

Database design revisited

Nikitas N. Karanikolas^a, Michael Gr. Vassilakopoulos^b

^a *Department of Informatics, Technological Educational Institute (TEI) of Athens, Athens, Greece*

^b *Department of Computer Science and Biomedical Informatics, University of Central Greece, Lamia, Greece*

Abstract: Information Systems design is affected by the simplicity of the relational model types, though these types do not correspond directly to entities of real applications. In this paper, we present another approach, where the Information System designers would be able to directly represent the real world in a database model, very close to the Entity-Relationship (ER) model. A query and manipulation language that can handle composite data types, close to ER diagrams, is the Conceptual Universal Database Language (CUDL). We demonstrate that a database modelled by ER diagrams can be directly expressed to the CUDL Abstraction Level (CAL), the data model corresponding to CUDL, by a set of rules for direct transformation of ER diagrams to CAL. This set consists of a basic set of five rules, which is extended with an extra rule that deals with specific situations appearing in practical applications. Consequently, the resulting more powerful and composite data can be directly maintained with the CUDL. In this way, the development process consisting in conceptual design (ER), transformation to Logical Relational Schema, usage of SQL for data manipulation/retrieval and the reverse steps to present the results in the conceptual level is simplified by conceptual design (ER), transformation to CAL and the usage of CUDL for direct manipulation/retrieval of real world (conceptual) structures.

Keywords: Information Systems Design, Database

1. Introduction

An assumption of a common methodology that is used up to today in the design of Information Systems based on relational databases is that we do not know the structure of information that we were called to model and manipulate. For this reason we begin with interviews of the persons involved in the operation of a non-computerized system and from this process we arrive in a series of fundamental information elements (attributes), grouped in (usually) one and universal relation. Then, the specification of a set of functional dependencies among these attributes and the decomposition of the set of attributes (through normalization) into smaller relations which consist of subsets of the original set of attributes, in order to eliminate update anomalies and reduce data redundancy, follow. During data manipulation, the combination of attributes through the relational join operator is needed.

The normalization is the process that aims to produce the best from a (by nature) weak data model: —the relational model is limited with respect to semantic content (i.e., expressive power) and there are many design problems which are not

naturally expressible in terms of relationships [20]; —The relational model is weak when showing many-to-one relationships [17].

However, the real world that we are called to model with an Information System (often with the use of a database) seldom incorporates repetitions of data (data redundancy). As an example, consider a non-computerized managed library, where we do not incorporate multiple series of books (copies) and a corresponding number of bookshelves, in order to place the first set of book copies sorted according to the title, the second set of book copies sorted according to the first writer, the third set of book copies sorted according to the theme category, etc. On the contrary, we do not use any ordering (more precisely we use the ordering of books according to the date of import into the library) or use some of the orderings that interest us (usually thematic) and in addition we create indexes (with cards) for every one of the ordering that interest us. Each card contains the key of the classification and a reference to the natural ordering (the location of book in the bookshelves). Thus, we claim that the Information System design should not decompose the real world in its fundamental characteristics, but the Information System designers should be able to portray directly the real world in a model that provides more powerful structures (through composite data types), as those of the real world.

To support composite data types, meta-models (metadata schemata and the corresponding metadata) are needed. Handling of such (composite data) types by the programmer by SQL-like languages and extensions requires knowledge of the internal organisation of both metadata and data, unless a data manipulation language able to manipulate directly the composite data types is employed. The Conceptual Universal Database Language (CUDL) is such a language. We demonstrate that a database modeled by Entity-Relationship (ER) diagrams can be directly expressed to the CUDL Abstraction Level (CAL,) the data model corresponding to CUDL, by a set of rules for direct transformation of ER diagrams to CAL. This set consists of a basic set of five rules (first presented in [12]), that is extended with an extra rule (first presented in [13]) that deals with specific situations appearing in practical applications. Consequently, the resulting more powerful and composite data can be directly maintained with the CUDL language. In this way, the development process consisting in conceptual design (ER), transformation to Logical Relational Schema, usage of SQL for data manipulation/retrieval and the reverse steps to present the results in the conceptual level is simplified by conceptual design (ER), transformation to CAL and the usage of CUDL for direct manipulation/retrieval of real world (conceptual) structures. In this paper, we present an enhanced revisiting of the whole database design process, based on the rules presented in [12, 13], using examples that demonstrate the benefits for the database designer of using real- world data types.

The organization of the paper is as follows. Section 2 provides a review of data models that were proposed to overcome weaknesses of the relational model. However, having more powerful database models, but still using data manipulation languages designed only for fundamental data attributes, is a waste of time and resources. Therefore, in Section 3, we present the basic ideas related to CUDL and its abstraction level, CAL [9, 10, 11, 22]. The transformation from CAL to FDB (meta-model, [21, 23]) has been studied in [9, 11]. In Section 4, we present the 5 Rule Set for Transformation from ER to CAL [12] (dealing with the most common cases that appear in real life applications, for the transformation from ER diagrams to CAL), while in Section 5, we present the addition of an extra rule that is needed for the transformation of a relation between a weak and a strong (not identifying)

entity that also has attributes. Section 6 concludes the paper and presents future research possibilities.

2. Related Work

2.1. Generic Data Modeling

The Generic Data Modeling [4,8] approach is an outcome mainly emanate from research in the Medical Informatics domain. The fact that, in case of Health Care data maintenance, the amount of information and the complexity of information lead to a huge (Daedal / mazy) conceptual schema, concerned the Medical Informatics scientists. Moreover, the fact that the direct production of a logical schema for a relational DBMS, from a given huge conceptual schema, obviously conserve this Daedal characteristic, gave raise for research for alternative data modeling approaches. Another inherent characteristic of relational logical schemata is the difficulty for supporting data evolution (changing information needs). The research results for both problems (Daedal conceptual and logical schema and difficulty for data evolution in relational data) reveal the generic data modeling approach. This approach defines two generic transformations (namely —flattening! and —relation merging!). The later transformation (—relation merging!) is also the basis for supporting data evolution. These transformations, when applied to the original conceptual schema, produce a generic logical schema consisted of a reduced number of tables. However, this process is not a very strict procedure and actually it is depended from (the personal perception and) the choices made by the person who guides the process and applies the mentioned generic transformations. The final number of tables, as the outgrowth of the transformation process, is dependent from the choices made and is not a concrete (predefined) set of tables (as happens in other cases, e.g. in the FDB data modeling).

The disadvantage of Generic Data Modeling is that querying the resulting generic logical schema with standard SQL requires multiple statements and considerable intellectual effort, especially when the queries are intended to retrieve data for feeding data analysis tasks (e.g. feeding data mining applications). To overcome such difficulties, researchers have defined the Extended Multi-Feature (EMF) SQL extension [7] that provides simpler to understand, more compact and more efficient query constructions.

2.2. EAV Data Modeling

The Entity Attribute Value (EAV) data modeling [15,16] is also an outcome from research in the Medical Informatics domain. The motivation for the research that revealed the EAV data modeling was that, in the medical domain, the number of parameters (facts) that potentially apply to any clinical study is vastly more than the parameters that actually apply to an individual clinical study. For example, the potential number of laboratory examinations that a patient could be submitted to is a huge superset of the actually submitted examinations in a specific medical case (e.g. a patient suffering from a bilestone). Another reason, that motivated the research that revealed the EAV data modeling, was that clinical studies are subject to evolution as a result of medical research. As a consequence, the number of clinical parameters related to a clinical study is always differentiated (and, in most cases, are increased). Thus, the data model should be able to host new clinical parameters for any clinical study, without the need for data (structure) reorganization. The research, motivated by the above-mentioned reasons, revealed

the EAV data modeling. According to EAV design, metadata and data tables compose the logical database schema. The facts (that actually apply to a clinical study) are recorded into the data tables, as a triplet: the Entity, the Attribute and the Value. The Attribute is the recorded fact (clinical parameter) and the Entity is a composition of the relevant patient's identifier and some timestamp. The metadata tables are used to define the data composition (which clinical parameter, i.e. Attribute, pertains to which clinical study).

There are three main versions of the EAV data modeling [1] but all of them share the same basic principle (the triplet: Entity / Object, Attribute, Value). Another interesting feature of the EAV models is that they permit mixture of EAV stored and conventionally stored data. However, the existence of this heterogeneity complicates significantly the task of data querying. We should also mention that the EAV data modeling support facts evolution (equivalent to the addition columns in a relational table without the need for any reorganization) but does not support table (entity) evolution.

2.3. FDB Data Modeling

In previous works (by Yannakoudakis et al.) [21,23] there has been an investigation of dynamically evolving database environments and corresponding schemata, allowing storage and manipulation of a variable number of fields per record, variable length of fields, composite fields (fields having subfields), multiple values per field (either atomic or composite), etc. The ultimate goal of the research work of Yannakoudakis et al. was to make the design and maintenance of a database a simpler task for database designers, so as that they will not have to put in a lot of effort to design the database and later they will not have to pay special attention and work for database changes. Their research proposed a new framework for the definition of a universal logical database schema that eliminates completely the need for reorganisation at both logical and internal levels, even when major modifications of the database requirements have occurred. This new framework was called Frame Database Model (FDB) [23].

This Universal logical database schema is based on well and strictly defined set of Metadata and data tables and it does not permit any mixture with conventionally stored data. All the entities that are available to the user, along with their attributes, are documented exclusively in the metadata tables and the facts concerning the instances of the entities are recorded exclusively in the data tables of the FDB Universal schema. Another noteworthy feature of FDB is that it supports Schema evolution, both for facts and entities.

Moreover the FDB model allocates ways of imprinting strictly connected (Hardly related) information with innate (inherent/native) mechanisms, in contrast to the relational model that compels the creation of artifact structures (tables) to represent strictly connected (Hardly related) data. As an example, the relational model requires the creation of new table to store data that relate of the form one-to-many (the addresses or the telephones of customer). In contradiction to the relational model, the FDB model can maintain the same information with a field that is accommodated in the side of the —one| and accepts multiple values. Even more complex forms of strictly connected (Hardly related) information, are impressed, in the FDB model, with innate (inherent/native) mechanisms. For example a correlation of information with a form many-to-many (as are the DVDs that have been rented to a member of a Video Club) is maintained in one of the two connected sides without requiring the creation of a new table to correlate the information. That is to say, in the many-to-many correlations we follow a mechanism that emanates

from the real world (in the example of the Video Club we maintain inside the card of a customer a table with his/her renting of DVDs).

The most important fact in the FDB model is that it organizes information without any repetition of values. In order to be more precise, not only it does not proceed in repetitions but it ensures that these cannot be created. In the example of the addresses (or alternatively the telephones) of a customer the basic data of a customer (let us say name, surname and code) are stored once and all the different addresses that the customer may have are stored in the field addresses. That is to say, the use of a single big (universal) table to store/repeat as many times the basic attributes of a customer (name, surname and code) as the number of his/hers addresses is avoided naturally (without any effort). Thus the FDB model provides as an inherent feature the no redundancy property.

2.4. Not First Normal Form (NF2) or Nested Relational Data Modeling

The motivation for developing the Nested Relational data model was the fact that the Relational model has difficulties of modeling the real world; It is also inconvenient for handling even simple data structures commonly used in Information Retrieval. To overcome these problems, Researchers have proposed a relational model where Non First Normal Form (NF2) relations are allowed [18]. This extension encompasses the classical 1st Normal Form (1NF) model and adds, to the relational algebra, two basic operators (namely “nest” and “unnest”). Based on the “nest” operator, this proposal allows sets (as the result of “one-attribute” nest operation) and sets of sets (as the result of “multi” attribute nest operation) as attribute values. NF2 sets are equivalent to simple FDB fields with repetitions and NF2 sets of sets are equivalent to composite FDB fields with repetitions. The researchers have also proposed a query language extension for NF2 table definition and manipulation. However, the NF2 presents some weak points:

- It does not support Entity or Fact (Attribute) evolution
- It does not have Universal logical schema
- The proposed query language extension only undertakes (be engaged in) Retrieval statements
- This Retrieval statements are rather suggestions or hypothetical statements and are not parts of a mature language that handles relational tables with non-atomic (sets and sets of sets as) attribute values
- The notion that governs the whole idea, which has passed through and is reflected by the proposed query language extension, is that the subfields (of composite fields) are not directly accessed by the user.
- Related to the previous point is that the proposed query language uses Nested Select statements, whenever a restriction over a subfield should be applied

Possibly, these weak points have the consequence that, after 29 years, Non First Normal Form does not seem to be implemented as a DBMS. However, some researchers are still interested [19] and define languages supporting the eXtended NF2 (XNF2) model.

2.5. Object Oriented and Object Relational Databases

The weakness of the relational model to manage complex, highly interrelated information motivated the research for Object Databases (ODB) and Object Relational Databases (ORDB). Both models are also described in textbooks (for example Elmasri and Navathe, 2000 [5]).

The portability and interoperability of ODBs is ensured by the Object Model suggested by the Object Database Management Group (ODMG). The ODMG Object Model provides also the definition for an Object Definition Language (ODL) and an Object Query Language (OQL). The ODL statements seem to (or are influenced by) the Java language statements used for class definitions, while the OQL statements seems to (or are influenced by) the SQL statements used for data retrieval. At mid of the 1990's decade there were a notable number of ODB solutions, but today only half of those remain active. Our personal opinion is that the data management professionals do not like to bother themselves with strange programming constructions of classes, inheritance, etc., and consequently they do not decide easily to use an ODB.

The other direction, the Object Relational Databases, aim to provide solutions for complex and highly interrelated information management, without imposing complicated programming constructions. For this reason, they provide Black Box Complex data types for various purposes (management of time series, geographic point manipulations, face recognition, content-based retrieval of digital audio, image watermarking, image search, full-text search), Opaque types for extending the repertoire of Black Box Complex data types and User Defined Complex data types. Black Box Complex data types are named as Data Blades in Informix Universal Server and are named as Cartridges in Oracle. The User Defined Complex data types have similar characteristics to the ODMG Objects. The composition of User Defined Complex data types is based on simpler structures (namely: the Collection types and the Row types). The SQL3 standard provides an extension to the previous SQL standards, for handling the most of the characteristics added with the ORDBs.

2.6. Approaches based on Semantic Web technology

The purpose of Semantic Web is to share and reuse knowledge. The associated concepts, Web standards and query languages allow storing, querying and reasoning. Reasoning is a common-sense query processing and it requires a small number of deductions, in order to answer some query. To support the common-sense query processing, any Semantic Web system must be provided with a set of rules of deduction. However, Semantic Web is not designed with the intention of handling information efficiently and effectively, as Databases are.

There exist approaches that combine the Semantic Web information technology and the Database information technology. One of the goals of these combinations is the interoperability of Semantic Web representations (Ontologies) and Databases [3]. Underneath of this interoperability is the translation of schemas from one information technology to the other [2]. The Relational data model and also the Functional data model have been used as the representatives of the Database technology, in these approaches of information technology combinations [3,14]. The Resource Description Framework (RDF) and its associated Query Language (SPARQL) are the most used representatives of the Semantic Web, in these approaches of information technology combinations.

One of the purposes of this paper is to support a richer, than the Relational, data model, while the application developer, that uses this model, does not face complications in analysis, design and data manipulations. Our belief is that, the freedom of more complex data, that the Semantic Web information technology (through Ontologies) offers, operates against the efficiency, maintainability and simplicity of an application development oriented system.

3. The Conceptual Universal Database Language

It is obvious from the plethora of models (presented in the previous section) that there is a need for a DBMS able to provide more powerful data types, as those of the real world complex data. In the first four discussed models, handling of such types by the programmer by SQL-like languages and extensions requires knowledge of the internal organization of both metadata and data [12]. Thus, the ultimate goal should be to provide a data manipulation language able to manipulate directly the composite data types (without bothering the programmer with internal organization details), while permitting the direct expression of restrictions over subfields.

In the approach that we have adopted, the FDB [21, 23] is the underlying model for implementing our goal for a data manipulation language able to manipulate directly composite data types. We preferred the FDB model [10], since it is more compact and well defined than the other models and also support schema evolution. Karanikolas et al. [9, 11], introduced the syntax and semantics of the CUDL, a database language (both DDL and DML for the FDB model) that supports composite data types, i.e. attributes (tags in CUDL terminology) that combine more than one other sub-attributes (subfields in CUDL terminology). It also supports the existence of multiple values (repetitions) for both tags and subfields. The key difference of CUDL to other approaches is that not only tags, but also subfields and repetitions can be addressed in the CUDL statements. That means that the users can express retrieval restrictions over a specific subfield or can express the modification of a specific subfield in a specific repetition of some tag, etc. In [9], the authors focused mainly in presenting and analyzing the statement of value retrieval (in the schema and the data). More recently [11], they focused mainly in presenting and analyzing the CUDL statement of value modifications in the schema (schema changes) and the data. They also have undertaken the need for relationship declarations [10]. This need becomes more significant for the FDB- CUDL model, because the relationships between entities, in most cases, are implemented without the introduction of new tables. Without having methods to declare relationships, the user would face a refuting stage, where the model is self-explained (the user can consult only tag attributes and subfield attributes and carry off the data model) but the data relationships are totally undocumented. To cope with this need, the FDB model introduced the Authority links set. They also use the Authority links set to declare authority controls and reduce variability of expressions used for the same instance of an identity. All of these (relationship declarations and authority control declarations) are provided through CUDL statements.

Apart from being a data manipulation language able to manipulate directly composite data types, while permitting the direct expression of restrictions over subfields, CUDL allows the application programmer/designer to model the structures of its application with composite data types that are closer to the ER diagrams and sometimes without any decomposition of the ER entities into simpler ones (an example is depicted in [12]).

On the other hand the logical database level of any CUDL based application is the underlying FDB model. Thus, instead of transforming from ER to simple relational tables to provide a logical model for manipulation through SQL, we are able to transform from ER to CUDL Abstraction Level (CAL) entities (namely, CUDL data sets with composite data types). In other words, the classic database design triplet (ER, logical and physical design) is replaced by the quadruple: ER, CAL, logical and physical design, with a fixed logical design (the FDB model).

CAL is a model that supports collections of uniform (congenous) instances (frame objects in CAL terminology). The structure of each instance (of a concrete

collection) is a set of attributes (tags in CAL terminology). In contrast to the relational model, attributes (tags) can have either simple (predefined) data types or can be composed by a set of sub-attributes (subfields in CAL terminology). CAL model also supports the existence of multiple values (repetitions) for both tags and subfields. Each CAL collection of uniform (congenous) instances is called entity (also called —abstract entity| for reasons explained in [11]). The CUDL language is a database language for defining (DDL) and handling (DML) CAL entities.

Given the above characteristics of CUDL and FDB, it is explained why somebody would prefer to use a CUDL-DB versus an ORDB. The main advantages of CUDL-DB are:

- schema evolution,
- not unused fields (zero repetitions make unneeded the use of null values),
- several real world constructions are neither so complex as to demand data blades (and other extended capabilities of ORDBs), nor natively supported by the relational model (e.g. more than one telephone numbers of a customer require two tables with an 1:M relationship); these —intermediate| (neither simple, nor complex) data constructions do not bother the programmer and are natively supported in CUDL and CAL.

It is obvious, from the above explanation, that the transformation from ER entity types to CAL entities is a direct mapping. However, the transformation from ER relationship types to CAL structures is a more complicated process and we are going to consider it in the next section.

4. The 5 Rule Set for Transformation from ER to CAL

4.1. A Transformation Example

In order to make more evident the need for designing in upper levels we will provide an example of a really complicated ER (corresponding to a real situation) and its transformation to CAL entities (CUDL data sets with composite data types). We are going to present an Electronic Patient Record (EPR) that could be maintained in a real Hospital Information System (HIS). A brief description of our HIS follows:

- One patient can have one or more incidents treated in the hospital.
- During an incident the patient can be subject of a series of laboratorial examinations and also can be subject of a series of radiological examinations.
- During an incident the patient can be subject of zero or more operations (surgeries).
- There are a number of doctors (servant physicians) for each incident.
- Moreover, there are a number of doctors (surgeons, anaesthesiologists, etc) that participate in each operation.
- Each incident is characterized by a Social Security Institute (that undertakes the hospital fees – cover the expenses) and a diagnosis (final diagnosis of the incident).

The ER diagram of Figure 1 is detailed conceptual depiction of the HIS – EPR database. According to the ER diagram of Figure 1, there is a binary relationship between the Incident and the “Laboratorial Examination” entity types with cardinality ratio M:N. A similar binary relationship between the Incident and the “Radiological Examination” entity types, with cardinality ratio M:N, is also depicted. A third binary relationship with cardinality ratio M:N exists for the Incident and the Doctors entity types. This relationship is responsible for the servant physicians of the incident. The unique ternary relationship of the ER diagram is an identifying relationship of the weak entity type “Incident Operation”. There are two owner entity

types that identify the weak entity type “Incident Operation”. These are the Incident and the Operation entity types. There is a binary relationship, with cardinality ratio M:N, between the weak entity type “Incident Operation” and the (strong) entity type Doctors. The later relationship expresses the set of doctors (surgeons, anesthesiologists, etc.) that participate in each “Incident Operation”.

Next, we will provide the CAL entities that the ER diagram of Figure 1 is transformed. The CAL entities are (for composite and multivalued attributes, () and {} are used, respectively, Subsection 3.3.1 of [5]):

Incident

Incident_code, Date_started, Date_ended, Patient_code, SSI_code,
 {Lab_examinations (LE_code, LE_datetime, LE_result)},
 {Rad_examinations (RE_code, RE_datetime, RE_FilePath)},
 {Incident_operations (Operation_code, IO_datetime_started, IO_datetime_ended,
 {IO_doctors})},
 {Incident_doctors}, Diag_code.

Doctor

Doctor_code, name, surname.

Operation

Operation_code, name, cost.

Rad_Examination

RE_code, description, type, cost.

Lab_Examination

LE_code, name, normal_values, cost.

Patient

Patient_code, name, surname, father_name, tax_registration_number, date_of_birth.

Social_Security_Institute

SSI_code, name, immediate_insured_rate, intermediate_insured_rate.

Diagnosis

Diag_code, description, type.

Figure 2 portrays a CAL frame object (instance) for the entity “Incident”, while Figure 3 portrays some CAL frame objects (instances) for the entity type “Doctor”.

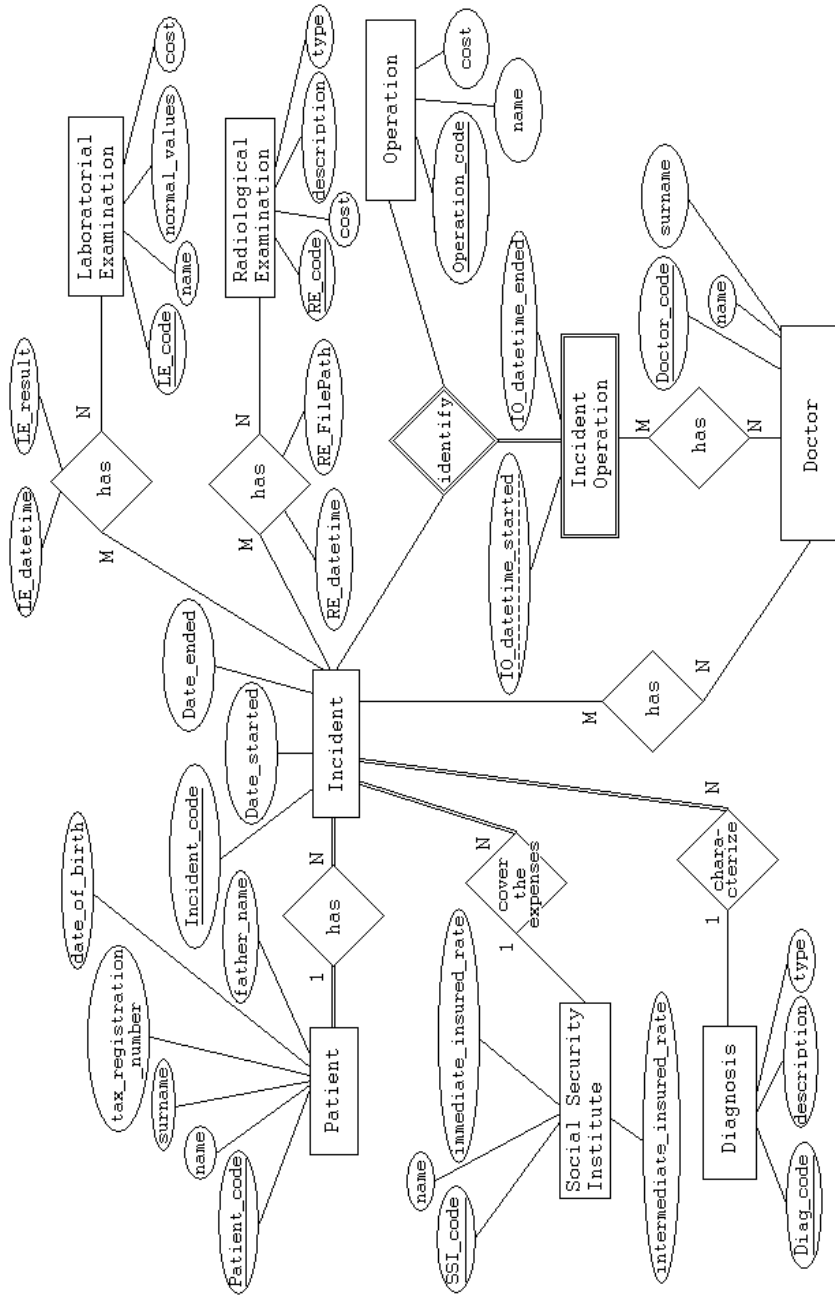


Figure 1. The ER diagram of the HIS-EPR.

Incident_code	S001			
Date_started	13/5/2007			
Date_ended	20/5/2007			
Patient_code	A001			
SSI_code	T001			
Incident_doctors	I001			
	I002			
	I079			
Diag_code	574			
Incident_operations	Operation_code	IO_datetime_started	IO_datetime_ended	IO_doctors
	E002	14/5/2007 13:35	14/5/2007 15:05	I001 I005 I100 I065
	E015	16/5/2007 12:00	16/5/2007 13:00	I012 I100 I032
Lab_Examinations	LE_code	LE_datetime	LE_result	
	UREA	15/5/2007 10:00	32,4 mg/dl	
	UREA	15/5/2007 14:30	32,5 mg/dl	
	UREA	16/5/2007 08:00	31,6 mg/dl	
	CREA	15/5/2007 10:00	1,17 mg/dl	
	CREA	16/5/2007 08:00	1,08 mg/dl	
	PROT	15/5/2007 10:00	7,19 g/dl	
	PROT	16/5/2007 08:00	6,95 g/dl	
Rad_Examinations	RE_code	RE_datetime	RE_FilePath	
	U/S Kidney	16/5/2007 12:00	\\FS1\RIS\Uaz34.tif	

Figure 2. An Incident CAL frame object.

Doctor_code	I001	Doctor_code	I032
name	Michael	name	...
surname	Garidopoulos	surname	...
	(a)		(d)
Doctor_code	I005	Doctor_code	I065
name	Paul	name	...
surname	Alexiou	surname	...
	(b)		(e)
Doctor_code	I012	Doctor_code	I100
name	...	name	...
surname	...	surname	...
	(c)		(f)

Figure 3. Some Doctor CAL frame objects.

4.2. The Basic Transformation Rule Set

The rules for transforming from (relationships types of) ER diagrams to CAL entities (at least for the cases that appear mostly in real applications, like those needed for transforming from the ER of Figure 1 to the CAL entity Incident) that were presented in [12] are the following:

Rule1. Binary relationships with cardinality ratio M:N can be hosted in one of the two related entities as a tag (field) with repetitions. In case that the relationship has no attributes, a simple (not composite) tag with repetitions is capable to store

the primary key values of the related (hosted) entity. Whenever the relationship has attributes, a composite tag (composed of the primary key values of the hosted entity and the relationship's attributes) with repetitions should be used. Obviously, the key of the hosting participant is not needed.

This rule applies for the relationship between the Incident and the "Laboratorial Examination" entity types and for the relationship between the Incident and the Doctor entity types. Since we have selected to host the relationships in the Incident entity type, the former relation is implemented with the following composite tag with repetitions:

{Lab_examinations (LE_code, LE_datetime, LE_result)}.

For the same reason, the later relation is implemented with the following simple tag with repetitions:

{Incident_doctors}.

Rule 2. Binary identifying relationships can be transformed to a composite tag with repetitions, hosted in the owner entity type. In this case, the composite tag should be composed by the partial key (of the weak entity) and the rest attributes of the relationship. Obviously, the key of the hosting owner entity is not needed. There is no application of this rule in the studied ER diagram.

Rule 3. The previous rule can be extended for ternary identifying relationships. Ternary identifying relationships can be transformed to a composite tag with repetitions, hosted in one of the (two) owner entity types. In this case, the composite tag should be composed by the partial key (of the weak entity type), the key of the hosted owner entity type and the rest attributes of the relationship. Obviously, the key of the hosting owner entity type is not needed.

This rule applies for the weak entity type □Incident Operation□. Since we have selected to host the relationships in the Incident owner entity type, the relation is implemented with the Incident operations composite tag with repetitions:

{Incident operations (Operation code, IO_datetime_started, IO_datetime_ended, {IO_doctors})}

The partial key of the weak entity type participates in the Incident_operations as the IO_datetime_started subfield. The key of the owner entity type Operation participates in the Incident operations as the Operation code subfield. The subfield IO_datetime_ended is an attribute of the relationship. The role of the IO_doctors subfield will be explained with the next rule.

Rule 4. Binary relationships with cardinality ratio M:N, relating a weak entity type with some other (strong) entity type, without having relationship attributes, can be transformed to an extra subfield with repetitions of the composite tag implementing the weak entity. This rule applies for the relationship between the weak entity type □Incident Operation□ and the entity type Doctor.

This rule explains the last constituent of the Incident operations tag, presented above. (The domain of {IO_doctors} is the power set of the Doctor code's domain.)

Rule 5. Binary relationships with cardinality ratio 1:N can be hosted in the N-side entity type, as a tag (field) without repetitions. In case that the relationship has no attributes, a simple (not composite) tag is enough for storing the primary key values of the opposite-side entity type. Otherwise, whenever the relationship has attributes, a composite tag should be used. The primary key of the opposite-side entity type and the relationship's attributes composes this composite tag. Obviously, the key of the hosting (N-side) participant is not needed.

This rule applies for the relationships of the Incident with the entity types Patient, "Social Security Institute" and Diagnosis, respectively. The tags

Patient_code, SSI_code and Diag_code of the CAL entity Incident implement these relationships.

Actually, the rules presented above are also responsible for updating the foreign-key constraints.

4.3. The Equivalent Relational Model

The relational model that corresponds to the ER diagram of Figure 1 is depicted in Figure 4. The entities of this ER diagram are directly mapped into relational model relations / tables, but the inter-entity relationships of the ER diagram are represented in the relational schema through foreign keys and extra relations / tables (in the case of M:N relationships). This representation is not obvious for the designer of the ER diagram and expressing queries of the real-world problem in the language that manipulates the relational model (e.g. SQL) is unnatural and not simple.

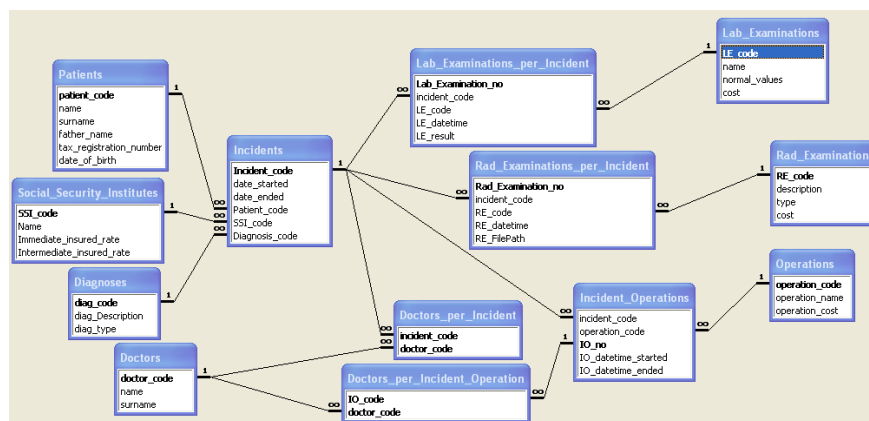


Figure 4. The relational model corresponding to the ER diagram of the HIS-EPR.

5. The Extra Transformation Rule

In the following we will provide an extension of the above set of rules. We are going to provide a rule that applies in a specific subcase of relation between entity types, because this subcase uncovers the barriers of the underlying Frame DataBase (FDB) model [23]. The FDB model permits the use of composite data types, but not in any (arbitrary) depth, as the Nested Relational model [6, 18] permits. This was a key decision, since the FDB model should be accompanied with a data manipulation language, CUDL [9, 11], able to directly manipulation of the provided composite data structures.

5.1. The Modified Conceptual Model

In the following conceptual (ER) diagram (Figure 5), we introduce two attributes that characterize the participation of a physician (doctor) to the weak entity “Incident Operation”.

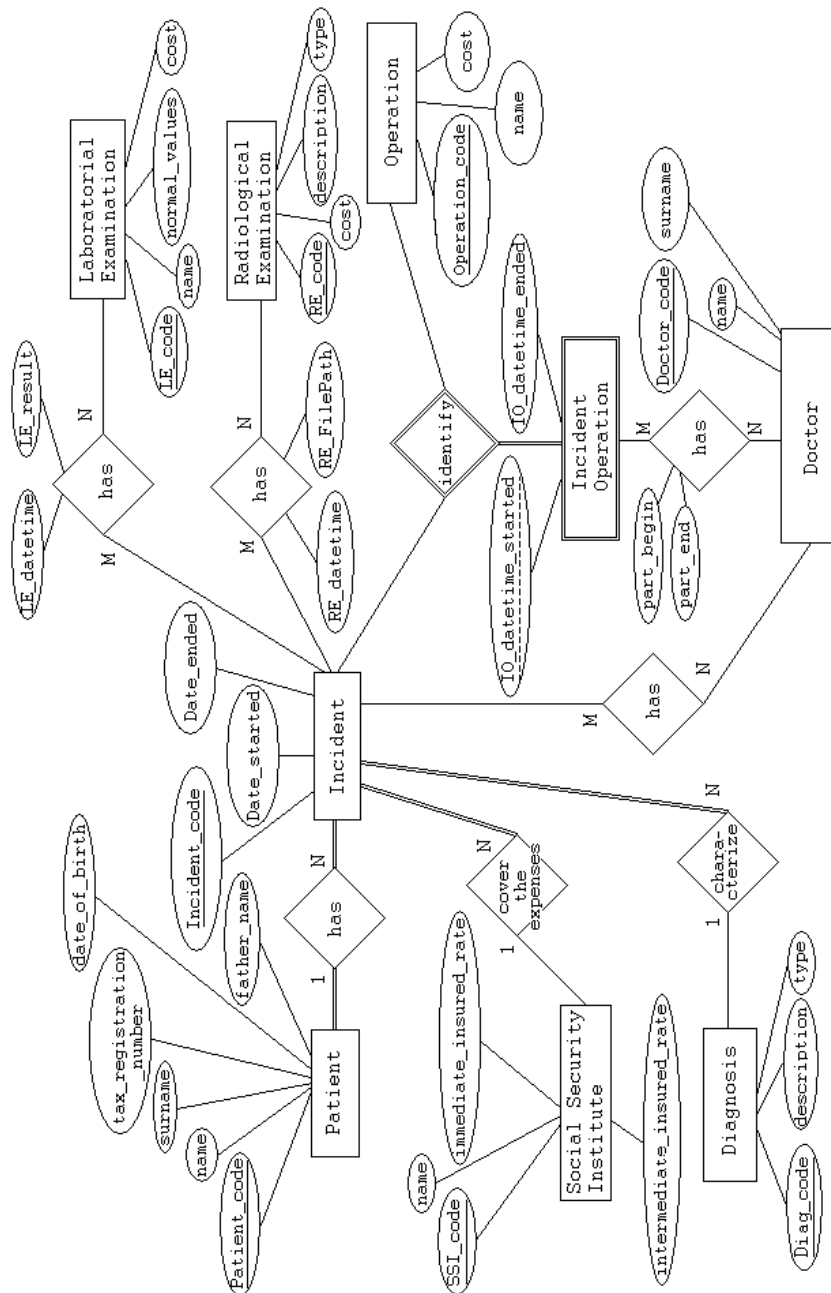


Figure 5. The modified ER diagram of the HIS-EPR.

The new attributes, introduced in figure 5, are:
 part_begin representing the date and time that the participation of the doctor to the incident operation started,
 part_end representing the date and time that the participation of the doctor to the incident operation finished.

These attributes define the duration of participation of a doctor in some specific operation that is accomplished during some incident. In other words, the user requirements demand not only to document the participation of doctors in incident operations, but also to document when the doctor gets involved and when he/she backs off the operation.

5.2. Handling attributes relating a weak with a strong entity type

The following rule handles the translation of a relation between a weak and a strong (not identifying) entity that has attributes:

Rule 6. Binary relationships with cardinality ratio M:N, relating a weak entity type with some other (strong) entity type, having relationship attributes (assume k attributes), can be transformed to an extra subfield of the composite tag implementing the weak entity. This extra subfield will not have repetitions and should have unique values in the range of the whole set of frames of the strong identifying entity. Moreover another composite tag with k+1 subfields and permitting repetitions should be introduced in the related (strong) entity type. The relationship (k) attributes, together with the unique value of the extra subfield of the composite tag that implements the weak entity will combined to form the (k+1) subfield values of some repetition of the new composite tag of the related (strong) entity.

This rule applies for the relationship between the weak entity type “Incident Operations” and the entity type Doctor.

5.3. The modified Incident and Doctor entity type

In contrast to Rule 4, Rule 6 requires the creation of identifiers for each instance of the weak entity of the relation and the storage of these identifiers into the frames (instances) of the strong (not identifying) related entity. This could be an alternative solution in case of a relation (between a weak entity type and a strong – not identifying – entity type) without attributes. This is the only solution when the relation (between the weak entity type and the strong – not identifying – entity type) has attributes.

The following are the CAL structures of the strong – identifying – entity type “Incident” that hosts the weak entity type “Incident Operation” and the strong – not identifying – entity type “Doctor”. The modifications (in relation to the previous versions of the same entity types) are emphasized.

Incident

Incident_code, Date_started, Date_ended, Patient_code, SSI_code,
 {Lab_examinations (LE_code, LE_datetime, LE_result)},
 {Rad_examinations (RE_code, RE_datetime, RE_FilePath)},
 {Incident_operations (Operation_code, IO_datetime_started,
 IO_datetime_ended, **IO_ID**)},
 {Incident_doctors}, Diag_code.

Doctor

Doctor_code, name, surname,
{Operations_participated (IO_ID, part_begin, part_end)}.

Figure 6 portrays a CAL frame object (instance) for the entity “Incident”, while Figure 7 portrays some CAL frame objects (instances) for the entity type “Doctor”.

5.4. Enforcement of data validation

The declaration of relationships between data structures is one of the most significant tools for the enforcement of data validation procedures in any database model. The relational model (through its language SQL) has introduced methods (constraint statements or sub-statements) that permit the declaration of relationships.

In CUDL the declaration of relationships between data structures is maintained in an FDB table (table of the FDB universal logical database schema) named Authority_Links [10]. The structure of this table is the following:

Authority_links (from_entity, from_tag, from_subfield, to_entity, to_tag, to_subfield, relationship_type)

For example, the M:N relationship between the entity type “Incident” and the entity type “Doctor” is declared in the following way:

(‘Incident’, ‘Incident_doctors’, null, ‘Doctor’, ‘Doctor_code’, null, 1)

The last attribute in the previous row of the Authority_links table is the integer value of one (1). It defines a simple M:N relationship where the data implementing the relationship are kept in the M side. The whole range of values of the relationship_type attribute is provided in table 1.

Incident_code	S001			
Date_started	13/5/2007			
Date_ended	20/5/2007			
Patient_code	A001			
SSI_code	T001			
Incident_doctors	I001			
	I002			
	I079			
Diag_code	574			
Incident_operations	Operation_code	IO_datetime_started	IO_Datetime_ended	IO_ID
	E002	14/5/2007 13:35	14/5/2007 15:05	317
	E015	16/5/2007 12:00	16/5/2007 13:00	318
Lab_Examinations	LE_code	LE_datetime	LE_result	
	UREA	15/5/2007 10:00	32,4 mg/dl	
	UREA	15/5/2007 14:30	32,5 mg/dl	
	UREA	16/5/2007 08:00	31,6 mg/dl	
	CREA	15/5/2007 10:00	1,17 mg/dl	
	CREA	16/5/2007 08:00	1,08 mg/dl	
	PROT	15/5/2007 10:00	7,19 g/dl	
Rad_Examinations	RE_code	RE_datetime	RE_FilePath	
	U/S Kidney	16/5/2007 12:00	\\FS1\RIS\ Uaz34.tif	

Figure 6. An Incident CAL frame object.

Doctor_code	I001		
name	Michel		
surname	Garidopoulos		
Operations_participated	IO_ID	part_begin	part_end
	317	14/5/2007 13:35	14/5/2007 15:05

(a)

Doctor_code	I005		
name	Paul		
surname	Alexiou		
Operations_participated	IO_ID	part_begin	part_end
	317	14/5/2007 13:35	14/5/2007 14:25

(b)

Doctor_code	I012		
name	...		
surname	...		
Operations_participated	IO_ID	part_begin	part_end
	318	16/5/2007 12:00	16/5/2007 13:00

(c)

Doctor_code	I032		
name	...		
surname	...		
Operations_participated	IO_ID	part_begin	part_end
	318	16/5/2007 12:00	16/5/2007 13:00

(d)

Doctor_code	I065		
name	...		
surname	...		
Operations_participated	IO_ID	part_begin	part_end
	317	14/5/2007 13:35	14/5/2007 15:05

(e)

Doctor_code	I100		
name	...		
surname	...		
Operations_participated	IO_ID	part_begin	part_end
	317	14/5/2007 14:25	14/5/2007 15:05
	318	16/5/2007 12:00	16/5/2007 13:00

(f)

Figure 7. Some Doctor CAL frame objects.

Table 1. The range of values for the attribute relationship_type.

type	Explanation	Symbol
1	Simple M:N, data implementing the relationship are kept in the M side	M:N/L
2	Simple M:N, data implementing the relationship are kept in the N side	M:N/R
3	S(M):N, data implementing the relationship are kept in the S(M) side	S(M):N/L
4	S(M):N, data implementing the relationship are kept in the N side	S(M):N/R
5	N:S(M), data implementing the relationship are kept in the N side	N:S(M)/L
6	N:S(M), data implementing the relationship are kept in the S(M) side	N:S(M)/R
7	S(M):S(N), data implementing the relationship are kept in the S(M) side	S(M):S(N)/L
8	S(M):S(N), data implementing the relationship are kept in the S(N) side	S(M):S(N)/R
11	Simple M:N, data implementing the relationship are kept in subfield in the M side	M:N/L _↓
12	Simple M:N, data implementing the relationship are kept in subfield in the N side	M:N/R _↓
14	S(M):N, data implementing the relationship are kept in subfield in the N side	S(M):N/R _↓
15	N:S(M), data implementing the relationship are kept in subfield in the N side	N:S(M)/L _↓
20	Authority control over Tag	A/T
21	Authority control over Subfield	A/S

According to table 1 and based on the early solution (ERD of figure 1) where there are no attributes for the relationship between the weak entity type —Incident Operation| (the operations that a patient undertook during the period of an incident) and the entity type —Doctor|, the relationship is declared as:

(‘Incident’, ‘Incident_operations’, ‘IO_doctors’, ‘Doctor’, ‘Doctor_code’, null, 3)

In case that the relationship between the weak entity type “Incident Operation” and the entity type “Doctor” has attributes (ERD of figure 5) and according to the modified CAL entity types, the relationship is declared as:

(‘Incident’, ‘Incident_operations’, ‘IO_ID’, ‘Doctor’, ‘Operations_participated’, ‘IO_ID’, 14)

5.5. The modified Relational Model

The relational model that corresponds to the modified ER diagram (Figure 5) is depicted in Figure 8. This relational schema is derived from the relational model of Figure 4 by a simple modification (addition of only two fields in an existing relation / table, the table Doctors_per_Incident_Operation). However, the deficiency of the non-easy expression of real-world problem queries in the manipulation language of the relational model still exists.

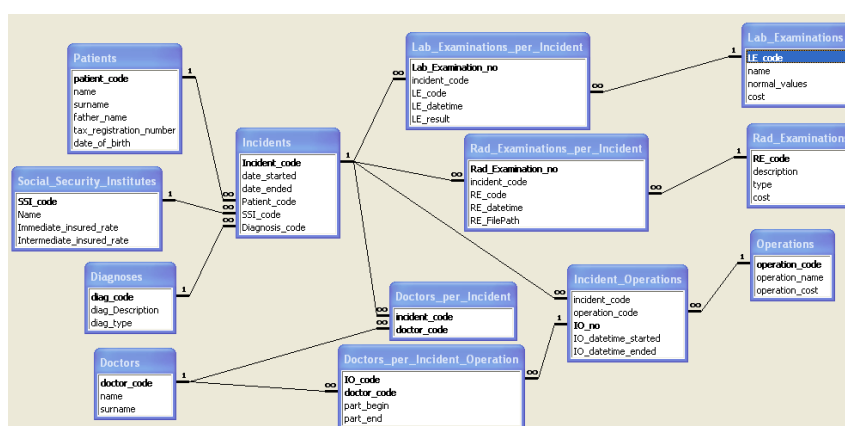


Figure 8. The relational model corresponding to the modified ER diagram of the HIS-EPR.

6. Conclusions

So far, the design of an application having a relational data repository mainly required the decomposition of the real world structures into very simple attributes, the composition of a logical schema with naive relational structures and the formation of query and manipulation (SQL) statements based on the logical schema. The need for a database query and manipulation language, like CUDL, able to handle directly composite entities of more abstract data models that offer composite / complex data types expressing real world entities (without transforming them to a relational logical schema) is apparent. The types used in CUDL allow a database designer / developer to express the structures of his application with types that are very close to the ER entity types, while ER entity types can be directly transformed to CAL entity types. Thus, the designers and developers can be concentrated with the business logic of their applications, instead of wasting time for the expression of statements that manipulate naive database structures.

The basic set of rules for transforming (the relationship types of) an ER diagram to CAL, so that the CUDL can be utilized to manipulate the resulting high level entities was supplemented with an extra rule about the translation of a relation between a weak and a strong (not identifying) entity that also has attributes. This rule applies in a specific subcase of relation between entity types that uncovers the barriers of the underlying Frame DataBase (FDB) model [23] (the FDB model permits the use of composite data types, but not in any depth).

Possible future research directions include:

- extending FDB and CUDL in order to explicitly document the foreign keys and other constraints,
- completing the analysis and design process that consists in the quadruple “ER, CAL, FDB, physical level” by providing a physical model and a transformation of the FDB to this model,
- implementing a database machine storing data based on this physical model and being able to process CUDL queries,
- designing algorithms for optimized processing of CUDL queries,
- evaluating this machine in terms of performance and suitability for development of complex applications.

References

- [1] Anhøj, J. (2003). Generic Design of Web-Based Clinical Databases. *Journal Medical Internet Research*, 4.
- [2] Atzeni, P., Cappellari, P., Torlone, R., Bernstein, P.A., Gianforme, G. (2008). Model-independent schema translation, *The VLDB Journal*, 17(6): 1347—1370.
- [3] Atzeni, P., Del Nostro, P., Paolozzi, S. (2008). Ontologies And Databases: Going Back And Forth. In: 4th International VLDB Workshop on Ontology-based Techniques for DataBases in Information Systems and Knowledge Systems (ODBIS 2008). Auckland, New Zealand.
- [4] Cai, J., Johnson, S., Hripcsak, G. (2000). Generic Data Modeling for Home Telemonitoring of Chronically Ill Patients. In: American Medical Informatics Association - Annual Symposium 2000 (AMIA 2000), pp. 116—120. Los Angeles, CA.
- [5] Elmasri, R., Navathe, S.B. (2000). *Fundamentals of Database Systems*, 3rd Edition. Addison Wesley Publishing Company: Reading, Mass.
- [6] Fischer P.C., Van Gucht D. (1985). Determining when a Structure is a Nested Relation. In: 11th Int. Conf. on Very Large DataBases (VLDB 1985), pp. 171—180. Morgan Kaufmann: Stockholm, Sweden.

- [7] Johnson S. B., Chatziantoniou D. (1999). Extended SQL for manipulating clinical warehouse data. In: American Medical Informatics Association Symposium (AMIA 1999), pp. 819—823.
- [8] Johnson, S.B. (1996). Generic data modeling for clinical repositories. *Journal of American Medical Informatics Association*, 3 (5): 328—339.
- [9] Karanikolas, N.N., Nitsiou, M., Yannakoudakis, E.J., Skourlas, C. (2007). CUDL language semantics, liven up the FDB data model. *In: 11th East-European Conf. on Advances in Databases and Information Systems (ADBIS 2007)*, local proceedings, pp. 1-16. Technical Univ. of Varna: Varna, Bulgaria.
- [10] Karanikolas, N.N., Nitsiou, M., Yannakoudakis, E.J. (2008). CUDL Language Semantics: Authority Links. *In: 12th East-European Conf. on Advances in Databases and Information Systems (ADBIS 2008)*, pp.123—139. Tampere Univ. of Technology: Pori, Finland.
- [11] Karanikolas, N. N., Nitsiou, M., Yannakoudakis, E. J., Skourlas, C. (2009). CUDL Language Semantics: Updating Data. *Journal of Systems and Software*, 82(6): 947—962. doi:10.1016/j.jss.2008.12.031
- [12] Karanikolas N. and Vassilakopoulos M. (2009). Conceptual Universal Database Language: Moving Up the Database Design levels. *In: Proc. Of ADBIS'09, LNCS 5739, Riga*, pp. 330—346.
- [13] Karanikolas N. and Vassilakopoulos M. (2010). Database Design with Real-World Structures. eRA-5. 5th Conference for the contribution of Information Technology to Science, Economy, Society and Education.
- [14] Martins, J., Nunes, R., Karjalainen, M., Kemp, G.J.L. (2008). A Functional Data Model Approach to Querying RDF/RDFS Data. *In: 25th British National Conference on Databases (BNCOD 2008)*, pp. 153—164. Cardiff, UK.
- [15] Nadkarni P. (2002). An introduction to entity-attribute-value design for generic clinical study data management systems. Presentation in: National GCRC Meeting. Baltimore, MD.
- [16] Nadkarni P.M. (2000). Clinical Patient Record Systems Architecture: An Overview. *Journal of Postgraduate Medicine*, 46 (3): 199—204.
- [17] Pavković, N., Storga, M., Pavlić, D. (2001). Two Examples of Database Structures in Management of Engineering Data. *In: 12th Int. Conf. on Design Tools and Methods in Industrial Engineering*, pp. 89—90. ADM-Associazione Nazionale Disegno di Macchine: Bologna.
- [18] Schek H.J., Pistor P. (1982). Data Structures for an Integrated Data Base Management and Information Retrieval System. *In: 8th Int. Conf. on Very Large DataBases (VLDB 1982)*, pp. 197—207. Morgan Kaufmann: Mexico City, Mexico
- [19] van Keulen, M., Vonk, J., de Vries, A.P., Flokstra, J., Blok, H.E. (2002). Moa: extensibility and efficiency in querying nested data. Technical Report TR-CTIT-02-19. Centre for Telematics and Information Technology, Univ. of Twente, The Netherlands.
- [20] Worboys, M.F., Hearnshaw, H.M., Maguire, D.J. (1990). Object-Oriented Data Modelling for Spatial Databases. *Int. Journal of Geographical Information Systems*, 4 (4): 369—383.
- [21] Yannakoudakis E.J., Diamantis I.K. (2001). Further improvements of the Framework for Dynamic Evolving of Database environments. *In: 5th Hellenic – European Conf. on Computer Mathematics and its Applications (HERCMA 2001)*, Athens, Greece
- [22] Yannakoudakis E.J., Nitsiou M., Skourlas, C., Karanikolas, N.N. (2007). Tarski algebraic operations on the frame database model (FDB). *In: 11th Panhellenic Conf. in Informatics (PCI 2007)*, pp. 207—216. New Technologies Publications: Patras, Greece.
- [23] Yannakoudakis E.J., Tsionos C.X., Kapetis C.A. (1999). A new framework for dynamically evolving database environments. *Journal of Documentation*, 55 (2): 144—158.