

A Context-based Name Resolution Approach for Semantic Schema Integration



Rosalie Barreto Belian

Supervisor: Profa. Ana Carolina Salgado

July, 2008

*Universidade Federal de Pernambuco
Centro de Informática
Pós-Graduação em Ciências da Computação*

Objective

- ◆ Specification and implementation of a name resolution approach for schema integration using contextual information
- ◆ Schema element names
 - Natural language words connected with their meanings
 - Meanings depending on the context to which they are related
 - Generic, domain-specific and contextual knowledge are necessary
 - various approaches use solely generic and domain -specific information
 - Are formed by Natural language words from any syntactical category
 - Are made-up by Multi-word terms
- ◆ *Context*
 - circumstantial elements that make a certain situation unique and comprehensible
 - a *context* contains metadata relating to its meaning, properties
 - its source, quality and precision, and organization

Objective

- ◆ Improve name resolution in schema integration using contextual information
 - Allowing context-bound interpretations
 - Providing richer semantic information
- ◆ To achieve this, we proposed
 - a general schema integration process
 - To validate the name resolution approach
 - an internal format to annotate schemas with lexical information
 - We are working at Integra
 - Integra uses X-Entity to represent schemas
 - the relevant contextual information for name resolution in schema integration
 - Represented as a context ontology including contextual and domain information

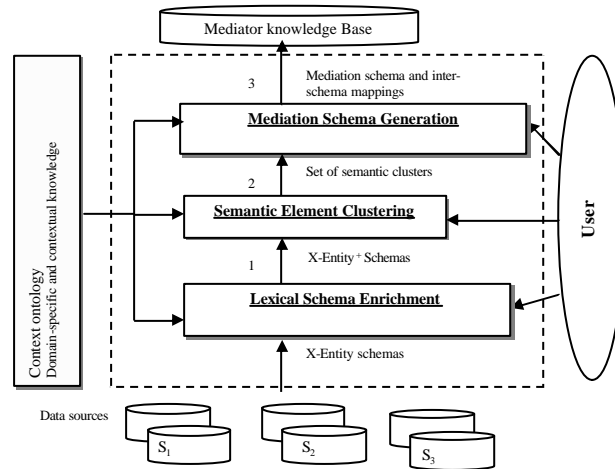
3

Name Resolution in Schema Integration

- ◆ Determine which real-world object a given schema element (entity, attribute or relationship) refers to
 - Spell-check and expansion of abbreviations and acronyms
 - Schema element sense disambiguation
 - Mediation schema element naming

4

The Schema Integration Process



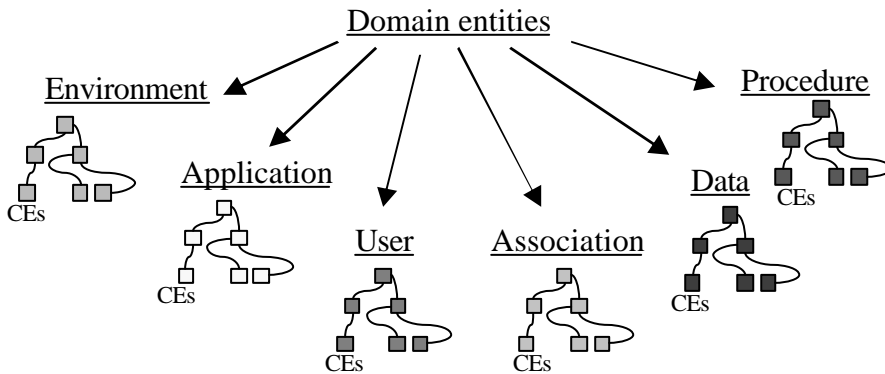
5

The Context Ontology

- ◆ Semantic information
 - Generic, domain-specific (UMLS)
 - Contextual information
- ◆ Contextual information necessary to schema integration/name resolution
 - Domain entities
 - Contextual elements (CEs)
- ◆ Context ontology
 - Reasoning, sharing of information
 - User-rules do declare specific situations that affect the meaning of terms such as defined in the domain
- ◆ Domain entities to DI (data integration)

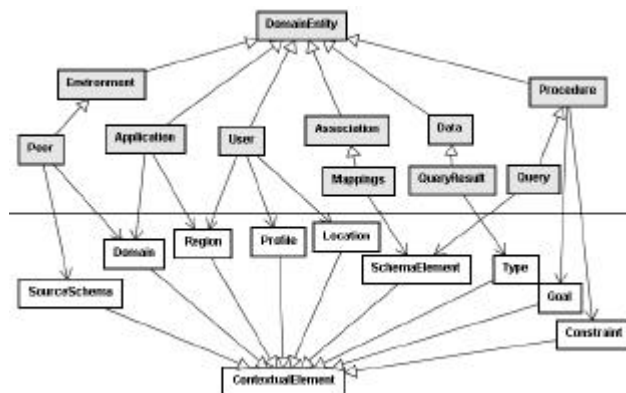
6

The Context Ontology



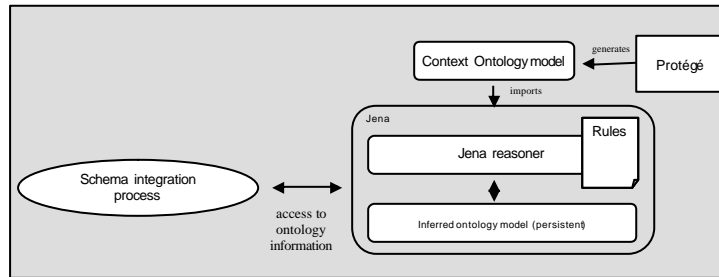
7

The Context Ontology



8

Prototype Implementation



- ◆ Context ontology
 - OWL DL using Protégé 3.2
 - Jena to reason on user-defined rules
 - W3C/SPARQL query language to access ontology information
 - UMLS 2007AC files
- ◆ Schema integration process
 - Java™ programming language version 1.5

9

Conclusion

- ◆ Context-based name resolution approach
 - Contextual information providing enriched semantic information for name resolution
 - Data source semantics (based on local semantics)
 - Mediation semantics (based on user and applications characteristics)
 - Ontology-based approach including domain -specific and contextual information
 - Augmenting the domain information with the contextual information required to name resolution
 - Model of relevant contextual information to name-resolution
 - May be used in other information integration issues

10

Future Work

- ◆ Inclusion of instance-level information in the Schema integration process
- ◆ Evaluation of the Name resolution approach (with and without context)
- ◆ Refinement of the Clustering process
- ◆ Evaluation of the generated Mediation schema
- ◆ Specification of a process to acquisition of generic and domain-specific knowledge
- ◆ Specification of a semantic similarity measure (Contextual? Feature-based?)

11

A Context-based Name Resolution Approach for Semantic Schema Integration



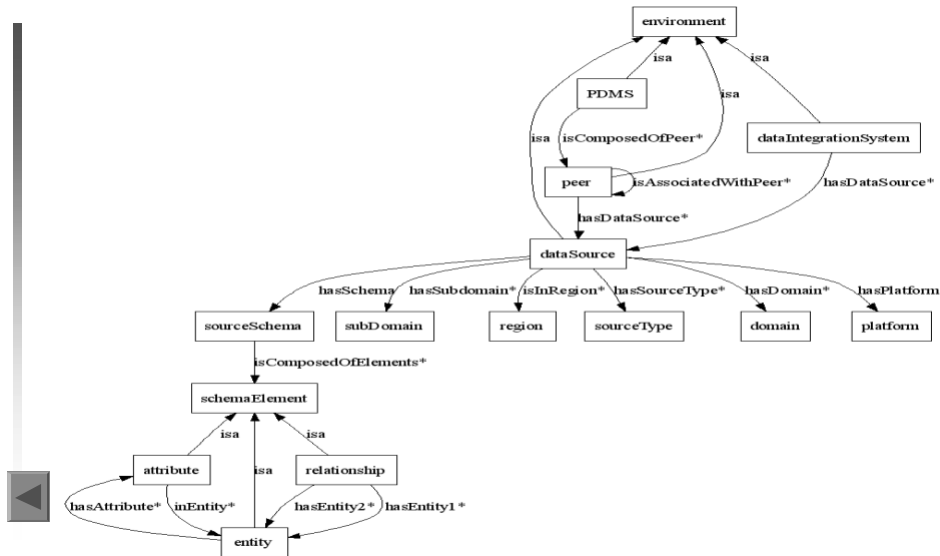
Rosalie Barreto Belian

Supervisor: Profa. Ana Carolina Salgado

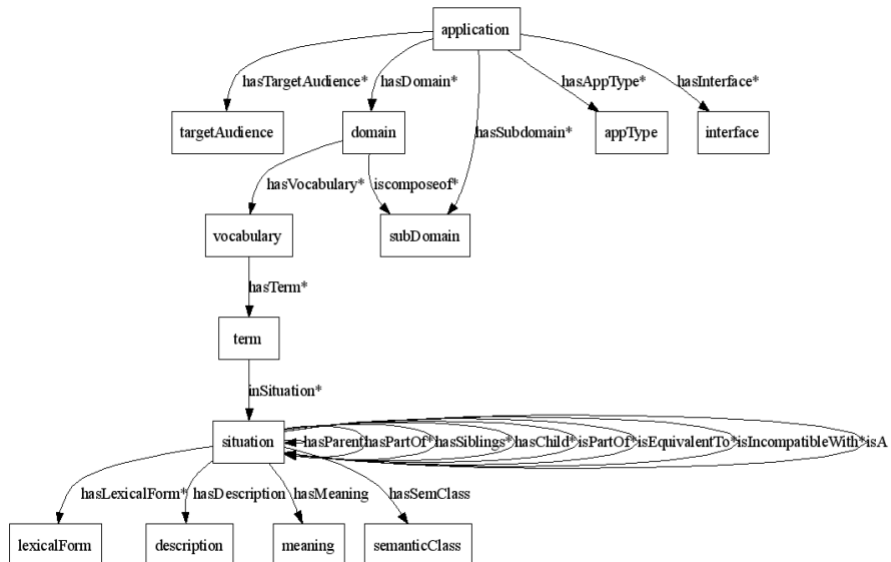
July, 2008

*Universidade Federal de Pernambuco
Centro de Informática
Pós-Graduação em Ciências da Computação*

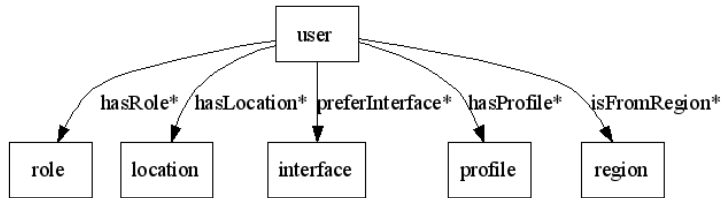
Environment CEs



Application CEs

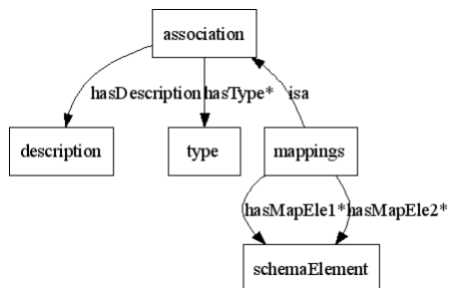


User CEs



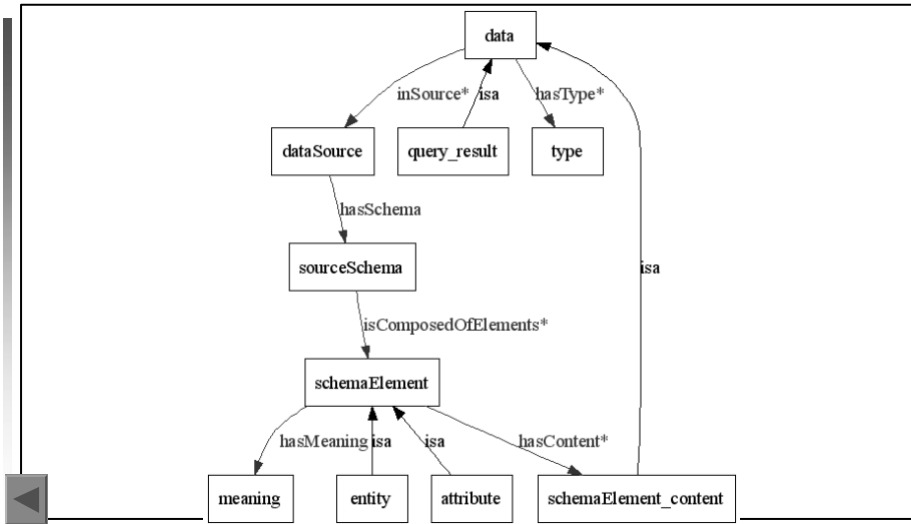
15

Association CEs

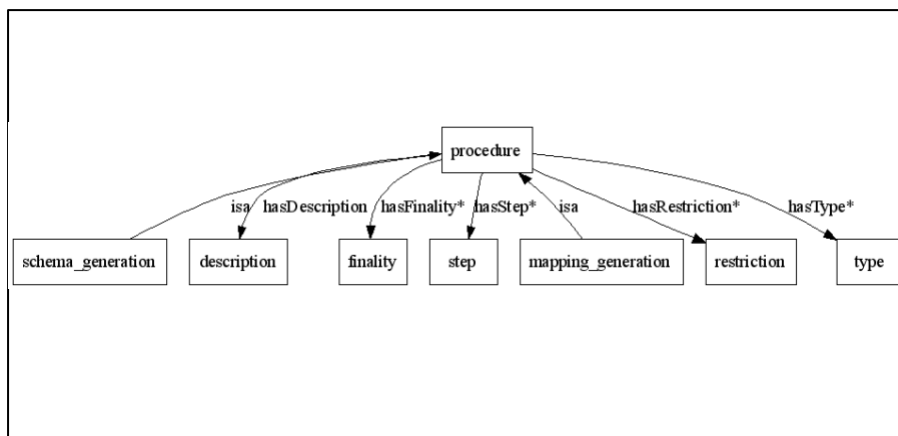


16

Data CEs



Procedure CEs



Lexical Schema Enrichment

- ◆ Step1 – preprocessing of element names from data source schemas
 - tokenization, unnecessary character removal, stop-words removal, case changes handling
- ◆ Step2 – searching schema element names in the ontology
 - Expansion of abbreviations and acronyms
 - Verification of the correct spelling
 - Element name sense disambiguation
 - Annotation of lexical information in X-Entity
- ◆ It is used syntactic matching methods and also, semantic information provided by the ontology to expand tokens and help to address syntactic issues

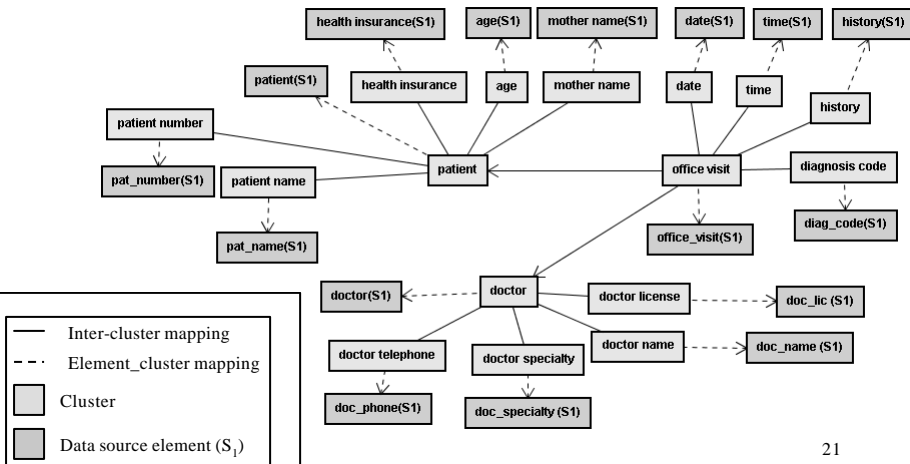
19

Semantic Element Clustering

- ◆ Linguistic-based approach
 - Groups similar schema elements using their names and intended meanings
 - Semantic information is retrieved from the ontology
- ◆ Entities and attributes are grouped into the most similar related semantic cluster
- ◆ Relationships are represented as inter-cluster relationships
 - Hypernymy/hyponymy (is-a), Meronymy/holonymy (part-of), Association (refers-to) and Attribution (has).

20

Set of Semantic Clusters (Snapshot₁ – S₁)



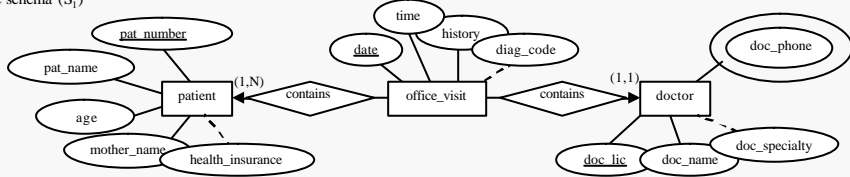
21

Mediation schema generation

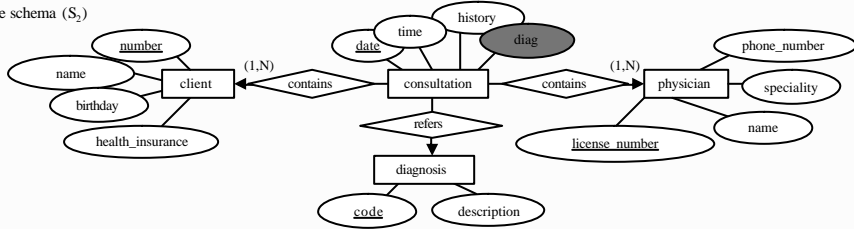
- ◆ Analyzes the set of clusters in order to
 - Generate the final mediation schema
 - Generate the inter-schema mappings (mediation vs. data source schema elements)
- ◆ Main issues
 - Definition of resulting relationships
 - Definition of inheritance relationships between mediation entities
 - Selection of final names for the mediation schema elements
 - Resolution of some structural differences
 - Mediation entity vs. data source attribute
 - Mediation attribute vs. data source entity

22

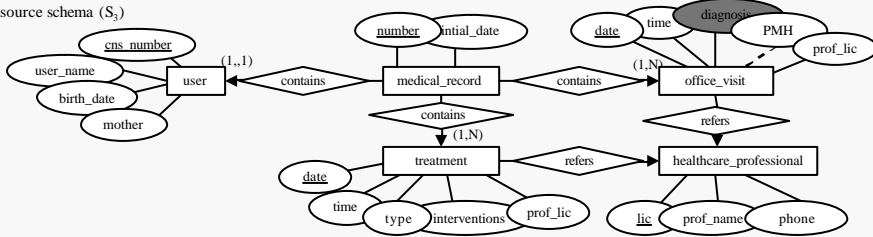
Data source schema (S_1)



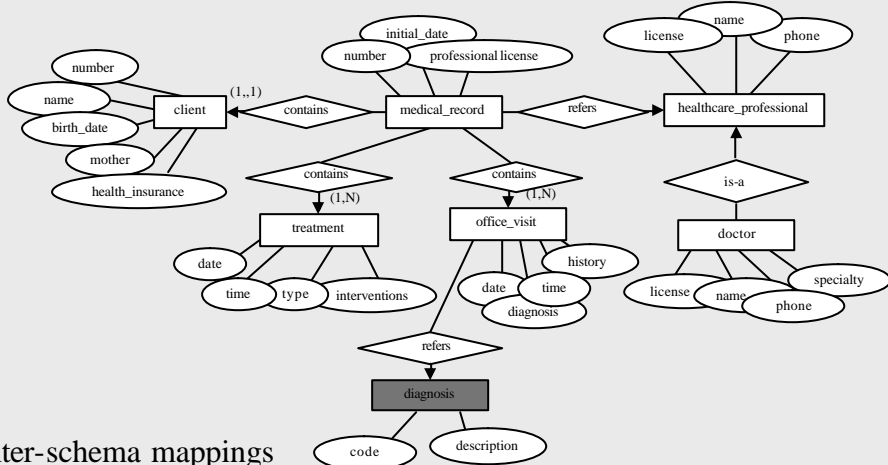
Data source schema (S_2)



Data source schema (S_3)



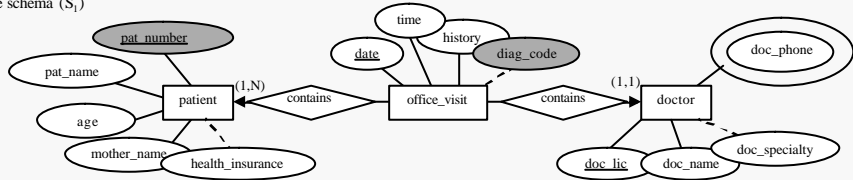
Mediation schema



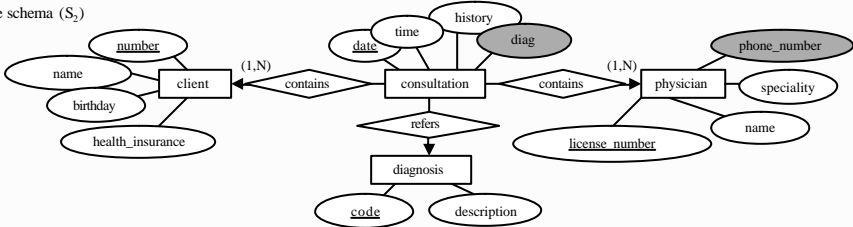
Inter-schema mappings

$\text{medical_record}_M, \text{medical_record}_M, \text{office_visit}_M, \text{office_visit}_M, \text{office_visit}_M, \text{diagnosis}_M, \text{diagnosis}_M, \text{code}_M @$
$\text{office_visit}_1, \text{diag_code}_1$
$\text{medical_record}_M, \text{medical_record}_M, \text{office_visit}_M, \text{office_visit}_M, \text{office_visit}_M, \text{office_visit}_M, \text{diagnosis}_M, \text{diagnosis}_M, \text{code}_M @$
$\text{medical_record}_1, \text{medical_record}_1, \text{office_visit}_1, \text{office_visit}_1, \text{diagnosis}_1$

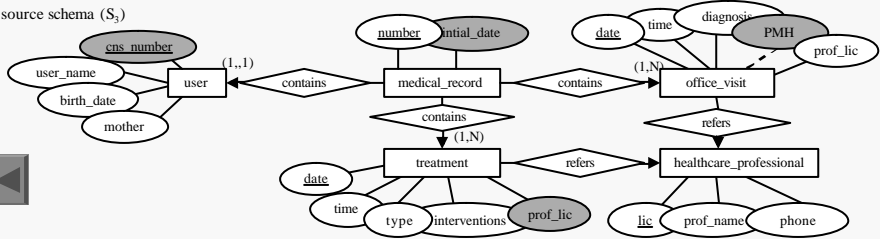
Data source schema (S_1)



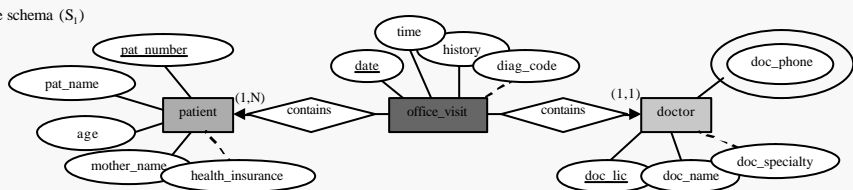
Data source schema (S_2)



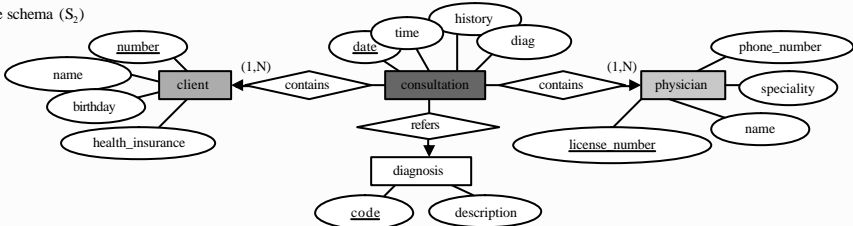
Data source schema (S_3)



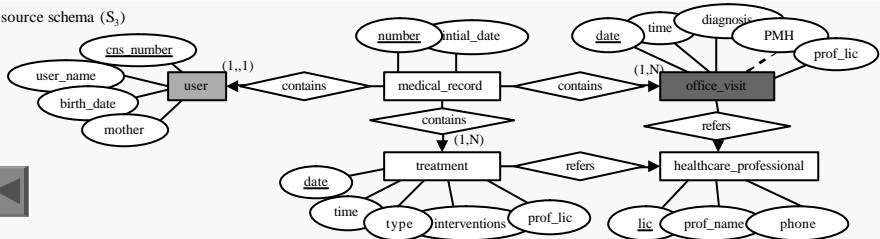
Data source schema (S_1)



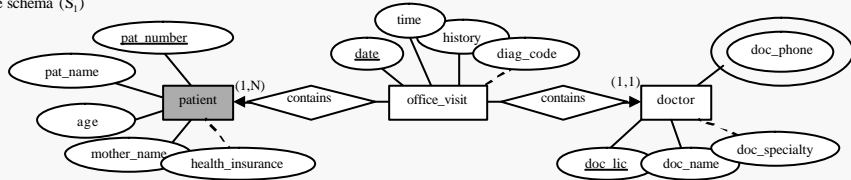
Data source schema (S_2)



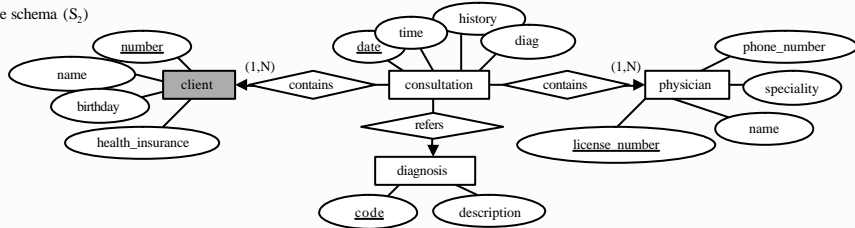
Data source schema (S_3)



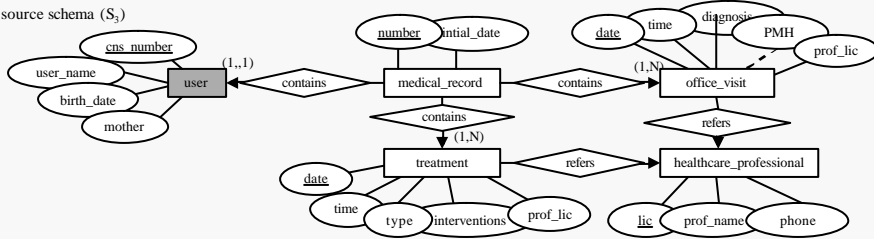
Data source schema (S_1)



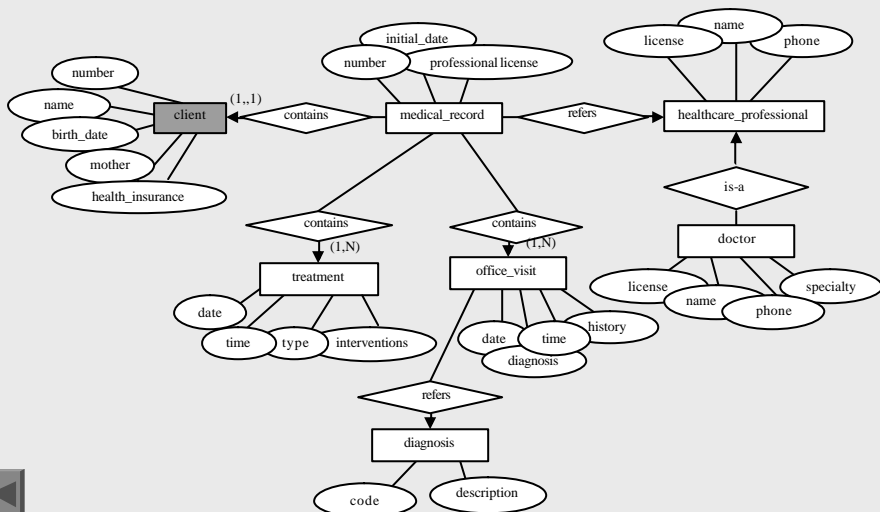
Data source schema (S_2)



Data source schema (S_3)



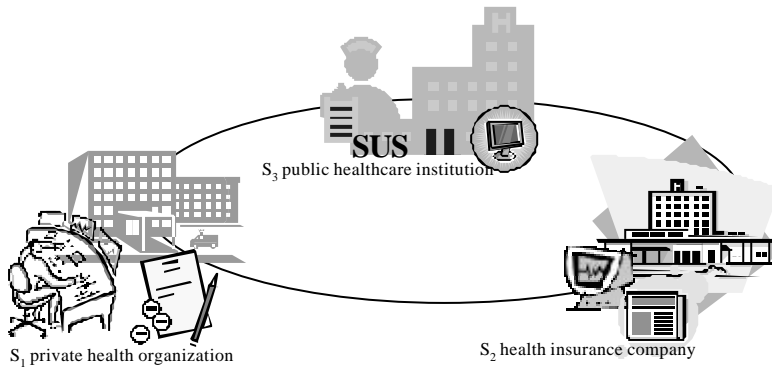
Mediation schema



The Motivating Scenario

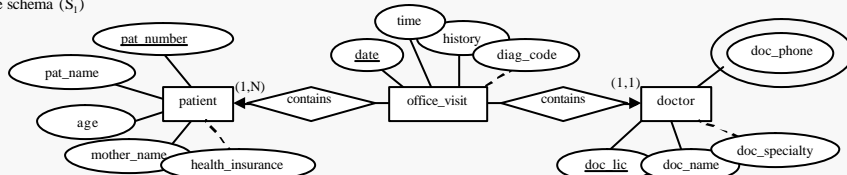
Data sources

- S_1 : medical data from a private health organization
- S_2 : health insurance company
- S_3 : public health institution

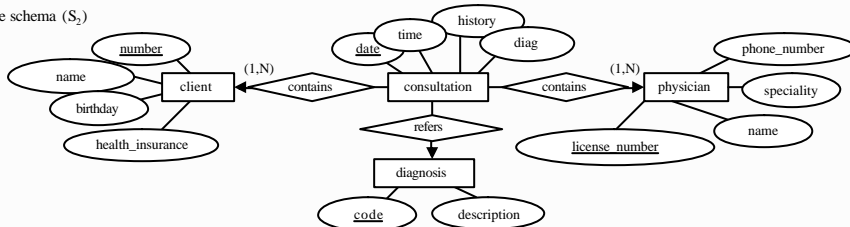


29

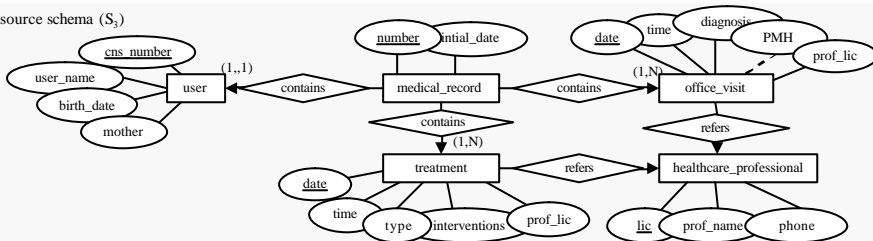
Data source schema (S_1)



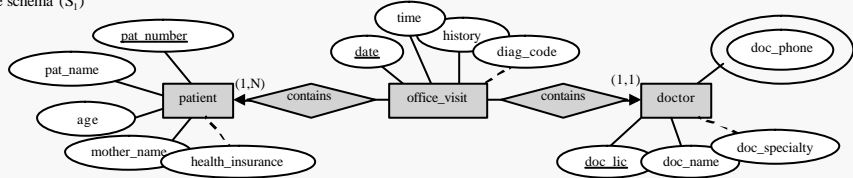
Data source schema (S_2)



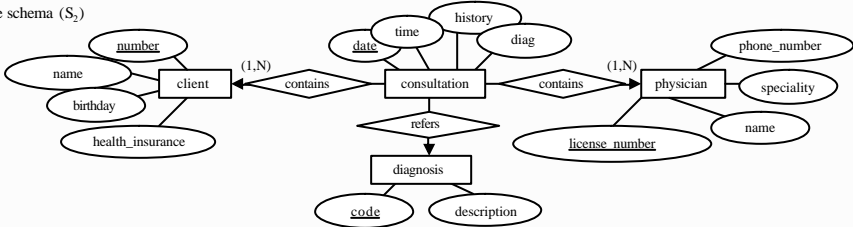
Data source schema (S_3)



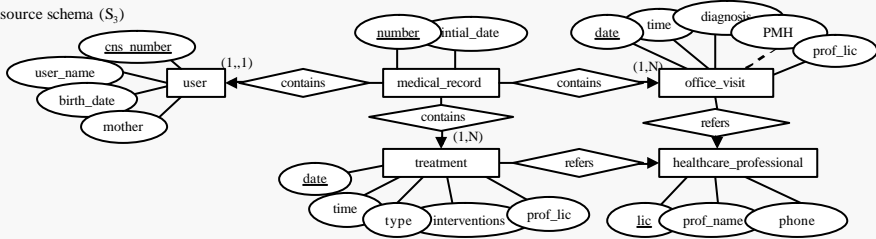
Data source schema (S_1)



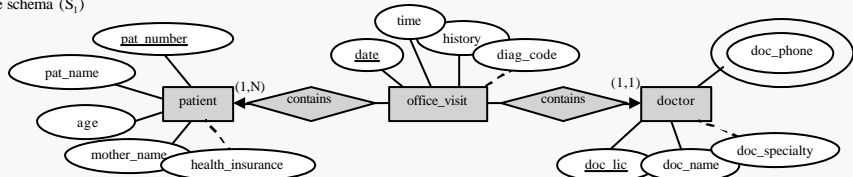
Data source schema (S_2)



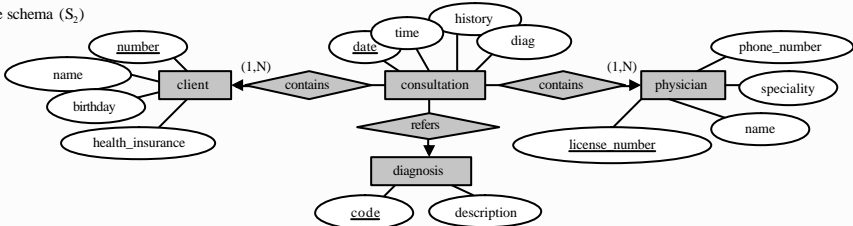
Data source schema (S_3)



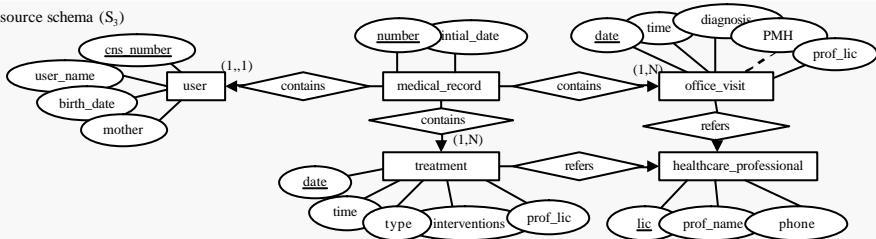
Data source schema (S_1)



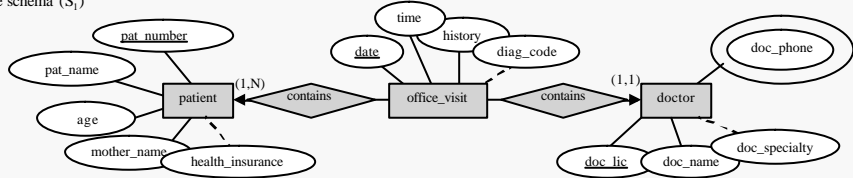
Data source schema (S_2)



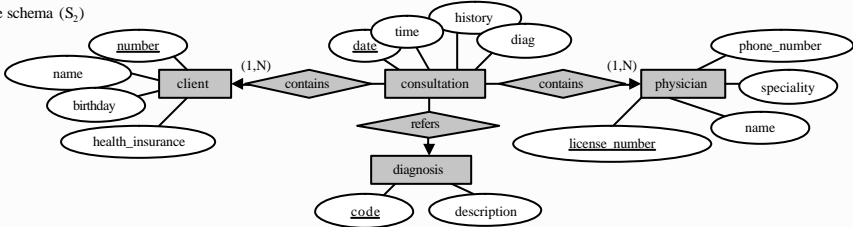
Data source schema (S_3)



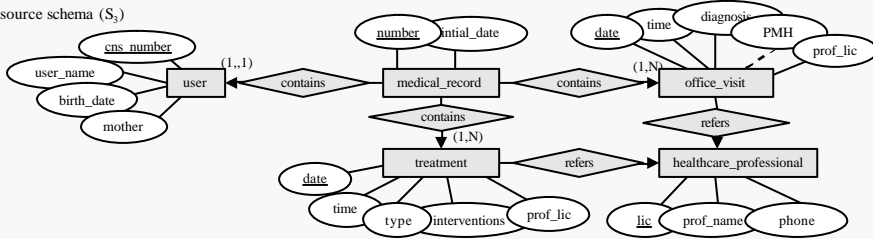
Data source schema (S_1)



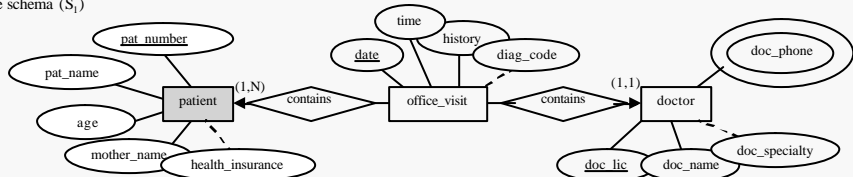
Data source schema (S_2)



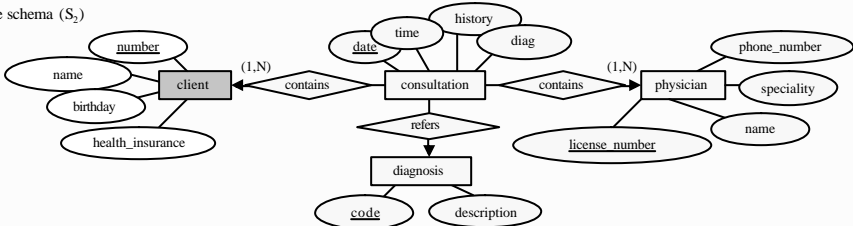
Data source schema (S_3)



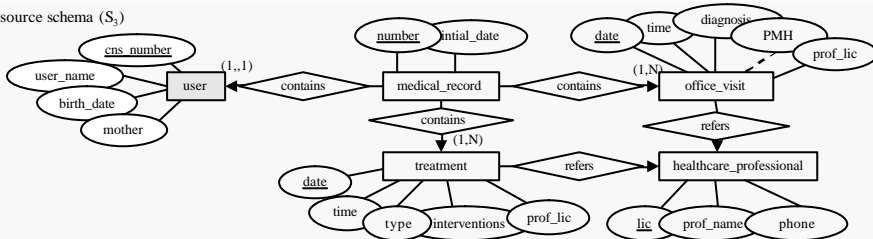
Data source schema (S_1)



Data source schema (S_2)



Data source schema (S_3)



Resulting mediation schema (S_1, S_2, S_3)

