



# European Symposium Review

**OHDSI Community Call**  
**June 28, 2022 • 11 am ET**



**OHDSI**  
OBSERVATIONAL HEALTH DATA SCIENCES AND INFORMATICS





# Upcoming OHDSI Community Calls

Date	Topic
July 5	NO MEETING
July 12	New Adopter Introductions and Q&A
July 19	Workgroup Updates
July 26	CDM Update Process



# Upcoming OHDSI Community Calls

Date	Topic
July 5	NO MEETING
July 12	New Adopter Introductions and Q&A
July 19	Workgroup Updates
July 26	CDM Update Process



# July 12: New Adopters & Community Members

**Our July 12 Community Call will be focused on new adopters of the OMOP CDM or new members of the OHDSI community.**

We are welcoming people to introduce themselves, share why they have joined the community and what impact they hope to make, and also ask a question to the broader community (if you wish). If you would like to take part in this event, please fill out this form to help us plan the session: <https://bit.ly/3A7JNkV>

Form in chat and on community calls page





# Three Stages of The Journey

**Where Have We Been?**

**Where Are We Now?**

**Where Are We Going?**





# A Record-Setting Submission Year!



Thank you to everybody who submitted brief reports to join our #OHDSI2022 Collaborator Showcase. We had a record amount (**more than 130!**) of submissions for poster presentations, software demos and oral presentations for the 2022 OHDSI Symposium, which will be held Oct. 14-16 in Bethesda, Md.

The scientific committee meets this week to begin the process of reviewing all submissions, and **selected presenters will be notified by August 1.**



# OHDSI Shoutouts!



**Any shoutouts from the community? Please share and help promote and celebrate OHDSI work!**

Have a study published? Please send to [sachson@ohdsi.org](mailto:sachson@ohdsi.org) so we can share during this call and on our social channels.  
Let's work together to promote the collaborative work happening in OHDSI!





# Three Stages of The Journey

**Where Have We Been?**

**Where Are We Now?**

**Where Are We Going?**





# Upcoming Workgroup Calls



Date	Time (ET)	Meeting
Tuesday	12 pm	Common Data Model Vocabulary Subgroup
Tuesday	3 pm	OMOP CDM Oncology Outreach/Research Subgroup
Wednesday	11 am	Open-Source Community
Wednesday	12 pm	FHIR and OMOP Terminologies Subgroup
Wednesday	7 pm	Medical Imaging
Thursday	10 am	Data Quality Dashboard
Thursday	12 pm	FHIR and OMOP Oncology Subgroup
Friday	9 am	GIS – Geographic Information Systems
Tuesday	10 am	Common Data Model

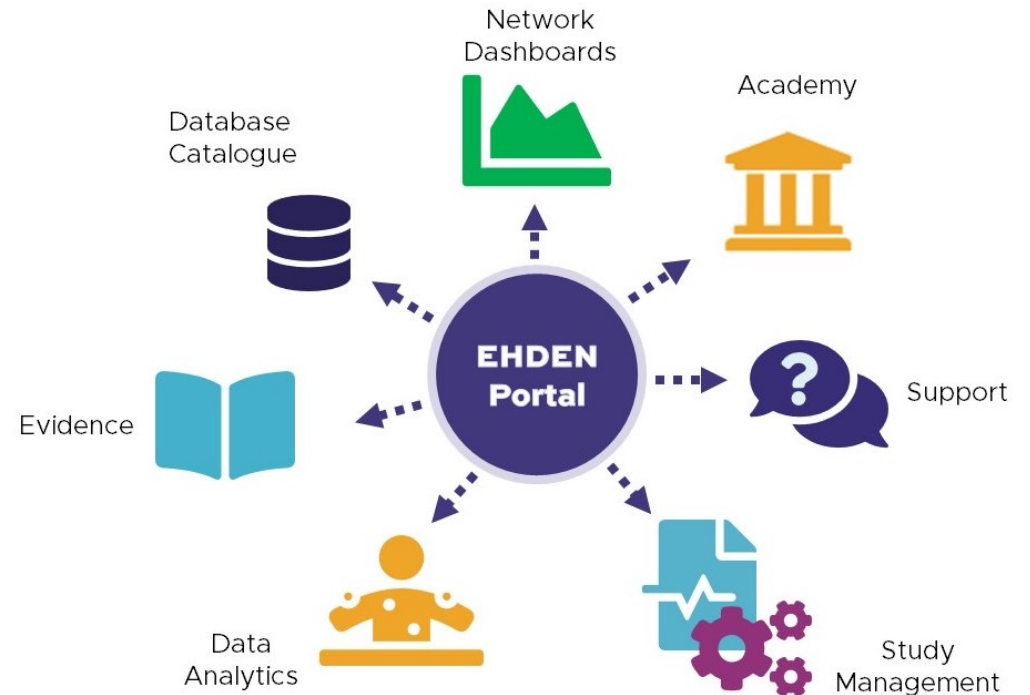
[www.ohdsi.org/upcoming-working-group-calls](http://www.ohdsi.org/upcoming-working-group-calls)





# EHDEN Portal/Data Catalogue

The EHDEN Portal, which provides free access to the research community, was launched at the 2022 OHDSI European Symposium. The Portal includes a Data Partner Catalogue (140 partners, >500M anonymous patient records) and Feasibility Dashboards that support data discoverability (findable under Findable, Accessible, Interoperable and Reusable (FAIR) principles).



[www.ohdsi.org/ohdsi-news-updates/](http://www.ohdsi.org/ohdsi-news-updates/)



# OHDSI EHR Data Survey



**MPhilofsky** Melanie Philofsky

6h

Hello friends with EHR data,

One of the Healthcare Systems group's objectives this year is "To provide support for transforming source EHR data to the CDM". Currently, we provide support through answering questions on the forums and during our regularly scheduled work group meetings. Another product we would like to provide to the community is a central repository of different OMOP sites, their underlying EHR system, and attributes. This will allow new OHDSI collaborators to find and reach out to sites with similar infrastructure, EHR systems, and/or research goals. Participating in this survey does NOT commit you to being a mentor, providing your ETL script, or even answering your email. However, we hope you embrace the spirit of our open source community and contribute to the cause. We all learn as we OMOP our data. I've been very active in the OHDSI community and digging deep into EHR data for 8 years, and I still learn something new every day. But I think all persons in any field of science continue to learn because science is continually evolving. Here's the [link](#) <sup>1</sup> to the google form.

    Reply



# Job Openings

Professor **Peter Rijnbeek** announced an opening for an epidemiologist to work with his team at Erasmus MC.

This position will be responsible for all aspects of observational research including protocol writing, input in the statistical analysis plan, study execution, interpretation of results and report/manuscript writing.

The application deadline is July 8, 2022.



## Epidemiologist

Published	Deadline	Location
9 Jun	7 Jul	Rotterdam



### JOB DESCRIPTION

This research will be performed in close collaboration with the [Observational Health Data Sciences and Informatics \(OHDSI\)](#) initiative, which is a global, multi-stakeholder, interdisciplinary collaborative to bring out the value of health data through large-scale analytics, and the EU-sponsored [European Health Data and Evidence Network \(EHDEN\)](#) which develop frameworks to generate reliable real-world evidence.

In your function as Epidemiologist you will be responsible for all aspects of observational research including protocol writing, input in the statistical analysis plan, study execution, interpretation of results and report/manuscript writing.



# Job Openings

**Odysseus Data Services** recently announced two openings, one for an epidemiologist and one for a data scientist.

Check out the links on the community calls page or reach out to a member of the Odysseus team to learn more!

Odysseus Data Services (Odysseus) has an exciting opening for an **Epidemiologist**. This role will be responsible for supporting the development, maintaining, and troubleshooting of the cutting-edge distributed solutions in the Real-World Evidence (RWE) area, utilized by the researchers in Pharmaceutical, Healthcare and Payer industries. Odysseus is looking for a self-driven individual who can hit the ground running, quick learner and wants to be a part of our dynamic global team.

## Responsibilities

- Lead and contribute to the design of observational database analysis, including authoring protocol, reviewing and providing relevant epidemiological and project-specific comments to statistical analysis plans and analysis output
- Participate in the design and development of standardized analytic tools to generate reliable and reproducible evidence in a network of observational data
- Contribute to the execution of observational database analyses using standardized analytical tools and writing statistical packages
- Contribute to the dissemination of scientific information through technical reports and publications in peer-reviewed literature.
- Work closely with healthcare and pharmaceutical customers to identify their needs
- Contribute to the development of complex phenotypes using advanced analytic approaches (i.e. machine learning, incorporating unstructured data sources using NLP, etc.)

## Qualifications

- Graduate degree (MS, PhD, MD, etc) in epidemiology, biostatistics, pharmacy, public health or related clinical discipline plus two years' experience in observational research. PhD preferred
- Experience in designing and conducting healthcare studies and in development and applications of advanced analytics solutions
- Strong epidemiology and biostatistics background
- Experience using OHDSI tools and analytical methods is a big plus

Odysseus Data Services (Odysseus) has an exciting opening for a **Healthcare/Clinical Data Scientist**. This role will be responsible for supporting the development, maintaining, and troubleshooting of the cutting-edge distributed solutions in the Real-World Evidence (RWE) area, utilized by the researchers in Pharmaceutical, Healthcare and Payer industries. Odysseus is looking for a self-driven individual who can hit the ground running, quick learner and wants to be a part of our dynamic global team.

## Responsibilities

- Lead and contribute to the design, development and documentation of standardized analytic tools that will be executed against a network of observational data
- Lead the execution of observational database analyses using standardized analytical tools and writing statistical packages
- Provide technical support for the data and analysis infrastructure and scientific support
- Contribute to writing of protocols and statistical analysis plans, methods development, conduct of simulation studies and statistical/mathematical modeling studies
- Lead and contribute to the development of complex phenotypes using advanced analytic approaches (i.e. machine learning, incorporating unstructured data sources using NLP, etc)
- Contribute to the dissemination of scientific information through technical reports and publications in peer-reviewed literature.
- Lead and contribute to the development of novel analytic tools and techniques to leverage the EHR data for rapid, reliable and reproducible evidence generation



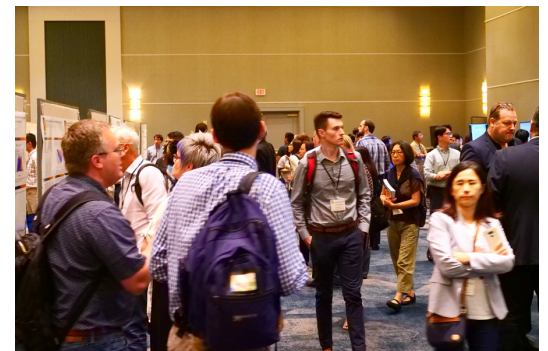


# 2022 OHDSI Symposium

Registration is OPEN for  
**#OHDSI2022!**

The 2022 OHDSI Symposium  
will be held Oct. 14-16 at the  
Bethesda North Marriott Hotel  
& Conference Center.

[www.ohdsi.org/ohdsi2022symposium](http://www.ohdsi.org/ohdsi2022symposium)







# An Introductory Journey From Data To Evidence

OHDSI2022 Tutorial • Saturday, Oct. 15 • Bethesda, Md.



**The OHDSI Journey:  
Where Are We Going?**

**Patrick Ryan**



**OMOP Common Data  
Model and Vocabulary**

**Clair Blacketer**



**ETL – A Source Database  
Into OMOP CDM**

**Melanie Philofsky**



**Creating Cohort  
Definitions**

**Asieh Golozar**



**Phenotype Evaluations**

**Gowtham Rao**



**Characterization**

**Kristin Kostka**



**Estimation**

**Martijn Schuemie**



**Prediction**

**Jenna Reps**



**The OHDSI Journey: Where  
Do We Go From Here?**

**George Hripcsak**



# Workgroup Activities

Saturday, Oct. 15, and Sunday, Oct. 16

Saturday, Oct 15					
Start Time (ET)	End Time (ET)				
800	900	Tutorial	HADES Hack-a-thon: Part 1	Oncology WG	FHIR-OMOP: Terminologies Subgroup, Part 1
900	1000				FHIR-OMOP: Increasing the Value of Data Through a Rich Set of Attributes
1000	1100				
1100	1200				
1200	1300		Lunch	Lunch	Lunch
1300	1400		Methods Research (PLE/PLP)	Oncology WG (continued)	FHIR-OMOP: Data Model Harmonization Subgroup
1400	1500			Natural Language Processing	FHIR-OMOP: Oncology Subgroup
1500	1600				
1600	1700				
1700	1800				FHIR-OMOP: Terminologies Subgroup, Part 2
1800	1900				
Sunday, Oct 16					
800	900	All-Hands Workgroup Meeting			
900	1000				
1000	1100				
1100	1200				
1200	1300	Lunch		Lunch	Lunch
1300	1400	Phenotype Evaluation	HADES Hack-a-thon: Part 2	Education	CDM and Data Quality
1400	1500			Health Equity	
1500	1600				
1600	1700				



# #OHDSISocialShowcase This Week

**Why predicting risk can't identify 'risk factors': empirical assessment of model stability in machine learning across observational health databases**

▲ PRESENTER: Aniek Markus  
(a.markus@erasmusmc.nl)  
CO-AUTHORS: Peter R. Rijnbeek,  
Jenna M. Reps

## INTRODUCTION:

- Some researchers incorrectly interpret prediction models:
  - Interpreting selected variables as factors that cause the outcome
  - Using selected variables for 'risk factor' detection (i.e. to identify variables associated with the outcome)
- We illustrate potential issues by investigating the stability of >450 prediction models in a large-scale experiment, investigating model changes across databases (care settings) and phenotype definitions.

## METHODS:

- We developed models using LASSO logistic regression for nine prediction tasks: predicting nine COVID-19 vaccine outcomes of interest (O) identified by the U.S. Food and Drug Administration for the general population (T) in the next 1 year (TAR).
- Measure model stability:
  - Q1: How many variables are selected across models?
  - Q2: Are the same or different variables included across models?
  - Q3: Is the direction of the effect of variables the same across models?

## RESULTS:

- Q1: A higher number of outcome cases generally leads to more variables being selected using (Fig 3).
- Q2: Overall model stability was poor, slightly better for top (i.e. most important) variables (Fig 4). The impact of different target/outcome phenotype definitions was limited, but the top 10 variables differed across databases (Fig 1).
- Q3: The sign of the coefficient can vary greatly even for the top variables (Fig 2), less selected variables seem more likely to switch sign.



Be careful interpreting prediction models as the identified 'risk factors' appear to depend on study design choices.

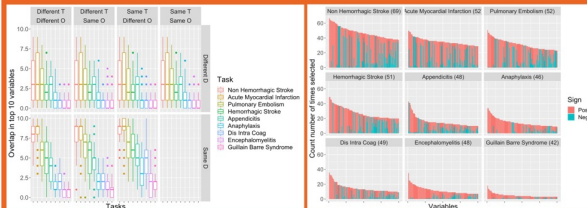


Figure 1. Overlap in the top 10 variables as defined by counting the number of common variables between each pair of models for same/different database (D), target population definition (T), outcome definition (O) across models

**TAKE AWAYS:**

- There is substantial variation in the selected variables across models.
- Different databases lead to different 'risk factors'.
- Interpreting the effect of 'risk factors' is problematic as the sign can differ across models.
- We recommend investigating model robustness across settings or using other techniques for 'risk factor' detection (e.g. univariate analysis).

T	O	Databases
General population	Acute myocardial infarction, Anaphylaxis, Appendicitis, Disseminated intravascular coagulation, Encephalomyelitis, Guillain-Barré syndrome, Hemorrhagic stroke, Non-hemorrhagic stroke, Pulmonary embolism	CCAE, Optum EHR, Optum DoD, MDCO, IQVIA, Germany, JMD, MDICR

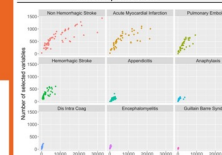


Figure 3. Number of outcomes vs the number of selected variables per prediction task.

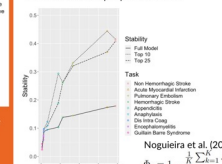


Figure 4. Model stability (stability estimator Nogueira et al. 2018) vs the average number of outcome cases across prediction tasks.



MONDAY

## Why predicting risk can't identify 'risk factors': empirical assessment of model stability in machine learning across observational health databases

Lead: Aniek Markus



# #OHDSISocialShowcase This Week

## OMOP Genomic mapping capacities in conversion of comprehensive genomic profiling results

Maria Rogozhkina<sup>1</sup>, Vlad Korsik<sup>1</sup>, Varvara Savitskaya<sup>1</sup>, Alexander Davydov<sup>1</sup>  
<sup>1</sup> - Odysseus Data Services

### Introduction

Omics data (genomics, proteomics, metabolomics and etc) tends to be the most important data currently because of its possibility to influence decisions in a variety of medical fields. The primary aim of this study was to evaluate the efficacy of representation of genetics data by OMOP Genomic vocabulary. We utilized a real world data personalized database (so that it can be applied to any enriched genomic data format such as FM (FoundationOne), VCF (Variant Call Format), GFF (General feature format) and others) to test the conversion capacity of the vocabulary.

The secondary aim was to define the best conversion strategy for real world database test results (it is one of the most commonly used comprehensive genomic profiling systems worldwide).

### Methods

We processed 3 source genomic data repositories from a real world databases including comprehensive genomic profiling systems.

- Source\_data\_type\_1 = copy number = CN (i.e. amplification or deletion),
- Source\_data\_type\_2 = short variant = SV (a single nucleotide variations or small insertion/deletion),
- Source\_data\_type\_3 = multiple myeloma genetic data = MMGD.

The conversion was performed on CDM 5.4 version with OMOP Vocabulary version: v5.0 09-APR-22.

Among the CN table we analyzed 6 033 de-identified records with name of gene, synonym name, type of mutation (amplification or deletion) and copy number:

source_table_name	gene_name	gene_syn	amplification/deletion	copynumber
CN	STK11	STK11	deletion	0
CN	CDK8	CDK8	amplification	8
CN	CDK4A	CDK4A	amplification	8

For mapping automation we applied the full-match approach with subsequent manual curation. Every distinct source name of the gene was a full-name counterpart in postcoordination approach. In precoordination approach distinct source name of the gene and type of mutation were a full name counterpart. The list of targets was aggregated by concept\_id list of preferred names and synonyms.

Single Nucleotide Variation table includes 233 793 records with information about the gene, substitution in DNA, RNA and protein with position, number of chromosome, type of alteration, sequencing coverage and many other (22 columns). Here are shown only the most important ones:

source_table_name	gene	RNA	protein	chromosome	codingtype	sequencingcoverage
SV_extended	DAKD	20816AC	P684A*752	17q	frameshift	1622
SV_extended	APC	20756AC	G693T	14q	missense	4384
SV_extended	TP53	4750T	R158	17p13	missense	6242

Both protein and RNA columns are filled well so source\_concept\_name was compiled as gene, RNA, protein. At first we made a conversion on the RNA column, if nothing was found, we did matching on the protein column, and if nothing was found again, we tried to do uphill mapping on the Genetic Variation class.

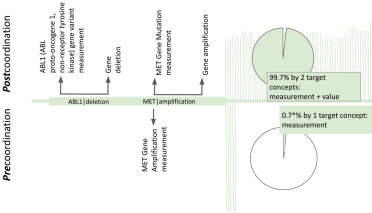
MMGD is a little table containing only 33 rows with information about specimen, method, gene and type of abnormality with all listed fields used to define source\_code.

source_table_name	specimen	method	gene
MMGD	Bone Marrow Aspiration	FISH	T134345
MMGD	Bone Marrow Biopsy	FISH	AMP15221



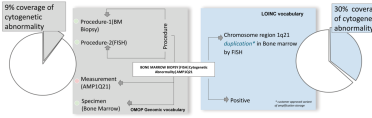
OHDSI ODYSSEUS DATA SERVICES

### Results

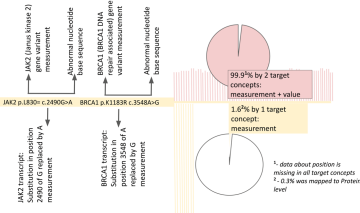


For copy number table the coverage of source data was performed by postcoordination as: meaningful coverage of both genetic variation and related copy number type may be achieved in 99.7% of cases with postcoordination with semantics equivaocal distribution between event and it's value.

In contrast, precoordination resulted in coverage of 0.7% of data, with absolute (one-to-one) semantics match coverage. I.e. postcoordination seems to be optimal strategy to define the changes in number of genes.



Explicit, 1-to-1 representation request is a MMGD scenario. Impossibility of OMOP Genomic to provide the way to reflect the cytogenetic abnormality, the need to store specimen and method attributes as separate facts resulted in targeting to non-genomic terminology.



Lack of appropriate target concepts results in need to uphill mapping in most cases of nucleotide variation tables.

Despite almost perfect (99.9%) SV table's records mapping to OMOP Genomic concepts the semantic coverage remains improper.

Full semantic match is attributed to 1.3% of codes were targeted to RNA Variant and 0.3% of codes were mapped to Protein Variant. Leftover concepts (97.7%) were targeted to Genetic Variation i.e. uphill mapping is the major storing strategy.

Table name	Modeling Approach	Mapping Rate	Gene coverage	Alteration type coverage	Coverage rate	Alteration coordinate coverage
SNV	Pre	~2%	→100%	→100%	→100%	→100%
	Post	~100%	→100%	0%	0%	0%
CN	Pre	~1%	→100%	~1%	NA	NA
	Post	~100%	→100%	→100%	NA	NA

Rationale for appropriate modeling approach selection for specific alteration types

### Conclusion

- OMOP Genomic vocabulary covers a large number of needs well, but some improvements are also needed. It is required to:
- perform deduplication and expand the list of synonyms for a better search using Clingen database
  - increase the number of concepts to cover a larger number of cases: transfer information from DoCM and Cancer Hotspots, take Clinically relevant variation from ICGC to the OMOP Genomic vocabulary.
  - make the genomic LOINC/SNOMED etc concepts non-standard, and then map them to the standard OMOP Genomic concepts.
  - prevent "combinatorial explosion" by allowing postcoordination at least for Copy Number changes
  - ratify the logic for storage of method and specimen

# TUESDAY OMOP Genomic mapping capacities in conversion of comprehensive genomic profiling results

## Lead: Maria Rogozhkina



1. Perseus - <https://github.com/SoftwareCountry/Perseus>
2. White Rabbit, and Rabbit-in-a-Hat <https://github.com/SoftwareCountry/WhiteRabbit>
3. Usagi - <https://github.com/OHDSI/Usagi>
4. Achilles - <https://github.com/OHDSI/Achilles>
5. DQM - <https://github.com/SoftwareCountry/DataQualityDashboard>
6. CDN Builder - <https://github.com/SoftwareCountry/ETL-CDMBuilder>

## Lead: Anton Ivanov







# #OHDSISocialShowcase This Week

*Using geospatial approaches and machine learning for asthma and COPD outcomes: a systematic review*  
Enriching OMOP CDM  
PRESENTER: **Daniel Jeannetot**  
d.jeannetot@erasmusmc.nl

**INTRO:**  
Asthma & COPD are major contributor to morbidity and mortality worldwide. OMOP CDM databases provide a unique opportunity to enrich Electronic health records with geospatial data and machine learning approaches to improve patient-level prediction. This systematic review shows that this is still an untapped approach which large potential for exploration

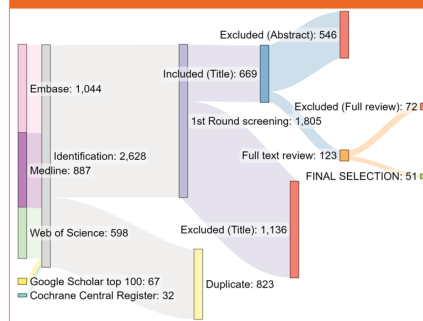
**METHODS**

1. Systematic review following PRISMA guideline
2. 4 databases queried
3. 3 reviewers involved in full text review
4. 12 specific characteristics for data extraction including type of models (ML/non-ML), spatial scale, spatial approach.

**RESULTS**

- 1805 papers screened.
- 123 Papers fully reviewed
- 51 Articles making use of geospatial approach and modelling to measure and predict asthma/related outcomes.

Asthma/COPD research has a lot of potential to benefit from machine learning algorithms and geospatial approaches, especially if combined with observational data.



Most papers approached spatial information in two steps: use of geospatial models to estimate an environmental factor exposure (e.g. specific air pollutants values at specific location) to then integrate values in non-spatial regressions models (e.g. linear, multivariate, etc.), removing specific geospatial and geographical processes information.

The type of scale used varied greatly, with most papers using a local administrative level (e.g. counties, neighbourhoods), thus local, but hard to compare or generalize. Only a few used grid-based spatial data, and even then the resolutions ranged from 5m grid to 1km grid and beyond, leading to widely disparate estimates and areas.

OHDSI provides a coherent and readily available infrastructure to help Asthma/COPD research leverage observational data, machine learning, and geospatial approaches for very large-scale analyses

List of included studies available in annex

## Inclusion criteria

- Has modelling/prediction methods
- Has geospatial/geostatistical approaches
- Explanatory variables include geographical/environmental (air pollution, green/blue space, etc.)
- Main outcome is COPD and/or Asthma related
- Population should be 18 years old or above.

## Search term categories

1. Asthma and/or COPD AND
2. Prediction models (OR, Modelling, Machine Learning etc.) AND
3. Spatial (OR geostatistical; geo\*, etc.) AND
4. ADULT (NOT children, etc.)

## Key points

- Population varied greatly in age groups and sample size (min= 105, max= +50000)
- Scale greatly varied but generally local
- <10 papers used Machine learning algorithms
- Most geospatial approaches are 2 steps
- < 10 papers used specific geostatistical tools
- Inconsistent quality and application of geospatial tools

Daniel Jeannetot, Johnmary Arinze, Victor Pera, Peter Rijnbeek, Katia Verhamme



## THURSDAY

## Using geospatial approaches and machine learning for asthma and COPD outcomes: a systematic review

Lead: **Daniel Jeannetot**



# #OHDSISocialShowcase This Week

A pilot study to evaluate the feasibility of using OHDSI analytical tools for supporting safety surveillance

Ceyda Tugba Pekmez Kristiansen<sup>1</sup>; Lasse Christensen<sup>1</sup>; Michael Stellfeld<sup>1</sup>; Atheline-Major Pedersen<sup>1</sup>; Ditte Mølgaard-Nielsen<sup>1</sup>; Mark White<sup>1</sup>; Peter Jørgensen<sup>1</sup>

<https://qarco.de/bd7bvT>

## OHDSI analytics tools have promising potential for utilising real world data sources to support validation of safety signals.

### Aim

- Real world data sources (RWD) can support validation of safety topics especially when the evidence from traditional safety data sources is scarce.
- Acute cholecystitis and acute cholelithiasis are known risks for Victoza® (liraglutide) and Saxenda® (liraglutide) (1).
- A known risk for liraglutide was chosen for the pilot study to evaluate the feasibility of implementing population level effect estimation into the safety surveillance process using the OHDSI analytics tools.

### Methods

- An observational new-user cohort was created for target drug exposure (liraglutide), comparator drug exposure (sulfamonomethoxazole or SGLT2 inhibitors), and the outcome of acute cholecystitis defined by the SNOMED code 65275005.
- The study cohorts were created using Truven MarketScan employer-based insurance claims data (2). Qualifying target and comparator cohort are shown in Figure 1.
- 1:3 propensity score (PS) matching was performed including age, gender, parity, body mass index, retropharyngeal lymphadenopathy, cardiovascular diseases, and obesity as covariates (Figure 2).
- Survival probabilities for acute cholecystitis were compared using HADES packages (3).

### Key results

**Figure 1: Comparative new user cohort definition**

**Figure 2: Propensity score distribution before and after the propensity score matching**

**Figure 3: Survival probability and the hazard ratio**

**Key result:**

- Cov-proportional HR: 2.26, CI [1.70 - 3.03]
- Minimum detectable relative risk: 1.62 ± 0.17 (SE)

### Summary

- A new-user comparative cohort study was conducted to evaluate the value of implementing population level effect estimation in a RWD setting.
- The application of the OHDSI analytics tools supports a previously validated safety signal of acute cholecystitis following the exposure of liraglutide.

### Conclusion

- Application of the CohortMethod R package supports a known risk of acute cholecystitis for liraglutide on a real-world data source.
- OHDSI analytics tools have promising potential for utilising real world data sources to support the validation of safety signals.
- Next steps will be a new test case for another therapeutic area including negative outcome controls and the data driven selection of covariates.

<sup>1</sup> Nordisk A/S, Global Safety, Safety Surveillance, Søborg, Denmark  
Presented at the European OHDSI Symposium, 2022.06.24, Rotterdam, The Netherlands

1. Nordisk A/S, Global Safety, Safety Surveillance, Søborg, Denmark  
Presented at the European OHDSI Symposium, 2022.06.24, Rotterdam, The Netherlands

2. Nordisk A/S, Global Safety, Safety Surveillance, Søborg, Denmark  
Presented at the European OHDSI Symposium, 2022.06.24, Rotterdam, The Netherlands

3. Schneeweid S, Glynn RJ, Schneeweid S, Glynn RJ, Schneeweid S, Glynn RJ. Effect of Liraglutide Compared With Placebo on Rates of Acute Cholecystitis in Patients With Type 2 Diabetes at High Risk for Cardiovascular Events in the LIRAGLUTIDE Trial. JAMA. 2019;321(10):955-963.

4. Schneeweid S, Glynn RJ, Schneeweid S, Glynn RJ, Schneeweid S, Glynn RJ. Effect of Liraglutide Compared With Placebo on Rates of Acute Cholecystitis in Patients With Type 2 Diabetes at High Risk for Cardiovascular Events in the LIRAGLUTIDE Trial. JAMA. 2019;321(10):955-963.

**FRIDAY**

**A pilot study to evaluate the feasibility of using Observational Health Data Sciences and Informatics analytics tools for supporting the validation of safety signals**

**Lead: Ceyda Pekmez**



4. We compare outcomes between drugs (and classes) for each pivotal RCT emulated using the packages included in HADES (formally known as the OHDSI Methods Library).

Scan QR code link to  
GitHub repository

 **ohdsi**



# #OHDSISocialShowcase This Week

## EHDEN Platform Roadmap

▲ PRESENTER: Michel Van Speybroeck

### INTRODUCTION

EHDEN has currently 140 Data Partners engaged. The EHDEN platform will now allow these 140 Data Partners and other stakeholders to participate and lead real world studies with a focus on reliability, robustness and ease-of-use.

### METHODS

The EHDEN platform is based on a set of core principles:

- Maximum use of available (OHDSI) components
- Additions will be open source as well
- Data privacy by Design
- Supports full study lifecycle: exploration, feasibility, execution, result collection, dissemination
- Study results through interactive dashboards
- Extensible Modular framework
- Robust and integrated security management
- Make the data FAIR

### RESULTS

The EHDEN platform will consist of the following components:

- The EHDEN portal (<https://portal.ehden.eu/>) as the point of entry for all EHDEN related evidence generation capabilities
- An EHDEN network study execution platform based on ARACHNE
- A study design component using Atlas (<https://ohdsi.org/software-tools/>)
- The EHDEN database catalogue contains metadata on all data sources
- The network dashboards offering summarized univariate statistics on the mapped data sources
- EHDEN Academy for the dissemination of training content (<https://academy.ehden.eu/>)
- The capability to share the generated evidence through the EHDEN evidence hub
- Single Sign On through Elixir (recently transitioned to Life Science Login)

### EHDEN Platform Capability Development



### Currently Available in the Data Catalogue

- Information on 67 data partners
- Description of the database
- Demographic coverage
- Type of database (hospital / registry / GP /...)
- Data Governance and Ethics
- Publications
- Database Dashboard (see below)

### Currently Available in Network Dashboards

- Info on 35 data sources covering 44 Million patients
- Filters on country, database type, data source
- Gender Distribution
- Age at first Observation
- Year of Birth
- Average number of records per person per entity
- Data Provenance
- Longitudinal (observation period) coverage
- Visit Types
- Number of record counts / descendant records per OMOP Concept / database

▲ Michel Van Speybroeck, Maxim Moinat, Julia Kurps, Sebastiaan Van Sandijk, José Luis Oliveira, Peter Rijnbeek



## The EHDEN Platform Roadmap

# TUESDAY

## Lead: Michel Van Speybroeck



# #OHDSISocialShowcase This Week

A standard ETL process from REDCap to OMOP

by Francesco Pozzoni

#### INTRO:

- Building an ETL process towards the CDM is a resource-consuming task
- REDCap is a worldwide used web application to manage and build eCRFs for nonprofit research studies and registries
- Aim: leverage REDCap data structure to build a **configurable** ETL procedure that can be adapted to different studies.

#### METHODS

- ETL process working with a **fixed procedural component** and a **study-specific configuration component**
- Mirth Connect, the ETL software built specifically for clinical data. It can perform multiple operations and provides high customizable options for the ETL development
- Two configuration files that provide the **mapping** between the information in the REDCap and OMOP CDM and the **filtering/transformation** procedures that need to be performed on the data

#### RESULTS

- Feature specific tests carried out in an environment that simulates a REDCap project
- Each basic functionality has been evaluated



## REDCAP → OMOP

REDCap studies



REDCap is a web-based eCRF platform, granted for free for non-profit scopes. It is widely adopted all around the world.

The mapping file translates the Rabbit in a Hat Arrows into a format readable by the ETL

Configuration files



The rule file provides filtering and advanced transformation options to be performed on REDCap fields



REDCap data is accessed via the REST API interface, a standard well-documented way to programmatically interact with the platform

Mirth Connect is an ETL software built specifically for clinical data. It can perform multiple operations and provides high customizable options for the ETL development



Data insert in the CDM database is done through SQL queries dynamically built by Mirth Connect

ETL

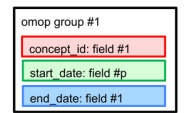
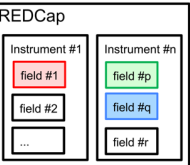
The target CDM version is the most up to date release 5.4 which includes the oncology extension



AMMO BAR

#### How to build the mapping file

- groups are collection of REDCap fields
- 3 types of groups: **person**, **visit** and **fact**
- group each field in REDCap that refers to a single "object" in OMOP
- link each **fact group** to a **visit group**
- translate** the grouping logic into a csv table



by Francesco Pozzoni, Matteo Gabetta, Mauro Bucalo, Nicola Barbarini



## A standard ETL process from REDCap to OMOP

## WEDNESDAY

### Lead: Francesco Pozzoni





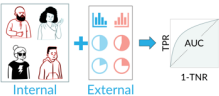
# #OHDSISocialShowcase This Week

## Learning Robust Models from Limited External Statistics

PRESENTER: Tal El Hay

### INTRO

- Model robustness is usually assessed by external validation
- In a previous work, we developed a method that estimates model performance on external data sources from their limited statistical characteristics



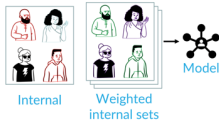
- Can we adopt a similar approach to train robust models, alleviating privacy concerns and communication costs?

### METHODS

- Search for weights that reproduce external statistics; generate a weighted copy of internal data with external characteristics.



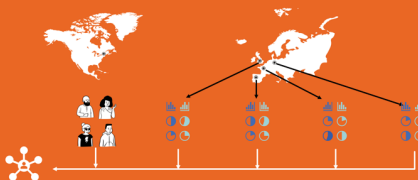
- Train model on internal and weighted sets



### AIM

- Train robust predictive models using:
  - various machine learning algorithms
  - patient-level internal data + population-level statistics from external sources
  - a single (or very few) communication round

Augmenting internal data with population-level statistics from external sources could improve model robustness to data-shift



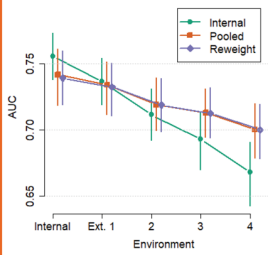
External performance (AUC) of a model trained on internal data degrades faster than for models trained on pooled data or using external statistics and reweighting



Scan for details about the reweighting algorithm

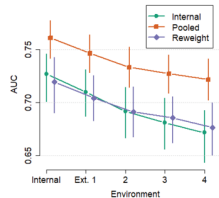
Estimating Model Performance on External Samples from Their Limited Statistical Characteristics, Conference on Health, Inference, and Learning (CHIL) 2022

Covariate shift, logistic regression

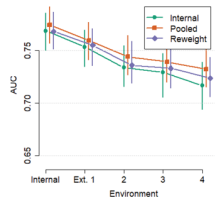


### ADDITIONAL RESULTS

- Model averaging may fail with linear models that combine environment-dependent and invariant predictors using Lasso.
- XGBoost: reweight model only slightly improves over internal one



- Neural network: reweight model outperforms internal but is not as good as pooled one.



### DISCUSSION

**Strengths.** requires only limited statistics (can use info from characterization studies); a single communication round

**Limitations.** may fail if insufficient statistics are used; suboptimal in comparison to pooled training.

**Future directions.** Adapt the method to non-linear models; optimize the choice of stats; introduce a distributionally robust objective.

Tal El Hay and Chen Yanover



Learning robust models from limited external statistics

THURSDAY

Lead: Tal El Hay



 **ohdsi**



# Where Are We Going?

**Any other announcements  
of upcoming work, events,  
deadlines, etc?**





# Three Stages of The Journey

**Where Have We Been?**

**Where Are We Now?**

**Where Are We Going?**







# June 28: The European Symposium

## Presenter:

Nigel Hughes • Director, Observational Health Data Analytics at Janssen Research & Development

