

Using ontologies for P2P data sharing and ETL design

Dimitrios Skoutas

(joint work with Alkis Simitsis, Verena Kantere, and Timos Sellis)

Dept. of Electrical and Computer Engineering
National Technical University of Athens, Hellas
dskoutas@dblaboratory.ece.ntua.gr
<http://www.dblaboratory.ece.ntua.gr>



Bertinoro Workshop on Information Integration (INFINT 2007)
Bertinoro, Italy

Outline

- Data sharing in PDMS
- Ontology-based ETL design

Peer Data Management Systems

- P2P networks
 - massive sharing of structured or unstructured data
- Structured vs. unstructured overlays
 - Structured
 - data is indexed and stored at peers with pre-specified characteristics, e.g. peer id
 - peers are connected to specific peers
 - Unstructured
 - no pre-specified connection of peers and storage/indexing of data

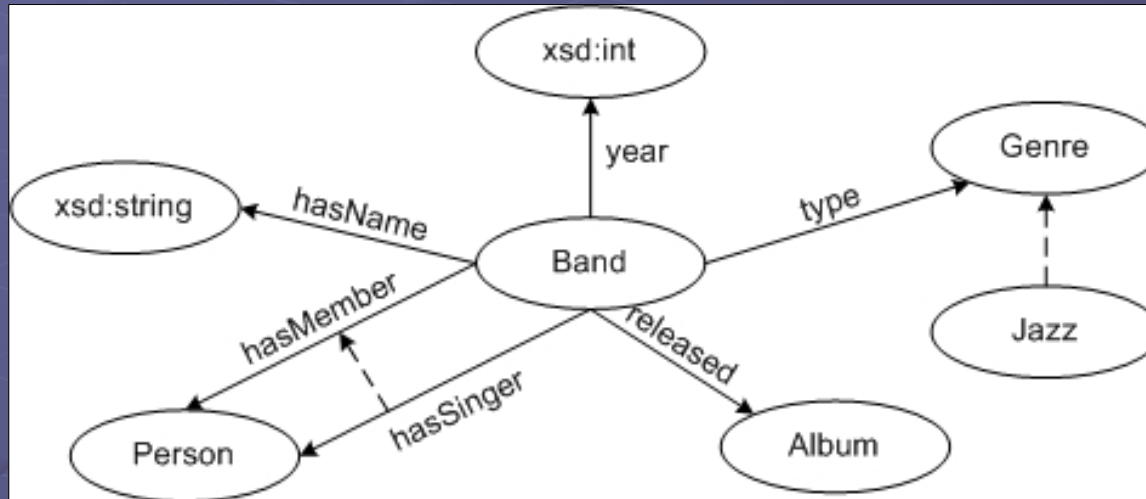
Problem description

- An unstructured PDMS
 - local relational database
 - pairwise mappings of the forms GAV/LAV/GLAV
 - SPJ queries
 - expressed on the local schema
 - reformulated according to the mappings
- An OWL ontology describing the application domain

Problem description (cont'd)

- Semantic peer $P = (R, O, A)$
 - R: local database schema
 - O: the domain ontology
 - A: semantic annotation
- Two issues
 - Semantic similarity between peer schemas
 - Semantic similarity between reformulated queries

Sample ontology and peer schemas



- P_1 : bands (name, members, year)
- P_2 : bands (name, singer, year)
- M_{P_1, P_2} : bands (name, members, year) :- bands (name, singer, year)

Semantic annotations

- P_1 : bands having at least one album
- P_2 : bands that play Jazz music, were formed before the year 2000, and have released at least 3 albums
- $\text{Band_}P_1 \equiv \text{Band } n \geq_1 \text{ released}$
- $\text{Band_}P_2 \equiv \text{Band } n \forall \text{ type.Jazz } n \geq_3 \text{ released}$
 $n \forall \text{ year.} (\leq 2000)$

Comparison of peer schemas

- Express the degree of relevance between the interests of peers (subset, superset, overlap, etc.)
- Semantic similarity is assessed by comparing the classes annotating the peer schemas
- The similarity function needs to be asymmetric

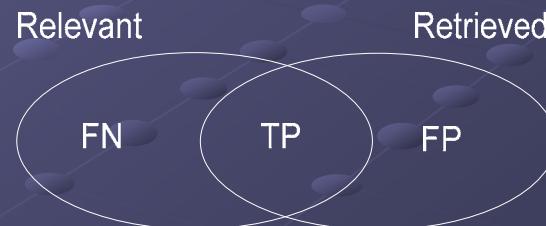
Recall and Precision

- Use notions recall and precision from Information Retrieval

$$recall = \frac{|relevant \cap retrieved|}{|relevant|}$$

$$precision = \frac{|relevant \cap retrieved|}{|retrieved|}$$

$$F_{\alpha} = \frac{(1 + \alpha) \cdot precision \cdot recall}{\alpha \cdot precision + recall}$$



TP: True Positive

FP: False Positive

FN: False Negative

Adapting recall and precision

- We adapt the notions of recall and precision as follows:

$$\text{recall}(C_1, C_2) = \frac{|\{x \mid x \in (C_1 \cap C_2)\}|}{|\{x \mid x \in C_1\}|}$$

$$\text{precision}(C_1, C_2) = \frac{|\{x \mid x \in (C_1 \cap C_2)\}|}{|\{x \mid x \in C_2\}|}$$

<i>Relationship</i>	<i>equivalent</i>	<i>subset</i>	<i>superset</i>	<i>overlap</i>	<i>disjoint</i>
<i>Recall</i>	1	1	<1	<1	0
<i>Precision</i>	1	<1	1	<1	0

Semantic similarity

- Similarity is assessed by the semantic information provided by the ontology
 - class hierarchy
 - properties
 - property hierarchy
 - value and cardinality restrictions

Calculating recall and precision

- Based on the class hierarchy
 - ratio of common ancestors
 - similarity increases with the hierarchy depth

$$\text{recall}_C(C_1, C_2) = \frac{|A(C_1) \cap A(C_2)|}{|A(C_2)|}$$

$$\text{precision}_C(C_1, C_2) = \frac{|A(C_1) \cap A(C_2)|}{|A(C_1)|}$$

Calculating recall and precision (cont'd)

- Based on the properties of classes
 - property hierarchy

$$\textit{recall}(p_1, p_2) = \frac{|A(p_1) \cap A(p_2)|}{|A(p_2)|}$$

$$\textit{precision}(p_1, p_2) = \frac{|A(p_1) \cap A(p_2)|}{|A(p_1)|}$$

Calculating recall and precision (cont'd)

- Based on the properties of classes (cont'd)

- property restrictions

- value restrictions on object properties: \forall type.Jazz
- value restrictions on datatype properties: \forall year.(≤ 2000)
- cardinality restrictions: \geq_3 released

<i>case</i>	<i>recall</i> (R_1, R_2)	<i>precision</i> (R_1, R_2)
$R_1 \equiv R_2$	1	1
$R_1 \subseteq R_2$	1	0.5
$R_1 \supseteq R_2$	0.5	1
$R_1 \cap R_2 \neq \emptyset$	0.5	0.5
$R_1 \cap R_2 = \emptyset$	0	0

Calculating recall and precision (cont'd)

- Based on the properties of classes (cont'd)
 - combining property hierarchy and restrictions

$$\text{recall}(p_1, p_2) = \frac{|A(p_1) \cap A(p_2)|}{|A(p_2)|} \cdot \prod_{R(p_2)} \text{recall}(R'_i(p_1), R_i(p_2))$$

$$\text{precision}(p_1, p_2) = \frac{|A(p_1) \cap A(p_2)|}{|A(p_1)|} \cdot \prod_{R(p_1)} \text{precision}(R_i(p_1), R'_i(p_2))$$

Calculating recall and precision (cont'd)

- Based on the properties of classes (cont'd)

$$\text{recall}(C_1, C_2) = \frac{\sum_{P(C_2)} \text{recall}(p'(C_1), p(C_2))}{|P(C_2)|}$$

$$\text{precision}(C_1, C_2) = \frac{\sum_{P(C_1)} \text{precision}(p(C_1), p'(C_2))}{|P(C_1)|}$$

Calculating recall and precision (cont'd)

- Extending to sets of classes

$$\text{recall}(C_1, C_2) = \frac{\sum_{C_i \in C_1} \text{recall}(C_i, C'_i)}{|C_1|}$$

$$\text{precision}(C_1, C_2) = \frac{\sum_{C_i \in C_2} \text{precision}(C'_i, C_i)}{|C_2|}$$

Query reformulation

- A query Q_o is forwarded from peer P_i to peer P_j , and is reformulated as Q_r
 - some attributes may not be rewritten
 - some attributes may be rewritten approximately
 - some conditions may be lost
 - some conditions may be inserted

Extending the similarity measure

● Rewritten attributes

■ t was rewritten to t'

- $\text{recall}(t, t') = \text{recall}(p_t, p_{t'})$
- $\text{precision}(t, t') = \text{precision}(p_t, p_{t'})$

■ t was not rewritten

- $\text{recall} = 0$
- precision is not affected

Extending the similarity measure (cont'd)

● Rewritten conditions

1. A query Q_0 is issued at peer P_i
2. $C_{Q_0} \leftarrow$ the set of classes annotating the relations in Q_0
3. $C_{Q_0,e} \leftarrow$ the classes in C_{Q_0} enhanced with additional value restrictions, according to the conditions specified in Q_0
4. Q_0 is forwarded to peer P_j and is rewritten as Q_r
5. $C_{Q_r,e} \leftarrow$ apply steps 2 and 3 for Q_r

Example

- Q_0 : SELECT name, members, year FROM bands
WHERE year \geq 1980 AND year < 1990
- $C_{Q_0,e}$: Band_P₁ \equiv Band $\pi \geq_1$ released $\pi \forall$ year.([1980,1990))
- Q_r : SELECT name, singer, year FROM bands
WHERE year \geq 1980 AND year < 1990
- $C_{Q_r,e}$: Band_P₂ \equiv Band $\pi \forall$ type.Jazz $\pi \geq_3$ released
 $\pi \forall$ year.([1980,1990))

Similarity measure for rewritten queries

- Combine results for attributes and conditions

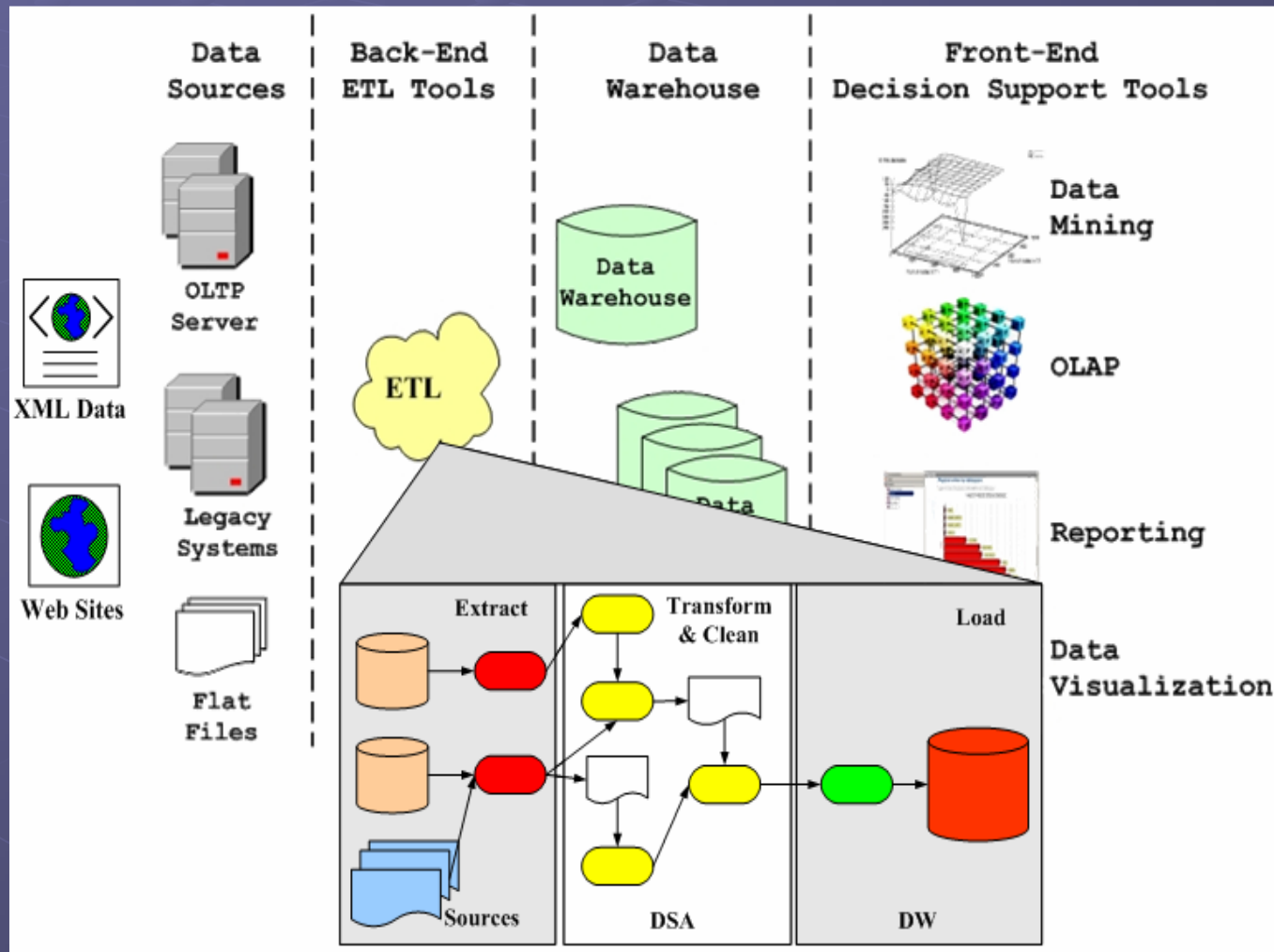
$$\text{recall}(C_{Q_{o,e}}, C_{Q_{r,e}}, Q_o, Q_r) = \frac{\sum_{t \in S(Q_o)} \text{recall}(p_t, p_{t'})}{|S(Q_o)|} \cdot \text{recall}(C_{Q_{o,e}}, C_{Q_{r,e}})$$

$$\text{precision}(C_{Q_{o,e}}, C_{Q_{r,e}}, Q_o, Q_r) = \frac{\sum_{t \in S(Q_r)} \text{precision}(p_{t'}, p_t)}{|S(Q_r)|} \cdot \text{precision}(C_{Q_{o,e}}, C_{Q_{r,e}})$$

Outline

- Data sharing in PDMS
- Ontology-based ETL design

Extract-Transform-Load (ETL)



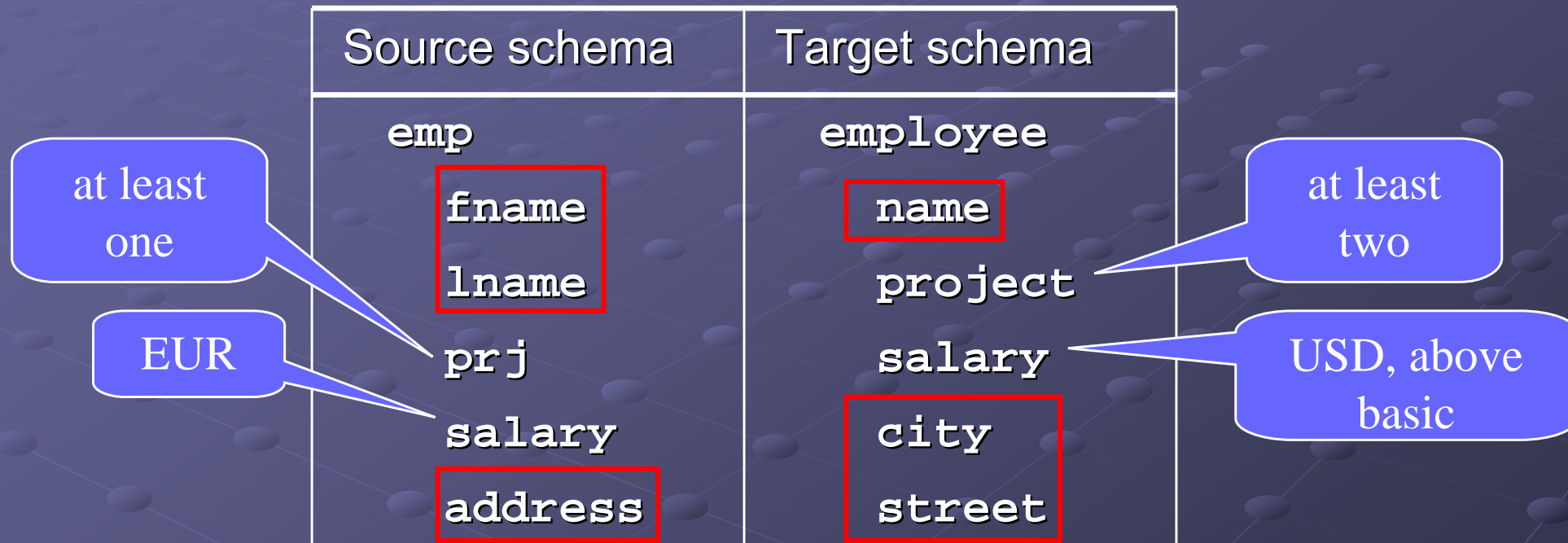
Problem description

- The problem of heterogeneity in data sources
 - structural heterogeneity
 - data stored under different schemata
 - semantic heterogeneity
 - different naming conventions
 - e.g., homonyms, synonyms
 - different representation formats
 - e.g., units of measurement, currencies, encodings
 - different ranges of values
- Two main goals
 - specify inter-schema mappings
 - identify appropriate transformations

Overview

- Construct a suitable application ontology
- Annotate the datastores
 - establish mappings between the datastore schemas and the ontology
- Apply reasoning techniques to
 - select relevant sources
 - to identify required transformations

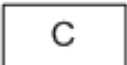
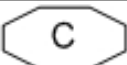

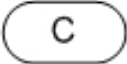
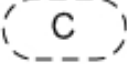
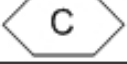






Reference example



The application ontology

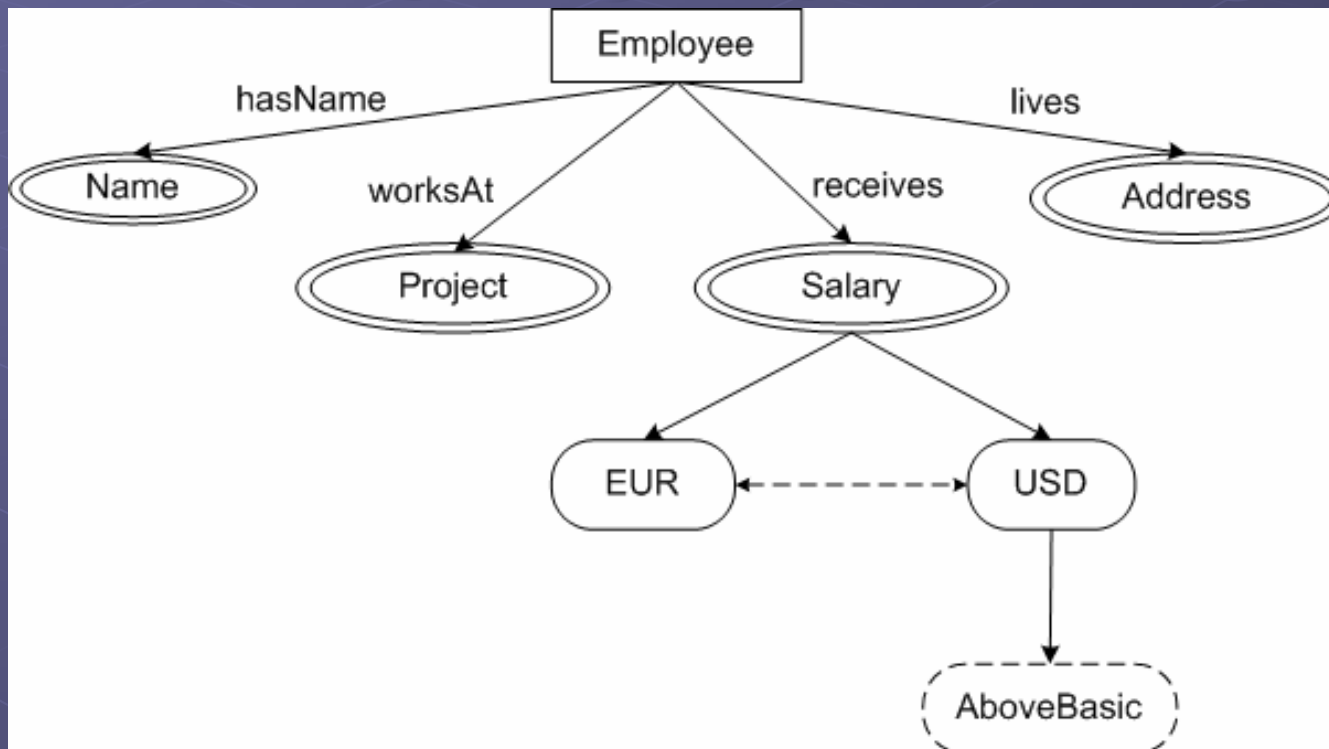
- A suitable application ontology is constructed to model
 - the concepts of the domain
 - the relationships between those concepts
 - the attributes characterizing each concept
 - the different representation formats or (ranges of) values for each attribute
- A graph representation specified for the ontology
 - graph nodes → classes
 - graph edges → properties
 - visual notation
 - different symbols used for each type of class or property

Ontology graph notation

	Type	Represents	Symbol
Nodes	concept-node	a class $c \in C_C$	
	aggregation-node	a class $c \in C_G$	
	type-node	a class $c \in C_{TP}$	
	format-node	a class $c \in C_{TF}$	
	range-node	a class $c \in C_{TR}$	
	aggregated-node	a class $c \in C_{TG}$	
Edges	property-edge	a property $p \in P_P$	
	convertsTo-edge	property convertsTo	
	aggregates-edge	property aggregates	
	groups-edge	property groups	
	subclass-edge	class hierarchy	
	disjoint-edge	disjointness of classes	

Reference example (cont'd)

- The application ontology graph

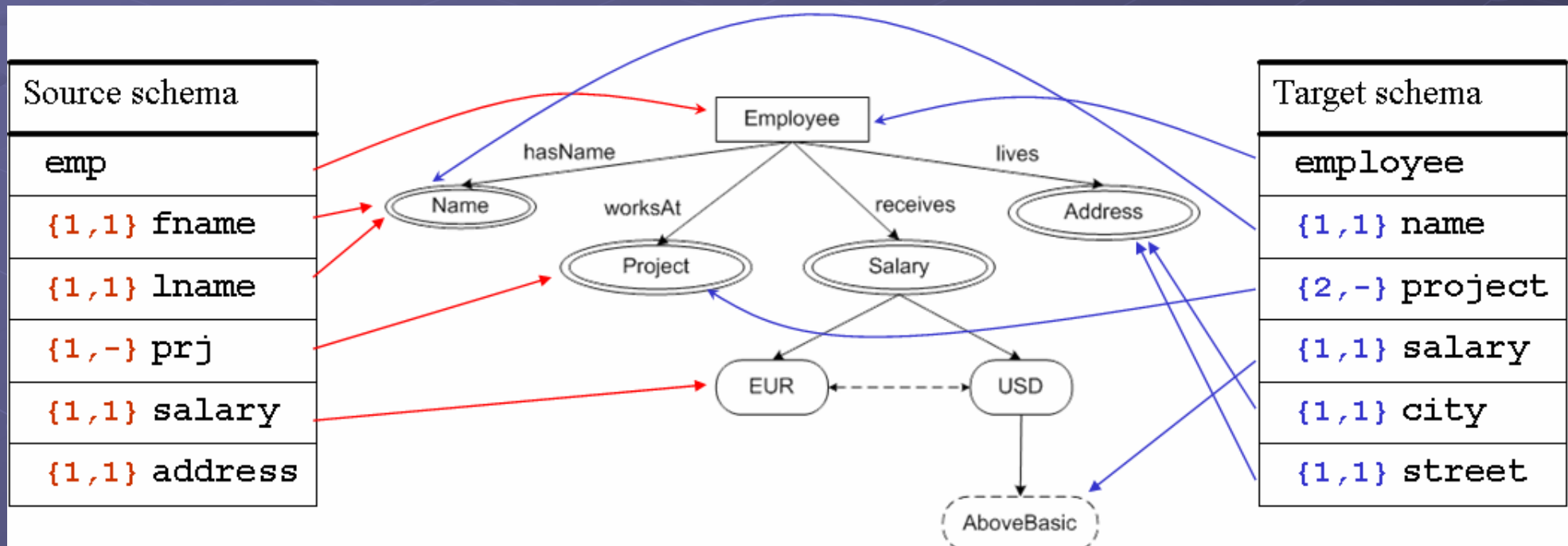


Datastore Annotation

- Semantic annotation → correspondences between the datastore schema and the ontology
- Relations are mapped to concept-nodes
- Attributes are mapped to nodes of the following types:
 - type-node
 - format-node
 - range-node
 - aggregated-node
- Defined classes are created to express the semantics of the schema elements

Reference example (cont'd)

- Datastore mappings



Reference example (cont'd)

- Datastore definitions

- $S.Emp \equiv Employee \sqcap \forall hasName.Name \sqcap =1hasName \sqcap \forall worksAt.Project \sqcap \geq 1worksAt \sqcap \forall receives.EUR \sqcap =1receives \sqcap \forall lives.Address \sqcap =1lives$
- $T.Employee \equiv Employee \sqcap \forall hasName.Name \sqcap =1hasName \sqcap \forall worksAt.Project \sqcap \geq 2worksAt \sqcap \forall receives.AboveBasic \sqcap =1receives \sqcap \forall lives.Address \sqcap =1lives$

ETL Transformations

- Generic types of ETL transformations

Operation	Description
RETRIEVE (n)	Retrieves recordsets from the underlying provider node n
EXTRACT (c)	Extracts, from incoming recordsets, the part denoted by c
MERGE	Merges recordsets from two or more provider nodes
FILTER (c)	Filters incoming recordsets, allowing only records with values of the template type specified by c
CONVERT (c_1, c_2)	Converts incoming recordsets from the template type denoted by c_1 to the template type denoted by c_2
AGGREGATE (f_g, g_1, \dots, g_n)	Aggregates incoming recordsets over the attributes g_1, \dots, g_n applying the aggregate function denoted by f_g
MINCARD (p, \min)	Filters out incoming recordsets having cardinality less than \min on property p
MAXCARD (p, \max)	Filters out incoming recordsets having cardinality more than \max on property p
UNION	Unites recordsets from two or more sources
JOIN	Joins incoming recordsets
STORE	Stores incoming recordsets to a target datastore

Generating ETL transformations

- Two main steps
 - select relevant sources to populate a target element
 - identify required data transformations
- Based on the use of a reasoner to infer subsumption relationships between the defined classes representing the datastores

Generating ETL transformations

● Selecting relevant sources

- a source node n_s , mapped to class c_s
- a target node n_T , mapped to class c_T
- n_s is provider for n_T , if
 - c_s and c_T have a common superclass
 - ensures that the integrated data records refer to the same kind of entity in the domain
 - c_s and c_T are not disjoint
 - prevents integration between datastores with conflicting constraints

Generating ETL transformations

- Identifying data transformations (I)
 - a RETRIEVE operation for each provider node
 - a MERGE operation to combine data from several provider nodes
 - an EXTRACT operation to extract a portion of data from a provider node

Generating ETL transformations

● Identifying data transformations (II)

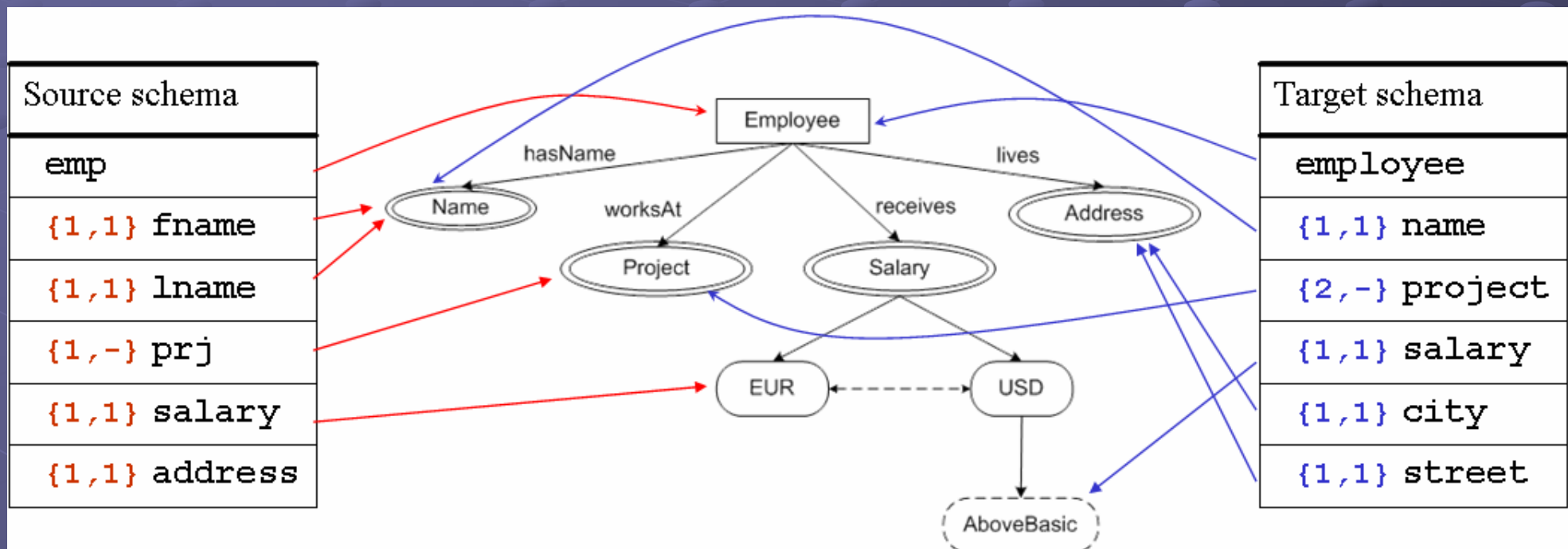
- if $C_S \equiv C_T$ or $C_S \sqsubset C_T$, no transformations are required
- if $C_T \sqsubset C_S$, AGGREGATE, FILTER and/or MINCARD/MAXCARD operations are required
 - choice of specific operation(s) is based on parsing the class definitions and analyzing the value and cardinality restrictions found
- else, as previous plus CONVERT operations

Generating ETL transformations

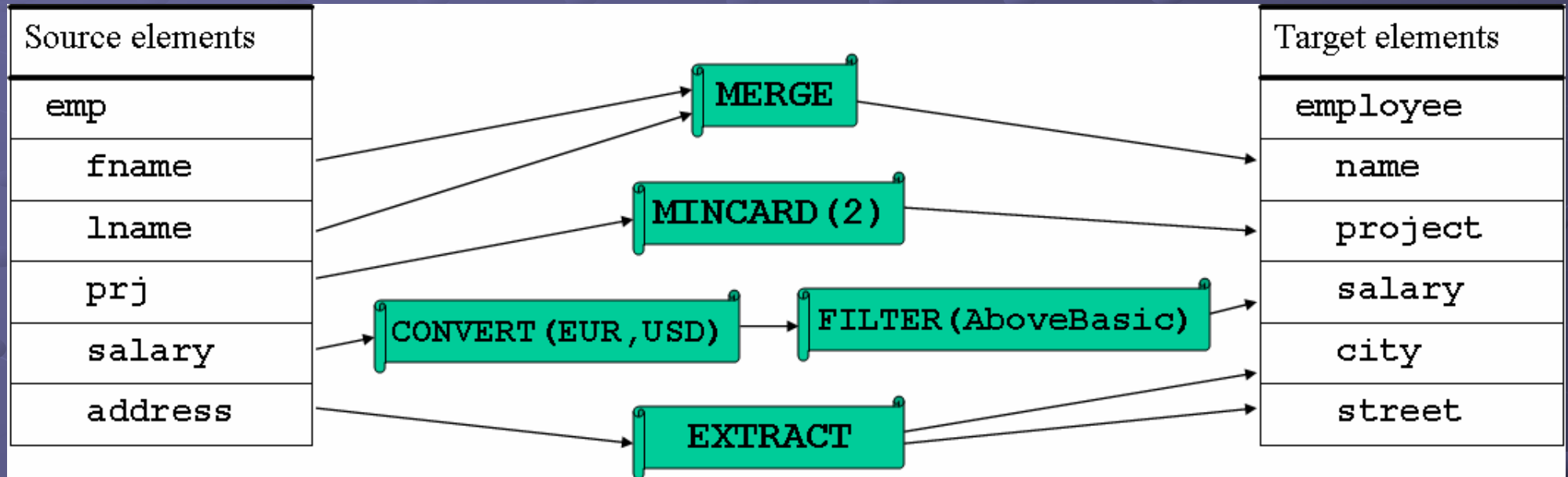
- Identifying data transformations (III)
 - a JOIN operation to combine recordsets from nodes, whose corresponding classes are related by a property
 - a UNION operation to combine recordsets from nodes, whose corresponding classes have a common superclass
 - a STORE operation to denote loading of data to the target datastore

Reference example (cont'd)

- Datastore mappings



Reference example (cont'd)



Current and Future Work

● PDMS

- route queries in the network
- apply the approach to social network applications

● ETL

- NLG techniques for textual representation of datastore descriptions and ETL operations
- impact of schema evolution

● Semantic Web services

- service ranking
- QoS aspects

Thank You

