

HEARTFAID'S ECRF: LESSONS LEARNT FROM USING A TWO-LEVEL DATA ACQUISITION AND STORAGE SYSTEM FOR KNOWLEDGE DISCOVERY TASKS WITHIN AN ELECTRONIC PLATFORM FOR MANAGING HEART FAILURE PATIENTS

¹ ANDRZEJ A. KONONOWICZ, ¹ KATARZYNA STYCZKIEWICZ, ¹ BOGUMIŁA BACIOR,
² MATKO BOŠNJAK, ² RAJKO HORVAT, ² MARIN PRCELA, ² DRAGAN GAMBERGER,
³ ANGELA SCIACQUA, ⁵ MARIA CONSUELO VALENTINI, ¹ KALINA KAWECKA-JASZCZ,
^{4,5} GIANFRANCO PARATI, ⁶ DOMENICO CONFORTI

¹ Jagiellonian University Medical College, Kraków, Poland

² Rudjer Boskovic Institute, Zagreb, Croatia

³ University "Magna Graecia" of Catanzaro, Department of Experimental and Clinical Medicine, Italy

⁴ University of Milan – Bicocca, Department of Clinical Medicine and Prevention, Milan, Italy

⁵ Department of Cardiology, S. Luca Hospital, Istituto Auxologico Italiano, Milano, Italy

⁶ University of Calabria, Department of Electronics, Informatics, Systems (DEIS), Italy

Abstract: Case report forms are important sources of medical knowledge in all clinical studies. Electronic versions of these forms have several advantages compared to traditional paper-based questionnaires, and they have been adopted in many contemporary research projects in medicine. This paper presents a framework for creating case report forms designed with a two-level approach. Data at the generic information model level is stored in EAV (entity-attribute-value) tables and extended by tables facilitating specification of the questionnaire layout. The second layer (knowledge model) specifies the domain specific concepts describing the field of application of the questionnaire. This framework has been applied and tested in the frame of an EU FP6 research project – HEARTFAID – the objective of which was to build a knowledge-based platform supporting the management of elderly patients suffering from heart failure. Data collected by the electronic case report form (eCRF) was used in the project's knowledge discovery and decision support tasks. The work presents a new way for effective extraction of the data necessary for the integration with the knowledge discovery process in a distributed, service oriented framework of the HEARTFAID platform. It is demonstrated that it is feasible to implement these tasks using the two-level EAV table design.

Keywords: Electronic Data Capture, Remote Data Entry, EAV, Two Layers Modelling, eCRF, Spring Framework

1. Introduction

Electronic Data Capture (EDC) techniques have been used in clinical trials for a long time [8]. The first EDC systems (also known under other names e.g. Remote Data Entry (RDE) Systems) date back to the early 1970s [7]. Since that time a huge amount of applications (either academic or commercial) for creating, managing and publishing medical on-line forms has been developed. The electronic versions of questionnaires seem to have lot of advantages in comparison to their paper-based counterparts [18]. Among the assets of EDC are cost savings, faster dissemination of forms and collection of data, built-in validation mechanisms, easy maintenance and export to statistical packages.

The conventional method of designing database schemes for questionnaires is to map a form to a single table or a set of tables in a relational database in which each attribute (question from the form) is stored into an individual column [10, 14]. Even though this technique works fine for many applications, it has become apparent that this method is not always effective [14], especially in bio-medical research or electronic health records. This problem pertains to databases with a large, heterogeneous list of fields from which many are optional and can be omitted. In such databases new fields are often added, altered or removed after the database has been deployed, and this introduces additional complication in its structure. Designing such databases with the conventional approach is possible but often troublesome and ineffective due to the limitations of tra-

ditional RDBM systems (e.g. a maximum limit of 255 columns in a single database table in some RDBM systems) or to the need to frequently update the database structure.

An alternative approach is to store records as association lists containing (*attribute name, attribute value*) pairs of variables [13]. A database that stores information in that form is called an entity-attribute-value (EAV) database. This storage method by itself is not new since it dates back to at least the time when the LISP programming language was created. However, its application in relational databases has not yet been very widespread. Classical EAV-databases contain one large table with just three columns: identifier of the described object, identifier of the attribute, value of the attribute. Additionally, dictionaries are required which contain metadata describing the attributes applied.

This simple design technique enables a very flexible method of space-efficient storage of heterogeneous data. However, it should be acknowledged that also this approach is not free from flaws. It is well suited for one object-at-a-time queries in which all information about a single object (e.g. patient) is returned, but it is less efficient in complex attribute-centric queries [14]. For such situations special frameworks facilitating more advanced searches in this model are implemented (as e.g. QAV: *querying entity-attribute* framework [13]), thus empowering the user to browse the data more easily. The overhead that is needed to organise the data in an EAV manner is often not worth the effort for the simple and static databases used in many business applications. The queries are also not time-effective, rendering them less suitable for commercial usage. These drawbacks are, however, not as obvious in research projects and clinical trials where more emphasis is put on the flexibility of the tool than its efficiency.

The EAV model can be considered the first generic layer of a database. This tier may be used in virtually any field of application, and can be extended by additional tables supporting more complex data design. An example of a model with such additions is represented by the EAV/CR by Nadkarni et al. [14], which enhances the EAV by structures for the representation of classes and relationships. Other approaches customize the EAV to store clinical forms [5]. The EAV model with its extensions represents an information model of the database which is domain independent. In a two-level approach to database design, a second layer (i.e. the knowledge model), is added on top of the first [11]. This model specifies the domain specific concepts describing the field of application of the questionnaire. It may consist of terminologies and ontologies related to a given specialization field. The values that can be entered into the information model can be constrained by knowledge model archetypes [2] – i.e. special templates that specify at runtime the way data can be entered. Archetypes may be specified autonomously by subject matter experts (e.g. clinicians) without the need to consult database specialists. A clear separation between the first and second level of the database makes the architecture flexible and reusable.

The aim of this paper is to report on the information obtained while implementing a vast two-level electronic case report form (eCRF) which was designed for the cardiology domain. The eCRF is part of a large knowledge-based platform called HEARTFAID supporting the management of elderly heart failure patients, developed in the frame of the EU FP6 research program. It was required by the HEARTFAID project

that the eCRF system implements insertion, modification and querying of large forms (containing over 700 attribute values for each patient). The system needed to be well integrated with the remaining services of the platform.

2. The HEARTFAID Platform

Heart failure (HF) occurs when the heart fails to pump enough blood to meet the metabolic needs of the body's tissues and/or organs. The prevalence of this pathological condition is very high – approximately 10 million patients suffering from HF in Europe. Chronic (C)HF is a disease of older people; the Framingham study noted a doubling of prevalence with each advancing decade, reaching a rate ranging from 7% to 10% in those aged 80 and older. The mortality of patients with severe HF is also high, approaching 50% in the course of one year in NYHA IV¹ class patients. However, it is believed that through regular monitoring and personalised management of patients affected by this condition, their survival rate and quality of life can be significantly improved.

The role of the HEARTFAID platform is to support physicians and healthcare personnel (e.g. nurses) in managing heart failure patients, while at the same time empowering patients to self-monitor their health condition [3, 4]. HEARTFAID is a web-based platform of services integrating several diverse modules (Fig. 1). Its basic function is to collect patient-related biomedical data from different sources (e.g. mobile and wearable measurement devices or medical imaging systems) and enable access to previously collated data from electronic health records. Part of the system includes declarative and procedural knowledge taken from evidence-based sources such as medical guidelines and carefully selected research papers [6, 16]. The users of the platform may securely access the data it contains in a standardised manner. The platform gives access to data taking into account the different roles and rights of the users. New knowledge can be discovered based on the data collected on the platform by employing newly developed artificial intelligence tools. The system has the potential to support physicians in making clinical decisions in the workplace also by alerting them if a dangerous situation is detected. All HEARTFAID services are integrated by an enterprise service bus (ESB). The system utilises a single sign-on mechanism. Users interact with the system through a customisable web portal. The anticipated results of integrating the platform into clinical practice include a reduction in the re-admission of HF patients to hospital, improvement in the quality of treatment and a decrease in management costs [3, 4].

A knowledge-based platform like HEARTFAID requires various forms of medical data acquired from different sources. Data collected from mobile and wearable devices are covered by the Aml service. The role of the HEARTFAID's electronic case report form (eCRF) is to handle all data required by clinicians that need to be inserted manually by the medical personnel. Beyond the scope of the eCRF is the storage of medical knowledge in the form of rules or ontologies which are used for inference in knowledge discovery and decision support systems. However, these modules exploit the data collected by the

1 NYHA – New York Heart Association Functional Classification – A four scale classification of heart failure extent

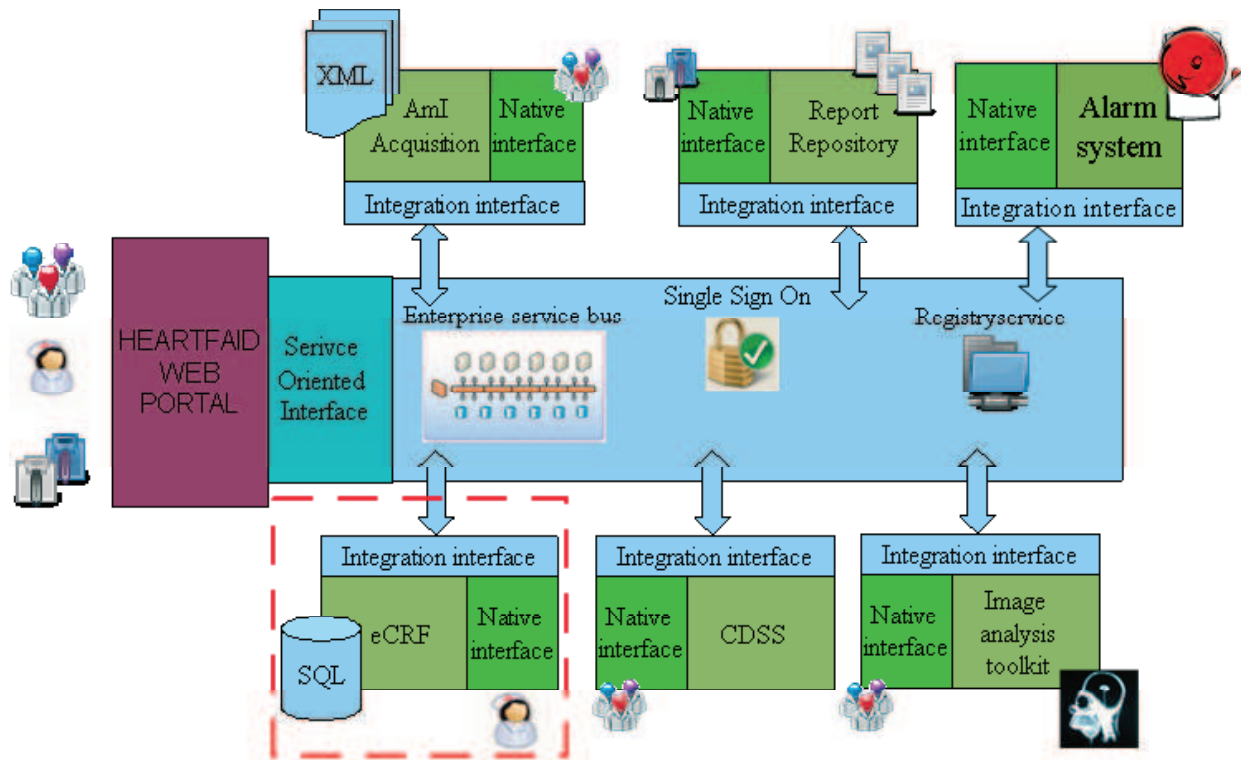


Fig. 1. General overview of the HEARTFAID services

eCRF. The eCRF is intended to be used by medical personnel in the hospital and is not accessible by patients. It plays the role of a specialised electronic health record, collecting heart failure data from a multitude of sources. From a medical (i.e. cardiographic) perspective the eCRF is useful because it gives easy access to the results of lab tests, to treatment schedules and to the prognostic assessment of HF patients.

The HEARTFAID eCRF comprises three parts: the baseline, additional visits and final evaluation forms. Each of these forms is uniquely assigned to a patient and can be filled out only once, with the exception of the additional visit form which may be repeatedly compiled without limitations. Questions in the eCRF questionnaire may be combined into groups. The activation of a group may be triggered in real-time by the input of a specific value by the user. Question groups may be nested to an unlimited depth. Most of the questions are of simple types: Boolean values, text strings, numerical values (integer, real numbers) and dates. However, there are also more complex types of questions which involve, for instance, the selection of a value from a controlled vocabulary, or the use of a special tool to specify a medication and its dosage from a hierarchical list of products (drug class, international name and generic name). It is also possible to add new drugs to the list. Some questions are grouped into matrices (tables) of values of simple types. The forms also contain rules for validating inserted values.

3. Method

While planning the implementation of the eCRF we looked for off-the-shelf products that were web-based, available free, open source, flexible enough to add new question types, able to

support large questionnaires with nested groups of questions, based on XML and J2EE technologies, and easy to integrate into the HEARTFAID platform. None of the existing tools we found for designing web-based questionnaires (e.g. ArchiMed [5], Form Handler [21], Instant Survey [24], Survey Monkey [26], WebEAV [15], Zoomerang [27]), fully met our demands. For that reason it was decided to implement the eCRF from scratch. Since the number of questions was large (more than 700), quite diverse, potentially changeable and the efficiency of the tool was not a critical factor, it was decided to employ a two level architecture. The idea of a two-level approach emerged in electronic health records development [11]. Following this approach database structures are divided into two separated models: information model and knowledge model. The information model represents stable and generic concepts, whereas the knowledge model depicts the dynamics of the problem field [11]. In the HEARTFAID's eCRF the information model expresses a generic database for storing clinical forms following the EAV paradigm. The classical EAV data model has been extended to facilitate the usage of complex web-based forms. In Fig. 2 the ERD (entity-relationship-diagram) of the information model underlying the HEARTFAID eCRF is presented. The model is generic – i.e. it does not contain any information specific to the heart failure domain and can be used in diverse multi-centred clinical trials. The EAV model was implemented in a RDBM system. Basic EAV tables were extended by additional tables for storing hierarchical question groups and for user management. Similar approach was taken in other EAV projects (as e.g. in EAV/CR representation by Nadkarni et al [14]). The *User Center* and *User* tables enable the separation of patients coming from different research institutions and enable access to the data only by entitled users. The *Patient*

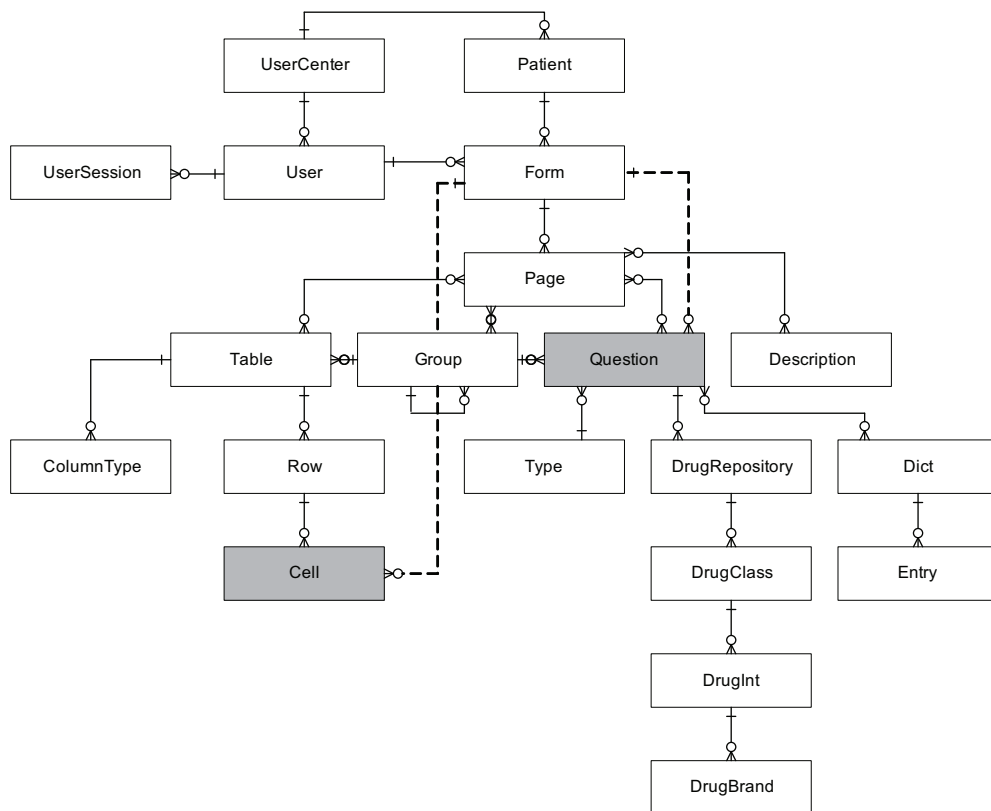


Fig. 2. ERD of the information model under laying the HEARTFAID eCRF

```

<bean id="physical_exam_systolic_blood_pressure" class="org.javs.ecrf.mvc.model.Type" singleton="true">
  <property name="type" value="integer"/>
  <property name="html">
    <value><![CDATA[Systolic blood pressure:]]></value>
  </property>
  <property name="cui" value="C1306620"/>
</bean>

<bean id="question_hfcrfsa_1_g1_1" class="org.javs.ecrf.mvc.model.Question" singleton="false">
  <property name="order" value="1"/>
  <property name="type" ref="physical_exam_systolic_blood_pressure"/>
  <property name="value" value="" />
</bean>

```

Fig. 3. XML archetypes specifying the values that can be inserted into the form

table contains basic patient data. Since the questions are assigned to pages and these pages may contain many levels of nested question groups or tables this structure is reflected by the *Page*, *Table* and *Group* entities. The grey-shaded *Question* and *Cell* tables are classic EAV tables containing a reference to the type of the question, the owning entity (i.e. *Page*, *Group*, or *Row*) and the value. Additionally, these tables contain information about the time of creation and last modification of the values, version number, as well as the identity of the user that modified the value. The dashed line between the *Question* or *Cell* tables and the *Form* table is a redundant connection added for efficiency reasons to accelerate queries with fields nested deep in many subgroups. The *Description* table contains textual information needed as additional description in

the forms. The *Drug[X]* ($X \in \{Repository, Class, Int, Brand\}$) and *Dict* tables represent respectively the pharmacological treatment and values from controlled vocabularies.

The archetypes (second layer of the model) constraining the values that can be inserted into the database are specified in XML syntax, compatible with the bean definition syntax of the Spring Application Framework [9]. An example of the specification of a question type and its instance is presented in Fig. 3. The first archetype bean example defines a type representing a patient's systolic blood pressure taken during a physical examination. This question has its description in HTML format (attribute *html*), stating that it accepts only integer values (attribute *type*) and a mapping to a concept in the UMLS ontology explaining its semantics (attribute *cui*). The second bean

is an instantiation of the previously mentioned question type (attribute *type*). The position at which the question is displayed in the question group is specified by the attribute order, and its default value may be specified by the attribute value. The archetypes often also contain lists of questions or subgroups aggregated by group type, or information about question groups being activated or deactivated based on specific values of the questionnaire fields inserted by the user.

Both archetype beans (i.e. type declaration bean and question instantiation bean) are mapped to POJO (plain old java objects) elements and are stored on demand in the eCRF database using the Hibernate Framework [12]. The way the archetypes are specified enables easy extension of the list of constraining rules (e.g. by information about the soft or hard limits for ranges of accepted values).

4. Results

The eCRF has been implemented in the course of the second year of the HEARTFAID project in the Java 5 programming language. The development has been accelerated by the usage of the Spring Application Framework [9] version 1.2 and Hibernate 3 [12]. The final knowledge model of the eCRF specified by XML archetypes included 735 question instances of 364 semantic types. Archetypes were created using a general purpose XML editor (Altova XMLSpy 2008 [20]). Data were stored in a MySQL 5.1 RDBM system. A simplified structure of the eCRF is presented in the figures included in the Appendices 1 and 2. In order to make the schemes legible, the number of fields for each object was limited to a maximum of 10 fields. The letters *b,a* and/or *f* denote in which eCRF type of form this question is located (i.e. *baseline*, *additional form* or *final visit*). The forms are presented online as HTML views created with

XSLT transformation of XML archetypes and data retrieved from the database (Fig. 4). The *top bar* contains the questionnaire's name, patient id and page number. The pages can be changed either through the list of pages in the *table of content* panel in the right part of the form or through the *backward* and *forward* buttons in the navigation bar. The form is automatically saved after changing a page or after clicking on the *submit* button. Activation of the *cancel* button rejects the last changes and exits the form. Question groups are marked by red boxes and activated by trigger questions (e.g. in Fig. 4 the group containing the *max. ST depression* question is activated by setting the value "yes" in the *ST depression* field). In Fig. 5 a 3x3 question table (matrix) of integers is presented. In addition, above the main form a list of detected validation errors is presented.

Communication with the eCRF with the HEARTFAID platform is established through an XML protocol implemented by one of the partners in the project (SYNOPSIS) including all the necessary information of an HL7 message [22] and following the transactions suggested by IHE [23]. The HEARTFAID middleware implements Patient Demographic Query HL7 V3 (PDQ) and Patient Identifier Cross-Reference HL7 V3 (PIX) profiles. In order to integrate the patient-related data into the platform a MPI (Master Patient Index) service is used which manages patient's demographic information and guarantees their unique identification in the environment. For instance, while registering a new patient on the platform, a message is sent from the HEARTFAID portal to the ESB which was implemented using the Mule open-source framework [25]. Mule descriptors for routing the information to a MIDA Graph (a workflow engine implemented by SYNOPSIS [19]) are read and transformed into information that is stored in the MPI and transmitted as HTTP XML messages to the eCRF service. The eCRF receives the messages, enrolls the patient and sends back a confirmation message [19].

The screenshot shows a web browser window displaying the HEARTFAID eCRF interface. The main content area is titled "10. 12-lead electrocardiography (25 mm/s)". It contains several input fields and radio buttons for clinical data. A "Table of content" panel on the right lists various sections, with "12) 12-lead ECG" highlighted. A "Navigation bar" at the bottom includes "Cancel" and "Submit" buttons. Red boxes highlight specific question groups, and blue callout boxes point to the "Question groups", "Table of content", and "Navigation bar" labels.

Fig. 4. User interface of the HEARTFAID eCRF

The screenshot shows a Mozilla Firefox browser window displaying the HEARTFAID eCRF form. At the top, a red banner reads "This page contains following errors". Below it, three validation errors are listed:

1. Question SVE: Field doesn't contain integer value
2. Question SDRM: Field doesn't contain a number
3. Question Table HR cells 1,2: Field doesn't contain integer value

A blue box labeled "Form validation" points to these errors. The main form is titled "12. 24 h Holter electrocardiography". It includes a date field set to "Jul 2 2009". Below the date is a table for HR (Heart Rate) data:

	mean	min	max
24h	test		
Day (6-22)			
Night (22-6)			

Below the table are fields for SVE (6.6 / 24h), SVTach (3 beats / 24h), and Atrial fibrillation/flutter (No/Yes). A blue box labeled "Question Table" points to the HR table. On the right, a red sidebar shows a navigation menu with items like Demographic, Inclusion, Exclusion, Card. history, etc.

Fig. 5. Question tables and form validation in eCRF HEARTFAID

The screenshot shows the HEARTFAID web interface in Internet Explorer. The URL is <http://www.staging.gp/heartfaid/index.php?page=4&patid=4127&test=204-405-9130-24492419203>. The patient ID is 097304231. The interface displays a list of diagnoses with associated reasoning:

- Diagnosis (6)
 - why? Heart failure signs
 - why? Heart failure symptoms
 - why? Suggested to perform BNP test
 - why? Diagnosed heart failure diastolic negative
 - why? Diagnosed heart failure systolic positive by echo
 - why? Patient has some signs or symptoms of acute decompensation
- Additive diagnosis (5)
 - Severity (4)
 - Prognosis (1)
 - Medication status (4)
 - Medication suggestions (2)
 - Medication warnings (4)
 - Management (2)
 - Other (4)

The interface also includes a sidebar with navigation options and a patient photo.

Fig. 6. The result of reasoning based on the data collected by the eCRF

The eCRF was deployed on the HEARTFAID platform in 2007 and since then it has been in constant use. Data from approximately 100 patients from four clinical centres [Università degli studi Magna Graecia, Catanzaro (Italy), Università degli studi di Milano Bicocca, Milan (Italy), Jagiellonian University Medical College, Kraków (Poland) and S. Luca Hospital, Istituto Auxologico Italiano, Milan (Italy)] have been collected.

The eCRF has been integrated with the HEARTFAID's Knowledge Discovery Service (KDS) and Decision Support

System (DSS) developed by Rudjer Boskovic Institute in Zagreb (Croatia). Both services require tight integration with the large amount of patient data collected by eCRF, however these services require substantially different data access types. DSS is always focussed on one patient while KDS requires information about all or most of available patient data that has been collected by the eCRF. Additionally, it must be noted that DSS requires effective access to the most recent information for all potentially relevant measurements regardless of when they

were collected and with a clear indication about when the data was acquired. In contrast to this, actual data collection time is not relevant for KDS, but it requires access to the data grouped according to the time of its collection, that data should be ordered by its historical order, and that it is identified by the time interval from previous measurements. Fig. 6 demonstrates a typical result from the decision support service while Fig. 7 and 8 illustrate the knowledge discovery service.

A unique property of the currently implemented KDS is that it integrates knowledge discovery algorithms with direct database access into one web-based service. This is not a simple task due to the complexity of the KD process [16]. The HEARTFAID service implements the modern random forest based machine learning algorithm [1] that has been reimplemented by Rudjer Boskovic Institute. The service has been built as a series of projects so that every project consists of different datasets with many tasks that can be performed for every dataset. Access to projects, datasets, and tasks is enabled through a web interface (Fig. 7).

Computationally, the most complex part of the service is the construction of the classifier and the preparation of a report (Fig. 8) based on the results of this process. The value of this newly implemented service is the realisation of direct access to the data in the eCRF and its automatic transformation into a form that can enter the KD process. Direct access to the relational database containing EAV tables by a traditional SQL interface is very laborious. This problem can be solved by implementing a special query module for external analytical services. An interface consisting of four generic functions has been implemented for the purpose of the knowledge discovery task. Table 1 contains a description of these functions:

These functions are available through a HTTP GET interface. In the following line a command is shown that starts a query involving the *getLastValue* function `http://localhost:8080/heartfaid/query.html?function=getLastValue&uuid=312&sid=physical_exam.weight`

After execution the interface searches in the eCRF database for all values of the *physical_exam.weight* attribute regarding the patient with the id 312. The result of the function is returned in simple XML syntax. This allows a clear separation of the data collection and query tool located at one centre (currently at Jagiellonian University Medical College, Kraków, Poland), from the knowledge discovery system located at a remote centre (currently located at Rudjer Boskovic Institute, Zagreb, Croatia).

Discussion

After the generic framework for designing questionnaires had been developed, the process of implementation of the HEARTFAID's eCRF knowledge model by specifying the XML archetypes took little time and did not cause any difficulties. The structure of the eCRF turned out to be more stable than initially anticipated, so the benefit of the flexibility of the architecture has not been fully used (with the exception of a few minor changes). On the other hand, the drawbacks of decreased database efficiency in this type of application are hardly noticeable. In the production database containing just a few users and approximately 100 forms installed on a Intel Core 2 Duo T5450 1,66Ghz,1GB RAM computer, loading a whole form from a database took in average 1360 ms, saving a modified

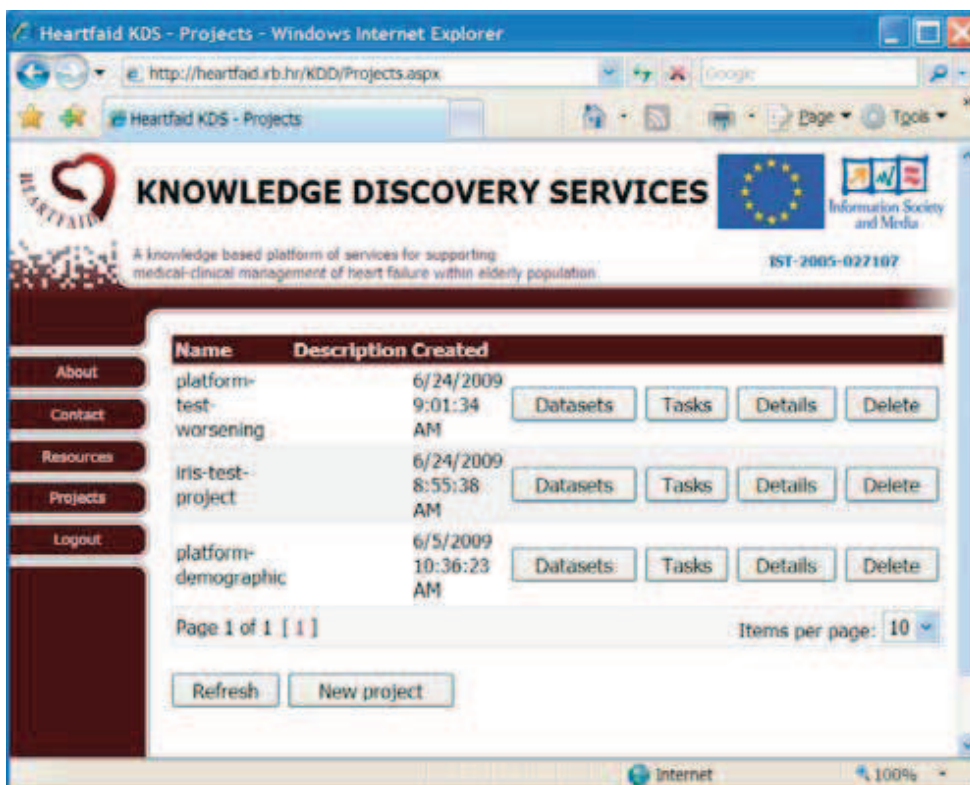


Fig. 7. The main page of the HEARTFAID knowledge discovery service with three current projects: "platform-test-worsening", "Iris-test-project", and "platform-demographic"

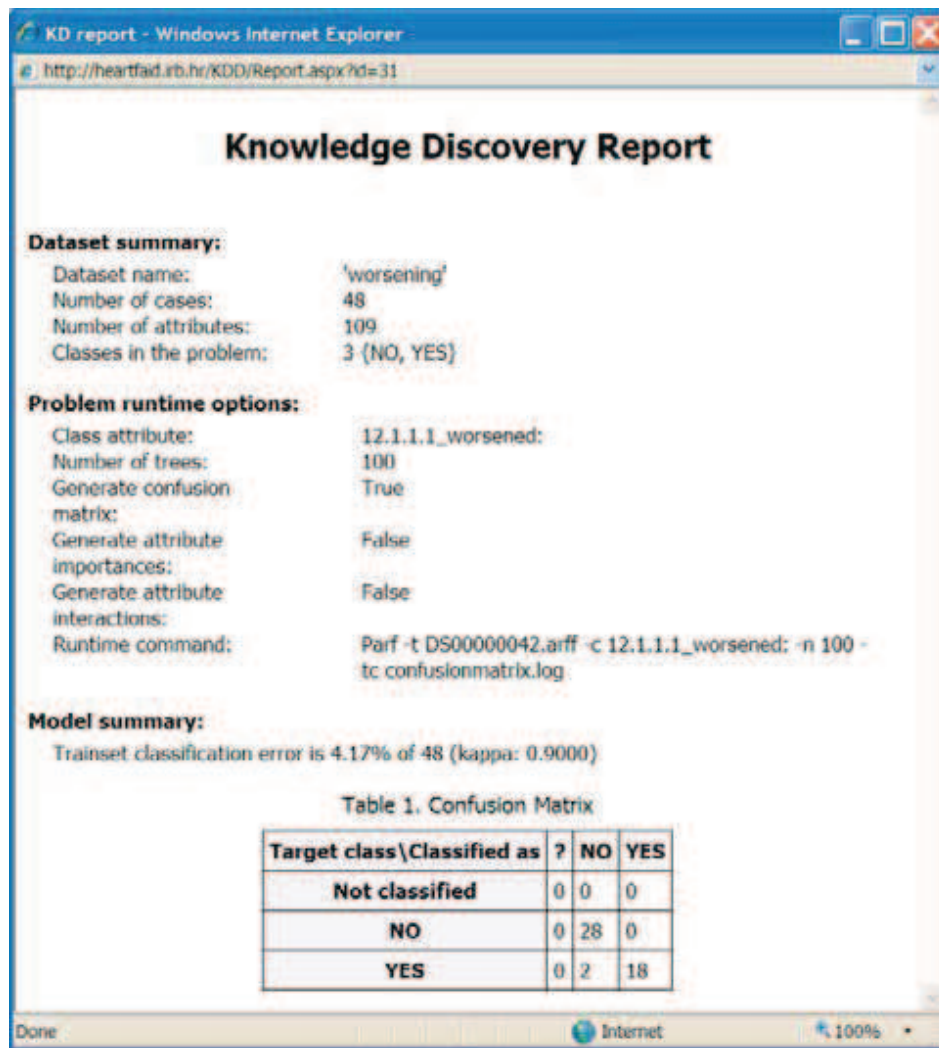


Fig. 8. The result of any KD task is a report. The figure presents a report for a two-class domain obtained after constructing a random forest with 100 trees. The main part of the report is the confusion matrix demonstrating the predictive accuracy measured by cross-validation on the training set.

form 600 ms, querying the last value of a selected parameter 84 ms. . Thanks to the application of XML technology the integration of the eCRF to the platform's enterprise service bus was easy and fulfilled the requirements of current medical informatics standards.

The future plan for the proposed architecture includes implementation of a graphical editor for the XML archetypes and extension of the list of constraints that can be used for the knowledge model's specification. Tighter integration of the eCRF with knowledge engineering and data mining tools through the proposed interface also seems to be important.

It is not easy to give definite advice about when to use EAV tables instead of traditional relational database design. If our highest priority is flexibility, and the number of collected attributes is very large and potentially often changeable, this suggests that a two-level EAV design should be used. In all other cases, a more traditional design would probably be more advantageous. When designing frameworks with EAV databases for knowledge discovery tasks it is imperative to also offer a special query module with an interface similar to that

presented in this paper, or to export the data to an external system with a different data model.

Conclusions

This paper presents a practical implementation of a two-level database system for a medical research project. The generic layer of this database uses EAV tables which are useful for designing large heterogeneous and frequently changeable database schemas, as are often found in research studies. In this system a method for implementing the concept of two-level architecture in a modern application framework (Spring Framework) has been demonstrated. The fact that the system has been in use for almost two years in the HEARTFAID project and has delivered useful data for other modules like a knowledge discovery module and decision support services proves the feasibility and the effectiveness of this solution. The significance of our work consists in the proposal of a new type of direct interface for accessing complex data structures with the

Tab. 1. eCRF interface for knowledge discovery tasks

Function Name	Description
<i>getLastValue</i>	Returns the last known descriptor value available for the patient. If all values are unknown the returned value is also unknown.
<i>getAnyValue</i>	Returns information concerning all previous visits. For numerical measurements it returns two values: minimum and maximum while for categorical attributes it returns most frequent (mode) value. If all values are unknown the returned value is also unknown.
<i>getDifference</i>	Returns the difference between the last available piece of data and the penultimate piece. If there are not two available entries the value is unknown. For numerical attributes (e.g. laboratory values) it is the difference (+/- value). For categorical attributes it is 0 (no change) and 1 (value changes) [or -1 improved, 0 no change, 1 worsening]
<i>getFlattenedTable</i>	For categorical values it returns the number of known values and the most frequent value. For numerical it returns mean, minimal and maximum value, range, standard deviation and slope.

output already prepared for artificial intelligence applications. Additionally, it is also clearly stated that this model is not appropriate for every database, especially not for large commercial databases, and therefore its adoption needs to be carefully considered.

References

- Breiman L., Random Forests, *Machine Learning* 45(1), pp. 5-32, 2001.
- Bird L., Goodchild A., Tun Z., Experiences with a Two-Level Modelling Approach to Electronic Health Records, *Journal of Research and Practice in Information Technology*, 35(2), pp. 121-138, 2003.
- Chiarugi F. et al., Support for the Medical-Clinical Management of Heart Failure within Elderly Population: the HEARTFAID Platform, Proc. of ITAB, Ioannina, Greece, 26-28 October 2006.
- Conforti D. et al., HEARTFAID: A Knowledge Based Platform for Supporting the Clinical Management of Elderly Patients with Heart Failure, *The Journal on Information Technology in Healthcare*, 4(5), pp. 283-300, 2006.
- Duftschmid G., Gall W., Eigenbauer E., Dorda W., Management of data from clinical trials using the ArchiMed system, *Med. Inform. Internet*, 27(2), pp. 85-98, 2002.
- Gamberger D., Prcela M., Jović A., Šmuc T., Parati G., Valentini M., Kawecka-Jaszcz K., Kononowicz A. A., Candelieri A., Conforti D., Guido R., Medical Knowledge Representation Within Heartfaid Platform, Healthinf, Funchal, Madeira – Portugal, 2008.
- Helms R. W., Entering Data from Remote Terminals in Clinical Centers using IBM's OS/TSO in the Kidney Transplant Histocompatibility Study, Technical Report 007, Chapel Hill, NC University of North Carolina, KTHS Statistics and Data Management Center, Department of Biostatistics, 1973.
- Helms R. W., Data Quality Issues in Electronic Data Capture, *Drug Information Journal*, 35, pp. 827-837, 2001.
- Johnson R., Hoeller J., Arendsen A., Risberg T., Sampaneau C., *Professional Java Development with the Spring Framework*, John Wiley & Sons, 2005.
- Merzweiler A., Weber R., Garde S., Haux R., Knaup-Gregori P., TERMTrial – terminology-based documentation systems for cooperative clinical trials, *Comput. Meth. Programs Biomed.*, 78, pp. 11-24, 2005.
- Michelsen L., Pedersen S. S., Tilma H. B., Andersen S. K., Comparing different approaches to two-level modelling of electronic health records., *Stud. Health Technol. Inform.*, 116, pp. 113-118, 2005.
- Minter D., Linwood J., *Hibernate From Novice to Professional*, Apress, 3 edition, 2006.
- Nadkarni P., QAV: querying entity – attribute – value metadata in a biomedical database, *Comput. Meth. Programs Biomed.*, 53, pp. 93-103, 1997.
- Nadkarni P. et al., Organization of Heterogeneous Scientific Data Using the EAV/CR Representation, *J. Am. Med. Inform. Assoc.*, 6(6), pp. 478-493, 1999.
- Nadkarni P., Brandt C., Marengo L., WebEAV: Automatic Metadata-driven Generation of Web Interfaces to Entity-Attribute-Value Databases, *J. Am. Med. Inform. Assoc.*, 7(4), pp. 343-356, 2000.
- Prcela M., Gamberger D., Bogunovic N., Developing Factual Knowledge from Medical Data by Composing Ontology Structures, MIPRO 2007, Opatija, Croatia.
- Sonicki Z., Gamberger D., Smuc T., Sonicki D., Kern J., Data mining server: On-line knowledge induction tool, in: Proc. of Medical Informatics Europe, IOS press, pp. 330-334, 2002.
- Wyatt J. C., When to Use Web-based Surveys, *J. Am. Med. Inform. Assoc.*, 7(4), pp. 426-430, 2000.
- HEARTFAID Consortium, D28 – Integration and Interoperability middleware prototype, 2008.
- Altova XMLSpy, <http://www.altova.com/xml-editor/>
- Form Handler, <http://www.formhandler.net>
- HL7, Health Level 7, <http://www.hl7.org>
- IHE, Integrating the Healthcare Enterprise, <http://www.ihe.net>
- Instant Survey, <http://www.instantsurvey.com>
- Mule, ESB <http://www.mulesoft.org/display/COMMUNITY/Home>
- Survey Monkey, <http://www.surveymonkey.com>
- Zoomerang, <http://www.zoomerang.com>

Appendix 1 – Knowledge model of the HEARTFAID eCRF, Simplified – Part 1 of 2

Anamnesis		
blood pressure change	(a.f)	enum[change]
bradycardia	(b)	boolean
bradycardia change	(a.f)	enum[change]
chest pain	(b)	boolean
chest pain change	(a.f)	enum[change]
chest pain remote	(b)	boolean
dyspnoea	(b)	boolean
dyspnoea change	(a.f)	enum[change]
dyspnoea remote	(b)	boolean
fatigue	(b)	boolean
<i>and 19 more fields</i>		

Echocardiography		
aorta ascending aorta diameter	(b.a.f)	double
aorta root diameter	(b.a.f)	double
contractility akinesis	(b.a.f)	boolean
left atrium anteroposterior diameter	(b.a.f)	double
left ventricle end-diastolic diameter	(b.a.f)	double
left ventricle end-diastolic volume	(b.a.f)	integer
mitral valve deceleration time	(b.a.f)	integer
mitral valve emax-amax	(b.a.f)	double
mitral valve mitral regurgitation	(b.a.f)	integer
pulmonary artery pressure	(b.a.f)	integer
<i>and 16 more fields</i>		

24 h Holter Electrocardiography		
atria fibrillation flutter	(b.a.f)	boolean
conduction abnormalities	(b.a.f)	boolean
conduction abnormalities details	(b.a.f)	textfield
date	(b.a.f)	date
heart rate HF	(b.a.f)	double
heart rate LF	(b.a.f)	double
heart rate pNN50	(b.a.f)	double
heart rate rMSSD	(b.a.f)	double
heart rate SDANN	(b.a.f)	double
heart_rate_total_power	(b.a.f)	double
<i>and 9 more fields</i>		

Final Visit		
date	(f)	double
required_hospitalization_date	(f)	double

Rehabilitation		
model required	(f)	textfield
time	(f)	integer

Chest X-ray		
cardio-thoracic_ratio	(b.a.f)	integer
comment	(b.a.f)	textarea
date	(b.a.f)	date
pulmonary_congestion_or_oedema	(b.a.f)	boolean

Quality of Life Questionnaire		
date	(b.f)	date
minnesota_total_score	(b.f)	integer
sf36_bodily_pain	(b.f)	integer
sf36_general_health	(b.f)	integer
sf36_mental_component_summary	(b.f)	integer
sf36_mental_health	(b.f)	integer
sf36_physical_component_summary	(b.f)	integer
sf36_role_emotional	(b.f)	integer
sf36_role_physical	(b.f)	integer
sf36_social_functioning	(b.f)	integer
<i>and 2 more fields</i>		

Family History		
primary_cardiomyopathy	(b)	boolean

Beat-to-beat Blood Pressure Monitoring		
baseline_finger_BP_SBP	(b.f)	integer
baseline_finger_HR	(b.f)	integer
comments	(b.f)	textarea
cuff_size	(b.f)	enum[cuff_size]
date	(b.f)	date
device	(b.f)	enum[device]
end_standing_CB_finger_BP_SBP	(b.f)	integer
end_standing_CB_finger_HR	(b.f)	integer
finger	(b.f)	enum[finger]
hand	(b.f)	enum[hand]
<i>and 13 more fields</i>		

Drug Therapy		
drug_therapy_change	(b.a.f)	drug
drug_therapy_change	(a)	drug

Laboratory Assessment		
ALT	(b.a.f)	double
AST	(b.a.f)	double
blood_samples_for_DNA-RNA	(b.a.f)	boolean
BNP	(b.a.f)	pmol_mg
creatinine	(b.a.f)	umol_mg
creatinine_clearance	(b.a.f)	double
date	(b.a.f)	date
glucose	(b.a.f)	mmol_mg
glycated_hb	(b.a.f)	double
hb	(b.a.f)	double
<i>and 12 more fields</i>		

Physical Examination		
body temperature	(b.a.f)	double
diastolic blood pressure	(a.f)	integer
heart murmurs	(b.a.f)	boolean
heart murmurs_apex	(b.a.f)	boolean
heart murmurs_base	(b.a.f)	boolean
heart murmurs_diastolic	(b.a.f)	boolean
heart murmurs_systolic	(b.a.f)	boolean
heart sounds	(a.f)	boolean
heart sounds_bilateral	(b.a.f)	boolean
heart sounds_fourth	(b.a.f)	boolean
<i>and 24 more fields</i>		

Cardiopulmonary Exercise Testing		
AT	(b.f)	double
BP baseline DBP	(b.f)	integer
BP baseline SBP	(b.f)	integer
BP end DBP	(b.f)	integer
BP end SBP	(b.f)	integer
BP_peak_ex_DBP	(b.f)	integer
BP_peak_ex_SBP	(b.f)	integer
data_recorded	(b.f)	boolean
O2_pulse	(b.f)	double
RQ	(b.f)	double
<i>and 10 more fields</i>		

Additional Visit		
date	(a)	date
next_scheduled_visit_date	(a)	date
other_than_chf_reasons_of_visit	(a)	boolean
required_advice	(a)	boolean
required_advice_details	(a)	boolean

Appendix 2 – Knowledge model of the HEARTFAID eCRF – Simplified – Part 2 of 2

Six-minute walking test	
BP baseline DBP	(a) integer
BP baseline SBP	(a) integer
BP end DBP	(a) integer
BP end SBP	(a) integer
date	(a) date
HR baseline	(a) integer
HR end	(a) integer
SpO2 baseline	(a) integer
walking_distance	(a) integer

Lifestyle Information	
alcohol use	(b) boolean
physical activity	(b) enum[ph_activity]
smoking	(b) boolean
smoking cessation	(a,f) boolean
smoking cessation date	(a,f) date
smoking duration	(b) integer
smoking_no_cigarettes	(b) integer

Non Cardiovascular Medical History	
anemia	(b) boolean
anemia worsening	(a,f) boolean
bronchial asthma	(b) boolean
connective tissue diseases	(b) boolean
diabetes	(b) boolean
diabetes type	(b) enum[type12]
diseases not related to hf	(b) boolean
diseases_potentially_related_to_hf	(b) boolean
endocrine disorders	(b) boolean
exposure_to_endemic_diseases	(b) boolean

and 29 more fields

Demographic Data	
birthday	(b) date
death	(a,f) boolean
death_cause	(a,f) textfield
death_date	(a,f) date
sex	(b) enum[sex]
status	(b) enum[pat_status]

Cardiovascular Status	
aortic regurgitation	(b) boolean
aortic stenosis	(b) boolean
CABG	(b,a,f) boolean
cardiovascular reason of death	(a,f) boolean
cerebrovascular events	(b,a,f) boolean
changes in therapy	(a) boolean
chf status improved	(a) boolean
chf status requires hospitalization	(a) boolean
congenital heart disease	(b) boolean
congestive heart failure	(b) boolean

and 41 more fields

12-Lead Electrocardiography	
conduction LBBB	(b,a,f) boolean
conduction PQ	(b,a,f) integer
conduction QRS	(b,a,f) integer
conduction QT	(b,a,f) integer
conduction RBBB	(b,a,f) boolean
date	(b,a,f) date
heart rate	(b,a,f) integer
heart rate 24h max	(b,a,f) integer
heart rate 24h mean	(b,a,f) integer
heart_rate_24h_min	(b,a,f) integer

and 14 more fields

Substudy 1 - Inclusion Criteria	
age_gt_65	(b) boolean
chf	(b) enum[diag_chf]
diastolic dysfunction	(b) boolean
ef_lt_40p	(b) boolean
functional capacity	(b) boolean
hypertension	(b) boolean
idcm	(b) boolean
ihd	(b) boolean
informed consent	(b) boolean
sinus_rhythm_presence	(b) boolean

and 2 more fields

Substudy 1 - Exclusion Criteria	
AIDS	(b) boolean
autoimmune disorders	(b) boolean
cardiac resynchronization therapy	(b) boolean
drug or alcohol abuse	(b) boolean
gfr_lt_30	(b) boolean
hepatic disease	(b) boolean
immunosuppressive therapy	(b) boolean
malignancy	(b) boolean
no_informed_consent	(b) boolean
pacemaker	(b) boolean

and 3 more fields

Patient	
id	long
uid	string
initials	string
usercenter	integer
createTime	date
updateTime	date
createUser	string
updateUser	string