

Reasoning on Data: Challenges and Applications

Marie-Christine Rousset

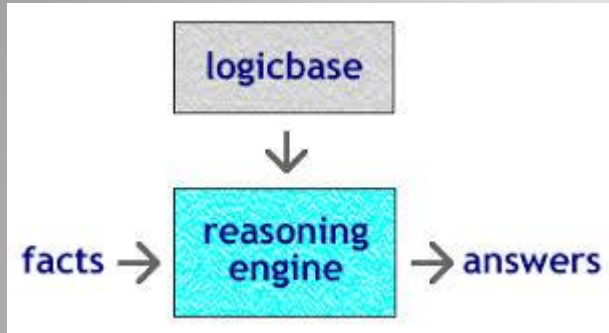
Laboratoire d'Informatique de Grenoble

Université Grenoble Alpes & Institut Universitaire de France

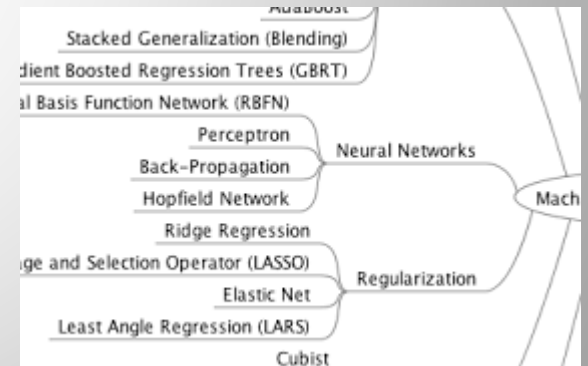
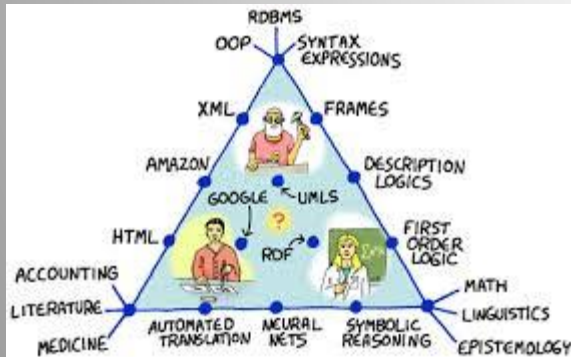
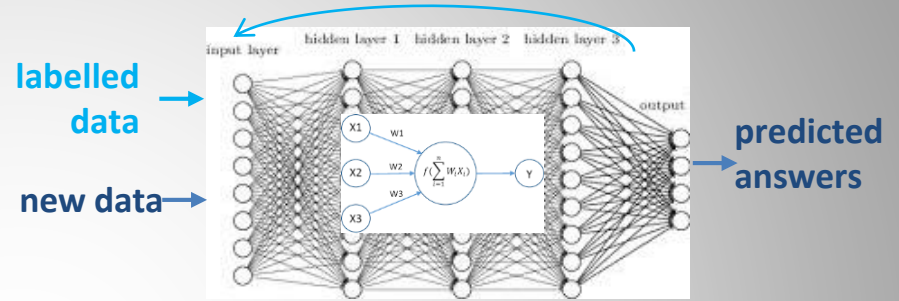


Two branches of Artificial Intelligence

Symbolic knowledge-driven approaches

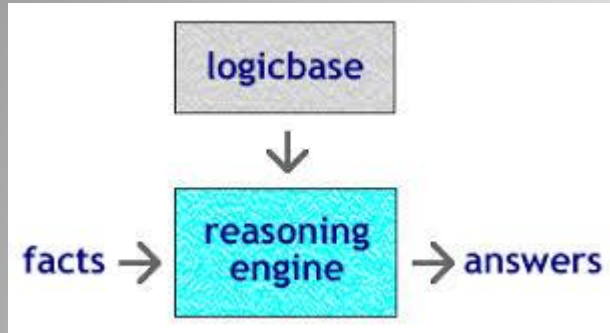


Numerical data-driven approaches

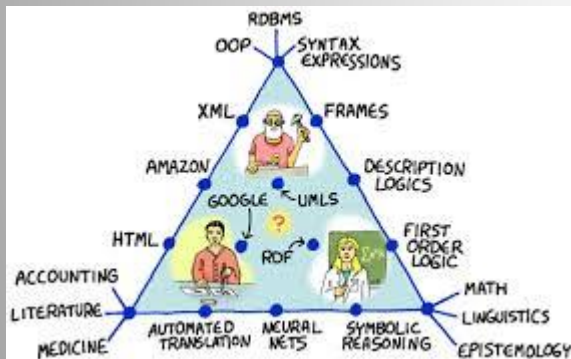


Two branches of Artificial Intelligence

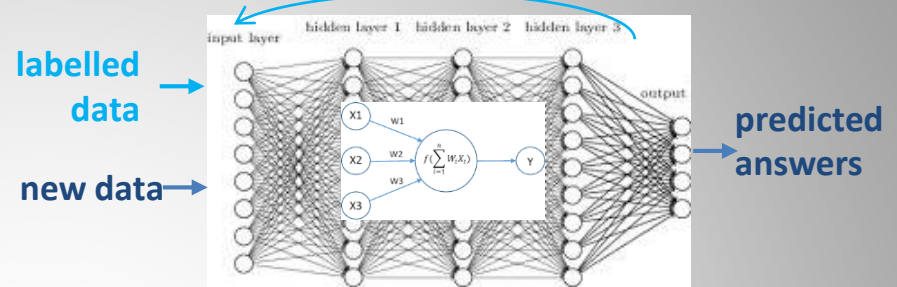
Symbolic knowledge-driven approaches



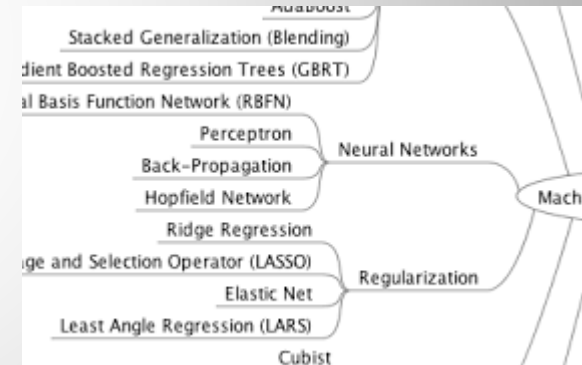
a.k.a Good-Old-Fashioned AI (GOFAI)



Numerical data-driven approaches



a.k.a Modern AI



Respective advantages and disadvantages

Explicability and transparency:

all reasoning steps to reach a conclusion are based on symbolic human readable representations

Robustness and scalability:

- the rules and knowledge have to be hand coded ... but more and more work on learning rules from data
- the generic reasoning algorithms may have a high computational complexity (atleast in the worst-case)

Automated Reasoning

- Problem studied in Mathematics, Logic and Informatics
 - Many decidability and complexity results coming from decades of research in the KR&R community
 - Several inference algorithms and implemented reasoners
 - The key point
 - first-order-logic is appropriate for knowledge representation
 - but **full first-order-logic is not decidable**
- ⇒ the game is to find restrictions to design:
- **decidable fragments** of first-order-logic
 - expressive enough for modeling useful knowledge or constraints

Key logic-based knowledge representation formalisms

- **Rules:** logical foundation of **expert systems**
 - the first successful and commercial AI systems (in the 1970s)
 - human expertise in a specific domain is captured as a **set of if-then rules**
 - given a **set of input facts**, the **inference engine** triggers **relevant rules** to build a chain of reasoning arriving to a particular conclusion
 - extended to fuzzy rules to deal with uncertain reasoning
- **Conceptual graphs:** a **graphical representation of logic**
 - logical formalism focused on representing individuals by their classes and relations (> mid-eighties)
 - originated from semantic networks (introduced to represent meaning of sentences in natural language)
 - reasoning algorithms based on **graph operations**
 - directly applicable to Linked Data for querying RDF knowledge bases (RDF graphs constrained by RDFS statements)
- **Description logics:** logical foundation of **ontologies** and the **Semantic Web**
 - (started in the early 1990s)

REASONING ON DATA: FOCUS ON

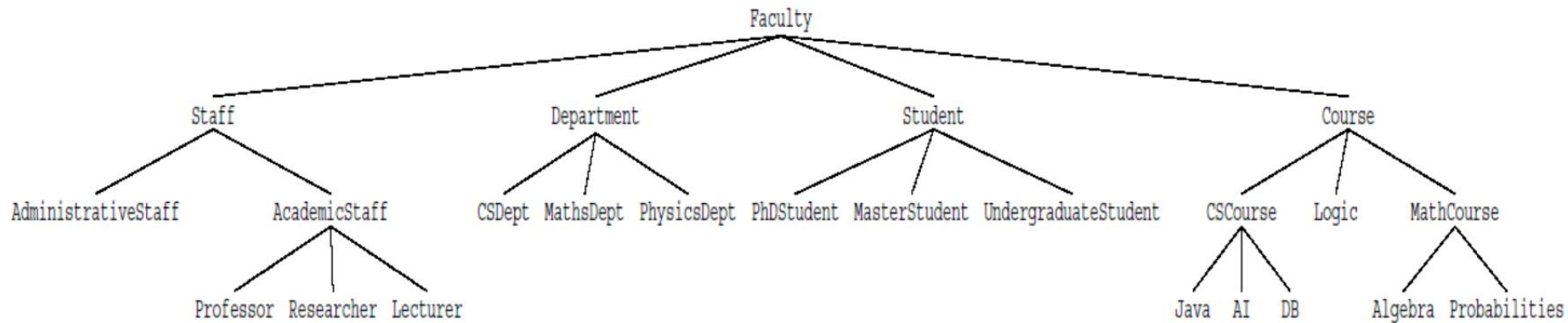
- **ONTOLOGY-BASED DATA QUERYING**
- **DATA INTEGRATION**
- **DATA LINKAGE (A.K.A KNOWLEDGE GRAPH COMPLETION)**

Ontologies

- A formal specification of a domain of interest
 - a vocabulary (classes and properties)
 - enriched with statements that constrain the meaning of the terms used in the vocabulary
 - *java* can be a *dance*, an *island*, a *programming language* or a *course*
 - the statement *java* is a subclass of *CS Courses* makes clear the corresponding meaning for *java*: it is a course
 - With a logical semantics
 - Ontological statements are axioms in logic
- ⇒ a conceptual yet computational model of a particular domain of interest.
- computer systems can then base decisions on reasoning about domain knowledge.
 - humans can express their data analysis needs using terms of a shared vocabulary in their domain of interest or of expertise

Example

A taxonomy (graphical representation of subclass constraints)



+ set of properties with constraints on their domain and range

TeachesIn (Academic Staff, Courses)

TeachesTo (Academic Staff, Students)

Manager (Staff, Departments)

+ additional constraints (not expressible in RDFS but in OWL)

Students disjoint from Staff

Only Professors or Lecturers may teach to Undergraduate Students

Every Department must have a unique Manager who must be a Professor

Query answering over data through ontologies

- **A reasoning problem**

- Ontological statements can be used to infer new facts and deduce answers that could not be obtained otherwise
- Subtlety: **some inferred facts can be partially known**

From the constraint “a professor teaches at least one master course”

$\forall x (\text{Professor}(x) \Rightarrow \exists y \text{Teaches}(x,y), \text{MasterCourse}(y))$

and the fact:

Professor(dupond) (RDF syntax: **<dupond, type, Professor>**)

it can be inferred the two following incomplete “facts” :

Teaches(dupond, v) , MasterCourse(v)

i.e, in RDF notation, two RDF triples with blank nodes:

<dupond, Teaches, _v> , <_v, type, MasterCourse>

Reasoning: a tool for checking data inconsistency

- Some ontological statements can be used as **integrity constraints**
 - “a professor cannot be a lecturer” ; “a course must have a responsible”
 $\forall x (\text{Professor}(x) \Rightarrow \neg \text{Lecturer}(x))$
 - $\forall x (\text{Course}(x) \Rightarrow \exists y \text{ResponsibleFor}(y,x))$
 - “a master course is taught by a single teacher”
 - “only professors can be responsible of courses that they have to teach”
 $\forall x \forall y (\text{Course}(x), \text{ResponsibleFor}(y,x) \Rightarrow \text{Professor}(y), \text{Teaches}(y,x))$
- Subtlety: **showing data inconsistency may require intricate reasoning** on different rules, constraints and facts
 - The facts: **Lecturer (jim), Teaches(jim, ue431) , MasterCourse(ue431)**
 - + the above integrity constraints
 - + the rule $\forall x (\text{MasterCourse}(x) \Rightarrow \text{Course}(x))$ leads to an inconsistency

Description Logics

- A family of class-based logical languages for which reasoning is decidable
 - Provides algorithms for reasoning on (possibly complex) logical constraints over unary and binary predicates
- This is exactly what is needed for handling ontologies
 - in fact, the OWL constructs come from Description Logics
- A fine-grained analysis of computational complexity with surprising complexity results
 - *ALC* is EXPTIME-complete
 - =>any sound and complete inference algorithm for reasoning on most of the subsets of constraints expressible in OWL may take an exponential time (in the worst-case)
 - “only professors or lecturers may teach to undergraduate students”
 - $\forall x \forall y (\text{TeachesTo}(x,y), \text{UndergraduateStudent}(y) \Rightarrow \text{Professor}(x) \vee \text{Lecturer}(x))$

$\exists \text{TeachesTo} . \text{UndergraduateStudent} \sqsubseteq \text{Professor} \sqcup \text{Lecturer}$

The same game again...

- Find restrictions on the logical constructs and/or the allowed axioms in order to:
 - design sublanguages for which reasoning is in P
 - EL, DL-Lite**
 - expressive enough for modeling useful constraints over data
- **DL-Lite: a good trade-off**
 - captures the main constraints used in databases and in software engineering
 - extends **RDFS** (the formal basis of OWL2 QL profile)
 - specially designed for answering queries over ontologies to be **reducible to answering queries over RDBMS with same data complexity** (at least for the fragment of union of conjunctive queries)

Reducibility to query reformulation

Query answering and data consistency checking can be performed in two separate steps:

- a **query reformulation step**
 - reasoning on the ontology (and the queries)
 - independent of the data
- ⇒ a set a queries: the reformulations of the input query

- an **evaluation step**
 - of the (SPARQL) query reformulations on the (RDF) data
 - independent of the ontology

⇒ Main advantage

- makes possible to use an SQL or SPARQL engine
- thus taking advantage of well-established query optimization strategies supported by standard relational DBMS

DL-Lite by example

Professor $\sqsubseteq \exists$ Teaches

$$\forall x (\text{Professor}(x) \Rightarrow \exists y \text{Teaches}(x,y))$$

\exists Teaches \sqsubseteq Course

$$\forall x \forall y (\text{Teaches}(x,y) \Rightarrow \text{Course}(y))$$

ResponsibleFor \sqsubseteq Teaches

$$\forall x \forall y (\text{ResponsibleFor}(x,y) \Rightarrow \text{Teaches}(x,y))$$

(funct ResponsibleFor)

$$\forall x \forall y \forall z (\text{ResponsibleFor}(y,x) \wedge \text{ResponsibleFor}(z,x) \Rightarrow y=z)$$

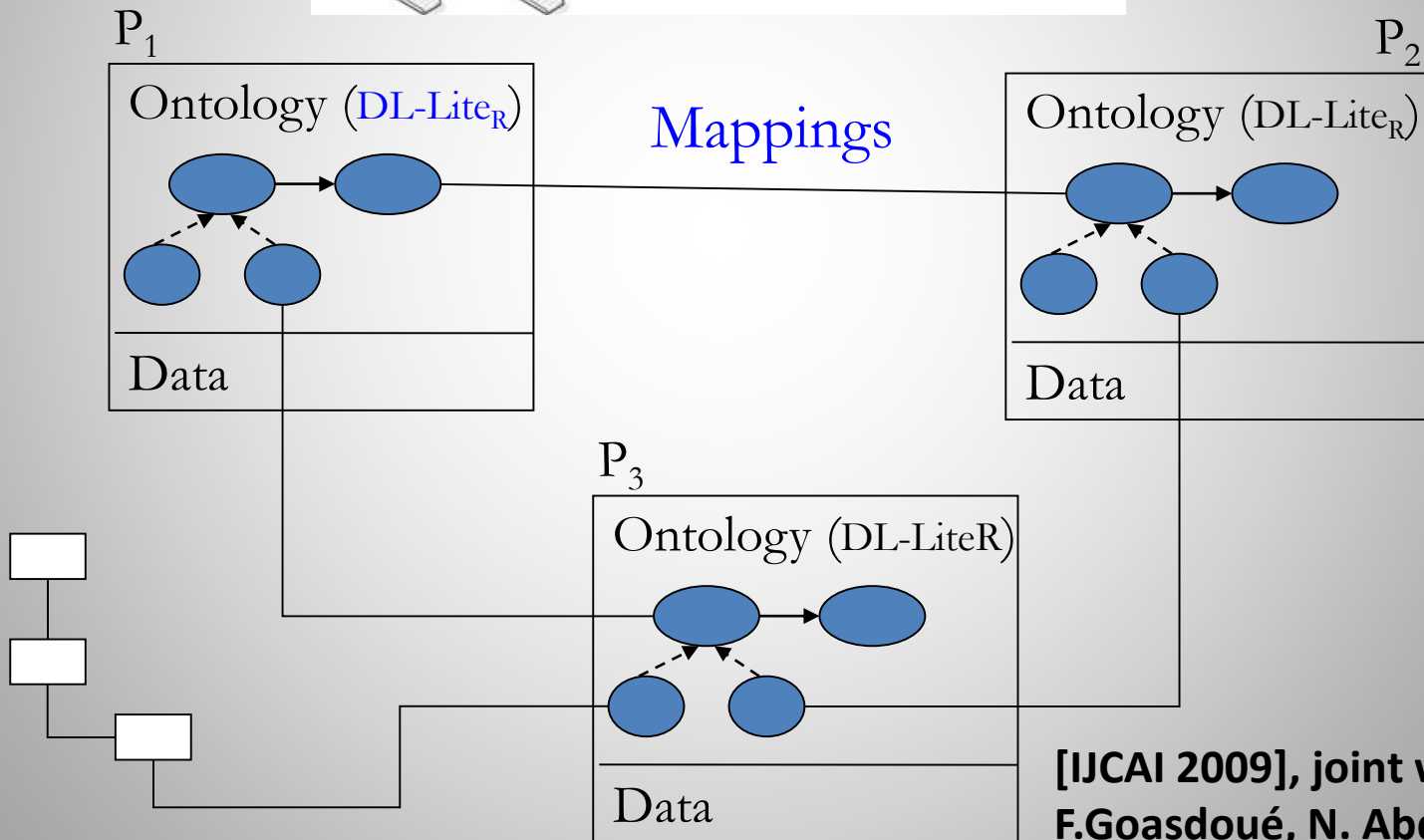
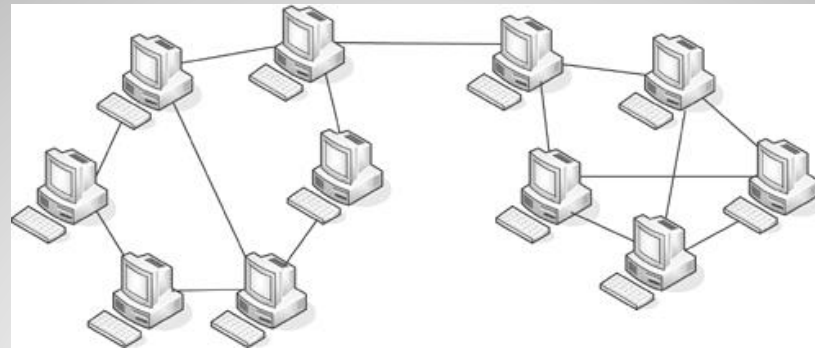
Lecturer $\sqsubseteq \neg (\exists \text{ResponsibleFor})$

$$\forall x \forall y (\text{Lecturer}(x) \wedge \text{ResponsibleFor}(x,y) \Rightarrow \perp)$$

DL-Lite: a frontier for CQ reducibility

- The **reasoning step** is **polynomial** in the size of the ontology
- The **evaluation step** has the same **data complexity** as standard evaluation of conjunctive queries over relational databases
 - in **ACo** (strictly contained in LogSpace and thus in P)
- The interaction between relation inclusion constraints and functionality constraints makes reasoning in DL-Lite **P-complete in data complexity**
 - **DL-Lite_R is CQ-reducible**
 - **full DL-Lite is not CQ reducible**
 - reformulating a query may require recursion (Datalog)

Decentralized ontology-based data querying



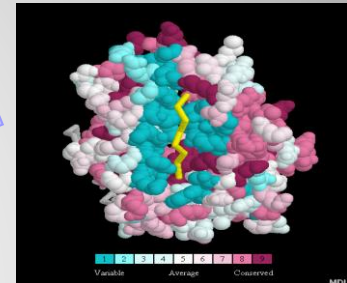
[IJCAI 2009], joint work with
F.Goasdoué, N. Abdallah

Data Integration

Web



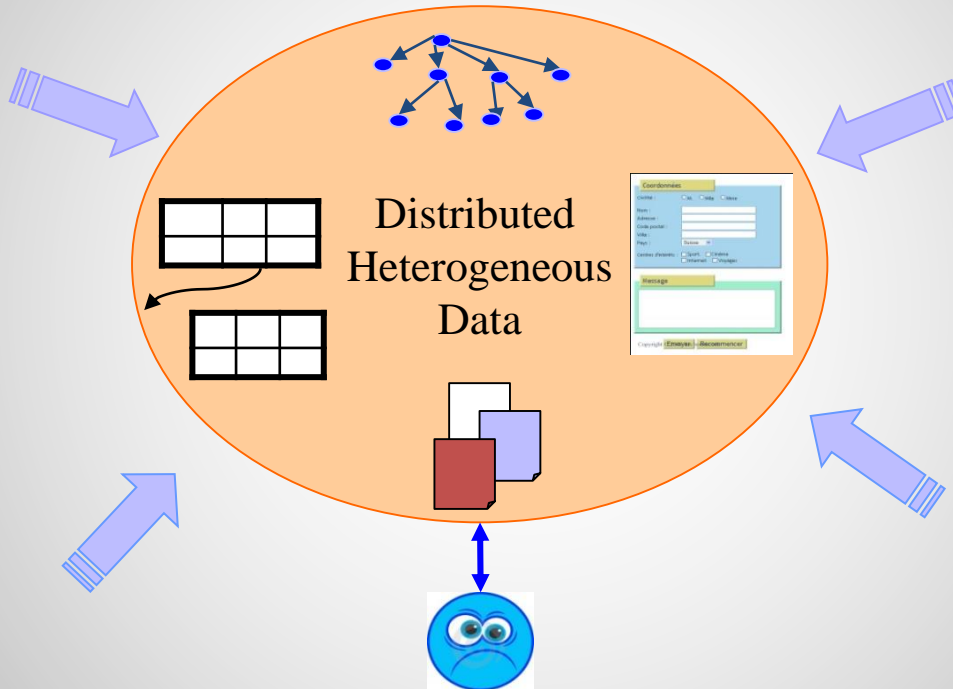
Sciences



Enterprise



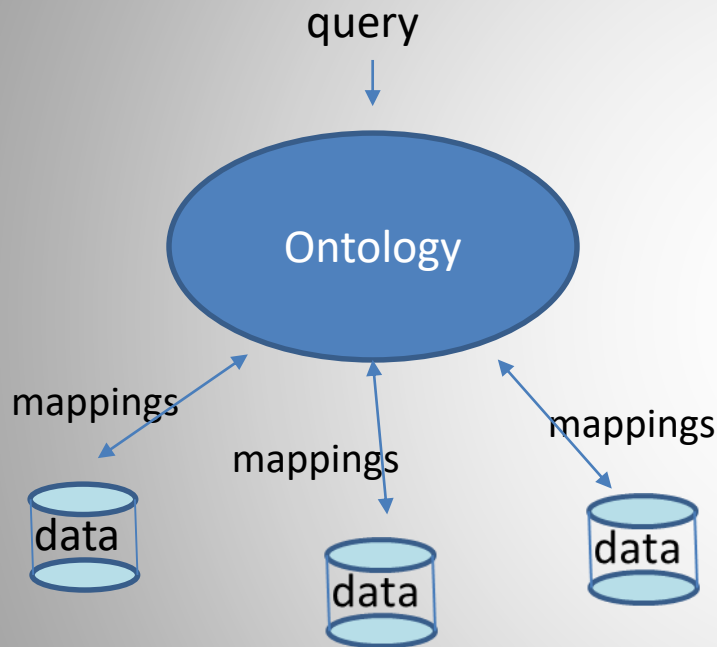
Administration



a difficult challenge !

Domain ontology + mappings:

the semantic glue between heterogeneous data sources



Two main algorithmic approaches

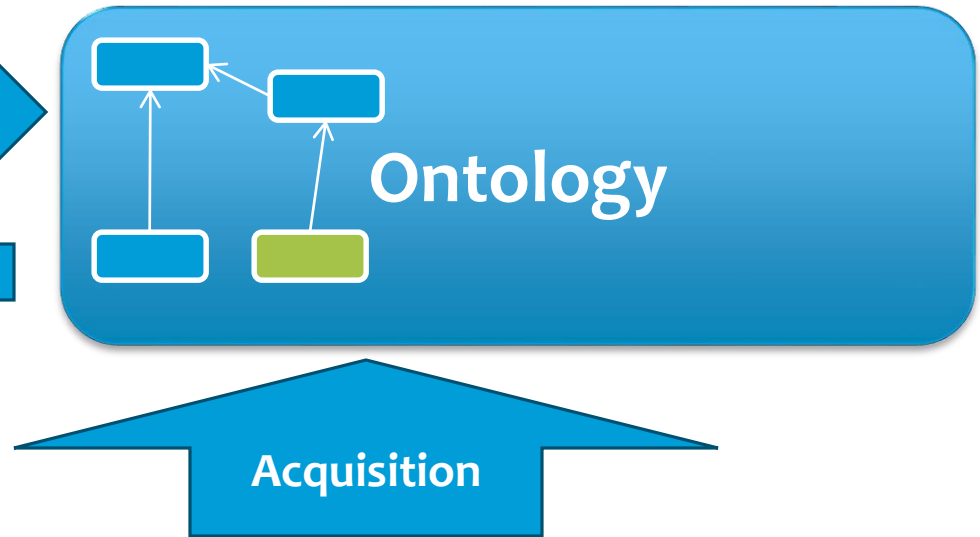
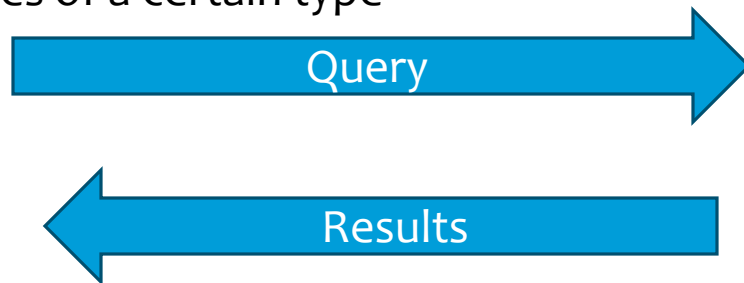
- 1. Answering queries by query rewriting :**
 - query reformulation using ontologies (backward reasoning)
 - query translation using mappings
- 2. Answering queries by data materialization:**
 - Data extraction and transformation using mappings (e.g., from relational to RDF)
 - Data saturation (forward reasoning on data and ontological statements)

The complexity and feasibility in practice depend on the languages used for expressing the queries, the mappings and the ontology

ANR project CONTINUUM (2008-2012)

CONTinuité de service en Informatique UbiqUitaire et Mobile
(joint work with F.Jouanot and J.Coutaz)

Find devices in the environment that offer services of a certain type



Environment

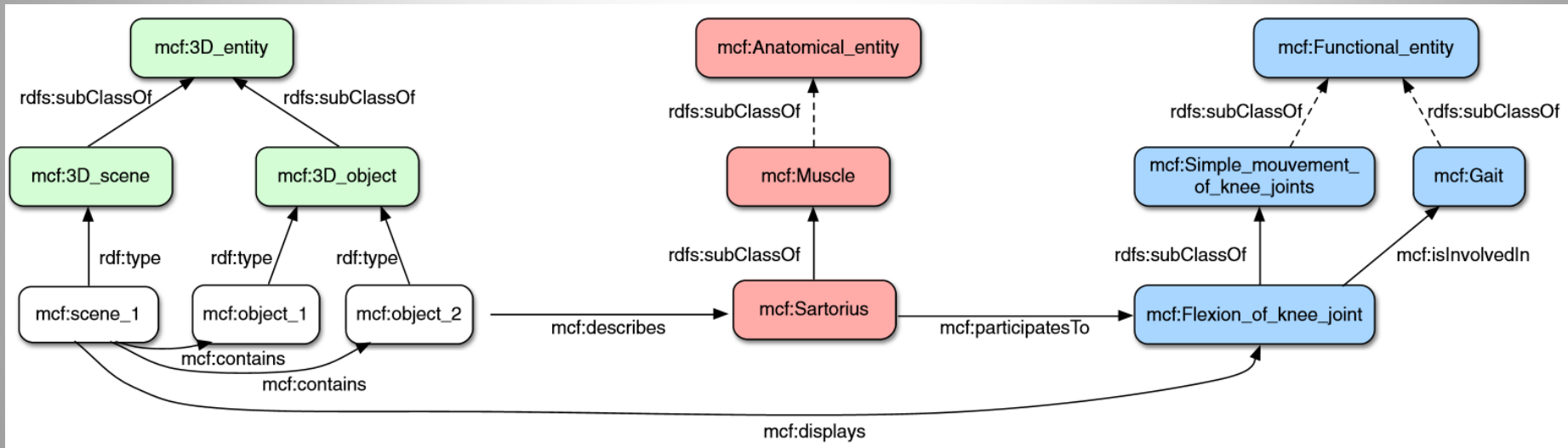


My Corporis Fabrica

(joint work with Olivier Palombi, LADAF, LJK)

[Journal of Biomedical Semantics, 2014]

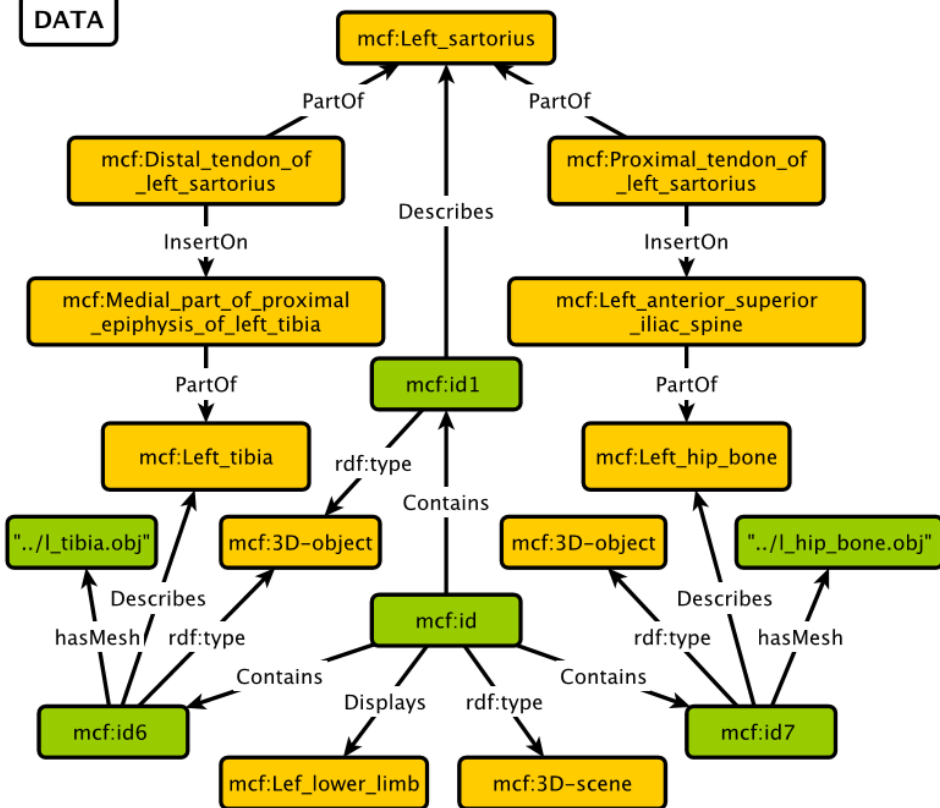
- Ontology-based integration of complex anatomical models
 - rules for mapping structural, functional, spatio-temporal and 3D models of anatomy



```
IF< ?x rdf:type mcf:3D-object >AND< ?x mcf:Describes ?y >  
AND< ?y rdfs:subClassOf mcf:Bone >  
THEN< ?x mcf:hasColour 'yellow' > (R12)
```

Support for interactive simulation and visualization

DATA



Corresponding triples :

```

< mcf:Distal_tendon_of_left_sartorius mcf:PartOf mcf:Left_sartorius >
< mcf:Distal_tendon_of_left_sartorius mcf:InsertOn mcf:Medial_part_of_proximal_epiphysis_of_left_tibia >
< mcf:Proximal_tendon_of_left_sartorius mcf:PartOf mcf:Left_sartorius >
< mcf:Proximal_tendon_of_left_sartorius mcf:InsertOn mcf:Left_anterior_superior_iliac_spine >
< mcf:Medial_part_of_proximal_epiphysis_of_left_tibia mcf:PartOf mcf:Left_tibia >
< mcf:Left_anterior_superior_iliac_spine mcf:PartOf mcf:Left_hip_bone >

< mcf:id rdf:type mcf:3D-scene > < mcf:id mcf:Displays mcf:Left_lower_limb >
< mcf:id1 rdf:type mcf:3D-object > < mcf:id6 rdf:type mcf:3D-object > < mcf:id7 rdf:type mcf:3D-object >
< mcf:id mcf:Contains mcf:id1 > < mcf:id mcf:Contains mcf:id6 > < mcf:id mcf:Contains mcf:id7 >

< mcf:id1 mcf:Describes mcf:Left_sartorius > < mcf:id1 mcf:hasMesh "../geometries/l_sartorius.obj" >
< mcf:id6 mcf:Describes mcf:Left_tibia > < mcf:id6 mcf:hasMesh "../geometries/l_tibia.obj" >
< mcf:id7 mcf:Describes mcf:Left_semimembranosus > < mcf:id7 mcf:hasMesh "../geometries/l_hip_bone.obj" >
    
```

QUERY

Query in English :

May I see, in the curent 3D scene, the bones on which the left sartorius muscle is inserted ?

Query in SPARQL

```

PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX mcf: <http://www.mycorporisfabrica.org/ontology/mcf.owl#>
    
```

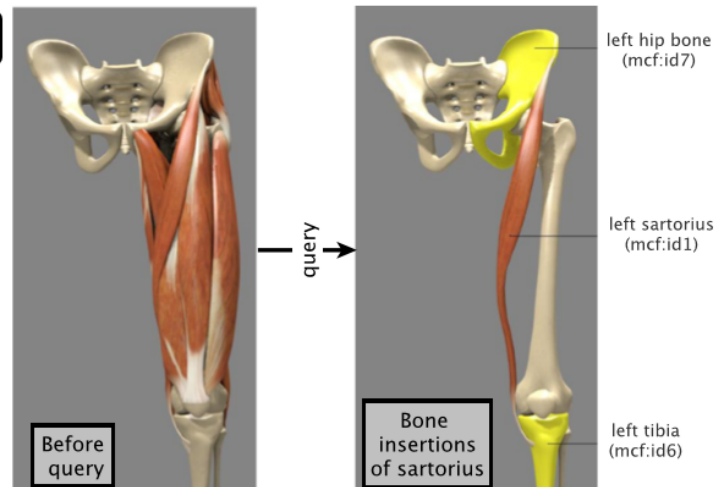
```

select ?bone ?individual ?mesh
where {?s mcf:PartOf mcf:Left_sartorius.
?s mcf:InsertOn ?z.
?z mcf:PartOf ?bone.
?individual rdfs:Describes ?bone.
?scene mcf:Contains ?individual.
?individual rdf:type mcf:3D-object.
?individual mcf:hasMesh ?mesh }
    
```

Answer

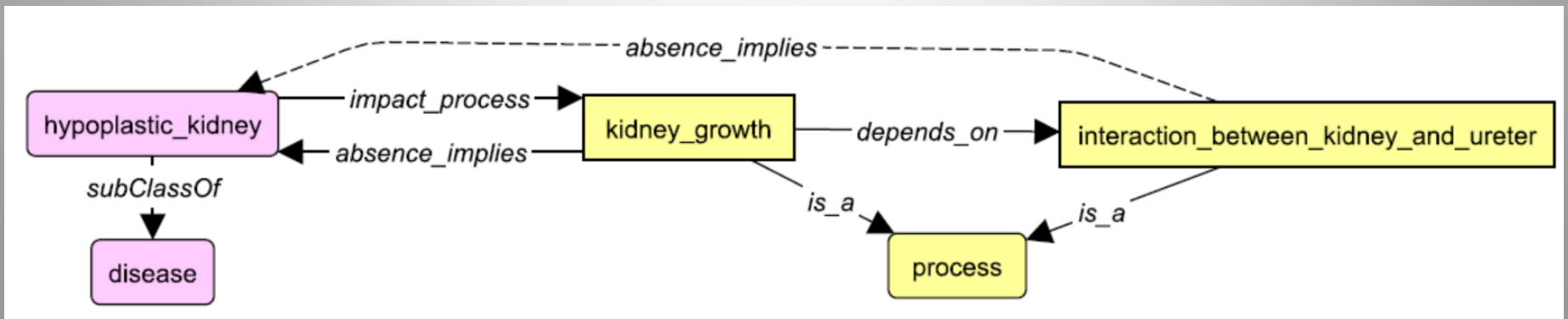
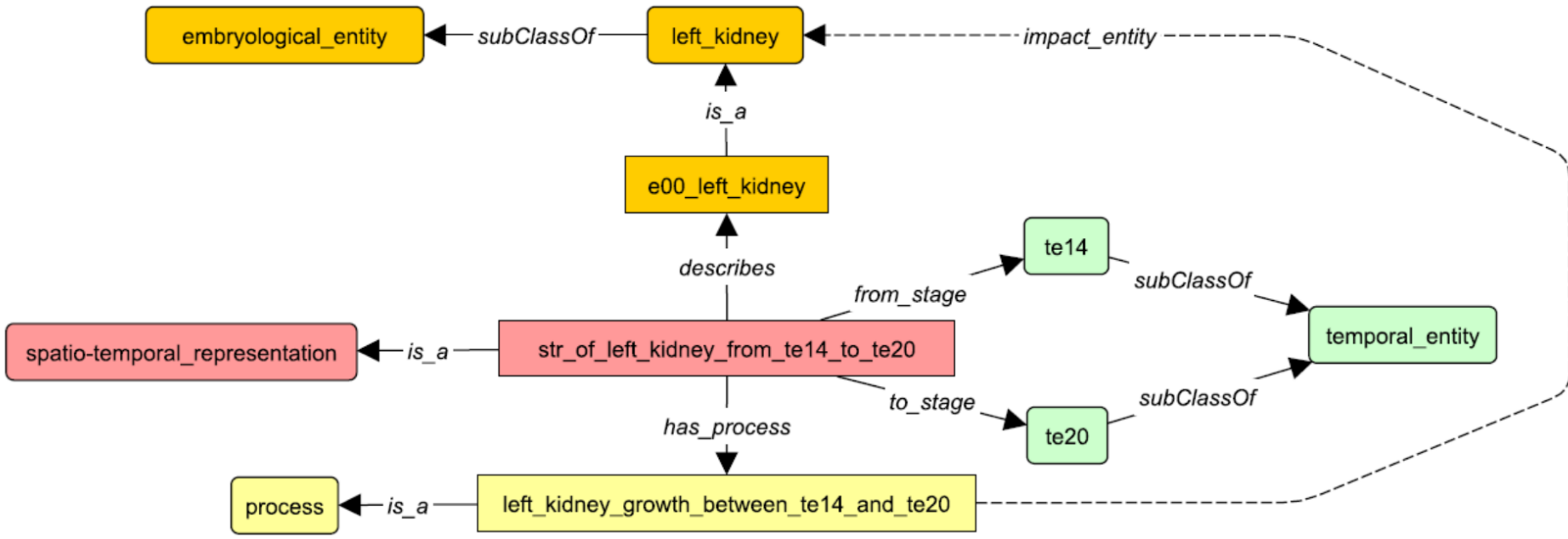
?bone	?individual	?mesh
mcf:Left_tibia	mcf:id6	../l_tibia.obj
mcf:Left_hip_bone	mcf:id7	../l_hip_bone.obj

3D view



My Corporis Fabrica Embryo

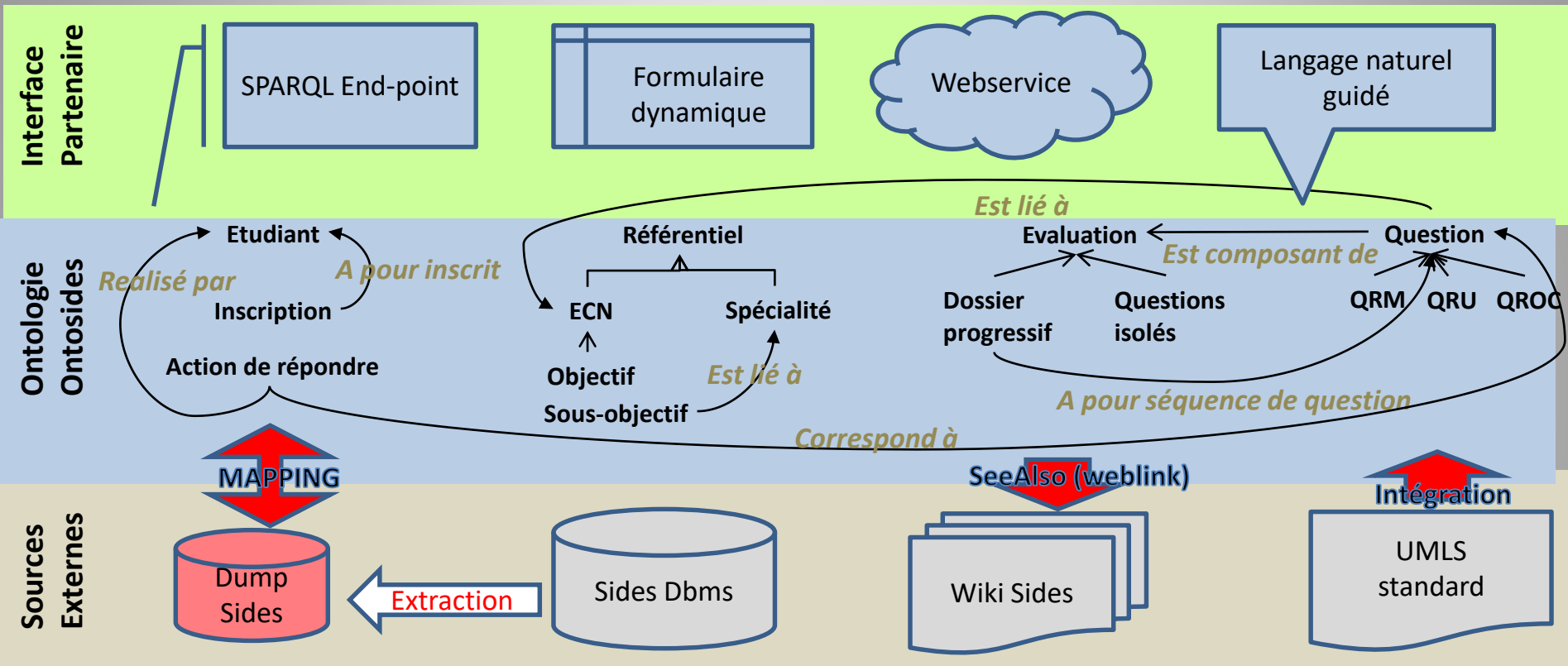
[Journal of Biomedical Semantics, 2015]



SIDES 3.0: e-learning in Medicine (2017-2020)



Ontology-based 3-layers architecture

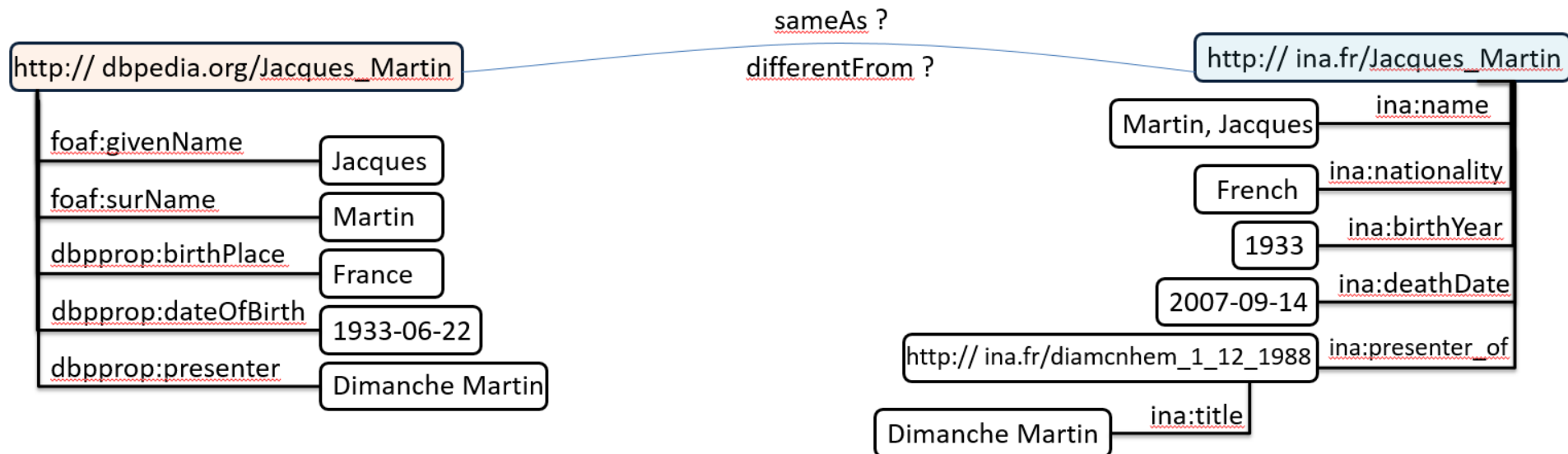


Materialization approach:

- small ontology (semi-automatically constructed)
- Instantiated by mappings with a Dump of SIDES (activities of 64 957 students over 3 years):
⇒ **1.5 Billions triples**
- Scales to complex SPARQL 1.1 queries (with aggregation and subqueries) for data analytics

Data linkage

- Deciding whether two URIs refer to the same real-world entity (within or across data sources)



- Crucial task for data fusion and enrichment
- A hot topic in Linked Open Data
- Also related to data privacy

Existing approaches

- **Numerical methods based** on aggregating similarities between values of some relevant properties
 - Specification through linkage rules (e.g., in Silk and LIMES) of:
 1. the properties to consider within the descriptions of individuals,
 2. the similarity functions to use for comparing their respective values,
 3. the functions for aggregating these similarity values
 - Linkage rules: defined manually or learned automatically
 - **Main weakness: no formal semantics and no rule chaining**
- **Symbolic methods** based on logical rules **equipped with full reasoning**
 - Translation of schema constraints into logical rules
 - Logical inference of sameAs facts
 - **Main weakness: not robust to incomplete and/or noisy data**
 - ⇒ 100% precision but risk of low recall

Probabilistic Datalog revisited to reason with rules and probabilities

- Joint work with M. Al-Bakri, M.Atencia, J.David and S.Lalande (**Qualinca ANR project with INA**)
 - [ECAI 2016] Uncertainty-Sensitive Reasoning for Inferring sameAs Facts in Linked Data
 - **ProbFR**: an inference algorithm that **computes the probability of inferred facts** as well as the uncertainty **provenance** of this computation
 - a series of experiments over real-world large RDF datasets showing the benefits and the scalability of our approach

Probabilistic Datalog (*)

- A simple extension of Datalog in which rules and facts are associated with **symbolic probabilistic events**
- Logical inference and probability computation are separated
 - **Step 1 (ProbFR)** : computation for each inferred fact of its **provenance** (the **boolean combination** of all the events associated with the input facts and rules involved in its derivation)
 - exponential in the worst-case
 - by-passed by a practical bound on the number of conjuncts in the provenances and a priority given to the most probable rules and facts
 - **Step 2:** computation of the probabilities of the inferred facts
 - from their provenances in which **each event of input facts and rules is assigned a probabilistic weight**
 - based on independence and disjointness assumptions to make it feasible

(*) N. Fuhr, Probabilistic models in information retrieval, The Computer Journal, 1992₃₀

Illustrative Example

Rules: uncertain rules are in red, certain rules are in blue

$r_1 : (?x \text{ sameName } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

$r_2 : (?x \text{ sameName } ?y), (?x \text{ sameBirthDate } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

$r_3 : (?x \text{ marriedTo } ?z), (?y \text{ marriedTo } ?z) \Rightarrow (?x \text{ sameAs } ?y)$

$r_4 : (?x \text{ sameAs } ?z), (?z \text{ sameAs } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

Facts: uncertain facts are in red, certain facts are in blue

$f_1 : (i_1 \text{ sameName } i_2)$ $f_2 : (i_1 \text{ sameBirthDate } i_2)$ $f_3 : (i_2 \text{ marriedTo } i_3)$

$f_4 : (i_4 \text{ marriedTo } i_3)$ $f_5 : (i_2 \text{ sameName } i_4)$

Provenance of inferred facts

Inferred facts	Provenance	Uncertainty Provenance
$(i_2 \text{ sameAs } i_4)$	$(e(r_1) \wedge e(f_5)) \vee (e(r_3) \wedge e(f_3) \wedge e(f_4))$	\top
$(i_1 \text{ sameAs } i_2)$	$(e(r_1) \wedge e(f_1)) \vee (e(r_2) \wedge e(f_1) \wedge e(f_2))$	$e(r_2) \wedge e(f_1)$
$(i_1 \text{ sameAs } i_4)$	$e(r_4) \wedge \text{Prov}((i_1 \text{ sameAs } i_2))$ $\wedge \text{Prov}((i_2 \text{ sameAs } i_4))$	$e(r_2) \wedge e(f_1)$

Illustrative Example (cont.)

Rules: uncertain rules are in red, certain rules are in blue

$$r_1 : (?x \text{ sameName } ?y) \Rightarrow (?x \text{ sameAs } ?y)$$

$$r_2 : (?x \text{ sameName } ?y), (?x \text{ sameBirthDate } ?y) \Rightarrow (?x \text{ sameAs } ?y)$$

$$r_3 : (?x \text{ marriedTo } ?z), (?y \text{ marriedTo } ?z) \Rightarrow (?x \text{ sameAs } ?y)$$

$$r_4 : (?x \text{ sameAs } ?z), (?z \text{ sameAs } ?y) \Rightarrow (?x \text{ sameAs } ?y)$$

Facts: uncertain facts are in red, certain facts are in blue

$$f_1 : (i_1 \text{ sameName } i_2) \quad f_2 : (i_1 \text{ sameBirthDate } i_2) \quad f_3 : (i_2 \text{ marriedTo } i_3)$$

$$f_4 : (i_4 \text{ marriedTo } i_3) \quad f_5 : (i_2 \text{ sameName } i_4)$$

Computation of the inferred facts probabilities

Inferred facts	Uncertainty Provenance	Probability
$(i_2 \text{ sameAs } i_4)$	\top	1
$(i_1 \text{ sameAs } i_2)$	$e(r_2) \wedge e(f_1)$	$Pr(e(r_2)) \times Pr(e(f_1))$
$(i_1 \text{ sameAs } i_4)$	$e(r_2) \wedge e(f_1)$	$Pr(e(r_2)) \times Pr(e(f_1))$

Illustrative Example (cont.)

Rules: uncertain rules are in red, certain rules are in blue

$r_1 : (?x \text{ sameName } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

$r_2 : (?x \text{ sameName } ?y), (?x \text{ sameBirthDate } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

$r_3 : (?x \text{ marriedTo } ?z), (?y \text{ marriedTo } ?z) \Rightarrow (?x \text{ sameAs } ?y)$

$r_4 : (?x \text{ sameAs } ?z), (?z \text{ sameAs } ?y) \Rightarrow (?x \text{ sameAs } ?y)$

Facts: uncertain facts are in red, certain facts are in blue

$f_1 : (i_1 \text{ sameName } i_2)$ $f_2 : (i_1 \text{ sameBirthDate } i_2)$ $f_3 : (i_2 \text{ marriedTo } i_3)$

$f_4 : (i_4 \text{ marriedTo } i_3)$ $f_5 : (i_2 \text{ sameName } i_4)$

Computation of the inferred facts probabilities

Inferred facts	Uncertainty Provenance	Probability
$(i_2 \text{ sameAs } i_4)$	\top	1
$(i_1 \text{ sameAs } i_2)$	$e(r_2) \wedge e(f_1)$	0.8×0.9
$(i_1 \text{ sameAs } i_4)$	$e(r_2) \wedge e(f_1)$	0.8×0.9

Experiments: interlinking DBpedia and MusicBrainz

Size and number of entities in the two datasets

Class	DBpedia	MusicBrainz
Person	1,445,773	385,662
Band	75,661	197,744
Song	52,565	448,835
Album	123,374	1,230,731
Number of RDF triples	73 millions	112 millions

86 rules from which 50 are certain and 36 are uncertain

ID	Rules
sameAsBirthDate	$(?x \text{ :solrPSimilarName } ?l), (?y \text{ skos:myLabel } ?l),$ $(?x \text{ dbo:birthDate } ?date), (?y \text{ mb:beginDateC } ?date)$ $\Rightarrow (?x \text{ :sameAsPerson } ?y)$
sameAsMemberOfBand	$(?x \text{ :solrPSimilarName } ?l), (?y \text{ skos:myLabel } ?l),$ $(?y \text{ mb:member_of_band } ?gr2), (?gr2 \text{ skos:myLabel } ?lg),$ $(?gr1 \text{ dbp:members } ?x), (?gr1 \text{ :solrGrSimilarName } ?lg)$ $\Rightarrow (?x \text{ :sameAsPerson } ?y)$

Experimental results

Gain of rule chaining

43,923 links not discovered by Silk among the **144,467 sameAs links discovered by ProbFR** between DBpedia and MusicBrainz

Gain of using uncertain rules for improving recall without losing much in precision (precision and recall estimated on samples)

	DBpedia and MusicBrainz					
	Only certain rules			All rules		
	P	R	F	P	R	F
Person	1.00	0.08	0.15	1.00	0.80	0.89
Band	1.00	0.12	0.21	0.94	0.84	0.89
Song	-	-	-	0.96	0.74	0.84
Album	-	-	-	1.00	0.53	0.69

Gain of exploiting probabilities to filter out wrong sameAs links

	P	R	F
Band ≥ 0.90	1.00	0.80	0.89
Song ≥ 0.60	1.00	0.54	0.72

Lessons learnt and perspectives

Probabilistic Datalog: a good trade-off for reasoning with uncertainty in Linked Data

Some restrictions compared to general probabilistic logical frameworks (e.g., Markov Logic)

- uncertain formulas restricted to Horn rules and ground facts
- probabilities computed for inferred facts only

Better scalability and more transparency

- explanations on probabilistic inference for end-users
- useful traces for experts to set-up the rules probabilities

Future work

ANR ELKER project

- A method to set up automatically the threshold for filtering the probabilistic sameAs facts to be retained
- A backward-reasoning algorithm on probabilistic rules for importing on demand useful data from external sources

Concluding message

- Semantic Web standards, data and applications are there, due to the simplicity and flexibility of the RDF data model
 - Promising applications are emerging for which reasoning on data is central
 - Fact checking
 - Interactive and personalized data exploration and analytics
 - Many challenges remain
 - to handle at large scale the incomplete and uncertain data
- => Combining numerical and symbolic AI is hard but worthwhile to investigate more deeply**

Joint work with many persons

Mustafa Al Bakri, Mohannad Almasri, Manuel Atencia, Shadi Baghernezhad, Jérôme David, Loic Druette (Univ. Lyon) , Fabrice Jouanot, Cyril Labbé, Steffen Lalande (INA), Behrooz Omidvar, Olivier Palombi (LADAF, LJK), Adam Sanchez, Federico Ulliana (LIRMM),...

THANKS