

Ontologies and Functional Dependencies for Data Integration and Reconciliation

Abdelghani BAKHTOUCHI
ESI, Algiers, Algeria
a_bakhtouchi@esi.dz

Ladjel BELLATRECHE
LISI/ENSMA, Poitiers, France
bellatreche@ensma.fr

Yamine AIT AMEUR
ENSEEIH, IRIT, Toulouse, France
yamine@enseeiht.fr

Domain ontology: definitions

What is an **ontology**?

« *An explicit specification of a conceptualization* » [Gruber 93]

Conceptualization of domain classes and properties [JoDS'08]

- Formal
- Consensual
- Capability to be referenced

« *a formal and consensual dictionary of categories and properties of entities of a domain and the relationships that hold among them* »

A Taxonomy of Domain Ontologies

- ❑ **DB: Canonical vocabulary of the concepts of a domain**
identifier + primitive concepts {classes / properties / datatypes}
→ **Canonical Conceptual Ontology (CCO)**

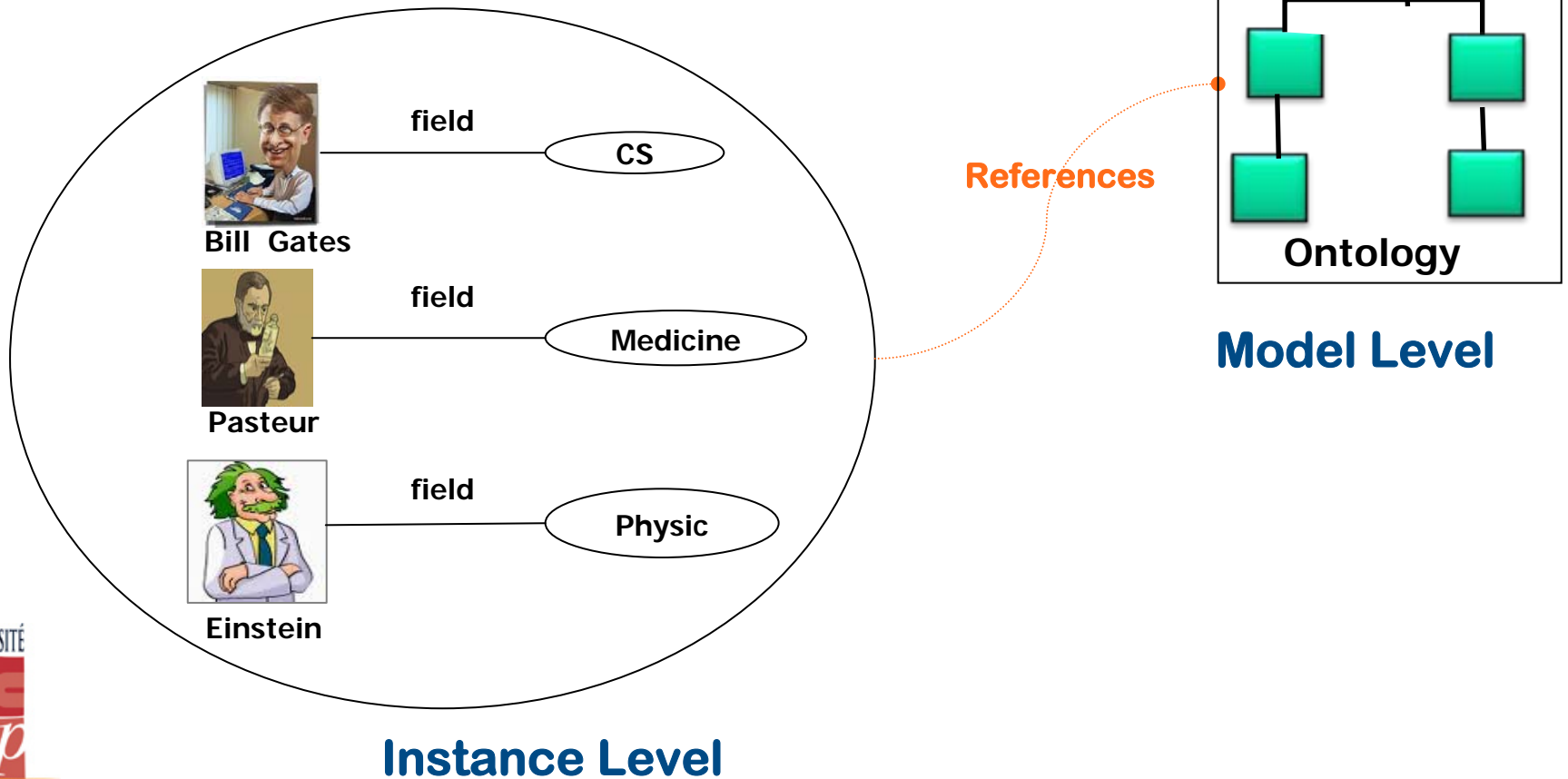
- ❑ **AI: Deductive capabilities**
Concept equivalence operators
→ **Non Canonical Conceptual Ontology (NCCO)**

- ❑ **Computational Linguistics: Terms of a domain**
{ words } + similarities and linguistic relationships
→ **Linguistic Ontology (LO)**

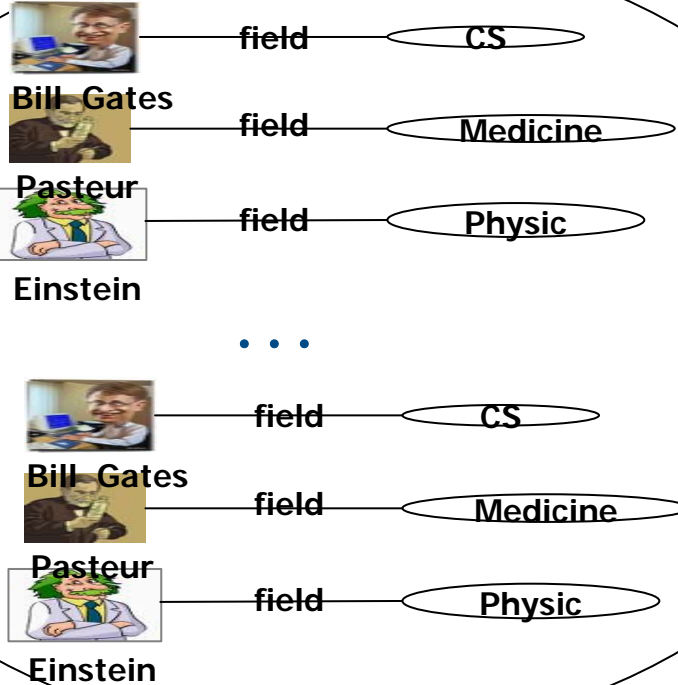
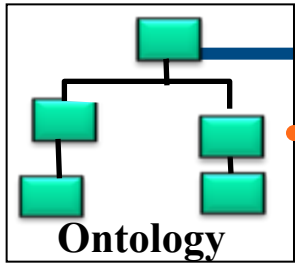
Three points of view = three complementary categories of Ontologies

Explicitation of Data Semantic

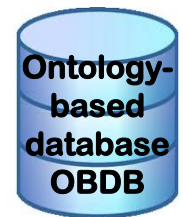
- Existence of ontological instances
 → RDF, RDF Schema, PLIB, OWL



Explosion of Ontological Instances



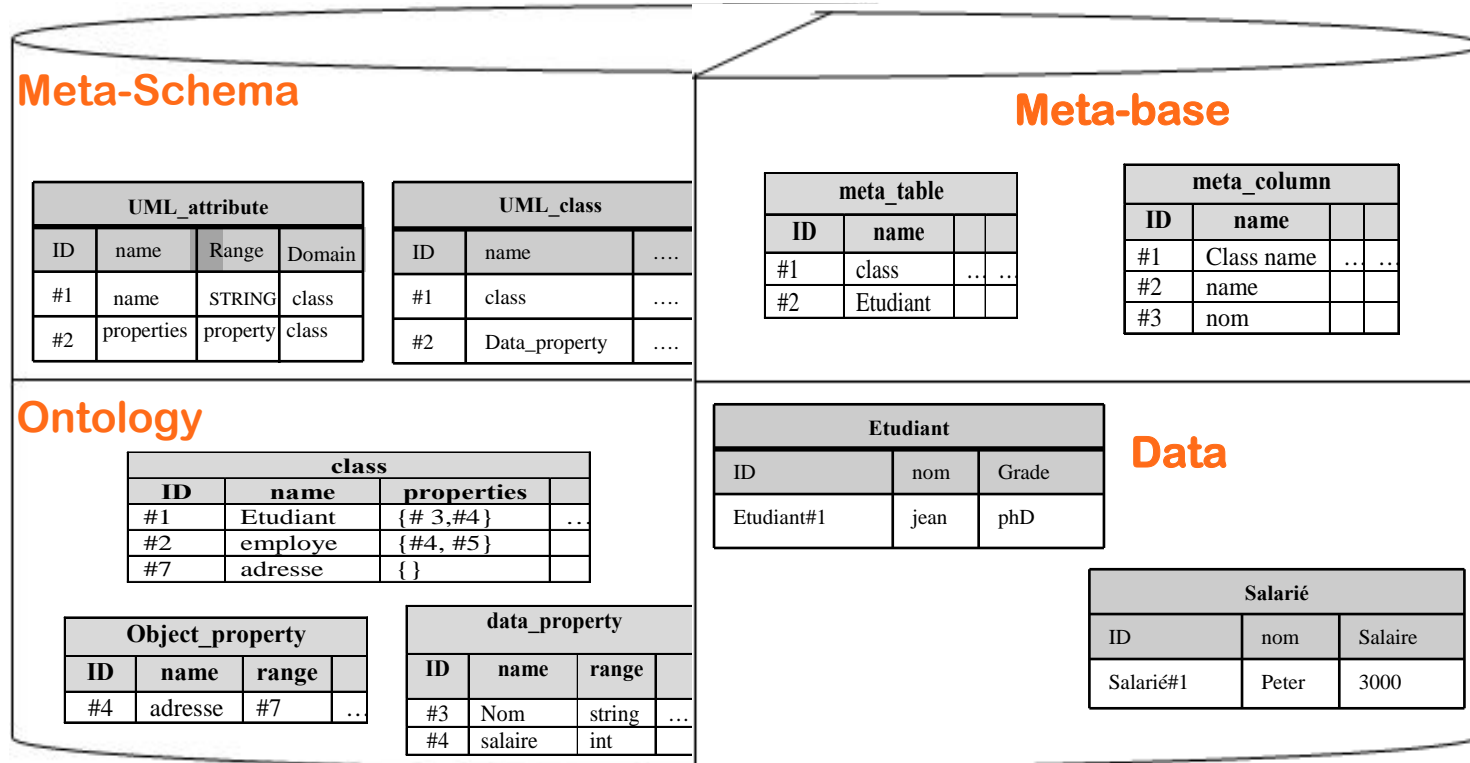
Persistence



Industrial: Oracle, IBM SOR
Academic: Sesame, Owingres, OntoDB

Ontology-based Database Sources: OntoDB

- **Schema**: meta-model of ontology model
- **Instances** :
 - Model of ontology +
 - meta-model of the used modeling language



- **Schema** : ontology model
- **Instances** : ontologies

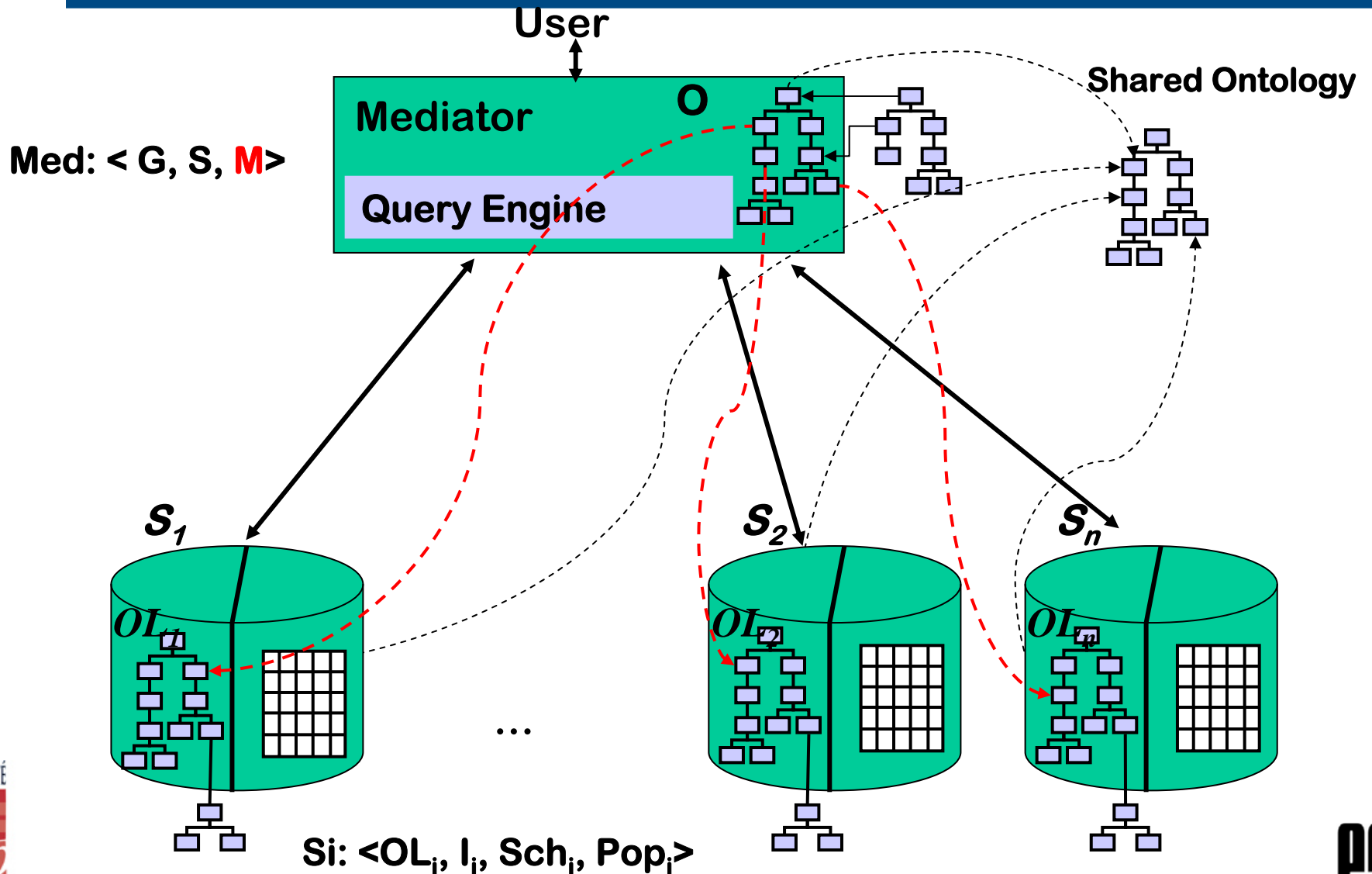
- **Schema** : subset of the ontology
- **Instances** : instances

Agenda

- 1. Ontology-based Integration System**
- 2. Data Reconciliation**
- 3. Components of our Integrated System**
- 4. Validation**
- 5. Conclusion & Perspectives**

OBDB Candidate for Integrated Systems

1. **Ontology-based Integration System**
2. Data Reconciliation
3. Components of our IS
4. Validation
5. Conclusion & Perspectives



Data Reconciliation

1. Ontology-based Integration System
2. **Data Reconciliation**
3. Components of our IS
4. Validation
5. Conclusion & Perspectives

❑ Existing Reconciliation Approaches:

1. Source Entities Representing the same Concept have the **Same Key**
Observer, PicseI2, COIN, etc.

→ **A Strong Hypothesis that Violates Source Autonomy**

2. Different Keys

- ❑ Use of Statistical Methods (Affinity Measures between Concepts)

→ **Not Really Suitable for Sensitive Applications (Bank, Engineering, ...)**

→ Data Integration and Reconciliation are Treated in **Isolated Way**

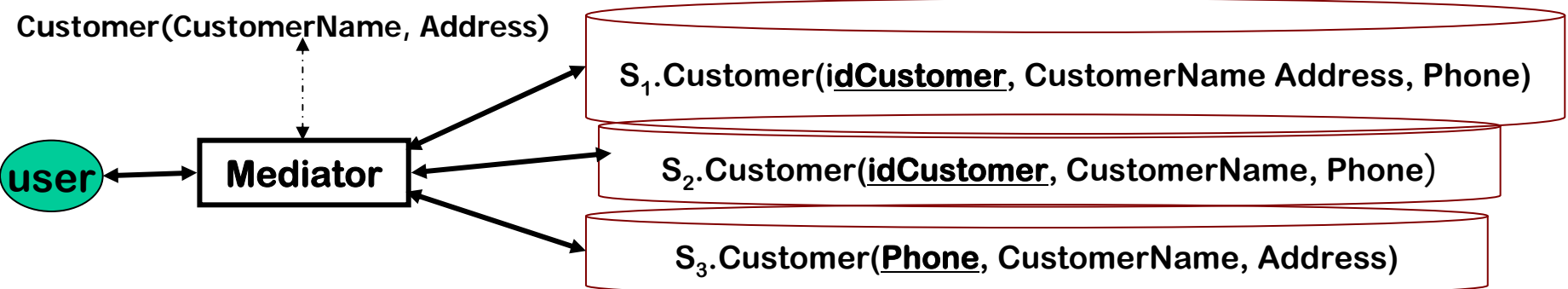
❑ Our Proposal: **Joint Treatment of Data Integration & Reconciliation**

→ Use of Functional Dependency (FD) Defined on Ontology Classes

→ Key Candidates become Ontological Concepts

Motivating Example

1. Ontology-based Integration System
2. **Data Reconciliation**
3. Components of our IS
4. Validation
5. Conclusion & Perspectives



1. Only results from S_1 and S_2 may be reconciled

→ lost of a part of results.

2. Reconciliation of Results from S_1 , S_2 and S_3 without Duplicate Elimination

→ Duplicates may cause errors (Fusion)

3. Duplicates elimination using Statistical Methods

→ Inexact Results

❑ **Functional Dependencies on Customer [IGPL'11, Romero et al. 09]**

- Fd₁: IdCustomer → CustomerName
- Fd₂: IdCustomer → Address
- FD₃: IdCustomer → Phone
- Fd₄: Phone → CustomerName
- Fd₅: Phone → Address.

❑ **Key Candidates:**

- IdCustomer (S_1, S_2)
- Phone (S_1, S_2, S_3): **Reconciliation Key**

Formal Definitions of an Ontology

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives

□ Traditional Definition [Jean'07, JoDS'08]

O: $\langle C, P, \text{Sub}, \text{Applic} \rangle$

- **C**: ontology classes;
- **P**: set of properties used to describe of ontology classes **C**;
- **Sub**: $C \rightarrow 2^C$ subsumption relationship;
- **Applic**: $C \rightarrow 2^P$ applicable properties for each class.

□ FD = (R, LP, RP): fd $R : LP \rightarrow RP$ [Calbimonte'09]

- **R**: root class
- **LP**: Left Part: list of properties $\{p_1 \dots p_n\}$
- **RP**: Right Part: a single property

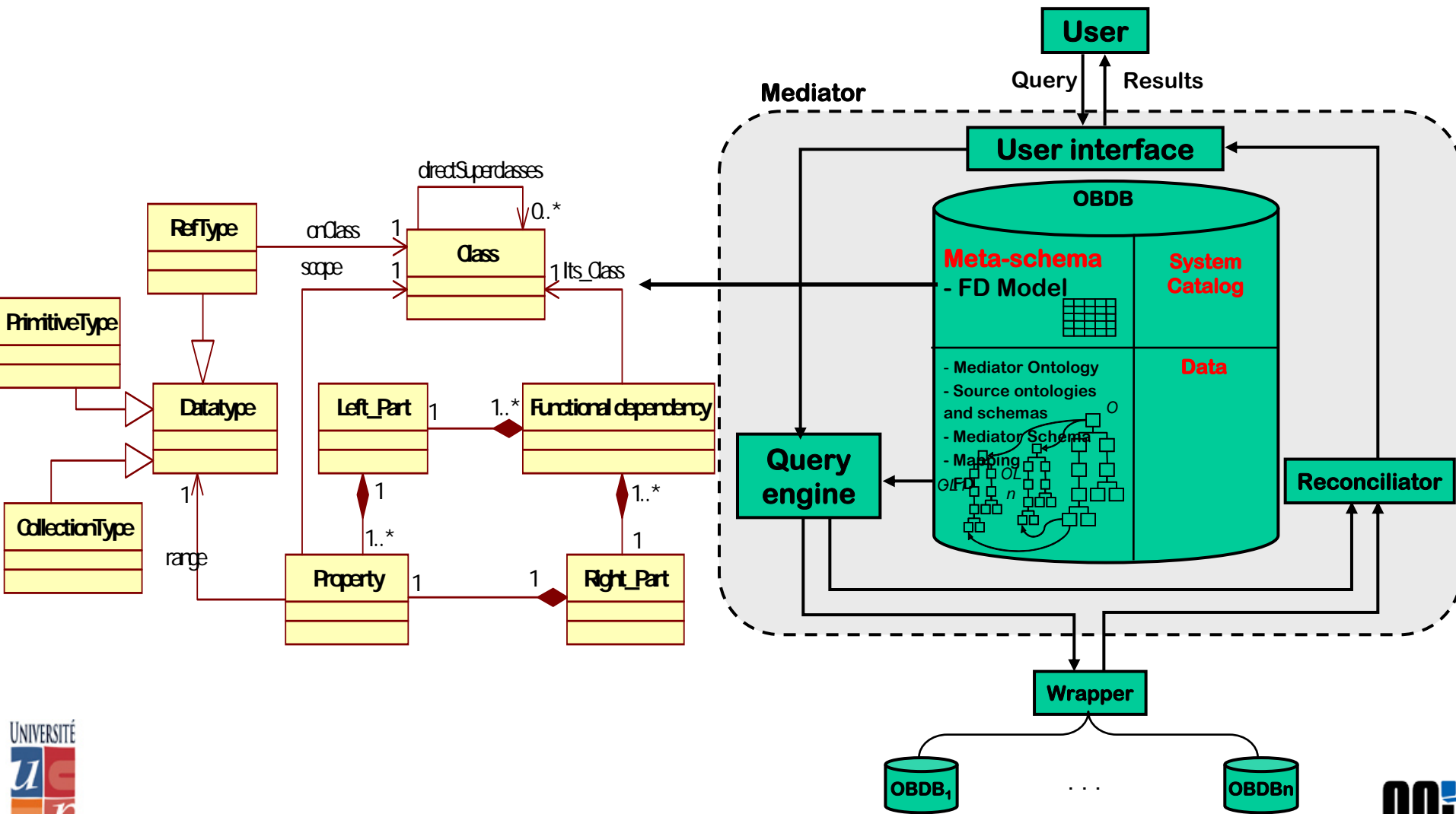
□ Formal Definition Enriched by FD

O: $\langle C, P, \text{Sub}, \text{Applic}, \text{FD} \rangle$ with **FD** : $\langle R, LP, RP \rangle$

$R \in C, \quad LP \in 2^P, \quad RP \in P$

Global Architecture

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives



Components of Mediator

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives

Med: $\langle G, S, M \rangle$

1. **G: $\langle O, Sch \rangle$**

- $O \langle C, P, Applic, Sub, FD \rangle$: mediator ontology
- $Sch: C \rightarrow 2^P$ associates to each class $c \in C$ the properties describing the instances of the class c , which are valuated in at least one integrated source.

2. **S = $\{S_1, S_2, \dots, S_n\}$**

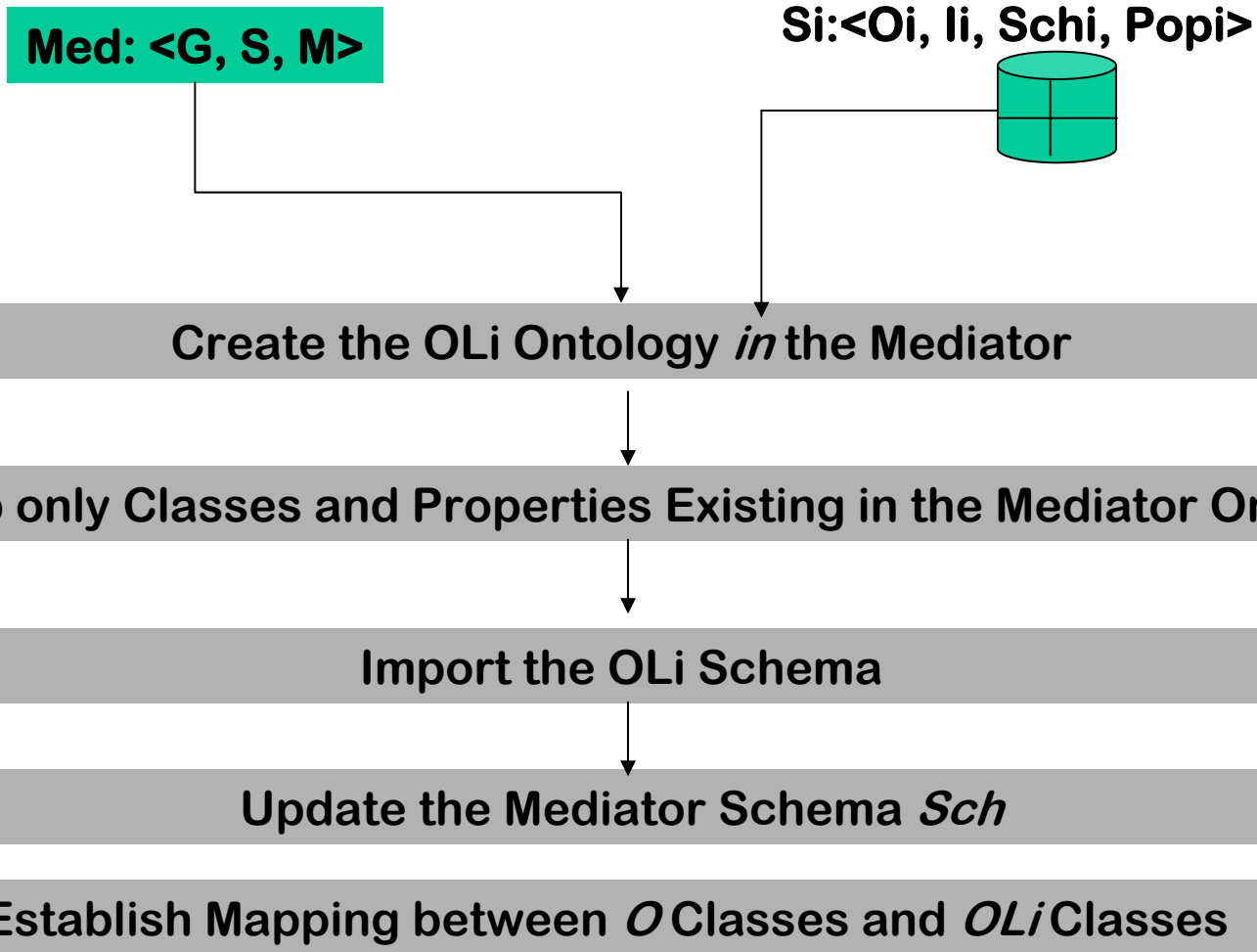
- $S_i: \langle OL_i, SchL_i \rangle$

3. **M: $C \rightarrow 2\{C_1 \cup \dots \cup C_n\}$**

- Mapping between the classes of mediator ontology O and the classes of source ontologies

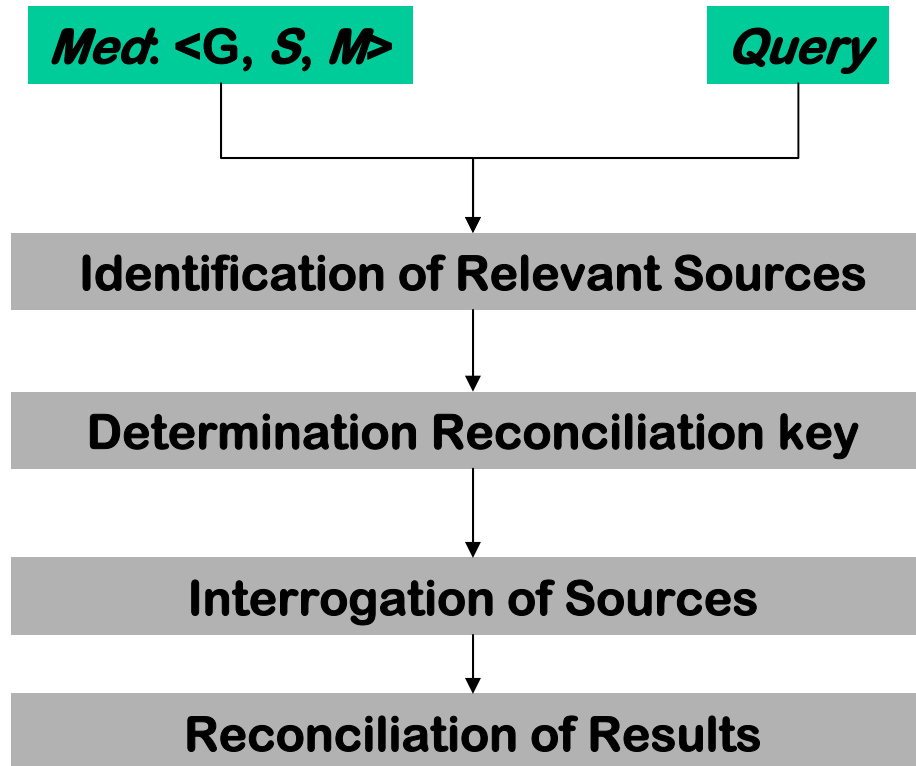
New source Integration

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives



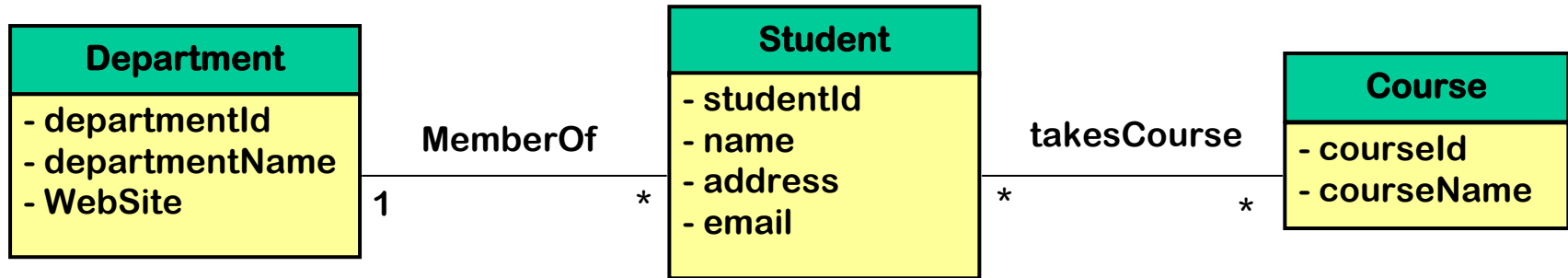
Query Processing

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives



Reconciliation Key Determination

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives



□ FD

- fd1..4 Student : studentId → name, address, email, MemberOf
- fd5..7 Student : email → name, address, MemberOf
- fd8,9 Department : departmentId → departmentName, WebSite
- fd10 Course : courseId → courseName

□ Query:

- $Q(x; y; z) :- \text{Student}(x1), \text{takesCourse}(x1, x2), \text{Course}(x2), \text{memberOf}(x1, x3), \text{Department}(x3), \text{name}(x1, x), \text{courseName}(x2, y), \text{departmentName}(x3, z)$

Reconciliation Key Determination

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives

❑ Projected Properties:

Proj = {Name, CourseName, DepartmentName}

❑ Direct FD

fd₁ .. fd₁₀

❑ Generated FD (by the means of **functional property**: MemberOf → departmentId; MemberOf → Email)

fd₁₁ : studentId → departmentId, departmentName, WebSite

fd₁₂ : email → departmentId, departmentName, WebSite

❑ Reconciliation Key (RK) Generation

– RK₀ contains all LP of DFs = {StudentId, Email, DepartmentId, CourseId}

– StudentId → DepartmentId ⇒ RK₁ = {StudentId, Email, CourseId}

– StudentId → Name, DepartmentName

– Email → Name, DepartmentName

1. If StudentId is valued in all sources ⇒ RK = {StudentId, CourseId}

2. If Email is valued in all sources ⇒ RK = {Email, CourseId}

Validation (I)

1. Ontology-based Integration System
2. Data Reconciliation
3. **Components of our IS**
4. Validation
5. Conclusion & Perspectives

❑ Dataset

- Lehigh University Benchmark
- Ontology: 45 classes, 32 properties (25 object properties)
- OBDB sources are generated
- 14 Queries

❑ Experiment Environment

- Intel Pentium IV, 3,2 GHz, 1 GB of Memory

❑ Execution Time of a Query includes:

1. Identifying Relevant Source(s)
2. Finding FD that hold for a query
3. Finding Reconciliation Key
4. Executing Query on each Relevant Source

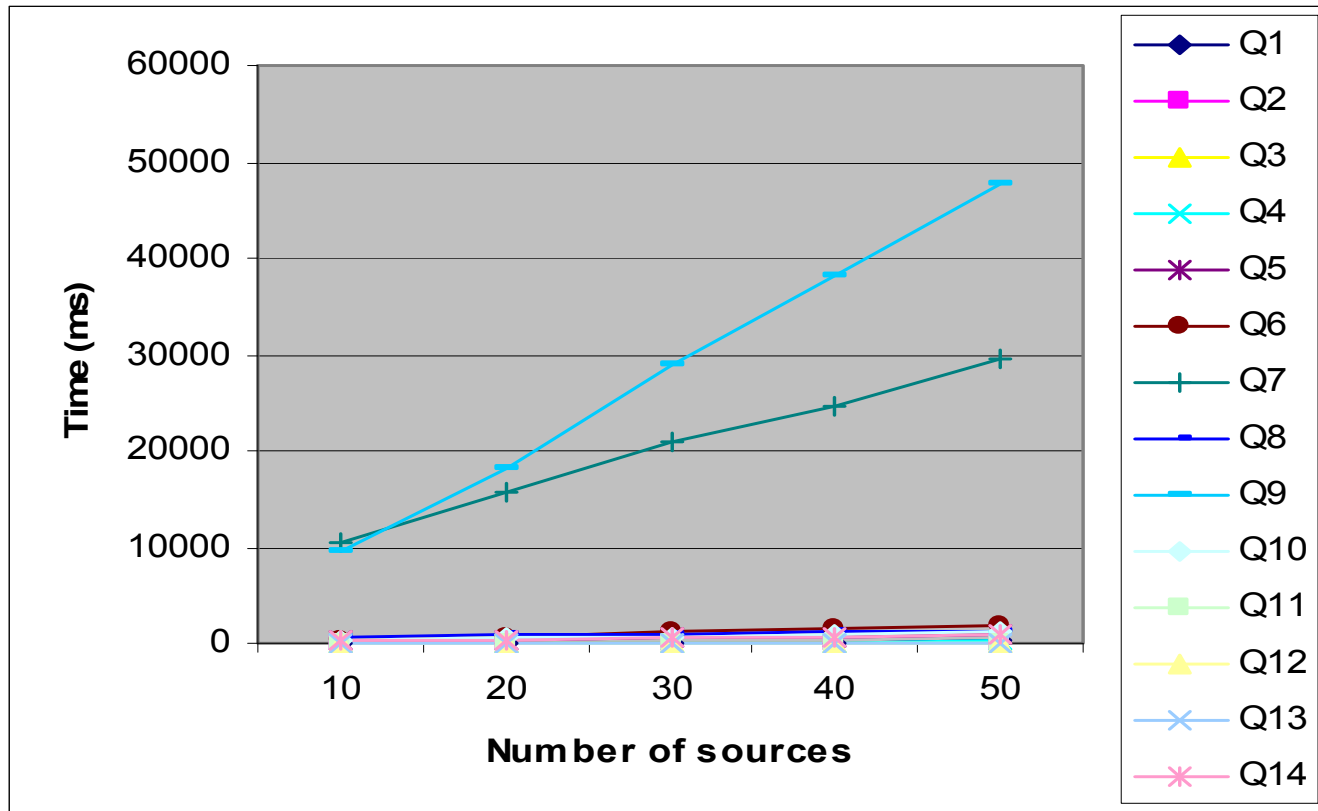
} **Mediator Tasks**

} **Source Task**

Validation (II)

1. Ontology-based Integration System
2. Data Reconciliation
3. Components of our IS
4. **Validation**
5. Conclusion & Perspectives

Query Response Time vs. Number of Sources



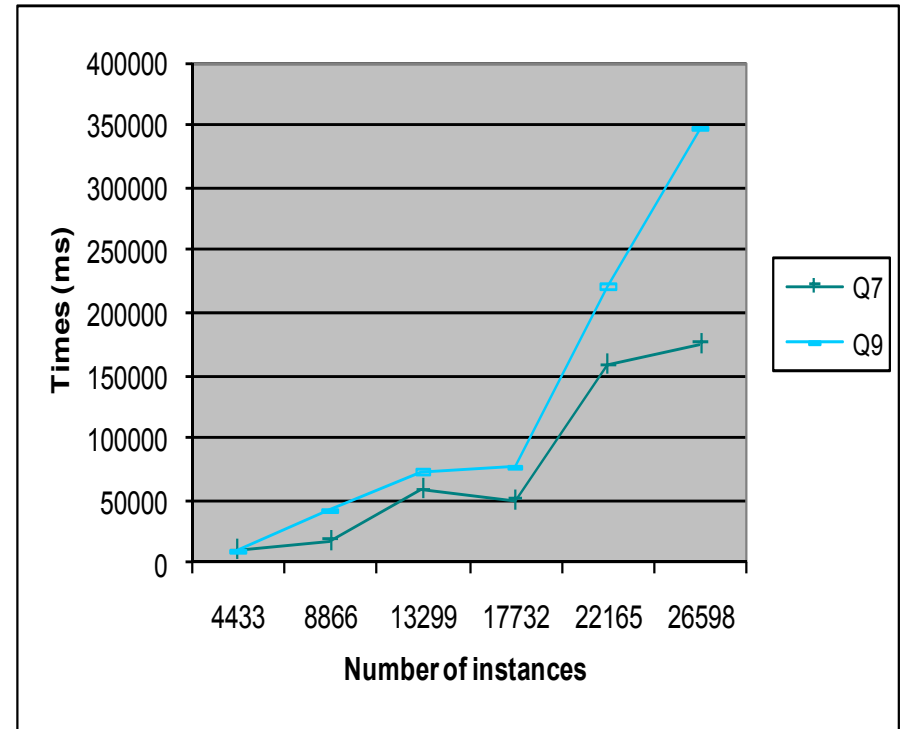
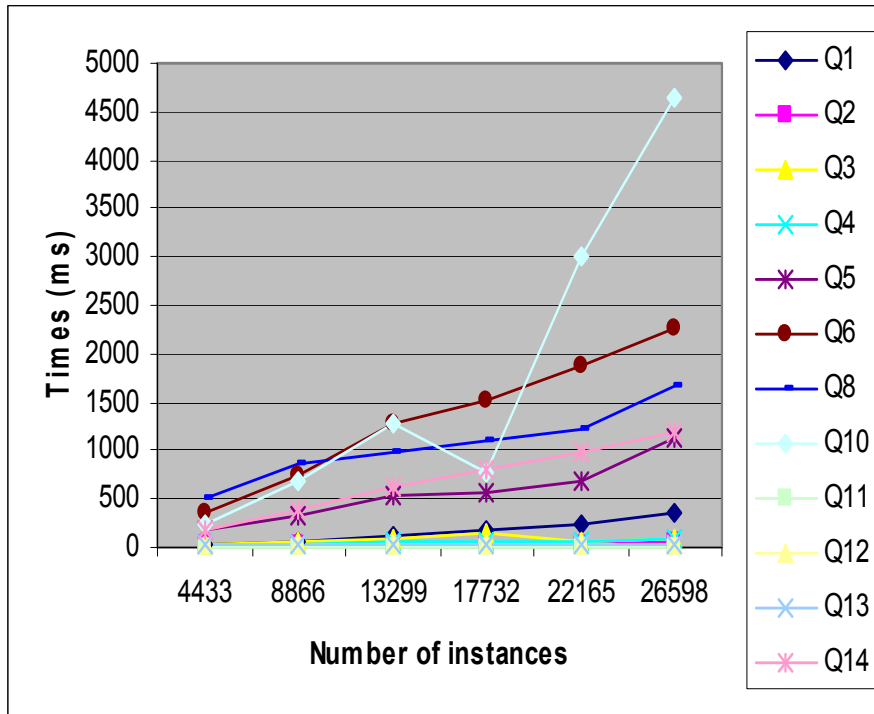
Lessons:

1. Mediator task processing time is **negligible**
2. Query processing time is **important**

Validation (III)

1. Ontology-based Integration System
2. Data Reconciliation
3. Components of our IS
4. **Validation**
5. Conclusion & Perspectives

Query Response Time vs. Number of Instances



Conclusion & Future Works

1. Ontology-based Integration System
2. Data Reconciliation
3. Components of our IS
4. Validation
5. **Conclusion & Perspectives**

- ❑ Enrichment of Ontology by FD
- ❑ Joint Solution for Data Integration and Reconciliation
- ❑ FD based Approach for Instance Reconciliation
- ❑ Validation using LUBM

- ❑ Intensive Experiments for Evaluating the Robustness of our System
- ❑ Development of Quality Metrics
- ❑ Ontology-based Data Warehouses