
x = a	y = m2
x = b	y = z0
x = c	y = x0

# APPROACH 2 TO RULES : BACKWARD CHAINING / QUERY REWRITING



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

« top-down »  
decomposition into  
2 steps [DL-Lite]

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

and more generally into a  
« first-order » query (core SQL query)

Query rewriting is independent from any factbase. For **any**  $F$ ,

$F, \mathcal{R} \models q$  iff  $F \models Q$  (i.e., if  $Q$  is a UCQ: there is  $q_i \in Q$  with  $F \models q_i$ )

**Pros:** independent from the data

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

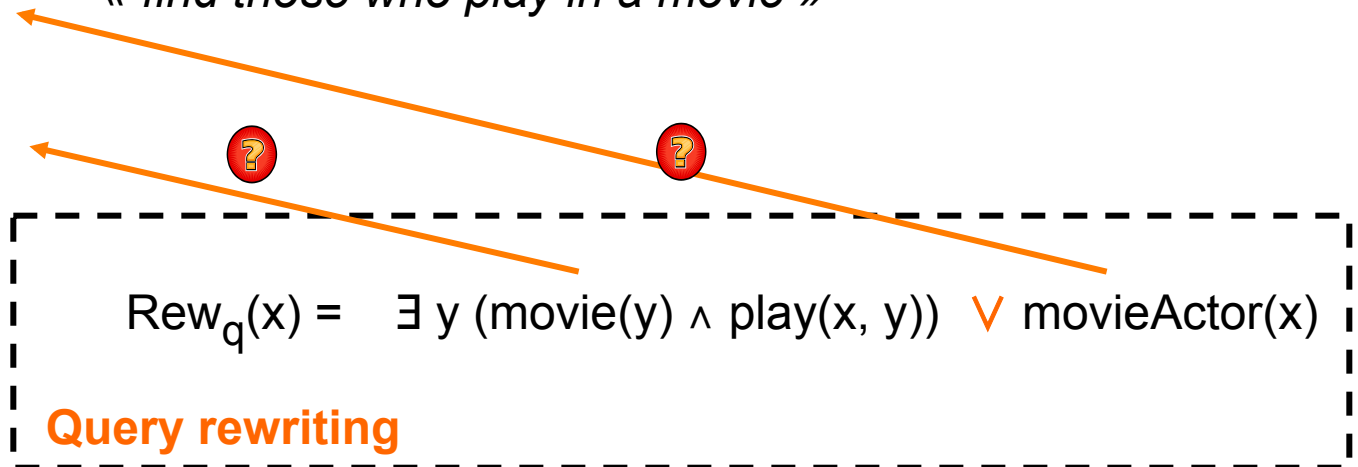
# EXAMPLE

$$\forall x (\text{movieActor}(x) \rightarrow \exists z (\text{movie}(z) \wedge \text{play}(x,z)))$$

movie(m1)  
movie(m2)  
movie(x0)  
movieActor(a)  
movieActor(b)  
play(a,m1)  
play(a,m2)  
play(c,x0)

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

« find those who play in a movie »

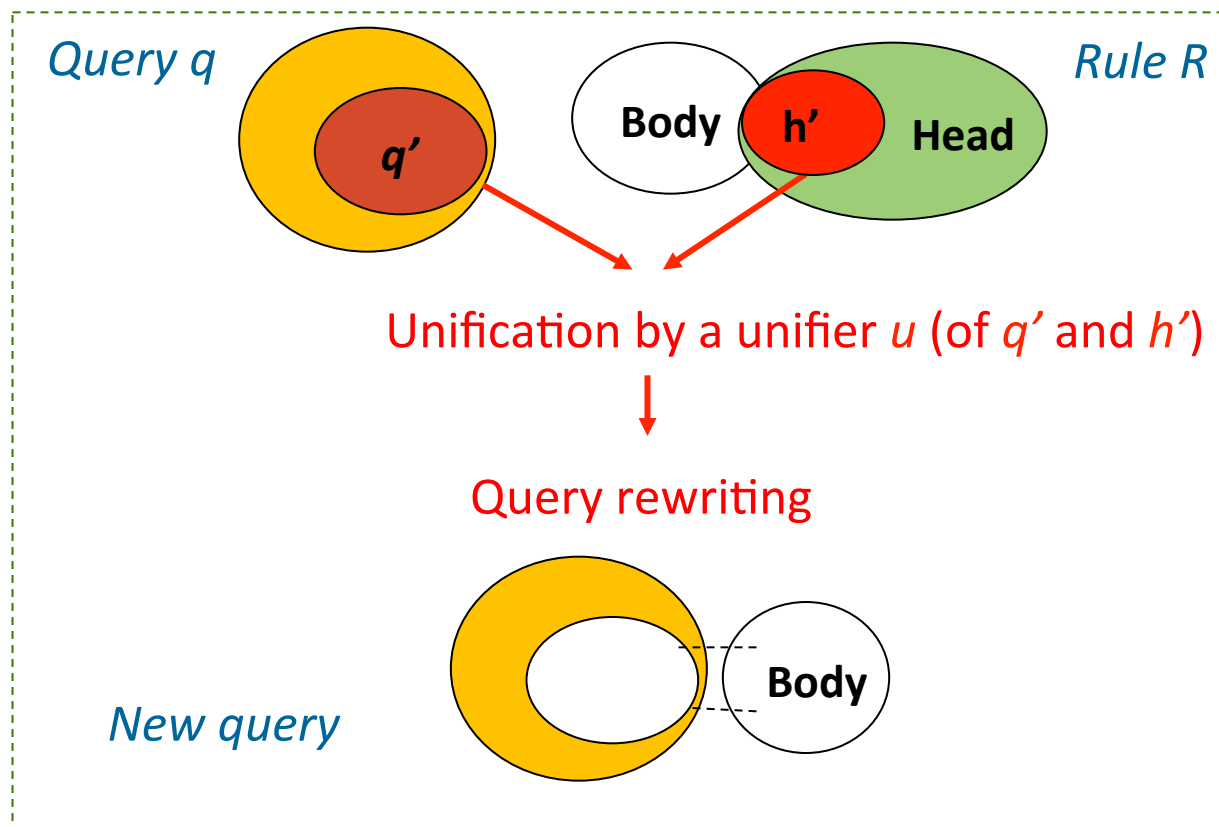


$x = a$	$y = m1$
$x = a$	$y = m2$
$x = c$	$y = x0$

$x = a$
$x = b$

# BACKWARD CHAINING SCHEME

Basic step:



Direct rewriting of  $q$  with  $R$  and  $u = u(q \setminus q') \cup u(\text{body}(R))$

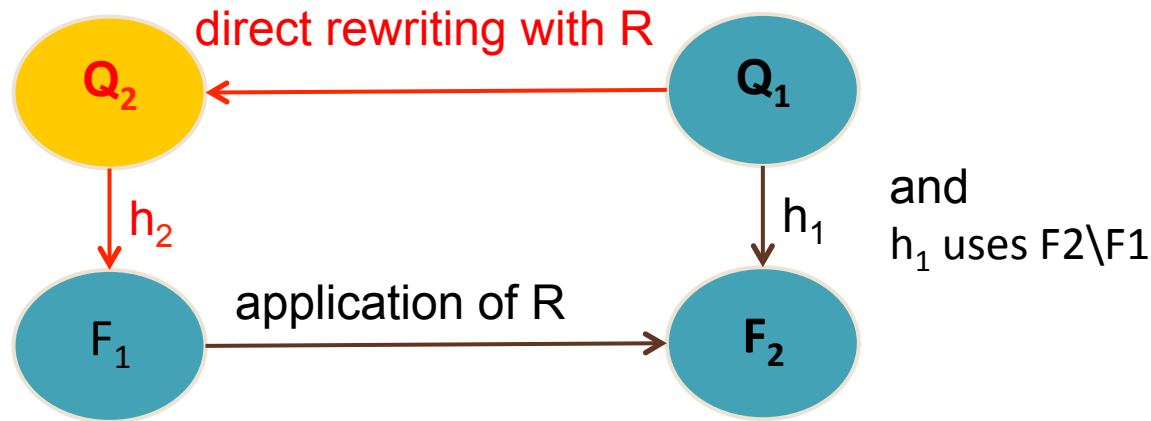
# BASIC PROPERTIES (1)

---

Let  $F_2$  be obtained from  $F_1$  by the application of Rule R

Let a query  $Q_1$  that maps to  $F_2$  by a homomorphism that uses at least one atom brought by R

Then there is  $Q_2$ , a direct rewriting of  $Q_1$  with R, such that  $Q_2$  maps to  $F_1$



The reciprocal property holds

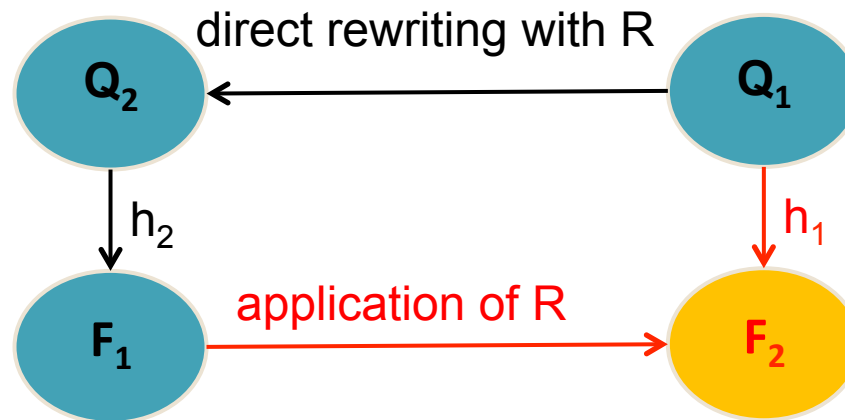
## BASIC PROPERTIES (2)

---

Let  $Q_2$  be a direct rewriting of  $Q_1$  with Rule R

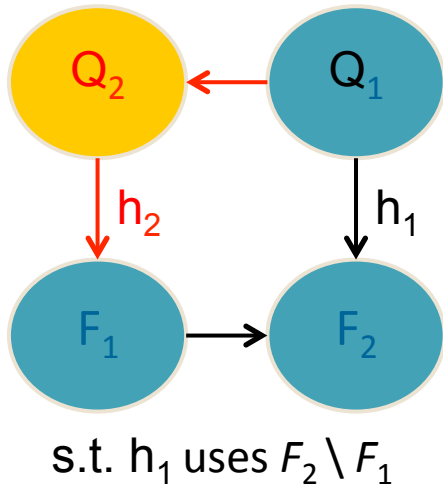
Let  $F_1$  be a factbase such that  $Q_2$  maps to  $F_2$

Then there is an application of R to  $F_1$  that produces  $F_2$  such that  $Q_2$  maps to  $F_1$





# EQUIVALENCE DERIVATION / REWRITING SEQUENCES

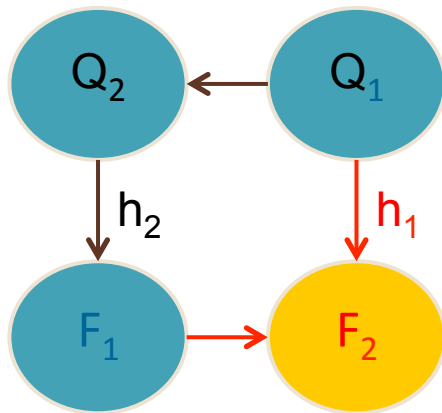


For any conjunctive query  $q$ , for any factbase  $F$ ,  
for any set of rules:

there is a homomorphism from  $q$  to  $F'$ , where  $F'$  is obtained  
from  $F$  by a **rule application sequence** of length  $\leq n$

**iff**

there is a homomorphism from  $q'$  to  $F$ , where  $q'$  is obtained  
from  $q$  by a **rewriting sequence** of length  $\leq n$



# TAKING INTO ACCOUNT EXISTENTIAL VARIABLES IN RULE HEADS (1)

- We want a complete set of **sound** rewritings (set of CQs):

$q_i$  s.t. for any  $F$ , if  $F \models q_i$  then  $F, \mathcal{R} \models q$

$R = \text{person}(x) \rightarrow \exists y \text{ hasParent}(x,y)$

$q = \text{hasParent}(v,w), \text{dentist}(w)$

$u = \{ x \mapsto v, y \mapsto w \}$

$\text{rew}(q,R,u) = q_i = \text{person}(v), \text{dentist}(w)$

$q_i$  is **unsound**:

$F = \text{person}(\text{Maria}), \text{dentist}(\text{Giorgos})$

$F \models q_i$  however  $(F, \{R\})$  does not entail  $q$

(1) If  $w$  in  $q$  is unified with an **existential variable** of  $R$ , then all atoms in which  $w$  occur must be part of the unification

## TAKING INTO ACCOUNT EXISTENTIAL VARIABLES IN RULE HEADS (2)

$R = p(x) \rightarrow \exists z1 \exists z2 r(x,z1), r(x,z2), s(z1,z2)$

$q = r(v,w), s(w,w)$

$u = \{x \mapsto v, z1 \mapsto w, z2 \mapsto w\}$

$\text{rew}(q,R,u) = q_i = p(v)$

$q_i$  is **unsound**:

$F = p(a)$

$F \models q_i$  however  $(F, \{R\})$  does not entail  $q$

(2) An **existential variable** of  $R$  cannot be unified with another term in  $\text{head}(R)$

# PIECE-UNIFIER (FOR BOOLEAN CQs)

A **piece-unifier**  $u$  of  $q' \subseteq q$  and  $h' \subseteq \text{head}(R)$

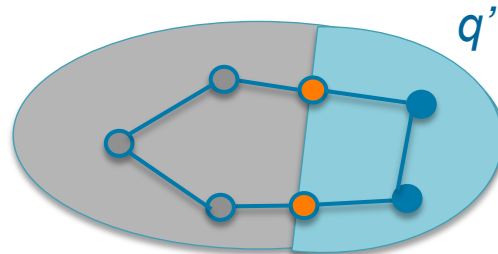
is a substitution of  $\text{var}(q' + h')$  by  $\text{terms}(q' + h')$  [if  $x$  is unchanged, we write  $u(x) = x$ ]

such that :

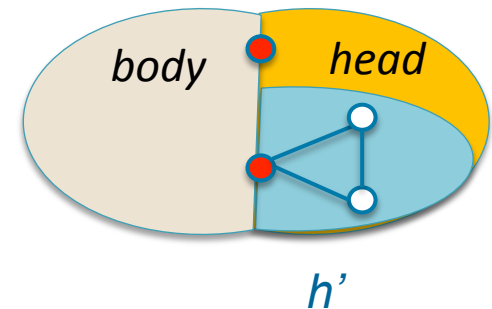
- $u(q') = u(h')$
- existential variables of  $h'$  are unified only with variables of  $q'$  that do not occur in  $(q \setminus q')$  (i.e., if  $x$  is existential and  $u(x) = u(t)$ , then  $t$  is a variable of  $q'$  and not of  $(q \setminus q')$ )

- variables shared by  $q'$  and  $(q \setminus q')$

Query  $q$



Rule  $R$



To extend the notion to general CQs:

universal variables cannot be unified with answer variables

# EXAMPLE

$$R = \text{twin}(x,y) \rightarrow \exists z \text{ motherOf}(z,x) \wedge \text{motherOf}(z,y)$$

$$q = \text{motherOf}(v,w) \wedge \text{motherOf}(v,t) \wedge \text{Female}(w) \wedge \text{Male}(t) ?$$

$$R = \text{twin}(x,y) \rightarrow \exists z \text{ motherOf}(z,x) \wedge \text{motherOf}(z,y)$$

$$q = \text{motherOf}(v,w) \wedge \text{motherOf}(v,t) \wedge \text{Female}(w) \wedge \text{Male}(t) ?$$

$$\text{piece-unifier } u_2 = \{z \mapsto v, x \mapsto w, y \mapsto t\}$$

$$\text{rewrite}(q,R,u_2) = \text{twin}(w,t) \wedge \text{Female}(w) \wedge \text{Male}(t)$$

$$R = \text{twin}(x,y) \rightarrow \exists z \text{ motherOf}(z,x) \wedge \text{motherOf}(z,y)$$

$$q = \text{motherOf}(v,w) \wedge \text{motherOf}(v,t) \wedge \text{Female}(w) \wedge \text{Male}(t) ?$$

$$\text{piece-unifier } u_1 = \{z \mapsto v, x \mapsto w, y \mapsto w\}$$

$$\text{rewrite}(q,R,u_1) = \text{motherOf}(v,t) \wedge \text{Female}(w) \wedge \text{Male}(w) \wedge \text{twin}(w,w)$$

If we rewrite again  
this query we could  
remove the first atom

# WHAT IF WE SKOLEMIZED RULES?

$R = \text{person}(x) \rightarrow \exists y \text{ hasParent}(x,y)$   
 $q = \text{hasParent}(v,w), \text{dentist}(w)$   
 $u = \{ x \mapsto v, y \mapsto w \}$   
 $\text{rew}(q,R,u) = q_i = \text{person}(v), \text{dentist}(w)$

$q_i$  is **unsound**:  
 $F = \text{person}(\text{Maria}), \text{dentist}(\text{Giorgos})$   
 $F \models q_i$  however  $(F, \{R\})$  does not entail  $q$

$\text{Skolem}(R) = \text{person}(x) \rightarrow \text{hasParent}(x, f(x))$

Classical most general unifier of  $\text{hasParent}(x, f(x))$  and  $\text{hasParent}(v, w)$ :  $v \mapsto x$  and  $w \mapsto f(x)$

$\text{rew}(q, R, u) = \text{dentist}(f(x)) \wedge \text{person}(x)$  which cannot be unified with a rule head  
(would not be kept in the output since it contains a skolem function)

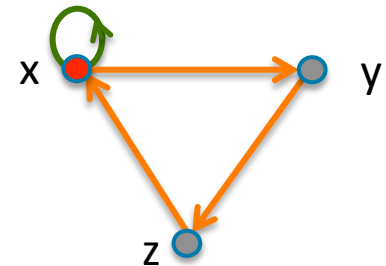
We could skolemize the rules and rely on usual m.g.u.  
then keep only rewritings without skolem function  
but this would create useless intermediate rewritings

# WHY « PIECES »?

A **piece** is a unit of knowledge brought by a rule:

- **Frontier variables** (and constants) act as **cutpoints** to decompose rule heads into pieces (« minimal non-empty subsets glued by existential variables »)

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



- A rule with  $k$  pieces can be decomposed into  $k$  rules, one for each piece, while keeping the same body

$$b(x) \rightarrow \exists y \exists z p(x,y) \wedge p(y,z) \wedge p(z,x)$$

$$b(x) \rightarrow q(x,x)$$

- It cannot be further decomposed (except by introducing new predicates)

# DECOMPOSITION OF RULES INTO ATOMIC HEAD RULES (1)

R:  $b(x) \rightarrow \exists y \exists z p(x,y) \wedge p(y,z) \wedge p(z,x)$

rule with single-piece head

Decomposition into rules with atomic head  
by introducing a **fresh predicate**

$R_0: b(x) \rightarrow \exists y \exists z p_R(x,y,z)$

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

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

$R_3: p_R(x,y,z) \rightarrow p(z,x)$

We lose the structure of the head

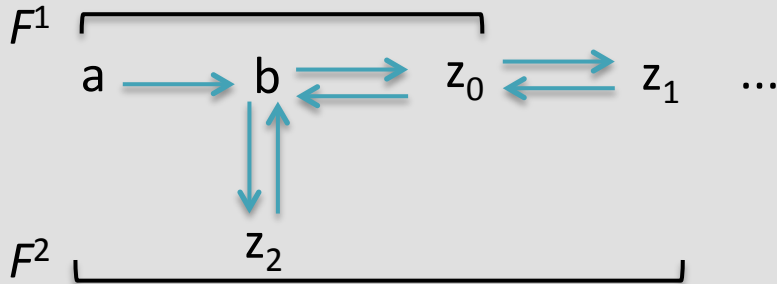
- much less efficient query rewriting
- may even lead to lose the property of having a **finite universal model** (if the set of rules has this property)



# DECOMPOSITION OF RULES INTO ATOMIC HEAD RULES (2)

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

$F^2 \equiv F^1$     ( $F^2$  maps to  $F^1$ )



hence  $F^* \equiv F^1$

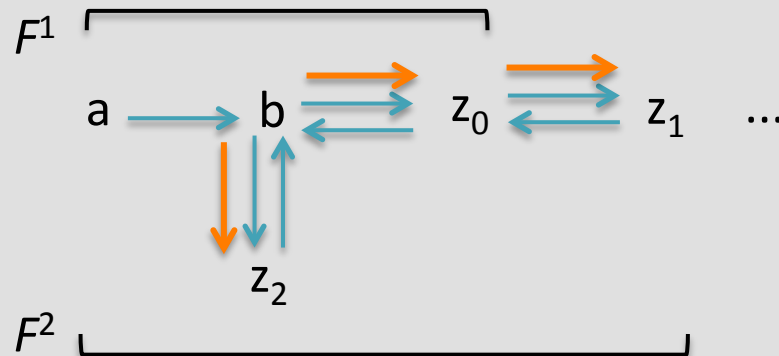
Finite universal model

After decomposition into atomic head rules:

$R_0 : p(x,y) \rightarrow \exists z p_R(y,z)$

$R_1 : p_R(y,z) \rightarrow p(y,z)$

$R_2 : p_R(y,z) \rightarrow p(z,y)$



$F^2 \not\equiv F^1$

No finite universal model

# OVERVIEW OF THE LECTURE

---

## **Part 1: Basics**

## **Part 2: KR formalisms and algorithmic approaches**

## **Part 3: Decidability issues in the existential rule framework**

Undecidability of the fundamental problem

Generic properties that ensure decidability

Main « concrete » decidable classes of existential rules

# SATURATION MAY NOT HALT

---

$R = \text{person}(x) \rightarrow \text{hasParent}(x,y) \wedge \text{person}(y)$

$F = \text{person}(a)$

$\wedge \text{person}(y_0) \wedge \text{hasParent}(a, y_0)$

$\wedge \text{person}(y_1) \wedge \text{hasParent}(y_0, y_1)$

No redundancies are added  
The KB has no finite universal model

However, here: query rewriting with  $R$  is **finite** for any  $q$

# QUERY REWRITING MAY NOT HALT

---

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

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

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

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

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

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

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

There is an infinite number of non-redundant rewritings

However, here: saturation with  $R$  is **finite** for any  $F$

There are cases where both processes do not halt  
(even if the factbase is known)

# UNDECIDABILITY OF THE FUNDAMENTAL PROBLEM

## Fundamental decision problem

Input:  $K = (F, \mathcal{R})$  knowledge base,  $q$  Boolean conjunctive query

Question: is  $q$  entailed by  $K$  ?

This problem is **undecidable** (only semi-decidable)

E.g. proof by **reduction from the word problem in a semi-Thue system**

**Input:** a set  $G$  of rules of the form  $w_i \rightarrow w_j$ , 2 words  $w_0$  and  $w_f$

**Question:** is it possible to derive (exactly)  $w_f$  from  $w_0$  using the rules in  $G$ ?

There is a *one-step derivation* from a word  $w$  to  $w'$  if

there is a rule  $w_i \rightarrow w_j$  in  $G$ , and  $w = w_1 w_i w_2$ ,  $w' = w_1 w_j w_2$

$w'$  is *derived from*  $w$  if

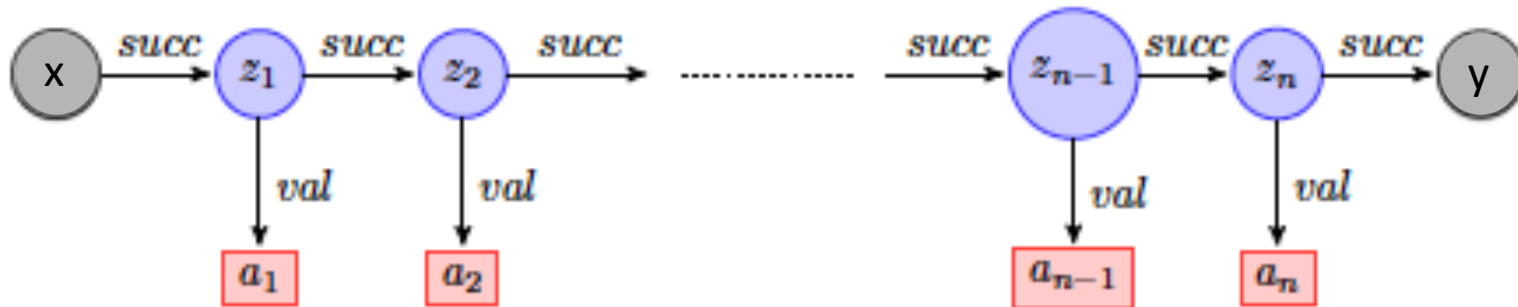
there is a (finite) sequence of one-step derivations from  $w$  to  $w'$

# REDUCTION FROM THE WORD PROBLEM

From  $G$ ,  $w_0$  and  $w_f$  we build a KB  $(F, \mathcal{R})$  and a Boolean CQ  $q$

**Vocabulary** constants: the **letters** occurring in  $G$ ,  $w_0$  and  $w_f$   
+ two special constants **B** and **E**  
binary predicates: **succ** and **val**

To a word  $w = a_1 \dots a_n$  we assign the following graph  $T(w, x, y)$   
where the  $z_i$  are existential variables and  $x, y$  are free

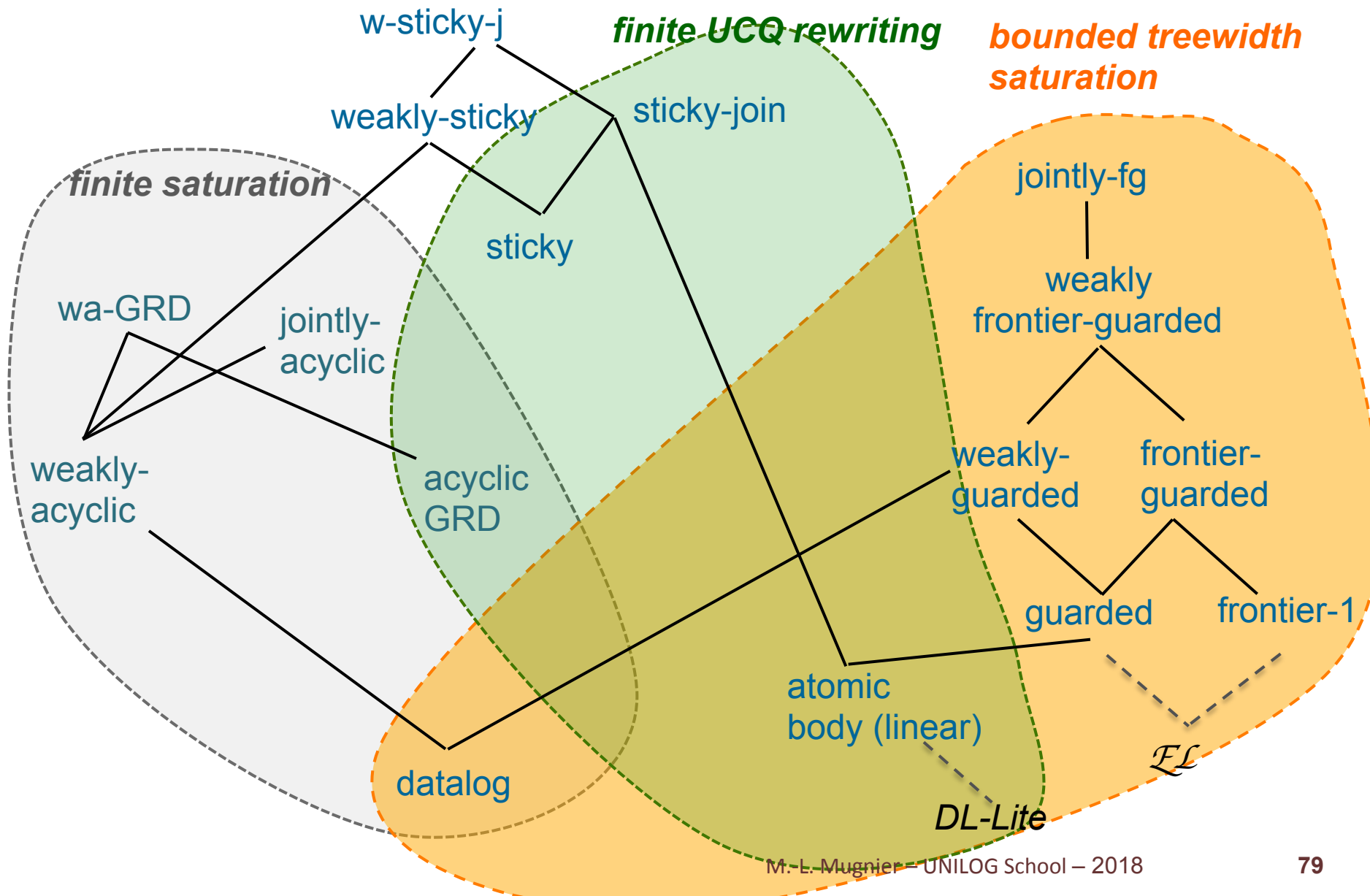


**Factbase**  $F = T(w_0, B, E)$       **Query**  $q = T(w_f, B, E)$

**Set of rules**  $\mathcal{R}$  is obtained by translating each rule  $w_i \rightarrow w_j$  into the existential rule  
 $\forall x \forall y (T(w_i, x, y) \rightarrow T(w_j, x, y))$

**Key:** any word  $w$  derivable from  $w_0$  with  $G$  corresponds to a path  $T(w, B, E)$   
in the saturation of  $F$  by  $\mathcal{R}$ , and reciprocally

# (PARTIAL) MAP OF DECIDABLE CASES



# GENERIC PROPERTIES THAT ENSURE DECIDABILITY

---

Three generic kinds of properties ensuring decidability:

- Saturation by Forward Chaining halts for any factbase (« finite expansion set », *fes*)
- Query rewriting halts for any conjunctive query (« finite unification set », *fus*, or UCQ-rewritability)
- Saturation by Forward Chaining may not halt *but* for any factbase the generated facts have a tree-like structure (« bounded treewidth set », *bts*)

None of these properties is *recognizable* [Baget+ KR 10]

but these properties provide *generic* algorithmic schemes



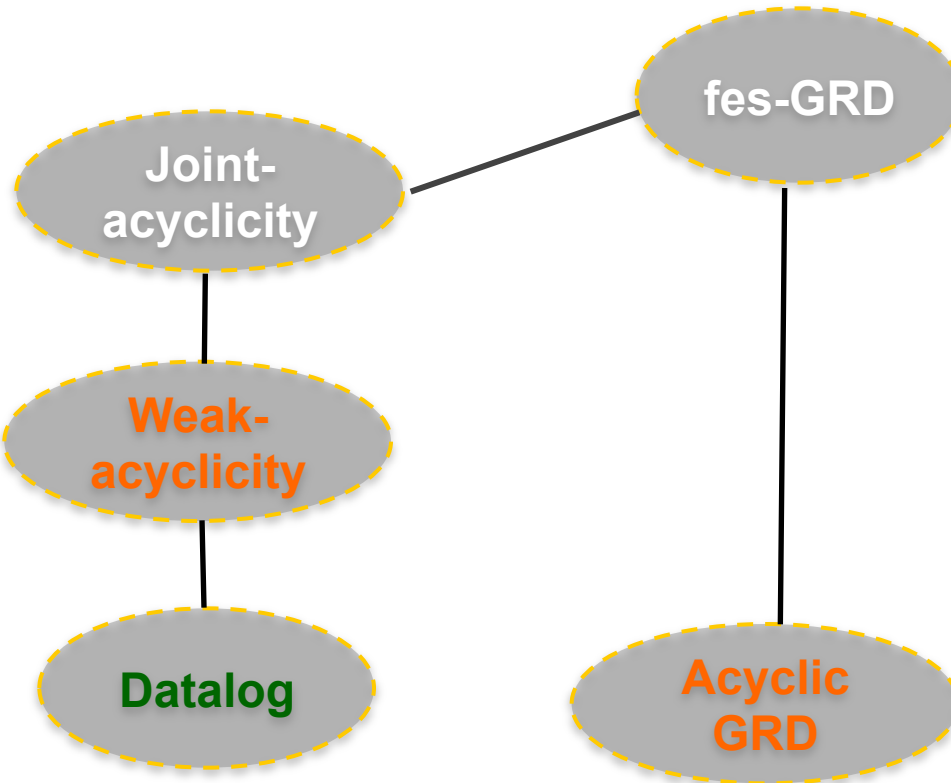
# Main Classes with **Finite Saturation** (*fes*)

Acyclic *existential*  
dependency graph

[Krötzsch+  
IJCAI' 11]

Acyclic *position*  
dependency graph

[Deutsch+ ICDT' 03]  
[Fagin+ ICDT 03]



GRD with *fes*  
strongly  
connected  
components

[Baget KR' 04]

Acyclic Graph of  
*Rule Dependencies*

[Baget KR' 04]

No existential variables

**Position dependency graph:** nodes are positions in predicates  
edges show how existential variables are propagated

**Graph of rule dependencies:** nodes are rules  
edges express that a rule may lead to trigger a rule

# WEAK-ACYCLICITY

## Position dependency graph

**nodes:** positions  $(p,i)$  in predicates

**edges:** for each frontier variable  $x$  in position  $(p,i)$  in a rule body

- an edge from  $(p,i)$  to each position  $(q,j)$  of  $x$  in the rule head

- a special edge from  $(p,i)$  to each position of an existential in the rule head

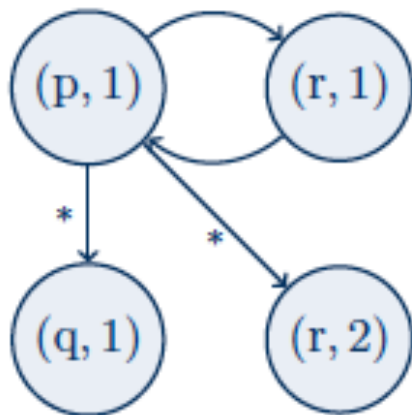
$\mathcal{R}$  is **weakly-acyclic** if its position graph contains no circuit with a special edge (\*)

$$R_1: p(x) \rightarrow \exists y r(x,y) \wedge q(y)$$

$$R_2: r(x,y) \rightarrow p(x)$$

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

$$R_2: r(x,y) \wedge r(y,x) \rightarrow p(x)$$



weakly acyclic

not weakly acyclic

special edge  $(p,1) \rightarrow (r,1)$  due to  $R_1$   
edge  $(r,1) \rightarrow (p,1)$  due to  $R_2$

# ACYCLIC GRAPH OF RULE DEPENDENCY

## Graph of Rule Dependencies

**nodes:** the rules

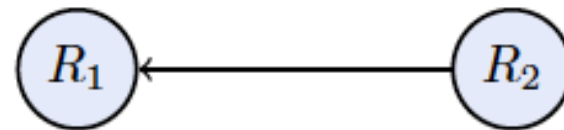
**edges:** an edge from  $R_i$  to  $R_j$  if an application of  $R_i$  may lead to trigger a new application of  $R_j$  («  $R_j$  depends on  $R_i$  »)

Dependency can be effectively computed by checking if there is a piece-unifier of  $\text{body}(R_j)$  and  $\text{head}(R_i)$

$$R_1: p(x) \rightarrow \exists y r(x,y) \wedge q(y)$$
$$R_2: r(x,y) \rightarrow p(x)$$

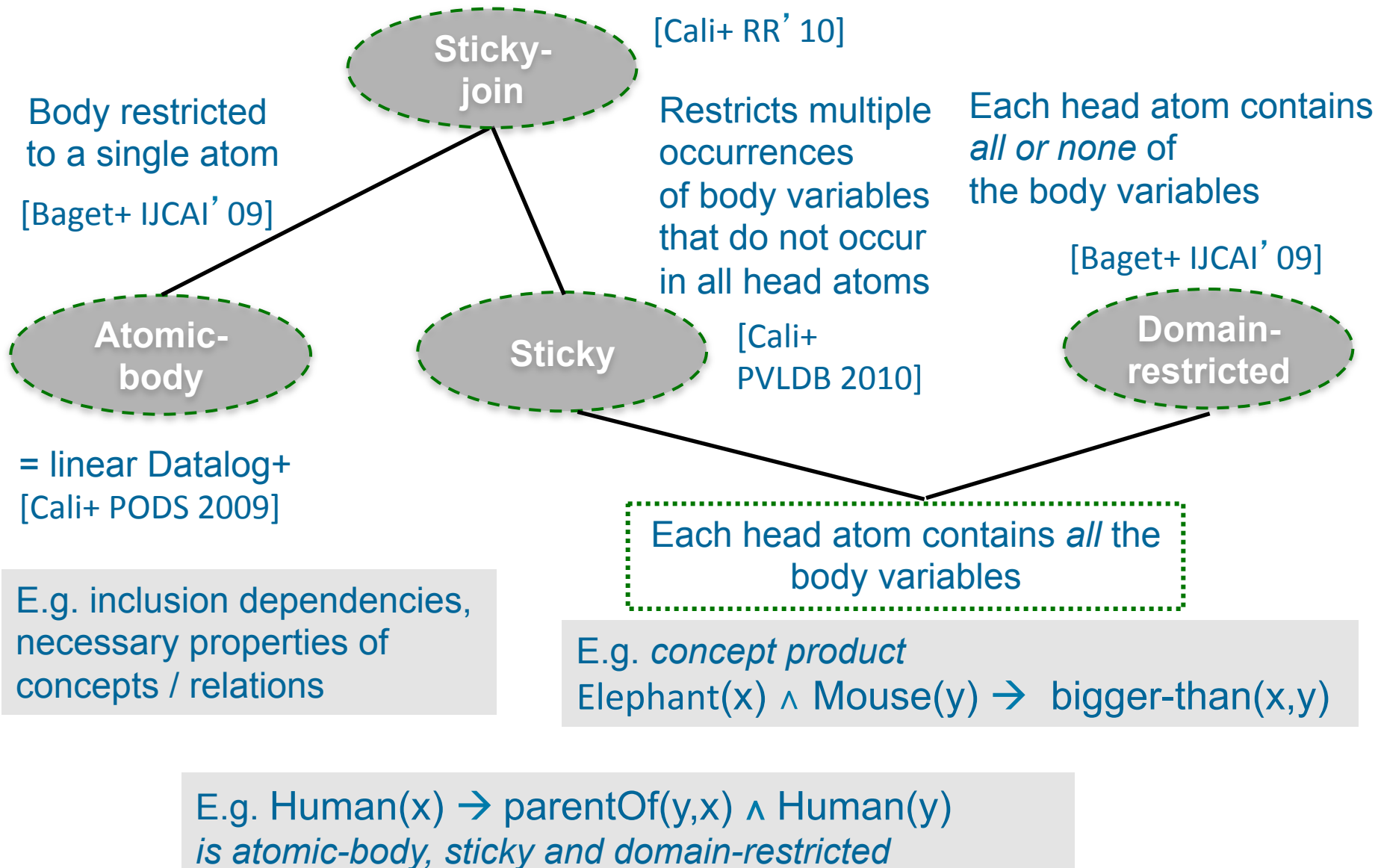
$$R_1: p(x) \rightarrow \exists y \exists z r(x,y) \wedge r(y,z) \wedge r(z,x)$$
$$R_2: r(x,y) \wedge r(y,x) \rightarrow p(x)$$

Cyclic GRD since  $R_1$  and  $R_2$  depend on each other



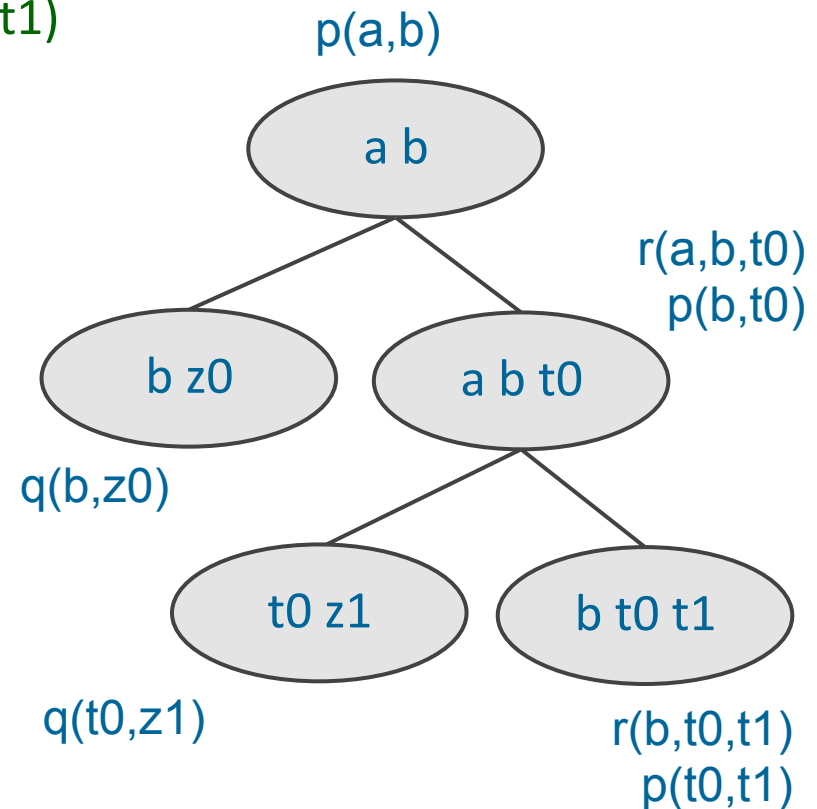
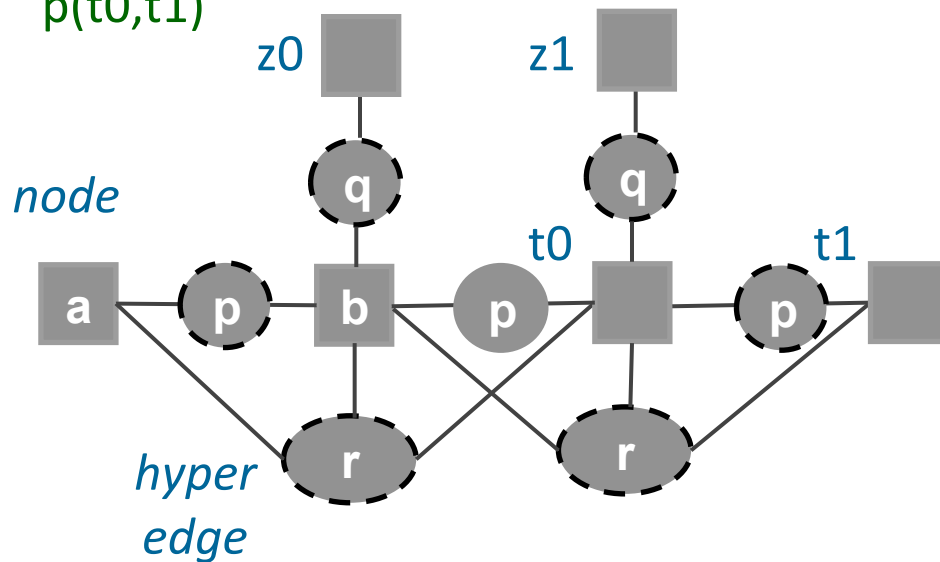
These examples show that weak-acyclicity and acyclic GRD are incomparable criteria  
Common generalizations of these two notions have been defined

# Main Classes with **Finite Query Rewriting** (*fus*)



# Decomposition Tree / Treewidth

$p(a,b)$   $q(b,z_0)$   $r(a,b,t_0)$   $p(b,t_0)$   $q(t_0,z_1)$   $r(b,t_0,t_1)$   
 $p(t_0,t_1)$



## Decomposition tree

- 1) each node (*term*) appears in a bag
- 2) each hyperedge (*atom*) has all its nodes in a bag
- 3) for each node  $x$ , the subgraph induced by the bags containing  $x$  is connected

**Width** of a tree decomposition = *max* number of nodes in a bag (minus 1)

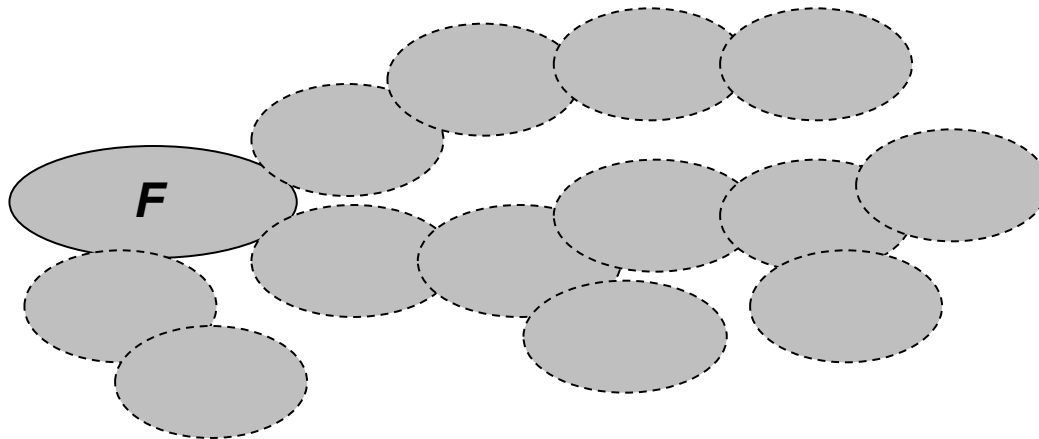
**Treewidth** of a graph = *min* width over all decomposition trees of this graph

# Bounded Treewidth of the Derived Facts (*bts*)

---

Essentially [Cali Gottlob Kifer KR'08]

$\mathcal{R}$  is *bts* if the forward chaining with  $\mathcal{R}$  generates facts with **bounded treewidth**:  
i.e., for any factbase  $F$ , there is an integer  $b$  s.t.  
any factbase  $\mathcal{R}$ -derived from  $F$  has **treewidth bounded by  $b$**



*fes* (*finite saturation*) is included in *bts*  
(bound given by the number of terms in the finite saturation)

The decidability proof does not provide a halting algorithm  
(relies on the bounded treewidth model property [Courcelle 90])

# Some Recognizable **bts** (and not fes) Classes of Rules

**Frontier:** variables shared by the body and the head

Guard only *affected* variables from the *frontier*

[Baget+ KR' 10]

Guard only the *frontier*

[Baget+ KR' 10]

$r(x,y) \wedge r(y,z) \rightarrow r(y,u) \wedge r(z,u)$

The *frontier* has size 1

[Baget+ IJCAI' 09]

Guard only *affected* variables (i.e. possibly mapped to new existentials)

[Cali+ KR' 08]

*datalog*

frontier  
1

guarded

An atom in the body *guards* all the body variables

[Cali+ KR' 08]

weakly  
frontier  
guarded

frontier  
guarded

weakly  
guarded

$r(x,y) \wedge r(y,z) \wedge r(x,z) \rightarrow \exists u r(z,u)$

$r(x,y) \wedge r(y,z) \wedge s(x,y,z) \rightarrow \exists u r(y,u) \wedge r(z,u)$

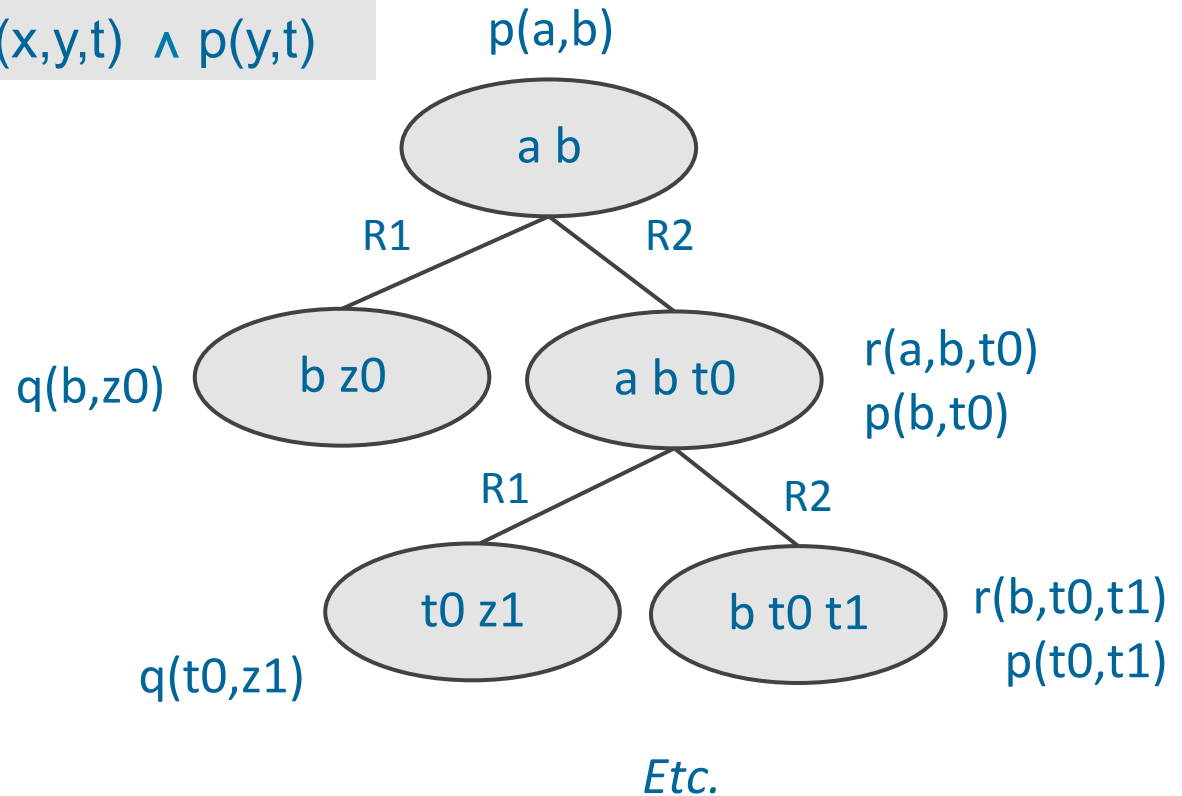
These classes are moreover « **greedy bts** » => a halting algorithm [Baget+ IJCAI' 11]

# Greedy *bts*

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

$$R_2 = p(x,y) \wedge q(x,z) \rightarrow \exists t r(x,y,t) \wedge p(y,t)$$

$$F = p(a,b)$$



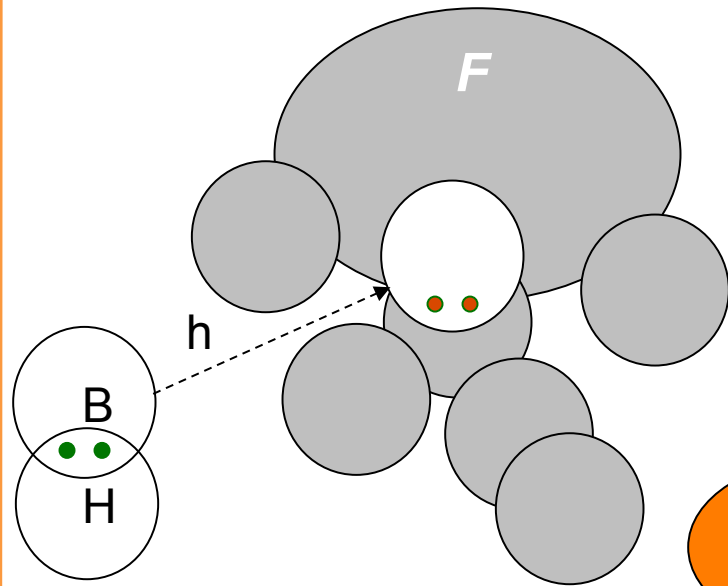
Greedy construction of a **decomposition tree** of derived facts  
with bounded width



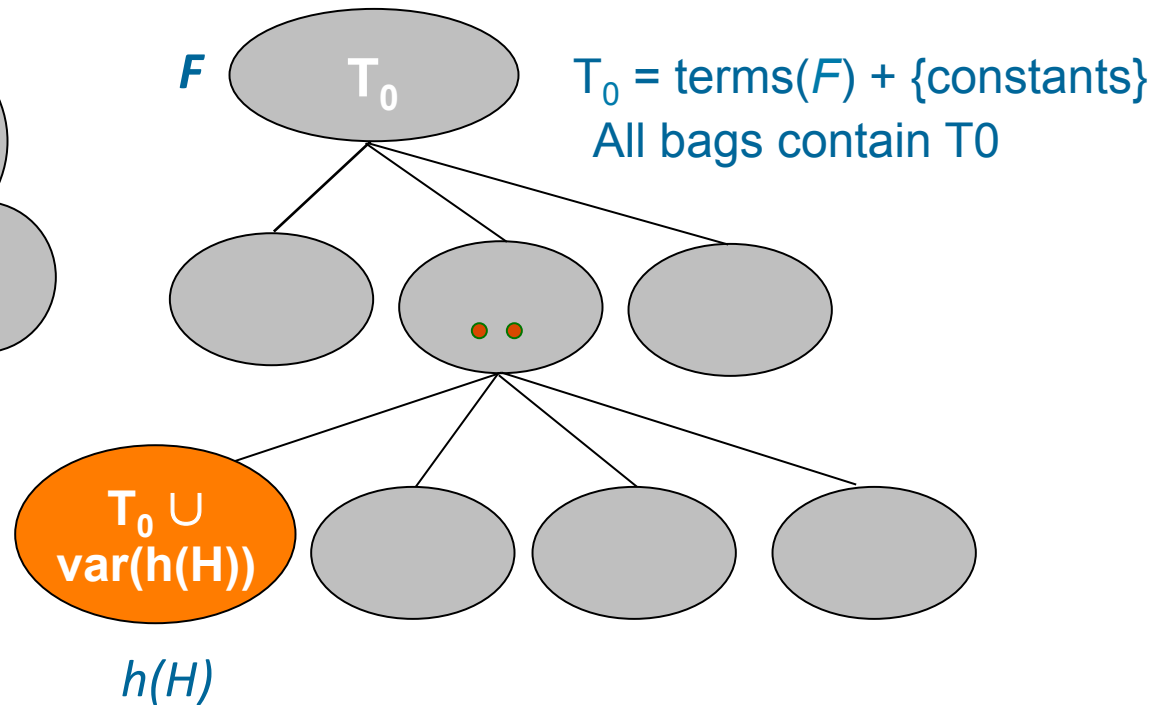
# The « Greedy bts » Property [Baget+ IJCAI' 11]

For any factbase, for each rule application,  
frontier variables not being mapped to initial terms are *jointly* mapped to variables occurring in atoms added by a single previous rule application

*Derived facts*



*Decomposition tree*



# Main Ideas of the Algorithm for *gbts* (1)

---

Build a **finite** decomposition tree that encodes a potentially infinite fact

1. Bag **pattern** = { *homomorphisms from part of a rule body to « current fact » that use some terms of the bag* }
  - A rule is applicable to the current factbase *iff* a bag pattern contains its body
  - FC can be performed on the decorated tree
2. **Equivalence relation** on bags

Only one bag per equivalence class is developed  
The other nodes are *blocked*

**Bounded number** of equivalence classes → finite « full blocked tree »  $T^*$

# Main Ideas of the Algorithm for *gbts* (2)

---

Query this finite decomposition tree

[Baget+ IJCAI 2011]  $q$  added as a rule «  $q \rightarrow match$  »

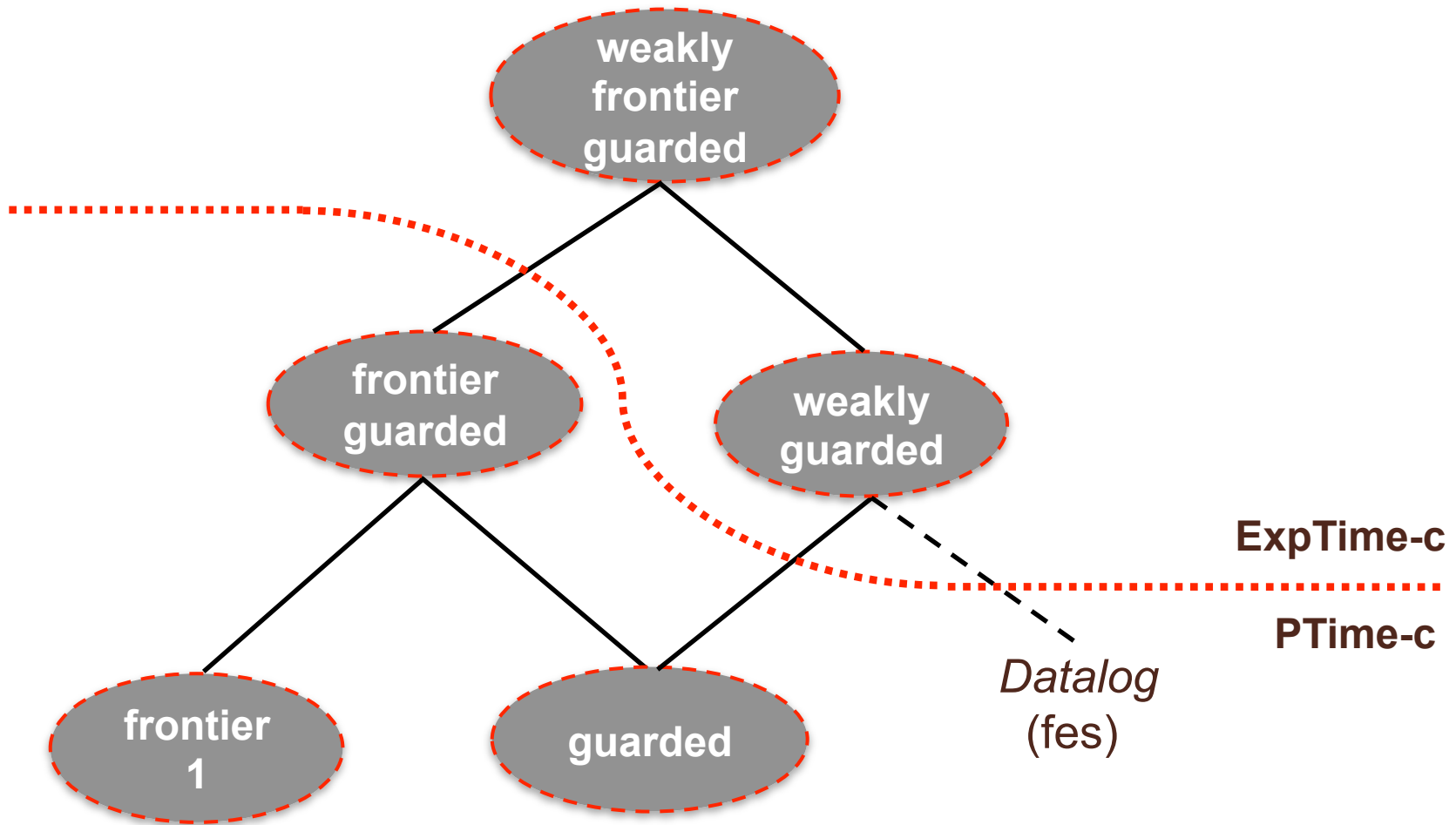
$q$  is entailed iff *match* occurs in a bag pattern  
i.e.,  $q$  maps by homomorphism to *atoms*( $T^*$ )

[Thomazo+ KR 2012] offline /online separation

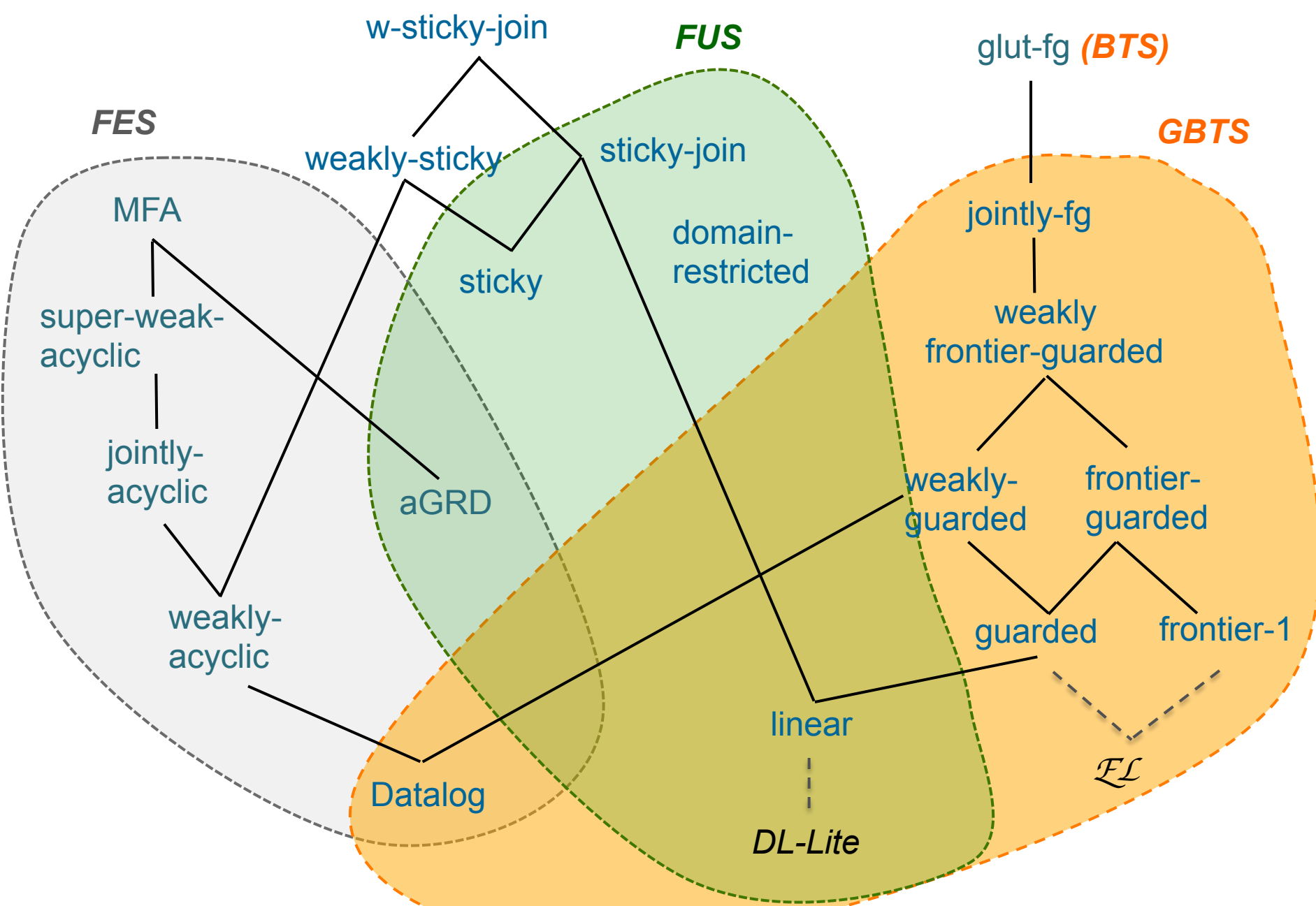
(1) compilation: tree  $T^*$  built independently from *any* query

(2) querying: *any*  $q$  is entailed iff it maps by *\*-homomorphism* to  $T^*$   
i.e.  $q$  maps by homomorphism to a *bounded* « *development* » of  $T^*$

# Data Complexity of gbts Classes



Previous algorithm is worst-case optimal on gbts  
for data / combined complexity.  
Can be specialized to be optimal on these gbts subclasses



# CONCLUSION

---

- Reasoning with ontologies is becoming central in many data-centric applications
- Solid theoretical foundations with a range of ontological formalisms that offer various tradeoff expressivity/complexity
- Ongoing research
  - Go beyond (unions of) conjunctive queries, e.g. combine them with navigational queries like regular path queries
  - New query rewriting techniques that target more powerful languages, e.g. Datalog
  - New query answering techniques that combine materialisation and query rewriting
  - Study the interaction of the ontology with mappings, which is key to efficient query answering over heterogeneous data
  - Representing and reasoning with temporal and spatial data
  - Dealing with data inconsistencies
  - . . . .

## (Small) Bibliography

Bienvenu M., Leclère M., Mugnier, M.-L. and Rousset, M.-C., [Reasoning with Ontologies](#), chapter 6, volume 1 in « A guided tour of artificial intelligence research », Springer, to appear.

**Introductions to several aspects of ontology-mediated query answering with description logics or existential rules in the Reasoning Web summer school books:**

in particular:

Bienvenu, M. and Ortiz, M. (2015). [Ontology-mediated query answering with data tractable description logics](#). 11th International Reasoning Web Summer School , volume 9203 of LNCS , pages 218–307. Springer.

Mugnier, M. and Thomazo, M. (2014). [An introduction to ontology-based query answering with existential rules](#). 10th International Reasoning Web Summer School, volume 8714 of LNCS, pages 245–278. Springer.

Gottlob, G., Orsi, G., Pieris, A., and Simkus, M. (2012). [Datalog and its extensions for semantic web databases](#). 10th International Reasoning Web Summer School , volume 7487 of LNCS, pages 54–77. Springer.

These syntheses provide further references

# APPENDIX: FURTHER DETAILS

---

- Fundamental definitions and properties for the  $FOL(\exists, \wedge)$  fragment
- Piece-unifiers



# INTERPRETATIONS / MODELS (1)

- **Vocabulary**  $\mathcal{V} = (\mathcal{P}, C)$ , where  $\mathcal{P}$  = finite set of predicates  
 $C$  = set of constants
- **Interpretation**  $I = (D_I, \cdot^I)$  of  $\mathcal{V}$ , where  
 $D_I \neq \emptyset$  (domain)  
for all  $c$  in  $C$ ,  $c^I$  in  $D_I$   
for all  $p$  in  $\mathcal{P}$  with arity  $k$ ,  $p^I \subseteq D_I^k$
- Furthermore, **unique name assumption**: for all  $c$  and  $d$  in  $C$ ,  $c^I \neq d^I$
- **Simplifying assumption** (in line with the unique name assumption):  
 $C \subseteq D_I$  and for all  $c$  in  $C$ ,  $c^I = c$

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

$$I: \quad D_I = \{a, b, d_1\} \quad p^I = \{ (b, a), (b, d_1), (d_1, b) \}$$
$$r^I = \{ (d_1, d_1, a) \}$$

- $I$  is a **model** of  $f$  (built on  $\mathcal{V}$ ) if  $f$  is true in  $I$

## INTERPRETATIONS / MODELS (2)

- Let  $f$  in  $\text{FOL}(\exists, \wedge)$ .  $I$  is a **model** of  $f$  iff

there is a **mapping**  $v$  from  $\text{terms}(f)$  to  $D^I$  such that  
for all  $p(e_1, \dots, e_k)$  in  $f$ ,  $(v(e_1), \dots, v(e_k))$  in  $p^I$

$$I: \quad D_I = \{a, b, d_1\} \quad p^I = \{(b, a), (b, d_1), (d_1, b)\}$$
$$r^I = \{(d_1, d_1, a)\}$$

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

$$v: \quad x \mapsto d_1 \quad y \mapsto b \quad z \mapsto d_1$$

- Interpretations can be seen as **sets of atoms**  
(with elements from  $D \setminus C$  seen as variables)

$$p(b, a), p(b, x_1), p(x_1, b), r(x_1, x_1, a)$$

- $I$  is a **model** of  $f$  iff there is a **homomorphism** from  $f$  to  $I$

# HOMOMORPHISMS AGAIN AND AGAIN

---

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

- We have:

If  $I_1$  maps  $I_2$  then, for any  $f$ ,  $I_1$  model of  $f \Rightarrow I_2$  model of  $f$

- To a formula  $f$  in  $\text{FOL}(\exists, \wedge)$ , we assign its **isomorphic model  $M(f)$**   
(also called **canonical model**)

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

$$M(f): \quad D = \{dx, dy, dz, a\}$$

$$p^{M(f)} = \{ (dx,dy), (dy,dz) \}$$

$$r^{M(f)} = \{ (dx, dz, da) \}$$

# NICE SEMANTIC PROPERTIES OF $\text{FOL}(\exists, \wedge)$

---

- The canonical model  $M(f)$  is **universal**, i.e., for all  $M'$  model of  $f$ ,  $M(f)$  maps to  $M'$

Proof: Let  $M'$  model of  $f$ . Then,  $f$  maps to  $M'$ . Since  $M(f)$  isomorphic to  $f$ ,  $M(f)$  maps to  $M'$

- $g \models f$  (i.e., every model of  $g$  is a model of  $f$ ) iff  
 $f$  maps by homomorphism to  $M(g)$  iff  
 $f$  maps by homomorphism to  $g$

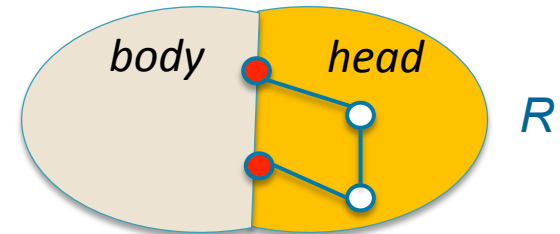
Proof:

$\Rightarrow$  Assume  $g \models f$ . In particular  $M(g)$  is a model of  $f$ , hence  $f$  maps to  $M(g)$

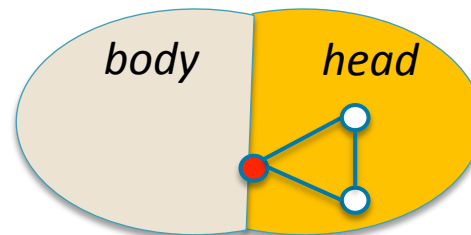
$\Leftarrow$  Assume  $f$  maps to  $M(g)$ . Since  $M(g)$  is universal: for any  $M'$  model of  $g$ ,  $f$  maps to  $M'$ , i.e.,  $M'$  is a model of  $f$ , hence  $g \models f$

# WHY « PIECES » ? (CONT'D) – PIECE-UNIFIERS

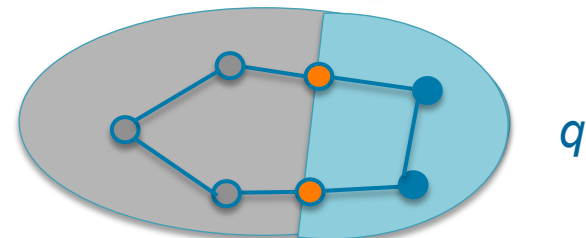
- Unification must « map » parts of  $q$  according to « pieces » that can be provided by a rule application  
(otherwise it is unsound or useless)



*specialization of the frontier*



*homomorphism*



The terms of  $q$  unified with the frontier  
(or with constants) cut  $q$  into « pieces »

⇒ entire « pieces » of  $q$  must be mapped  
to the pieces of the rule

# PIECE-UNIFIER: ALTERNATIVE DEFINITION

Let  $u_1: \text{frontier}(R) \rightarrow \text{frontier}(R) \cup \text{constants}$   
( $u_1$  is a specialization of the frontier of  $R$ )

Let  $u_2$  be a homomorphism from  $q' \subseteq q$  to  $u_1(\text{head}(R))$

**Cutpoints:** terms of  $q'$  mapped to  $u_1(\text{frontier}(R))$   
or to constants

The cutpoints cut  $q$  into « **pieces** »

$u_1+u_2$  is a piece-unifier if  $q'$  is composed of *pieces*

