

# IBM Research Report

## Model Driven Approach for Biomedical Data Integration

**David Carlson**

Veterans Health Administration  
Kalispell, MT USA

**Ariel Farkash**

IBM Research Division  
Haifa Research Laboratory  
Mt. Carmel 31905  
Haifa, Israel

**John T. E. Timm**

IBM Research Division  
Almaden Research Center  
650 Harry Road  
San Jose, CA 95120-6099 USA



Research Division

Almaden - Austin - Beijing - Cambridge - Haifa - India - T. J. Watson - Tokyo - Zurich

# Model Driven Approach for Biomedical Data Integration

David Carlson<sup>a</sup>, Ariel Farkash<sup>b</sup>, John T.E. Timm<sup>c</sup> (authorship ordered alphabetically)

<sup>a</sup> Contractor to the Veterans Health Administration, Kalispell, MT, United States

<sup>b</sup> IBM Haifa Research Lab, Haifa University Campus, Mount Carmel, Haifa, Israel

<sup>c</sup> IBM Almaden Research Center, San Jose, CA, United States

## Abstract

A core challenge in biomedical data integration is to enable semantic interoperability between its various stakeholders as well as other interested parties. Promoting the adoption of worldwide accepted information standards along with common controlled terminologies is the right path to achieve that. This paper describes a solution to this fundamental problem by proposing an approach of semantic data integration based on information models serving as a common language to represent health data, coupled with a technology that is able to represent the data semantics. We used the HL7 v3 Reference Information Model (RIM) [1] to derive a specific data model for the integrated data, the Web Ontology Language (OWL) [2] to build an ontology that harmonizes the metadata from the disparate data sources, the Unified Modeling Language (UML) [3] to model the data representation, and the Object Constraint Language (OCL) [4] to specify UML model constraints. The Essential Hypertension Summary CDA document and related models are derived from Hypergenes, a European Commission funded project [5] exploring the Essential Hypertension disease model.

**Keywords:** CDA, Ontology, OWL, Modeling, UML, OCL

## Introduction

Biomedical information repositories typically contain data related to a specific clinical domain with semantics unique to the originating systems [6]. These disparate data sources pose a challenge for data integration [7] that is paramount for improved patient-centric care [8], as well as for secondary use of the data for analysis of aggregated data in context of clinical research, public health surveillance, and decision support [9].

In this paper we depict a complete solution to this fundamental problem by proposing an approach of semantic data integration based on information models that serve as a common language to represent health data, ontology based metadata harmonization, technology used for creating and constraining data models, and an engine for instance generation.

Our workflow, as depicted in Figure 1, commences in the clinical domain where a clinical expert must identify the information elements or variables of interest needed for a par-

ticular study. This activity may be partly based on what data is available and how it is collected, but a common theme is that the practitioner does not care about the data format or explicit representation of the complexity of the data, only that certain data elements should be available for further analysis.

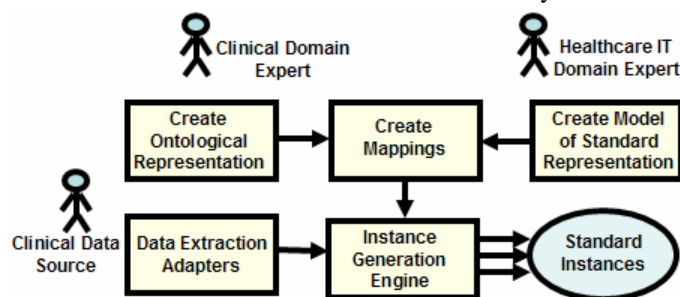


Figure 1 – Workflow

Our approach is intended to work over multiple, heterogeneous data sources, thus we have chosen to represent data using widely adopted international standards for healthcare semantics and interoperability. These standard exchange formats, along with a set of constraints, serve to unify the data into a single semantically unambiguous format that makes operations on the data straightforward from a technological standpoint. However, a clinician is most likely unfamiliar with use of standard healthcare interoperability representation, which naturally falls into the healthcare IT domain expertise.

We have chosen to use a set of industry standard modeling languages to bridge between these two fundamentally different domains and roles required for proper integration of the data. In our approach, the clinical domain expert works with an ontology-based approach using semantic web technologies to represent the metadata needed for harmonization, while the healthcare IT domain expert uses software modeling languages to: create model-based representations of the standard format, apply constraints to this format for a domain of interest, and, in collaboration with the clinical domain expert, create mappings between the ontological representation of the variables of interest and the standard-based information models. The annotated model created by the healthcare IT domain expert at design time, is then used by the instance generation engine at runtime in order to transform the data to the standard format that conforms to the constrained model.

## Background & Related Work

The HL7 v3 Reference Information Model (RIM) is used to derive consistent health information standards such as laboratory, problem and goal-oriented care, public health, and clinical research. It is an ANSI and ISO-approved standard that provides a unified health data ‘language’ to represent associations between entities who play roles that participate in acts. For example, an organization entity plays a role of laboratory that participates in an observation act. Or, a person entity plays a role of a surgeon who participates in a procedure act, and so forth. Acts may relate to other acts through “act relationships”, thus providing a mechanism to describe complex actions.

Clinical Document Architecture (CDA) [10] is a constrained subset of the RIM that specifies terminology encoded structure and semantics of clinical documents. The CDA standard is used to serialize clinical data as XML that is described by a published W3C XML Schema. In most applications, the general CDA structure is further constrained by a set of templates that are standardized and published in an implementation guide, such as the Continuity of Care Document (CCD) [11]. A CDA document instance includes template identifiers (e.g. “2.16.840.1.113883.10.20.1.28”) that determine the template specifications to which the data must conform.

Most CDA template specifications, such as CCD, are written using structured English expressions that are based on the XML schema element relationships. These conformance statements are usually implemented using Schematron rules to augment the CDA XML schema. Our work, however, includes methods and open source software tools for representing CDA documents and template constraints using the Unified Modeling Language (UML) and the Object Constraint Language (OCL). Details and examples of this approach are described in the Methods and Results sections of this paper.

Healthcare applications that produce or consume XML instances for CDA must include the appropriate template identifiers, as specified in the implementation guide. An additional capability of our model-driven software tools is to generate Java runtime libraries that support a registry of CDA templates and enables instance validation.

The UML modeling language is dominant among IT domain users, whereas clinical domain experts often work with formal ontology definitions. The Web Ontology Language (OWL) is a semantic markup language for publishing and sharing ontologies on the World Wide Web. It is endorsed by the World Wide Web Consortium (W3C) [12]. OWL is often used as the framework for converging distinctive terminologies into one coherent ontology; many successful examples exist in clinical research and medical informatics domains [13, 14, 15].

There has been some prior work in both using OWL ontologies in conjunction with instance generation [16], and in using OWL to add semantic annotations to UML information models [17]. These methods are applied and extended to support ontological mapping, representation modeling, formal constraining, and instance generation in our research.

## Methods

### Users

The use case diagram in Figure 2 illustrates the primary activities involved in our approach and the user roles required to perform these activities.

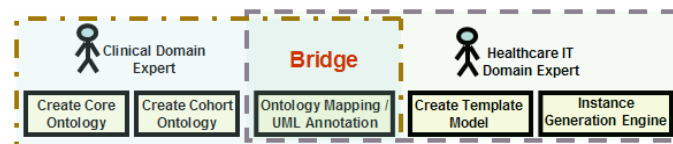


Figure 2 – Use Case Diagram

The clinical domain expert is responsible for creating the core ontology. The core ontology contains conceptual abstractions for a given clinical research domain and includes all the data elements required for secondary use by clinical researchers. The cohort ontology contains data elements specific to a cohort using cohort terminology. The cohort ontology is created by the clinical domain expert for each cohort that wishes to participate in the data integration. Using common ontology development tools such as Protégé [18], mappings are created between these cohort ontologies and the core ontology. This process is described in greater detail in the next sub section.

The healthcare IT domain expert is responsible for creating the CDA template model using a UML tool. The CDA template model contains classes, attributes, and relationships that are used to further constrain the CDA model to a particular clinical research domain. There are implicit relationships between classes in the template model and concepts in the core ontology. These relationships are made explicit by creating mappings on the CDA template model as UML annotations, providing the basis for generating the annotated template model.

The artifacts produced by these different users and the relationships between them are captured in Figure 3 below.

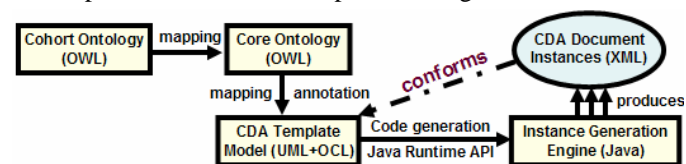


Figure 3 – Artifact Relationships

### Data Integration

Healthcare data integration involves harmonization, validation, normalization, and transformation into standard structures that are accepted by the healthcare and medical research communities. Relationships between data items are often defined implicitly, e.g., in documentation or as tacit knowledge of experts. These implicit relationships must be expressed in an explicit and standard way so that analysis algorithms not aware of the implicit semantics could use them effectively.

## Harmonization

Integration of data from dissimilar data sources must first undergo a process of conceptual harmonization, i.e. convergence of the sources metadata to a single and agreed-upon terminology. For example, blood pressure measurements from three different cohorts of essential hypertension are outlined in Figure 4. This outline depicts the underlying data model for the blood pressure measurements taken by the three cohorts.

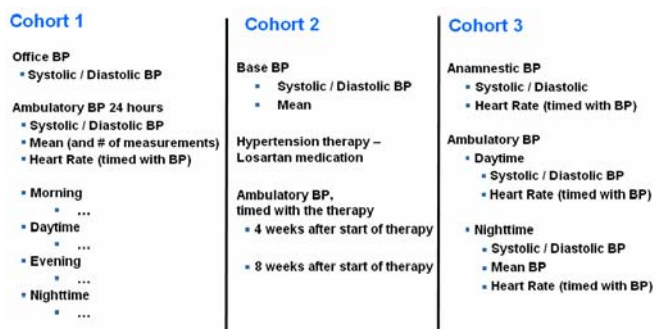


Figure 4 – Various blood pressure measurement schemes

Comparing data between the different cohorts is not a trivial task. The metadata is named differently, so how can one deduce that: Cohort 1 “Office BP”, Cohort 2 “Base BP”, and Cohort 3 “Anamnestic BP” all refer to the same conceptual data? Furthermore, looking at Ambulatory Blood Pressure findings one can see that Cohort 1 temporal divisions are to “Morning, Daytime, Evening, and Nighttime”, whereas in Cohort 3 we find “Daytime and Nighttime” only; Cohort 2 blood pressure observations relate to four and eight weeks after start of therapy, thus completely incomparable to the above data.

In order to be able to compare data of different cohorts, one should first converge to a core terminology. Using OWL, we leveraged technology used for semantic web representation, to map all cohort variables to a core ontology able to represent the base conceptual terms for the target domain, e.g. Essential Hypertension. A schematic diagram is shown in Figure 5.

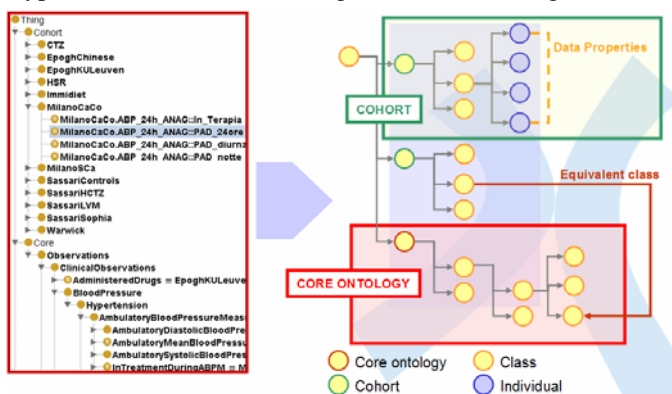


Figure 5- Ontology schematic diagram, left side is a screen capture of ontology using Protégé

The process starts by creating a cohort class (OWL class) for each metadata variable, thus each cohort contains a flat list of cohort classes. We then map each cohort variable in accor-

dance with harmonization effort to a core ontology class by specifying an equivalent class relationship. In case of n:1 mapping, cohort instances (OWL individuals) are created, allowing the class to maintain 1:1 mapping, and additional parameters (as Data Properties) are added to capture the instance disparities. Thus, following the example shown in Figure 4, Cohort 2 Ambulatory Blood Pressure would contain one class with two individuals, having a temporal parameter to specify for four or eight weeks after therapy.

## Normalization & Validation

Having crossed the hurdle of defining metadata in comparable terms, one is left with the challenges of deducing and validating data values for each metadata variable under the cohort’s data model, as well as normalizing values in correspondence to harmonized standard units. This task is a complex one due to differences in: units of measurement, classifications, and diversity of protocols. We will not elaborate on efforts made in this area in order not to diverge from the focus of this paper.

## Transformation to Intermediate Data Representation

Data is first extracted via a suitable adapter from data source proprietary formats, such as an excel file or MySQL database, and copied into a generic data container. The container is conceptually a map where the *key* is a cohort variable and the *value* is the matching value. It is important to note that only part of the variable values is mapped to the core ontology. Other parts may be individuals, belonging to a mapped class, or parameters. The inference is performed by the instance generation engine receiving both the data container and the ontology as input; this will be described in the Methods section.

## Data Representation

### Capturing Richness of Data

Having similar sets of metadata represented in an agreed-upon terminology provides the basis for syntactic interoperability [19], i.e. the ability to compare two orphan items of data. However, biomedical data is typically complex, consisting of associations and dependencies between discrete data items as well as between common structures. Consider the example in Figure 4: in Cohort 2, the Ambulatory Blood Pressure is measured while the subject is treated by a medication called Losartan. This calls for associating the act of observing the blood pressure with the act of administering the drug so that semantics is explicitly represented. This information may be crucial to physician, i.e. high blood pressure while under Losartan regimen has a completely different meaning than without such intervention. Therefore, in order to capture the full richness of the data, these kinds of associations should be established during the data integration process when the experts responsible for the data source can provide the implicit semantics often hidden in unstructured documentation or in their minds. As described in the background, the HL7 v3 RIM provides a unified ‘language’ to represent health acts such as observations, procedures and substance administrations. Using CDA as a RIM derived domain specific standard



facilitates the explicit representation of the rich semantics of these data. Referring back to the examples discussed above, the blood pressure measurements are represented as CDA observations and, when appropriate, these observations are associated with a substance administration of Losartan.

## CDA Model

The CDA UML model was created as an implementation model that is primarily based on two artifacts: (1) the CDA Refined Message Information Model (R-MIM) from HL7 and (2) the CDA XML Schema. This implementation model was developed to support the existing code generation and serialization mechanisms present in the Eclipse Modeling Framework (EMF). The model was imported into an EMF model and ultimately transformed into a set of Java classes. The Java classes in conjunction with a set of additional utility classes make up the base runtime API that can be used to produce, consume and validate instances of CDA.

## Template Modeling & Annotation

The template model is a domain-specific model that constrains the CDA model. Classes in a template model extend those in the CDA model. Constraints are modeled using directed associations, property redefinitions, and OCL expressions. The CDA Profile for UML is used to capture additional metadata needed during model transformation and at runtime. Annotations on template model elements including UML classes and properties are used to describe all core ontology variables and their possible parameterizations, each appearing at a unique location in the template model. Annotations are used to map between the core ontology and the CDA template model. After a template model has been created, it is transformed into an implementation model which leads to the generation of a domain-specific API for constructing and validating instances.

## Instance Generation Engine

The instance generation engine takes a data container that contains data values corresponding to variables in the cohort ontology as input and produces CDA document instances that conform to the template model. Using the ontology mappings, which were specified by the clinical domain expert at design-time, it resolves each variable in the data container to a corresponding variable in the core ontology. Annotations from the template model are then used to map core ontology variables to unique paths in the output tree and store data values in the leaves of the tree. Values that were specified as default or fixed in the template model such as template identifiers and coded attributes are also generated automatically.

## Results & Discussion

In the frame of Hypergenes, an FP7 European Commission funded project exploring the Essential Hypertension disease model, we had to deal with 18 historical cohort data sources with diverse clinical and environmental data. We chose HL7 v3 RIM meta-model and data types for data representation and CDA as our data model. Additionally we needed to apply a

template to constrain CDA to a document specialized for describing an Essential Hypertension Summary document (EH-CDA). Needless to say it was a perfect opportunity to put theory to test. In this section we will describe how the technology was used as well as illustrate a concrete example based on work done for Hypergenes project.

## Essential Hypertension Ontology

Hypergenes project assimilated clinical data from 18 cohort data sources. The harmonization process involved consulting with scientific experts in order to elucidate exact intention in each data element. The metadata was discussed at length in order to identify the list of variables, their meaning, variable associations, value ranges, and additional parameterization. The core ontology taxonomical structure was built based on data analysis of preliminary results and the macro-classes of intermediate phenotypes and environmental risk factors defined for Essential Hypertension. The core ontology was used as a reference for mapping the variables in each of the cohorts.

## Essential Hypertension Template Model

Once the metadata was fully accounted for, we created a template model hence constraining CDA to an Essential Hypertension Summary document. Figure 6 depicts a part of this model pertaining to a Blood Pressure Finding.

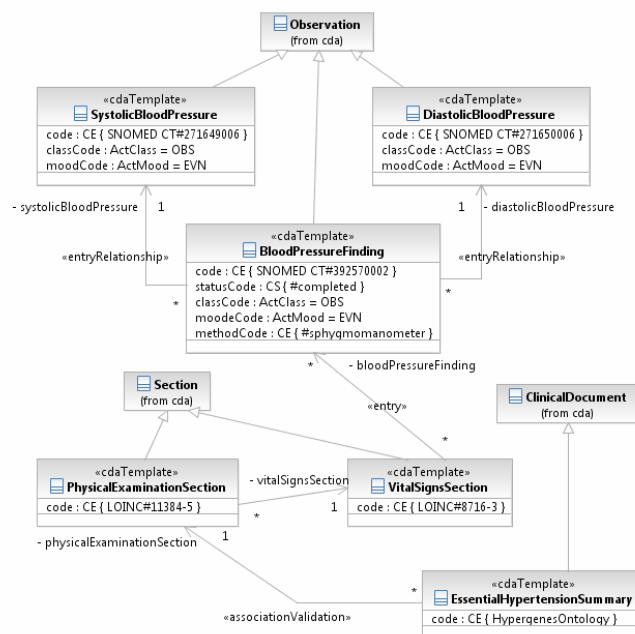


Figure 6 – Blood Pressure Observation in template model

The BloodPressureFinding class in the EH-CDA template model extends Observation class from the CDA model. The template identifier was specified in a property of the <<cdaTemplate>> stereotype. Additionally, the code attribute was used to capture metadata about the specific code in SNOMED-CT. This gives the template precise semantics from a clinical perspective. The directed associations in the diagram (e.g. VitalSignsSection to BloodPressureReading), which

were used to capture the required structure of the EH-CDA document, are converted into equivalent OCL constraints during the model-to-model transformation.

### Instance Generation for Essential Hypertension

Figure 7 depicts a CDA instance snippet of a Blood Pressure finding Observation.

```

<observation classCode="OBS" moodCode="EVN">
  <templateId root="2.16.840.1.113883.3.18.99.1.1.8.1"/>
  <code code="..." codeSystem="..." codeSysName="SNOMED CT" displayName="Blood pressure finding"/>
  <status code="completed"/>
  <methodCode displayName="sphygmomanometer"/>
  <entryRelationship typeCode="COMP">
    <observation classCode="OBS" moodCode="EVN">
      <templateId root="2.16.840.1.113883.3.18.99.1.1.8.2"/>
      <code code="..." codeSystem="..." codeSysName="SNOMED CT" displayName="Systolic BP"/>
      <value unit="mmHg" value="171" xsi:type="PQ"/>
    </observation>
  </entryRelationship>
  <entryRelationship typeCode="COMP">
    <observation classCode="OBS" moodCode="EVN">
      <templateId root="2.16.840.1.113883.3.18.99.1.1.8.3"/>
      <code code="..." codeSystem="..." codeSysName="SNOMED CT" displayName="Diastolic BP"/>
      <value unit="mmHg" value="109" xsi:type="PQ"/>
    </observation>
  </entryRelationship>
</observation>

```

Figure 7 – EH-CDA Blood Pressure finding Observation

Several things should be noted here. The CDA XML structure giving context and capturing data intra relationships explicitly; the use of RIM’s capability to describe the relationships accurately, thus Systolic BP is a component (typeCode=“COMP”) of Blood Pressure finding; and the use of standard healthcare terminologies, specifically SNOMED CT.

One of the major challenges of implementing the instance generation engine was in creating an algorithm that analyzes the annotations specified in the template model and traverses the model to generate path expressions. This is partly due to the fact that the template model is actually a logical model that hides some of the underlying structure of the base CDA model.

### Conclusion

In this paper we discussed a model-driven approach for integrating biomedical data using three complementary technologies. We used semantic technology in the form of an ontology definition language (namely OWL) to describe data elements of interest for a particular clinical research domain. We discussed the use of XML-based healthcare interoperability standards for clinical documents and the role they play in semantic interoperability across multiple data sources. Finally, we discussed the use of UML to bridge between the clinical domain expert and the healthcare interoperability expert and to facilitate generation of a runtime to produce conforming instances.

### Acknowledgments

The research leading to these results has received funding from the European Community’s Seventh Framework Program FP7/2007-2013 under grant agreement n° 201550.

### References

- [1] HL7 Reference Information Model, Health Level Seven, <http://www.hl7.org/v3ballot/html/infrastructure/rim/rim.htm>
- [2] Web Ontology Language, <http://www.w3.org/TR/owl-features/>
- [3] Unified Modeling Language at OMG, <http://www.omg.org/spec/UML/2.0/>
- [4] Object Constraining Language specification, <http://www.omg.org/technology/documents/formal/ocl.htm>
- [5] Hypergenes FP7 European Commission Project, <http://www.Hypergenes.eu/>
- [6] Stroetmann V. et al. “Semantic Interoperability for Better Health and Safer Healthcare”. SemanticHEALTH Project Report, 2009. Published by the European Commission, [http://ec.europa.eu/information\\_society/ehealth](http://ec.europa.eu/information_society/ehealth)
- [7] Heiler S: Semantic interoperability. ACM Computing Surveys 27(2):pp271-273, 1995.
- [8] Gold J. D., Ball M. J. “The Health Record Banking imperative: A conceptual model”, IBM Systems Journal, Vol 46, No 1, 2007
- [9] Bock B.J. et al. “The Data Warehouse as a Foundation for Population-Based Reference Intervals”. American Journal of Clinical Pathology, 120; pp662-670, 2003.
- [10] Dolin R. H. et al, “HL7 Clinical Document Architecture, Release 2”, JAMIA 2006;13:pp30-39
- [11] Ferranti et al, “The Clinical Document Architecture and the Continuity of Care Record: A Critical Analysis”, JAMIA 2006; 13:pp245-252.
- [12] Smith M. K., Welty C., McGuinness D. L. <http://www.w3.org/TR/owl-guide/>, 2004
- [13] McGonigle D., Mastrian K., “Nursing Informatics and the foundations of knowledge”. p.97 <http://nursing.jbpub.com/informatics>
- [14] Schultz S., Boeker M., Stenzhorn H., “How Granularity Issues Concern Biomedical Ontology Integration“. MIE, p. 863, 2008.
- [15] Golbreich C., Zhang S., Bodenreider O., “The foundational model of anatomy in OWL: Experience and perspectives”, Web Semantics: Science, Services and Agents on World Wide Web, Vol 4, Issue 3:pp181-195, 2006.
- [16] Farkash A. et al. “[Biomedical Data Integration – Capturing Similarities While Preserving Disparities](#)”, proceeding of IEEE EMBC 2009.
- [17] Carlson, D., “Semantic Models for XML Schema with UML Tooling,” proceeding of SWESE 2006.
- [18] Protégé, a [free, open source](#) ontology editor and knowledge-base framework. <http://protege.stanford.edu/>

[19]Heiler S: Semantic interoperability. ACM Computing Surveys 27(2):271-273, 1995.

**Address for correspondence**

Ariel Farkash  
IT for Healthcare & Life Sciences  
IBM Haifa Research Lab.  
E-mail: [arielf@il.ibm.com](mailto:arielf@il.ibm.com)