



UNIVERSITY OF TARTU

# Transforming Estonian health data to the Observational Medical Outcomes Partnership (OMOP) Common Data Model: lessons learned



**Sulev  
Reisberg, Ph.D.**



# Estonia & University of Tartu

- Small Eastern European country
- Member of European Union and NATO
- Population 1.4 M
- “Estonia is the world's most digitally advanced society” (WIRED)
- University of Tartu - the oldest and largest university in Estonia



# Research group of health informatics



# Our recent paper




JAMIA Open, 2023, 6(4), ooad100  
<https://doi.org/10.1093/jamiaopen/ooad100>  
Research and Applications

AMIA  
INFORMATICS PROFESSIONALS. LEADING THE WAY.

OXFORD

## Research and Applications

# Transforming Estonian health data to the Observational Medical Outcomes Partnership (OMOP) Common Data Model: lessons learned

Marek Oja , PhD<sup>1,\*</sup>, Sirli Tamm, MSc<sup>1</sup>, Kerli Mooses, PhD<sup>1</sup>, Maarja Pajusalu, MSc<sup>1</sup>, Harry-Anton Talvik, MSc<sup>1,2</sup>, Anne Ott, MSc<sup>1</sup>, Marianna Laht, MD<sup>1</sup>, Maria Malk, MSc<sup>1</sup>, Marcus Lõo, MSc<sup>1</sup>, Johannes Holm, MSc<sup>1</sup>, Markus Haug, MSc<sup>1</sup>, Hendrik Šuvalov, MSc<sup>1</sup>, Dage Särg, MSc<sup>1,2</sup>, Jaak Vilo , PhD<sup>1,2</sup>, Sven Laur, PhD<sup>1</sup>, Raivo Kolde, PhD<sup>1</sup>, Sulev Reisberg , PhD<sup>1,2</sup>

<sup>1</sup>Institute of Computer Science, University of Tartu, 51009 Tartu, Estonia, <sup>2</sup>STACC, 51009 Tartu, Estonia

\*Corresponding author: Marek Oja, PhD, Institute of Computer Science, University of Tartu, Narva mnt 18, 51009 Tartu, Estonia (marek.oja@ut.ee)

**Author Contributions:** Dr M. Oja and S. Tamm are considered co-first authors and Dr R Kolde and Dr S Reisberg are considered co-last authors of this work. In addition, they had full access to all the data in the study and take responsibility for the integrity of data and accuracy of the data analysis.

### Abstract

**Objective:** To describe the reusable transformation process of electronic health records (EHR), claims, and prescriptions data into Observational Medical Outcome Partnership (OMOP) Common Data Model (CDM), together with challenges faced and solutions implemented.

**Materials and Methods:** We used Estonian national health databases that store almost all residents' claims, prescriptions, and EHR records. To develop and demonstrate the transformation process of Estonian health data to OMOP CDM, we used a 10% random sample of the Estonian population ( $n = 150\,824$  patients) from 2012 to 2019 (MAITT dataset). For the sample, complete information from all 3 databases was converted to OMOP CDM version 5.3. The validation was performed using open-source tools.

**Results:** In total, we transformed over 100 million entries to standard concepts using standard OMOP vocabularies with the average mapping rate 95%. For conditions, observations, drugs, and measurements, the mapping rate was over 90%. In most cases, SNOMED Clinical Terms were used as the target vocabulary.

**Discussion:** During the transformation process, we encountered several challenges, which are described in detail with concrete examples and solutions.

**Conclusion:** For a representative 10% random sample, we successfully transferred complete records from 3 national health databases to OMOP CDM and created a reusable transformation process. Our work helps future researchers to transform linked databases into OMOP CDM more efficiently, ultimately leading to better real-world evidence.

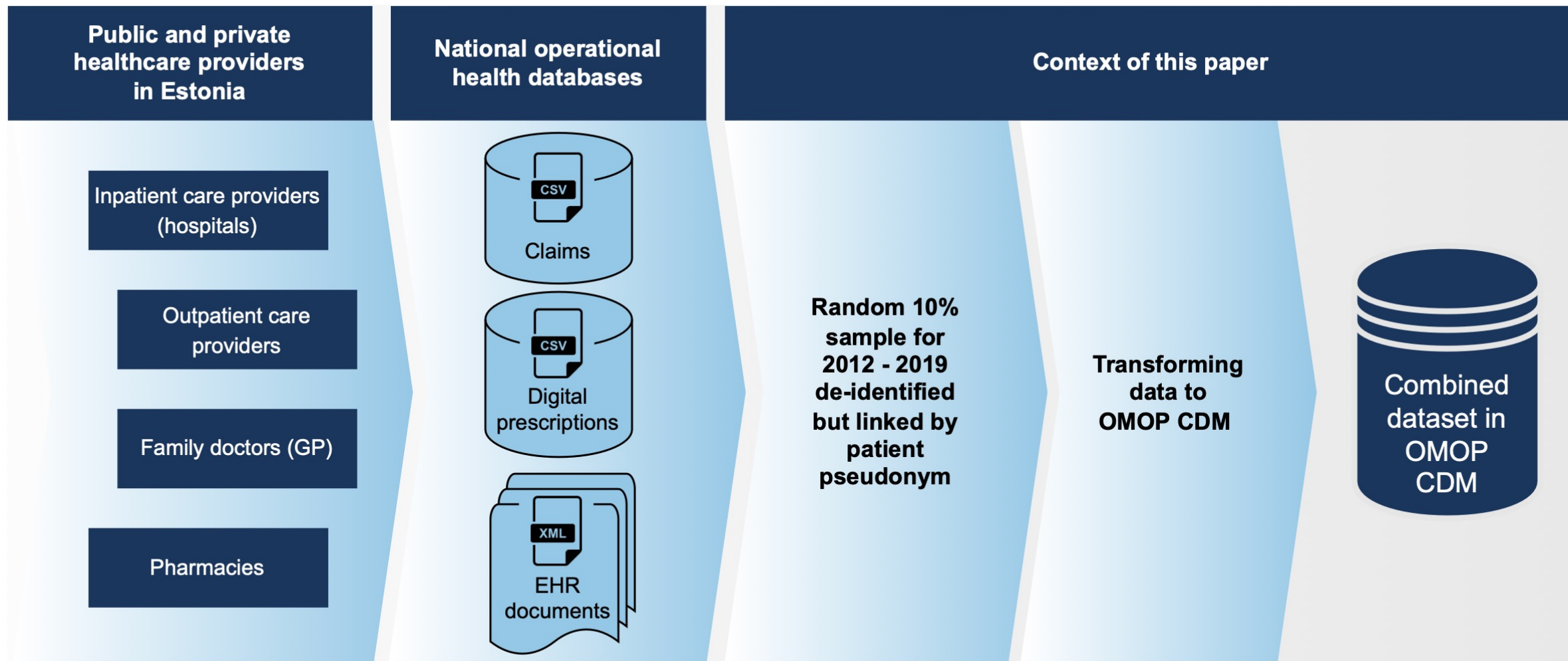
### Lay Summary

Health data can be found in various sources and formats, making it challenging for researchers. To address this issue, one possible approach is

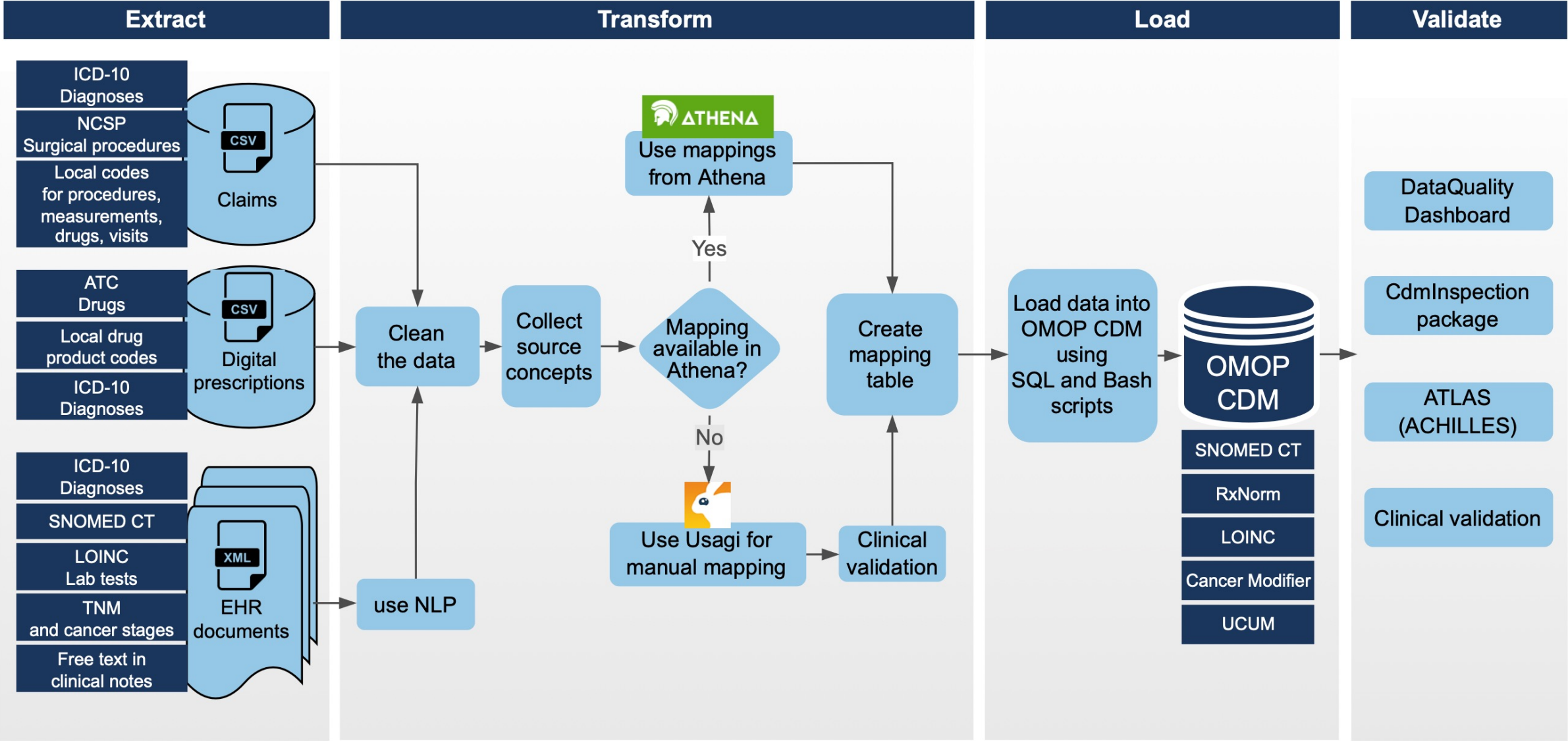


<http://doi.org/10.1093/jamiaopen/ooad100>

# Data sources & context of the paper



# Overview of ETL



## Reviewer 2

“ Starting with just 10% makes sense, but it'd be great to have stats for the other 90% too.

- Extensive table of challenges + examples + solutions
- Dataset used in a number of health studies
- A separate paper will be submitted that focuses on data extraction from source documents



**Table 4.** Main challenges and solutions of the current work.

**Challenge**

The same health event is represented in several source datasets without a clear link between them, potentially leading to duplicates.

No clear guidelines for choosing target vocabulary when multiple standard OMOP vocabularies are available. Additionally, there is no roadmap indicating which standard vocabularies may no longer be considered standard for OMOP CDM in the near future.

Hard to keep manual mapping files up to date as the standard target concepts change over time when updating the vocabularies.

**Example**

The same diagnosis code for a patient may be recorded in an EHR, claim, and prescription files. However, it may be difficult to link these documents to a single event due to the absence of a unique identifier for the case.

Physician Current Procedural Terminology Fourth Edition (CPT4) and SNOMED CT are both standard OMOP vocabularies for procedures; similarly, LOINC and SNOMED CT are for lab tests. The National Cancer Institute Thesaurus (NCIt) was a standard OMOP vocabulary at the beginning of our study, but not standard anymore.

Local code "9124," which is used for vaccination against diphtheria and tetanus, was mapped to SNOMED CT code "73152006" (administration of diphtheria and tetanus vaccine). That target concept changed from standard to nonstandard at some point in time. Thus, we had to remap it to the concept code "1657590" from RxNorm vocabulary (diphtheria toxoid vaccine, inactivated/tetanus toxoid vaccine, inactivated injection).

Atypical squamous cells of undetermined significance (ASC-US) result of the Papanicolaou test recorded in our datasets by SNOMED code 73152006. SNOMED codes

**Solution**

Transform each record as they are (even if duplicates) but add the provenance information to the record so one can use it when making cohorts.

Use the target vocabulary you are more familiar with. Keep in mind that what constitutes a standard OMOP vocabulary may change over time.

Whenever updating the vocabulary, recheck the mappings in Usagi before running the transformation. Usagi automatically creates the list of nonstandard mappings so one can fix them before the actual data transformation.

When working with historical codes always check the most recent target code for this event to reuse the same code.

# Highlights



# THANK YOU!

<https://health-informatics.cs.ut.ee>



ut.ee



info@ut.ee



tartuylkool  
tartuuniversity



unitartu  
unitartuscience  
unitartutiksu

