



UNIVERSITY
OF TRENTO - Italy



DIPARTIMENTO DI INGEGNERIA E SCIENZA DELL'INFORMAZIONE

– KNOWDIVE GROUP –

Fast Healthcare Interoperability Resources

Document Data:

- date -

Reference Persons:

Shaun McNaughton, Jacopo Mocellin, Sander Martins
Gonçalves, Zuhairia Ibnat

© 2020 University of Trento
Trento, Italy

KnowDive (internal) reports are for internal only use within the KnowDive Group. They describe preliminary or instrumental work which should not be disclosed outside the group. KnowDive reports cannot be mentioned or cited by documents which are not KnowDive reports. KnowDive reports are the result of the collaborative work of members of the KnowDive group. The people whose names are in this page cannot be taken to be the authors of this report, but only the people who can better provide detailed information about its contents. Official, citable material produced by the KnowDive group may take any of the official Academic forms, for instance: Master and PhD theses, DISI technical reports, papers in conferences and journals, or books.



Contents

1	Knowledge Graph Codebook	1
1.1	Knowledge Graph general description	1
1.2	Data level	1
1.2.1	Datasets general details	1
1.2.2	Datasets metadata documentation	1
1.3	Ontology level	1
1.3.1	Ontology general details	1
1.3.2	Ontology metadata documentation	1
1.4	Knowledge Graph Evaluation	2
2	Knowledge Graph Development Process	3
2.1	Scope Definition	3
2.2	Personas	4
2.3	Scenarios	5
2.4	Inception	6
2.4.1	CQs definition	6
2.4.2	Initial Datasets description	8
2.4.3	Datasets generation method	9
2.4.4	Datasets metadata documentation	10
2.4.5	Datasets collection process	13
2.4.6	Inception level evaluation	13
2.5	Informal Modeling	14
2.5.1	Schema level	14
2.5.2	Data level	17
2.5.3	Informal Modeling Evaluation	20
2.6	Formal Modelling	21
2.6.1	Schema level	21
2.6.2	Data level	24
2.6.3	Formal Modeling Evaluation	25
2.7	Data integration	25
2.7.1	Data integration operations and tool	25
2.7.2	Variance respect Formal Modeling datasets	26

Revision History:

Revision	Date	Author	Description of Changes
1.0	19/10/2020	Jacopo	Scope definition section written
1.1	21/10/2020	Jacopo	Inception section written
1.2	21/10/2020	Sander	Datasets generation method
1.3	21/10/2020	Shaun	General layout and editing
1.4	22/10/2020	Jacopo	Scope and inception Edits
1.5	22/10/2020	Zuhairia	Datasets Metadata Documentation written
1.6	22/10/2020	Shaun	Initial competency queries and scenarios added
1.7	23/10/2020	Sander	Metadata Synthea and Smart
1.8	23/10/2020	Zuhairia	Query edits
1.9	23/10/2020	Zuhairia	Synthea attribute table formation
2.1	05/11/2020	Jacopo	IM Schema level documentation draft
2.2	06/11/2020	Zuhairia	Metadata table formation
2.3	06/11/2020	Shaun	EER Diagram added
2.4	06/11/2020	Jacopo	IM Schema level documentation finalized+overall improvements
2.5	06/11/2020	Sander	Datasets management process
3.1	26/11/2020	Jacopo	Formal modelling schema section written
3.2	27/11/2020	Shaun	Added ontology visualisation, revised EER and lexical information
3.3	27/11/2020	Sander	Added Formal Modelling datasets management
3.4	27/11/2020	Sander	Added Variance respect Informal Modeling datasets

1 Knowledge Graph Codebook

The first of the two sections, in the current document, contains the codebook of the whole KG (Knowledge Graph), including the description of all the data and information that it contains.

1.1 Knowledge Graph general description

This sub section aims to give a general description of the KG, reporting:

- the context/domain in which the KG lives and works;
- *The Problem* the KG aims to solve;
- How the KG can solve *The Problem*

1.2 Data level

The data level section aims to describe in details the (final version of) datasets collected and managed by the KG, with a description of each variable involved.

1.2.1 Datasets general details

In this section are reported the metadata at datasets level, so the metadata regarding the sources, the authors, the collection methods, and so on.

1.2.2 Datasets metadata documentation

In this section are reported the metadata at dataset attribute level, through a description of each variable involved in the datasets collected, specifying the variable types, meanings, value-set (possible values), and every other meaningful variable information.

1.3 Ontology level

The ontology level section aims to describe the underlying KG ontology, through the description of its elements at each level, reporting so the language, conceptual and schema resources used within it.

1.3.1 Ontology general details

This first sub section of the ontology level description, report the general details such as authors, sources and the description of external ontology eventually adopted to generate the final one.

1.3.2 Ontology metadata documentation

In this section instead, are reported the more specific metadata describing the single elements of the ontology (terms, concepts, ETypes and relations).

1.4 Knowledge Graph Evaluation

In the final section of this first chapter, the KG Evaluation is reported. It aims to describe, through specific metrics, the quality of the overall KG on different aspect, like domain coverage, usability, domain representation, and other meaningful aspects.

2 Knowledge Graph Development Process

The second chapter of this document aims to describe, in a detailed way, the KG development process. The sections below describe each phase of the KG building project, reporting for each phase, the description of the datasets and their evolution respect the previous phases, the schema construction which will generate the KG ontology in the end, as well as the description of the procedures adopted to manage the data and finally achieve those results. Moreover for each phase is reported an evaluation section, which aims to evaluate the quality of the results achieved at the end of each phase.

2.1 Scope Definition

In this part of the project report, we provide the foundational specifications for the further development of our data integration project. In particular, we are going to outline what is the problem our solution is going to address, why we chose that problem, what is the context in which the problem exists in, and why the iTelos methodology is the most suitable technical approach to solving the problem.

EHRs (Electronic Health Records) are systems that allow to store all medical data concerning a patient in a digital format. While innovation could allow increased portability across different health organizations, current EHR implementations offer limited transferability to other systems because of several layers of obstacles. This problem is formalized as a problem of data interoperability.

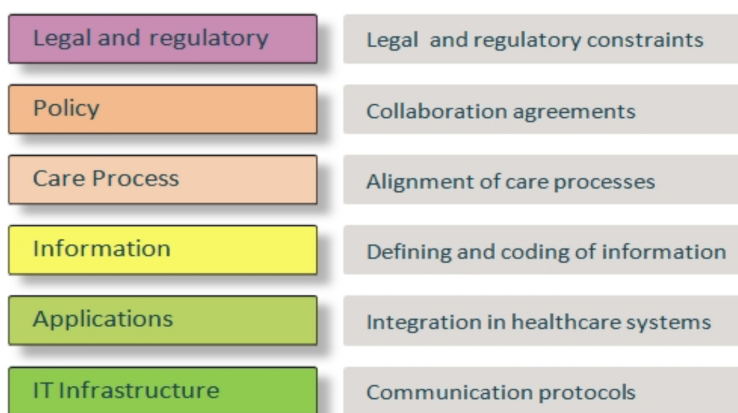


Figure 1: (from eHN's Refined eHealth European Interoperability Framework)

While several layers of complexity need to be addressed in order to achieve complete interoperability, we are going to focus our efforts to implementing a solution that aligns semantic differences across different EHR systems, and therefore our intervention can be ascribed to the “Information” layer of this scheme. This layer is comprised of all aspects of the data model, including coding terminology and the formatting of the medium for transportation of the information.

The problem of incompatible data representations is particularly taxing in Europe, where, on one hand, free circulation of EU citizens across the 27 member states is fostered by a friendly regulatory scenario, but, on the other hand, seamless circulation of citizen's health data has been hindered by the lack of complete interoperability.

In the current European medical context, the eHDSI (e-Health Digital Service Infrastructure) connects eHealth national contact points allowing them to exchange two sets of health data: patient summaries and ePrescriptions. This system serves as reference in the field of interoperable EHRs.

Recently, the European Commission produced and adopted a recommendation on a framework for the further development of a European EHR exchange format that integrates what already accomplished by the eHDSI project. This consists of a set of common technical specifications and principles that will enable citizens to securely access and exchange their health data across borders in the EU.

In the annex to the Commission Recommendation on a European Electronic Health Record exchange format, a total of five types of data are defined as target of a European cross-border exchange solution:

- Patient Summaries
- ePrescriptions/eDispensations
- Laboratory reports
- Medical images and reports
- Hospital discharge reports

In line with this proposed framework, while selecting the types of data that our data integration system should be able to interpret, we opted for restricting our focus to those mentioned above in order to adhere with the requirements of a realistic EHR solution existing in the current European context described.

In this problem space, the iTelos methodology proves to be the most appropriate because it allows the alignment of information that in our case, coming from different source datasets, is scattered and represented in a diversified manner.

2.2 Personas

As a top level view, a solution that could help foster data interoperability could benefit several stakeholders in the healthcare panorama:

- Healthcare Facilities: could reduce costs and inefficiencies by reducing the amount of testing needed
- Healthcare professionals: could save time by not having to enter information by hand in a system. Could become able to interpret information in different languages. Could see the history of a patient: reaction to medications, previous operations, allergies, and other intolerances
- Patients: enables to consult a specialist or receive emergency treatment in another EU country
- Researchers: Could enable to recruit a diverse population for clinical studies while being able to access accurate retrospective reporting information about the patients.

Among these, we chose to focus on developing a solution that could fit primarily the needs of patients and health-care professionals taken as a unique dyad of subjects cooperating to foster better outcomes in the therapeutic process

We will now illustrate some proposed personas and relative hypothetical use cases.

- Giannantonio, aged 35, is an Italian sales representative for an important clothing brand who is asked to travel frequently from one country to the other for work. His agenda is made up by appointments fixed with very little advance, which doesn't allow him to have much time to prepare luggage for his business trips or to bring extra paperwork along with him. He has good understanding of English, but not in a medical domain.
- Rita is an Italian woman, aged 73, who lives in a rural town near Caltanissetta. Once a year she travels to the Netherlands to visit her son and his family who emigrated there. She has no knowledge of foreign languages.
- Pedro is a Spanish man aged 45 who was recently diagnosed with Parkinson disease. He is travelling to Germany constantly in order to receive life-saving treatments. However, his clinical situation is complex since he suffers several co-occurring chronic conditions.
- Jean is a 52 year old, General Practitioner from France. He practices in a small clinic near Rennes, but is also working as a Forensic Medical Officer for the French police. While he loves both of his jobs, the additional stress from COVID-19 has been pressuring him to focus on a single job.

2.3 Scenarios

As the exchange of health information is an inherently complex and has widely varying needs, we use the personas and expand on them by creating a hypothetical scenarios to which having integrated health systems could be used.

Giannantonio: During one of his frequent trips to Greece, he is struck down by a heart attack. Not having any form of medical documentation with him, the medical team helping him needs to retrieve information elsewhere in order to know which treatments he might not respond well to. [Access to patient history]

Rita: During one of her visits to the Netherlands, Rita falls to the ground and breaks an arm. Even though her condition is not severe, the physician recommends close supervision by a colleague upon hospital discharge in order to avoid incorrect calcification of her fragile bones. Once back home, Rita is aware of the need to consult with her trusted physician in order to avoid complications, but would not be able to report to him any of the indications given by the Dutch doctor. On the other hand, her trusted doctor is unable to understand languages other than Italian. [Access to discharge reports data]

Pedro: Due to his complex medical conditions, the different medical teams who treat him in the two different countries need to be constantly updated on prior medications administered in order to inform their future clinical decisions. [Access to patient history]

Jean: Before seeing a new patient, Jean spends time reviewing that patient's medical history. He does this to better understand the patient at both a physical and emotional level. This in turn, allows him to mentally prepare for the consultation and to act in a professional manner. In doing so, he wishes to review at least any basic health information on the patient, for example age, blood type, and past medical history before seeing them. Now, he

wishes to access medical records from an Italian hospital to help diagnose a patient who has recently arrived in France and has a history of hospitalisations. [Access to patient history]

2.4 Inception

This section is dedicated to the Inception phase description. Here are reported the initial definitions for CQs (Competency Queries), initial datasets collected and the relative metadata. For each of those elements the procedures and the tools adopted to achieve the results, have to be reported in the sections below.

2.4.1 CQs definition

In order to create the competency queries, we use the scenarios and persona descriptions provided in the previous section, to come up with possible queries that a person could want from the system.

The scenarios provide a set of actions for which a competency query can be modelled from. These actions and the mapping to the relevant competence query for each scenario is given below.

In future iterations, we might consider shrinking our focus to a narrower set of pathologies (and therefore queries) since our current dataset provide a data about patients with a broad and diverse set of conditions that could prove to be too extensive to represent in the final DKG.

Scenario	Question	Action	Description
Giannantonio	1.1	Access to known allergies	A list of known allergies and intolerances is returned
Giannantonio	1.2	Access to history of circulatory diseases	A list of current and previous conditions is returned.
Giannantonio	1.3	Access to lab reports	Results from latest blood analysis are returned
Giannantonio	1.4	Access to lab reports	Complete results from all blood analyses are returned
Giannantonio	1.5	Access to specific lab data	Latest measurement of MCHC RBC Auto-mCn is returned
Rita	2.1	Access to patient name	Name and last name of the patient are returned
Rita	2.2	Access to patient language	Desired spoken language is displayed
Rita	2.4	Access to known conditions	The system returns whether the patient suffers osteoporosis
Rita	2.5	Access to conditions history, other language	Previous Complete list of conditions is shown, in Dutch
Pedro	3.1	Access to patient allergies	The system shows whether the patient is allergic to sulfonamides
Pedro	3.2	Access to patient allergies	The system returns the complete list of recorded allergies
Pedro	3.3	Patient allergies, in different language	The system returns a complete list of allergies, in Spanish
Pedro	3.4	List of encounters	The system returns
			by date and time

Scenario	Question	Action	Description
Pedro	3.4	Access to patient allergies	The system shows whether the patient is allergic to sulfonamides
Pedro	3.5	Access to lab reports metadata	Date and time of the latest blood exams is shown
Pedro	3.6	Data source	The database name from in which patient data is sourced is displayed
Jean	4.1	Access to demographics	The system returns the place of birth of the patient.
Jean	4.2	Immunizations history	The complete list of immunizations is shown.
Jean	4.3	Specific immunization	The system returns whether the patient has been vaccinated for Pneumococcal (PCV).
Jean	4.4	Specific immunization	The system returns whether the patient has been vaccinated for Hepatitis A (HepA).

2.4.2 Initial Datasets description

In this section are reported the metadata at datasets level involved in the Inception phase, so those metadata regarding the sources, the authors, the collection methods, and other meaningful information.

For our data integration process, we had to rely completely on synthetic data. The collection of genuine, real-world data is almost impossible because of the special privacy measures that are in place to protect this kind of personal data. This means that real-world data is often protected and needs some kind of authentication to be accessed. Nonetheless, we were able to find satisfying patient data generators by looking into the solutions developed to test other existing EHR systems. Specifically, we are going to provide a list of found datasets and their metadata here.

- EMRBOTS dataset (<http://www.emrbots.org>): The dataset was created by Uri Kartoun, PhD, in 2015. We are using the 100000-patients dataset, which comprises in total data about 100.000 patients, 361.760 admissions, and 107.535.387 lab observation. The dataset was generated according to the literature on patient

population studies in order to create a realistic set of observations. More details about the data can be found on table 1 below.

- Synthea patient generator DBs (<https://github.com/synthetichealth/synthea>): This tool was developed by the MITRE corporation and released freely under an Apache license. The datasets generated with this tool can include data about over 90 different conditions thanks to its modular structure. Data generated includes: Conditions, Allergies, Medications, Vaccinations, Observations/Vitals, Labs, Procedures, CarePlans, Primary Care Encounters, Emergency Room Encounters, and Symptom-Driven Encounters. This data is matched with Configuration-based statistics and demographics. The data can be formatted in unstructured formats or structured (JSON) ones. The latter also allows for formatting according to HL7 standards. The generator can be queried thanks to a Java-based interface.
- Smart on FHIR data generator (<https://github.com/smart-on-fhir/sample-patients>): This tool is provided as a test suite for SMART on FHIR applications by SMART itself. We are going to exploit it to generate structured data that fits the FHIR specifications about a population of patients.

[Ultimately, our goal is to coherently aggregate this data in order to form the population of patients that is going to be represented in our DKG.]

The possibilities are either to keep the datasets separate or to build a single dataset by merging the ones mentioned above, assigning to single patients from the EMRBOTS dataset (or a subset of it) data from the other datasets. We are aware that these datasets currently do not refer to the same population, this we are going to focus on retaining validity in the case of performing this task. In this case, the final result is going to be a single dataset containing very diverse and overlapping data on a set of patients. This data is going to be in different formats, therefore constituting the ground for usage of a DKG to align diversity.]

2.4.3 Datasets generation method

EMRBOTS

EMRBOTS provided the method used for the creation of the dataset. The process starts by setting the desired configuration of the population. This means that is possible to set the percentage of gender, marital status, major language, ethnicity, date of birth and income level.

Once this pre configuration is made the next step is the generation of the individuals. For each generated patient, will be created a record that will be its lifetime electronic medical record (EMR). This record is created by randomly assigning the number of admissions, the length of each admission, the result of laboratory exams and the code for the main symptoms of the patient according to the list of International Classification of Diseases (ICD-10-CM). Laboratory results means a list with the measurements of levels of sodium, creatinine, or platelet counts, for example.

Synthea

Synthea is based on the PADARSER framework. This framework assumes that access to EHR real data is not available and therefore relies on public data to compose the population to be used for the generation of synthetic EHR. These public data are composed by aggregation of health incident statistics, clinical practice guidelines (used to develop caremaps) and medical coding dictionaries. These public will then feed the information and knowledge

model in order to generate synthetic data that is composed by patient profiles, diseases and their prevalences, caremaps and coded narratives. Synthea uses Health Level-7 (HL7) the main codification system.

Variable and Category	Patients (n=100,000)
Mean age as of 1/1/2015, years (SD)	57.8 (17.3)
Gender (%)	
Female	52.0
Ethnicity (%)	
White	49.0
Asian	23.0
African American	15.0
Unknown	13.0
Mean number of admissions per patient, days (SD)	3.6 (1.5)
Mean length of stay (SD)	11.0 (5.2)
% Population with length of follow-up (years)	
0 - 9	13.1
10 - 15	9.3
> 15	77.6
Population below poverty (%)	21.6
Comorbidities; Prevalence (%)	
Malignant neoplasm	41.4
Rheumatoid arthritis	25.6
Diabetes (type I or II)	24.4
Renal complications	17.0
Coronary artery disease	7.0
Laboratory values (Mean; SD)	
Blood urea nitrogen (mg/dL)	17.5; 7.2
Platelets (k/cumm)	284.9; 95.3
Creatinine (mg/dL)	0.9; 0.2
Albumin (gm/dL)	4.2; 1.0
Lymphocytes (%)	25; 5.8

Table 1: Data setting - size 100.000 - EMRBOTS

2.4.4 Datasets metadata documentation

In this section are reported the new metadata at attribute level, describing the types, meanings, value-set, validity, and so on, of each variable within the initial dataset collected.

The dataset created by EMRBOTS has the following attributes:

- Patient ID: A unique number representing a particular patient
- Gender: The sex of a patient; either male or female
- Birth: The date of the birth of the patient
- Race: This describes the ethnicity of the patient. The categories are white, Asian African-American, unknown
- Marital Status: The categories of marital status are married, single, divorced, separated, widowed, and unknown

-
- Language: The preferable language of the patients are mentioned here
 - Population percentage below poverty: This category is represented in percentage
 - Admission ID: An admission ID for the patient
 - Admission start date: The date the patient was admitted
 - Admission end date: The date the patient was released
 - Primary diagnosis code: This is an ICD10 code for the primary disease diagnosis
 - Primary diagnosis description: The description of the diagnosis is mentioned
 - Lab name: This category includes the tests of the patients. The CBC tests are: White blood cell count, Red blood cell count, Haemoglobin, Hematocrit, Mean corpuscular volume, MCH, MCHC, RDW, Platelet count, Absolute Neutrophils, Absolute lymphocytes, Neutrophils, Lymphocytes, Monocytes, Eosonophils, Basophils. The metabolic tests are: Sodium, Potassium, Chloride, Carbon dioxide, Anion gap, glucose, BUN, creatinine, total protein, Albumin, Calcium, Bili total, AST/SGOT, ALT/SGPT, ALK PHOS. The urinalysis tests include: pH, specific gravity, red blood cells, white blood cells
 - Lab value: This represents the numeric value of the tests
 - Lab units: This attribute shows the unit of the lab tests
 - Lab date: The date of the lab tests is shown in this attribute

The dataset created by Synthea has the following attributes: Information related to the patient:

- Patient ID: A unique number representing a particular patient
- Name: Corresponding name of the patient
- Gender: The sex of a patient; either male or female
- Birth: The date of the birth of the patient
- Race: This describes the ethnicity of the patient. The categories are white, Asian African-American, unknown
- Marital Status: The categories of marital status are married, single, divorced, separated, widowed, and unknown
- Address: Information related to the residence of the person. It is composed by Country and City
- Language: The primary language spoken by the patient
- Education level: Categorical variable indicating the education level of the patient
- Identification: Documents of the patient, could be passport, driver's id or other document
- Allergy: Information if the patient has any kind of allergy of food or medication
- Diagnostic: Information related to the identified disease with its related code (SNOMED CT format)

-
- Treatment: The treatment given to the patient for treating its condition
 - Dosage: Information of the quantity of medicine was prescribed
 - Period: Information related to the beginning and the ending of the admissions
 - Immunizations: The code of the vaccine applied or NA in case it was not prescribed

Other variables related to the diagnostics or other information can be found in the repository:

DCAT (insert link here)

The dataset created by Smart has the following attributes:

- Patient ID: A unique number representing a particular patient
- Name: Corresponding name of the patient
- Birth: The date of the birth of the patient
- Diagnostic: Information related to the identified disease
- Treatment: The treatment given to the patient for treating its condition
- Period: Information related to the beginning and the ending of the admissions

2.4.5 Datasets collection process

The collection process section aims to describe the collection methods adopted to obtain the initial datasets in this phase, as well as the tools used.

The collection has proven to be challenging for reasons mentioned earlier: Health data constitutes sensitive information and is therefore always protected with special provisions. Access to real-world health data most of the times is granted only via an accreditation process (As is the case for the MIMIC dataset, including data about real-world intensive care unit stays). After clashing with this bureaucratic obstacles, we started looking into free-to-use, open solutions available on the internet. We were finally able to find resources given by communities active in both academia and industry. The community behind the development of current EMR solutions has implemented testing datasets or dataset-generation tools that we eventually opted for using in our solution.

2.4.6 Inception level evaluation

We initially struggled with the initial concept behind what we were required, first to find data sources that were relevant and secondly to create a set of personas for which we could then create queries from. However, we can evaluate our performance of this stage of this phase by considering the following specific KPIs:

1. Number of datasets found: The current number of data sources amounts to 3, which is in line with the given indications.
2. Diversity of data collected: The number of datasets creates sufficient diversity in the data gathered. In fact, data collected is diverse in format (structured vs unstructured) and in representations (e.g. measurement units)
3. Coherence of data with purpose: Given the found data sources represent information in the medical domain, we are reassured that it fits coherently with the purpose of integrating medical information.

2.5 Informal Modeling

This section is dedicated to the Informal Modeling phase description. The Section is divided in Schema and Data level in order to report the details of the elements involved in the generation of the schema, as well as the description of the datasets evolution in this phase. Moreover each specific section, one for each level, reports the difference between the elements defined in this phase and the definitions in the previous phase, analyzing in this way the variance in the different phases.

2.5.1 Schema level

The schema level in this phase reports the first informal definition of the ETypes and of the Enhanced Entity Relationship (EER) model constructed using them. We further refine the set of competence queries in order to guide the construction of the EER and subsequent data filtering.

2.5.1.1 ETypes and EER Model definition

This section reports an informal definition of the ETypes involved in the datasets collected in the previous phase. This section includes a list of metadata/attributes associated to each of the elements generated and how they are related to each other. At this stage we are focussed on the relationships between our core entities required for our previously identified competency queries.

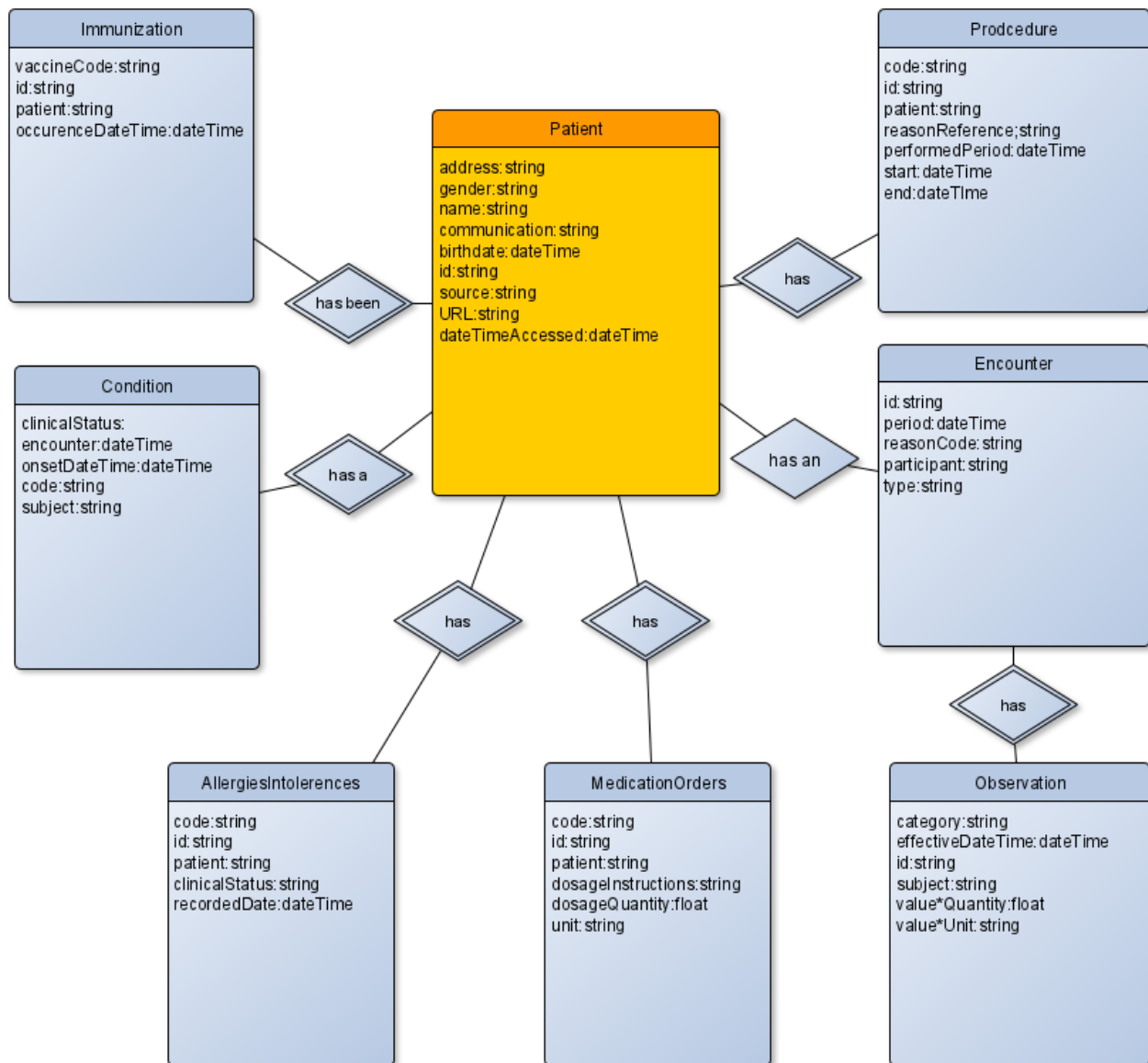


Figure 2: (FHIR Extended Entity Relationship Model)

In this first, informal version of our EER model, we have taken care of including the fundamental entities that could represent the variables in our selection of the data and to address our competency queries. Particularly, our initial core and contextual entities are:

- Patient:[core] This entity defines the concept of a human patient. In our schema, it is the concept for which all competency queries are drawn through.
- Immunization: [core] An entity that defines an vaccination instance. This is a core concept for queries that are focused on retrieving a patient's past medical history.
- Condition: [core] defines a current or past condition. This is a core concepts for queries that are focused on

retrieving a patient's past medical history.

- Allergies:[core] Defines current allergy or intolerance.This is a core concepts for queries that are focused on retrieving a patient's past medical history.
- MedicationOrder: [core] Defines a drug prescription administered by a physician. This is a core concept for queries involving past recommendations made by health professionals.
- Observation: [core] Defines lab observations and subsequent results and their associated units of measurements from testing.
- Encounter: [contextual] Defines an encounter with a health professional. The encounter may represent something like a consultation which may have an associated observation (e.g. a blood test).

Furthermore, we defined three more contextual entities in order to eventually represent the metadata about our sources of data within the EER. In this initial phase, we only limited ourselves to defining tentatively these entities and not the attributes contained.

- EMRBOTS patient
- Synthea Patient
- SMART patient

The contextual entities defined here are those that are required specifically for our particular competency queries. On the other hand, core Etypes are those collected from our reference ontology, that is the FHIR standard specification. Along with these, we sourced complemented our EER with Etypes from general-purpose reference ontologies in order to define our common entities and therefore ensure maximum compatibility with other schemas and ontologies. These entities have also been temporarily defined in order to account for the integration with other projects from the KDI course and will be expanded on demand.

- Doctor: defined as a Person entity sourced from schema.org (<https://schema.org/Person>)
- Hospital: as sourced from <https://schema.org/Hospital>. This entity inherits several properties from the entity Place.

Furthermore, the final version of our schema is going to include rdf label tags for each attribute and entity in order to represent a concept in a different languages. This is fundamental in order to satisfy the requirements from our competency queries.

2.5.1.2 Variance respect CQs definition

This section aims to define the variance between the schema elements produced in this phase, and the definition of the CQs reported in the previous phase. This a way to define the quality of the outcomes for the current phase as well as the alignment of the overall project development process.

In order to keep our work focused and manageable, we opted for restricting the entities under consideration to those related to the categories of Laboratory Reports and Patient Summaries that should be comprised in a EHR

system. We did not initially consider this option, since our goal was to define entities also regarding ePrescriptions/eDispensations, Medical images and reports and Hospital discharge reports.

For this reason, a great improvement on the previously defined CQs was needed, since we had to delete those not regarding strictly to the Patient Summaries and Laboratory Reports. At the same time, the total number of queries has more than tripled with respect to the previous phase, since a clearer understanding of the potential of the data collected has been gained.

2.5.2 Data level

The data level section in this phase reports the evolution of the datasets collected previously, reporting the metadata information for each new data, or new version of data, obtained.

2.5.2.1 Datasets management process

During the Informal Modeling phase the datasets collected in the previous phase are filtered and managed in order to obtain more suitable sets of data. In this section are described the procedures adopted to obtain that result.

After checking the variables contained on the datasets, was analysed the possibility of removing the ones that would not benefit the scope of our project and would not integrate with the initial set of entities defined. For each dataset will be presented the table correspondent to the variables that will be kept for our project.

Table 2: Filtered EMRBOTS

Dataset	Fields to be removed
EMRBOTS	Admission Id, Admission start date, Admission end date

Table 3: Filtered SMART

Dataset	Fields to be kept
SMART	Patient, Condition, Immunization, List, Observation

Table 4: Filtered SYNTHEA

Dataset	Keep only
SYNTHEA	“Resource type”: Patient, Allergy Intolerance, Immunization, Patient, Condition, Procedure, Diagnostic Report, Observation

2.5.2.2 Datasets metadata documentation

In this section is reported a list of new metadata in order to describe the modification performed on each datasets

and attribute, to achieve the new version of the datasets. The metadata from the EMRBOTS dataset is described below:

Field Name	Description
Title	Health records, Diagnosis, Testing, Test results
Description	All the health records were created to analyze realistic health related data. The data was generated to project primary diagnosis, required medical tests, lab observation of the patients in both macro and individual scale. The data focuses on the aforementioned attributes of the patient health record.
Category	Electronic Health Record
Keywords	Electronic Medical Records, Patient Simulation, FHIR, Electronic Health Record
Source of dataset	EMRBOTS
Link of the website	http://www.emrbots.org/terms.html
Total number of observations	100000
Spatial coverage	United States of America
Privacy Level	No privacy restriction
Authority level	View, use, copy or modify for non-commercial use
Last date of modification	06/11/20

Table: Metadata description of EMRBOTS dataset

The metadata description for the Synthea dataset is given below:

Field Name	Description
Title	Health records, Prevalent conditions, Diagnosis, Procedure, Testing, Test results, Observation
Description	All the health records were created to analyze realistic health related data. The data was generated to project primary diagnosis, required medical tests, procedures, lab observation, immunization of the patients in both macro and individual scale.
Category	Electronic Health Record
Keywords	Electronic Medical Records, Patient Simulation, FHIR, Electronic Health Record
Source of dataset	Synthea
Link of the website	https://github.com/synthetichealth/synthea
Total number of observations	2000
Spatial coverage	United Kingdom, Finland
Privacy Level	No privacy restriction
Authority level	View, use, copy or modify under the terms and conditions of Apache license 2.0
Last date of modification	06/11/20

Table: Metadata description of Synthea dataset

The metadata description for the SMART on FHIR dataset is given below:

Field Name	Description
Title	Health records, Prevalent conditions, Diagnosis, Procedure, Observation
Description	All the health records were created to analyze realistic health related data. The data was generated to project primary diagnosis, procedures, lab observation, immunization of the patients in both macro and individual scale.
Category	Electronic Health Record
Keywords	Electronic Medical Records, Patient Simulation, FHIR, Electronic Health Record
Source of dataset	SMART
Link of the website	https://smarthealthit.org/
Total number of observations	67
Spatial coverage	United States of America
Privacy Level	No privacy restriction
Authority level	View or modify the data using the SMART apps only
Last date of modification	06/11/20

Table: Metadata description of SMART on FHIR dataset

2.5.2.3 Variance respect Inception datasets

This section aims to define the variance between the data elements (datasets and attributes within them) produced in this phase, and the initial datasets collected in the previous phase. This a way to define the quality of the outcomes for the current phase as well as the alignment of the overall project development process. The data collection ended in the inception phase, therefore no new dataset has been added. Our already rich set of variables and medical concepts connected to the data suggested that adding more diversity would not have been required. At the same time, given each one of our datasets consists of a separate set of patients, collecting more information about one of these sets of patients would have been impossible given the constraints behind the collection of EHR data. For all these reasons we opted to work on the data already collected during the inception phase.

2.5.3 Informal Modeling Evaluation

The last section of the Informal Modeling phase report the evaluation of the outcomes obtained in this phase, through specific evaluation metrics.

For the time being, development of this section is paused as we await further indications.

What could be said for now is that the EER drafted adheres greatly to the concepts (entities) defined in the FHIR standard, given this is already a far-reaching framework specification that effectively represents medical data.

2.6 Formal Modelling

In this section, we describe the creation of the ontology in Protege, preparation of the data for integration and lastly the creation of the metadata in the DCAT standard.

2.6.1 Schema level

The schema level section in the current phase, reports the detailed description of the ontology generation.

2.6.1.1 Ontology definition

This section reports in details how the ontology is generated starting from the informal schema of the previous phase, which tools are used to do that, as well as usage of external ontology resources adopted to obtain the final KG ontology. Moreover a list of metadata is reported in this section, in order to describe all the elements of the ontology defined.

The informal EER was used as a starting point to define the entities constituting our final ontology, even though some changes (further detailed in the appropriate section of the report) had to be made with respect to the informally defined EER.

In order to define entities that would be in accordance with the FHIR standard (our reference ontology) we took advantage of the `fhir.ttl` RDF ontology officially released on [this documentation page](#). We used this resource just to perform a check of our understanding of how the FHIR standard works, since the actual definitions of etypes and attributes included in this ontology are not compatible with the very simplified representation we have devised during the informal modelling. We considered the possibility of pruning the FHIR ontology, but the fact that the official FHIR ontology defines all concepts as entities (while we define only few key entities and then add attributes on those).

Therefore, when it came to defining the actual ontology, we needed to work from scratch adding entities on top of a framework ontology provided by the UKC (iTelos repository). Through the provided KOS tool and related APIs we were able to first download this framework ontology, and then, after editing it, uploading it back to the UKC knowledge base.

The editing phase of this process has been performed using [Protégé](#). This tool allowed us to first define a set of entities and then to endow each of these of the desired attributes (Data properties). Below here follows a list of entities and relative attributes.

-
- Patient (GID "55936")
 - address: string
 - birthdate: dateTime
 - communication: string
 - gender: string
 - id: string
 - name: string
 - source: string
 - * URL: string
 - * dateTimeAccessed: dateTime
 - Condition (GID "77328")
 - ClinicalStatus: string
 - encounter: dateTime
 - onsetDateTime: dateTime
 - code: string
 - subject: string
 - Encounter (GID "110379")
 - id: string
 - period: dateTime
 - reasonCode: string
 - participant: string
 - type: string
 - Observation (GID "31728")
 - category: string
 - effectiveDateTime: dateTime
 - id: string
 - subject: string
 - valueQuantity: float
 - Procedure (GID "5269")
 - * unit: string
 - code: string
 - id: string
 - patient: string
 - reasonReference: string
 - performedPeriod: dateTime
 - * start: dateTime
 - * end: dateTime
 - Immunization (GID "4274")
 - vaccineCode: string
 - id: string
 - patient: string
 - occurrenceDateTime: dateTime
 - AllergyIntolerance (GID "77258")
 - code: string
 - id: string
 - patient: string
 - clinicalStatus: string
 - recordedDate: dateTime
 - MedicationOrder (GID "3440")
 - code: string
 - id: string
 - patient: string
 - dosageInstructions: string
 - dosageQuantity: float
 - * unit: string

The following are basic information about the metadata regarding the SKG generated:

- Author: FHIR team
- Version: 1.0
- Date Created: 26/11/2020 17:00 CET
- Class count: 32
- Data property count: 74

The first version of our ontology can be seen in Figure 3. It shows that everything is an entity and has a singular super entity. This is different to our EER model, because it is not linked through keys such as ID. For example, in our EER, an observation can only be integrated if it is linked to an encounter through an Encounter ID.

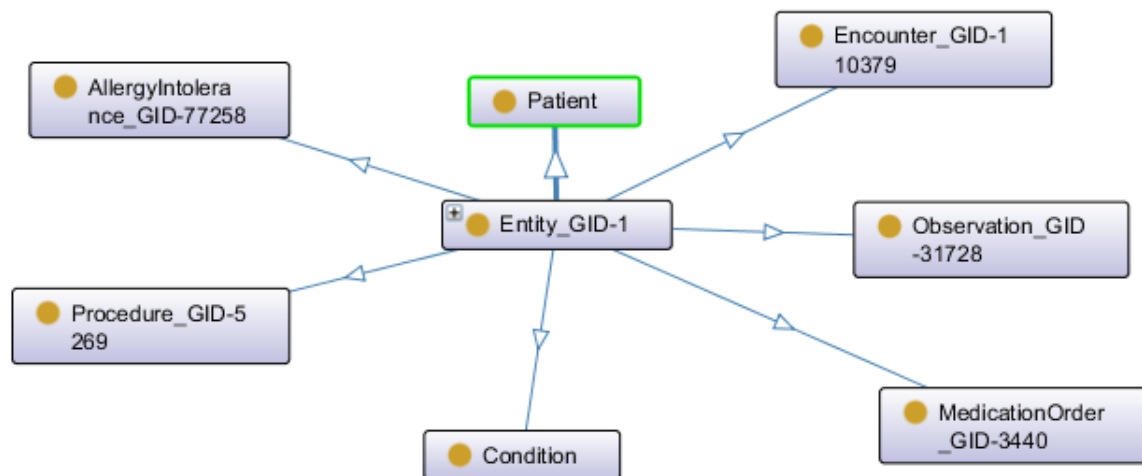


Figure 3: First version of our ontology

2.6.1.2 Variance respect to the EER Model

The EER was revised and refined according to the new ontology, this included implementing the indications we were given during the previous phase assessment, we needed to modify the total number of entities defined. The Patient entity has been unified, therefore no separate "EMRBOTS", "SMART", "SYNTHEA" patient entities were created. Instead, the defined patient entity type encompasses attributes that encode the source and some metadata information about the provenance of the medical record the patient belongs to.

Furthermore, Doctor, Hospital and Laboratory entities were not created since we followed the advice to focus closely on our project's scope.

2.6.1.3 Lexical information

As many of the entities have terms which can be ambiguous and our context requires interoperability between different countries with different languages, it is essential that the entities are clearly defined to a concept. We used concept definitions derived from the UKC (iTelos repository) to clearly identify the entities and it allows us to align our project with UKC concepts.

The following entities used the following definitions, where the format goes, concept, comment, and global ID;

- **Patient:** a person who requires medical care, 55936
- **Encounter:** come together, 110379
- **Observation:** facts learnt through observing, 31728
- **Condition:** the state of (good) health (especially in the phrases 'in condition' or 'in shape' or 'out of condition' or 'out of shape'), 77328
- **Procedure:** a procedure employed by medical or dental practitioners, 5269
- **Medication:** the act of treating with medicines or remedies, 3440
- **Allergy:** hypersensitivity to a particular allergen, 77258

2.6.2 Data level

2.6.2.1 Formal Modelling datasets management

In order to have a better way to align the datasets with the ontology's schema, was decided to unify the files for each data source. This unification wasn't straightforward since the files have unique features, each dataset case will be described bellow.

EMRBots

- Tool used: KNIME
- Input type:
- Output type: CSV
- 'LabsCorePopulatedTable.txt' - is a huge file. Not having enough memory to load it on Python (even using lazy load). It was decided to just filter this file using KNIME and keep it separated (even using KNIME the process of join was taking too long)
- the other files were filtered and joined into one output

Synthea

- Tool used: Python
- Input type: JSON

-
- Output type: JSON
 - The challenge in this case was to join the a large amount (1k) of JSON files into one while filtering. The result was one large JSON containing all the information.

SMART

- Tool used: Python
- Input type: XML
- Output type: JSON
- The XML file from this dataset was not formatted. Commonly used tools/libraries to convert it to JSON didn't work. It was necessary to search through all tags and elements, remove inconsistencies, filters the ones desired and add them to an empty JSON.

2.6.2.2 Datasets metadata documentation

In this section eventually new metadata information are added in order to describe the evolution of the datasets.

- The metadata generated during the SKG, LKG were explained by creating separate files. The link to the skg and dkg files is: (<https://github.com/UNITN-KDI-2020/kdi-fhir/tree/master/dataset/Formal>)
- The JSON file for the metadata documentation was created using the DCAT format.

2.6.2.3 Variance respect Informal Modeling datasets

After the filtering process, it was checked if the variables type in the dataset match the type of the variables of the ontology. For this step both types match. We can foresee though that on the data integration phase probably will be there a need to change how the variables are "allocated" on our dataset. One example could be the variable "Address", although it will remain as string, in our dataset the zip and street name are in different location. If the data integration phase needs to compile "Address" as one element, there will be need for changes.

2.6.3 Formal Modeling Evaluation

Similar to the informal modelling, evaluation metrics to this phase have not been released. This section will be modified as soon as this is completed.

2.7 Data integration

This section is dedicated to the Data Integration phase description.

2.7.1 Data integration operations and tool

This section is dedicated to the description of the usage of the data integration tool that allows to map the datasets generated and well formatted in the previous phases, with the final ontology generated. The last datasets adaptation performed using the tool, as well as the mapping operation are here detailed.

2.7.2 Variance respect Formal Modeling datasets

The last section of the data integration phase aims to describe the variance, analyzing the differences, between the datasets integrated with the ontology, in the data integration platform which contain the KG, and the datasets collected in the previous phase. This analysis can highlight the results of the operations performed during the final phase of the data integration process.