

# Mapping OMOP-CDM to RDF/OWL: Real World Data to the Semantic Paradigm

Achilleas CHYTAS<sup>a,b,1</sup>, Nick BASSILIADES<sup>b</sup>, Pantelis NATSIAVAS<sup>a</sup>

<sup>a</sup>*Institute of Applied Biosciences, Centre for Research and Technology Hellas*

<sup>b</sup>*School of Informatics, Aristotle University of Thessaloniki*

ORCID ID: Achilleas CHYTAS <https://orcid.org/0000-0001-8486-011X>

NICK BASSILIADES <https://orcid.org/0000-0001-6035-1038>

Pantelis NATSIAVAS <https://orcid.org/0000-0002-4061-9815>

**Abstract.** Real-world data (RWD) (e.g. data from Electronic Healthcare Records – EHRs, ePrescription systems, patient registries, etc.) gain increasing attention as they could support observational studies on a large scale. OHDSI is one of the most prominent initiatives aiming to support the harmonization of RWD and the development of relevant tools via the use of a common data model, OMOP-CDM. OMOP-CDM utilizes a certain degree of syntactic and semantic interoperability, still, it is based on a typical relational database format, and thus, the vision of a fully connected semantically enriched model is not fully realized. This work presents an effort to map the OMOP-CDM model and the data it hosts, to an ontological model using RDF to support the FAIRness of RWD and their interlinking with Linked Open Data (LOD) towards the vision of the Semantic Web.

**Keywords.** Semantic Web, Real-World Data, OMOP-CDM, Knowledge Graphs

## 1. Introduction

Real-world data (RWD) - e.g. data from Electronic Healthcare Records – EHRs, ePrescription systems, patient registries, etc. - are underutilized for uses other than its primary [1]. This can be attributed to semantic and syntactic interoperability issues, the fragmentation across institutions, as well as legal and ethical barriers. Additionally, it is common for RWD to face limitations regarding data quality and completeness due to errors and missing information that reduce their reliability and further impede their utilization. Privacy and security issues, such as compliance with regulations like HIPAA and GDPR [2] impose additional limitations on their usage, hindering data sharing among healthcare organizations even within the same country/region.

To address these challenges, Observational Health Data Sciences and Informatics (OHDSI) [3] is a global community maintaining and building open-source tools upon the OMOP Common Data Model (CDM). OMOP-CDM aspires to play a pivotal role as a de-facto standard supporting the harmonization of RWD to facilitate multi-site observational studies. Besides providing a common data framework, OMOP-CDM also is enriched with a certain degree of semantic interoperability which is achieved via the use of a plethora of interlinked vocabularies and widely accepted terminologies that

---

<sup>1</sup> Corresponding Author: Achilleas Chytas, email: [achytas@certh.gr](mailto:achytas@certh.gr)

cover a wide spectrum of all health-related concepts (e.g., diseases, laboratory exams, drugs, etc.) and even non-health-specific concepts (e.g. such as geographical location, ethnicities, etc.). However, the OMOP-CDM is built as a relational database schema and therefore it is not aligned with the Semantic Web vision as this was outlined by the use of RDF or OWL W3C recommendations.

There have been attempts to use RDF-based knowledge structures to support the OHDSI ecosystem, e.g. LAERTES [4] a knowledge base using RDF, or an effort to map the OMOP-CDM vocabularies to RDF [5]. This work presents an attempt to map OMOP-CDM to the RDF realm to bridge the gap between the world of OMOP-CDM and the Semantic Web ecosystem. Jean-Baptiste et. al [6] presented previously a related work for mapping the latest, but not currently adopted, version of OMOP-CDM (v6.0).

## 2. Methodology

The conversion of an OMOP-CDM database to a functional RDF Knowledge Graph (KG) was performed in a neutral/‘naïve’ manner, meaning the authors did not create a specialized ontological schema but rather kept the OMOP data structure and converted it to RDF as is. OMOP-CDM is a vastly generic schema and we argue that an ontological model in order to be functional, it should be created based on a more specific use-case scenario, (i.e. insurance claims, transmission patterns of infectious diseases, drug safety). As such, it was decided to create a plain representation of the database that can also be used by the rest of the OHDSI community as is. The idea is that this model could be semantically enriched to support specific study scenarios via the use of relevant ontologies integrated when needed, on top of the presented model.

R2RML is a language for expressing customized mappings from relational databases to RDF datasets [7]. R2RML is used to map the OMOP-CDM in RDF and convert relational data to a KG. The R2RML mappings are RDF graphs in Turtle syntax and can be used to map the relational OMOP-CDM data tables to relevant RDF concepts.

MIMIC-IV (Medical Information Mart for Intensive Care IV) is a large, and available upon-request relational database that with anonymized health data for over 40,000 Intensive Care Unit (ICU) patients [8] that is commonly used for exploring research questions and testing HC algorithms. This dataset was already converted to OMOP-CDM format<sup>2</sup> [9] and it was a prime candidate for the proposed data conversion pipeline. This dataset had certain changes from the current v5.4 due to updates in the CDM schema and included 100 patients. The most significant of those newly added tables are ‘*event*’ and ‘*event\_episode*’ which have been created to better model patients suffering from certain diseases, such as cancer, that require a specific and occasionally long-term treatment that is usually formulated in therapeutic protocols. Those changes were transferred to MIMIC-IV OMOP-CDM schema conversion to align it with the current CDM version.

A supplementary script written in R programming language was created to fine-tune the TURTLE file in order to create a proper and consistent instance of a KG. Indicatively:

- Locate and correct artifacts during mapping
- Set proper namespaces
- Set object properties

---

<sup>2</sup> <https://github.com/OHDSI/MIMIC>

- Set the domain and range of the object properties

Finally, in order to provide a quality assurance method to ensure proper conversion of OMOP-CDM data to an RDF KG, we built also a testing framework based on the rational of OHDSI HADES framework. HADES<sup>3</sup> is a suite of R packages that provide a variety of functionalities varying from statistical analysis, cohort creation, machine learning algorithms, functions that facilitate data sharing, and also tools that act as quality assurance which can assist in the validation the contents of an OMOP-CDM database. Along this line of thinking, a set of functions in R was created to validate the data conversion in RDF regarding both completeness and accuracy.

This testing framework validated the proposed transformation mechanism of MIMIC-IV data to RDF using via the execution of queries upon both the OMOP-CDM dataset and the KG and afterward a comparison of the results. These queries can be summarized as follows:

1. Count the number of patient records (rows in relational OMOP-CDM format, patient individuals in the case of the KG) and compare them.
2. For each table/class get the IDs of the rows/individuals and compare them.
3. For each non-empty OMOP-CDM relational table get a number of random rows, also get the corresponding individuals from the KG and compare each database field to the corresponding property.

The 3<sup>rd</sup> step of comparing a random number of rows instead of using the table as a whole was preferred since in a real-world scenario certain tables might be too large for a direct 1 to 1 comparison, although such functionality might be added on a later version as an option. Along with the tests, a reporting mechanism was developed that informs the user of the test results and can be used to identify the errors that might appear.

```
map:drug_exposure rr:logicalTable [rr:tableName "mimicredux.drug_exposure"];
rr:predicateObjectMap [
  rr:objectMap [
    rr:template "{drug_exposure_id}";
    rr:termType rr:Literal
  ];
  rr:predicate rdfs:label
], [
  rr:objectMap [
    rr:column "dose_unit_source_value";
    rr:datatype xsd:string
  ];
  rr:predicate <drug_exposure#dose_unit_source_value>
], [
  rr:objectMap [rr:template "http://omop_cdm/concept/{drug_concept_id}"];
  rr:predicate <drug_exposure#drug_concept_id>
], [
  rr:objectMap [rr:template "http://omop_cdm/person/{person_id}"];
  rr:predicate <drug_exposure#person_id>
],
...
```

**Figure 1.** Snippet of the R2RML mapping script for the OMOP-CDM table “drug\_exposure”

<sup>3</sup> <https://ohdsi.github.io/Hades/>

### 3. Results

The resulted RDF KG built using MIMIC-IV OMOP-CDM instance as the data source has over 6.4M axioms, 475K individuals, 22 classes and 62 object properties. Figure 2 presents the main classes of the ontological model and their interconnections.

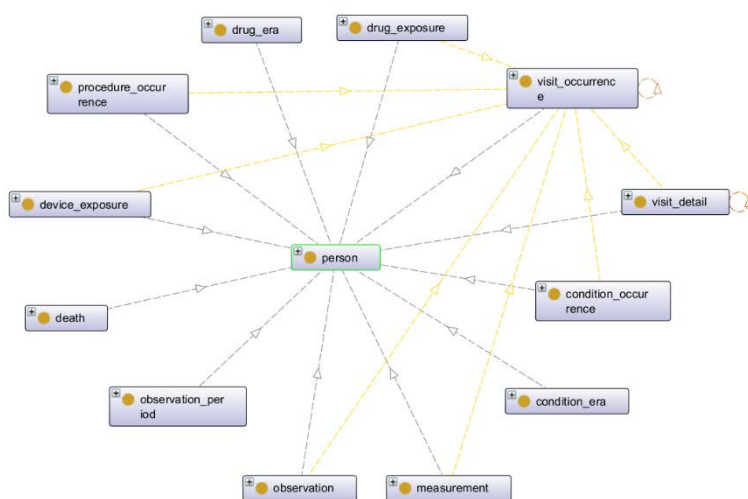


Figure 2. Important nodes of the ontological model

The first 2 tests resulted in a 100% compliance rate, while the 3<sup>rd</sup> test reported a few errors. Upon inspection of the reports, it was established that all errors on the testing dataset were due to data formatting reasons. More specifically, the errors were located in OMOP-CDM columns that contained verbatim the information from the initial data source (MIMIC-IV in our case), typically named 'X\_source\_value' where X stands for drug, condition, device, etc, and these inconsistencies occurred since R2RML tried to 'tidy up' the data fields. The error categories that were identified:

- Space padding, i.e. 'R89.4 ' instead of 'R89.4'
- Strings that depict numeric values that during the conversion were transformed to numeric and back to string, dropping zero padding in the process, i.e. '00160' to '160'.
- Encoding errors, i.e. for Greek or Chinese.

The R2RML mapping along with the validation scripts are accessible online in a Github repository<sup>4</sup>.

### 4. Discussion

In terms of potential limitations, we should mention that the MIMIC-IV relational OMOP-CDM dataset did not contain information in all the all tables or columns of the OMOP-CDM. This is a common occurrence as OMOP-CDM is an extensive model, and it is not unusual for a single data source to lack information for specific tables. This means that although the entirety of OMOP-CDM was mapped using R2RML, only the

<sup>4</sup> <https://github.com/achillec/omop2ttl>

tables contained data would be present in the final ontology as class, properties, and individuals and thus only these were tested in the context of this study.

Regarding future work, the first step is to use a dataset that contains data in all of the tables or use a synthetic dataset data in order to have a complete version of the RDF model. Then, other biomedical ontologies could be used to enrich the produced model and execute an observational study in order to evaluate potential benefits and challenges when using the RDF version of OMOP-CDM.

## 5. Conclusion

RWD modelled following the Semantic Web paradigm could in principle offer substantial benefits for healthcare research. By integrating RWD data with ontological models, interlinking with other knowledge structures could be facilitated and automatic reasoning could also be employed to support data processing.

By providing a semantic data model for a schema such as OMOP-CDM, with an ever-increasing adoption rate, will enable a large amount of RWD data to benefit from the interoperability and the extra computational methods that semantic technologies offer.

## References

- [1] J. Corrigan-Curay, L. Sacks, and J. Woodcock, "Real-World Evidence and Real-World Data for Evaluating Drug Safety and Effectiveness," *JAMA*, vol. 320, no. 9, pp. 867–868, Sep. 2018, doi: 10.1001/jama.2018.10136.
- [2] P. Voigt and A. Von Dem Bussche, *The EU General Data Protection Regulation (GDPR)*. Cham: Springer International Publishing, 2017. doi: 10.1007/978-3-319-57959-7.
- [3] G. Hripcsak *et al.*, "Observational Health Data Sciences and Informatics (OHDSI): Opportunities for Observational Researchers," *MEDINFO 2015: eHealth-enabled Health*, pp. 574–578, 2015, doi: 10.3233/978-1-61499-564-7-574.
- [4] R. D. Boyce *et al.*, "Large-scale adverse effects related to treatment evidence standardization (LAERTES): an open scalable system for linking pharmacovigilance evidence sources with clinical data," *Journal of Biomedical Semantics*, vol. 8, no. 1, p. 11, Mar. 2017, doi: 10.1186/s13326-017-0115-3.
- [5] J. M. Banda, "Fully connecting the Observational Health Data Science and Informatics (OHDSI) initiative with the world of linked open data," *Genomics Inform*, vol. 17, no. 2, p. e13, Jun. 2019, doi: 10.5808/GI.2019.17.2.e13.
- [6] J.-B. Lamy, A. Mouazer, K. Sedki, and R. Tsopra, "Translating the Observational Medical Outcomes Partnership - Common Data Model (OMOP-CDM) electronic health records to an OWL ontology," in *MEDINFO 2021 - 18th World Congress of Medical and Health Informatics*, online, France, Oct. 2021. Accessed: Mar. 15, 2024. [Online]. Available: <https://hal.science/hal-03479322>
- [7] "R2RML: RDB to RDF Mapping Language." Accessed: Dec. 08, 2023. [Online]. Available: <https://www.w3.org/TR/r2rml/>
- [8] "PhysioBank, PhysioToolkit, and PhysioNet | Circulation." Accessed: Dec. 08, 2023. [Online]. Available: <https://www.ahajournals.org/doi/full/10.1161/01.cir.101.23.e215>
- [9] M. Kallfelz *et al.*, "MIMIC-IV demo data in the OMOP Common Data Model." [object Object]. doi: 10.13026/P1F5-7X35.