# Welcome to OHDSI/ How To Get Started

## OHDSI Community Call
## Oct. 14, 2025 • 11 am ET

# Upcoming Community Calls

| Date | Topic |
| --- | --- |
| Oct. 14 | Welcome to OHDSI |
| Oct. 21 | Tribute to Andrew Williams/The Power of Collaboration |
| Oct. 28 | Meet the Titans |
| Nov. 4 | Collaborator Showcase Honorees |
| Nov. 11 | TBA |
| Nov. 18 | DARWIN EU 2025 Update |
| Nov. 25 | TBA |
| Dec. 2 | OHDSI/OMOP Research Spotlight |
| Dec. 9 | How Did OHDSI Do This Year? |
| Dec. 16 | Holiday Farewell To 2025 |

# Three Stages of The Journey

# Where Have We Been?

## Where Are We Now?

## Where Are We Going?

# OHDSI Shoutouts!

👏

Congratulations to the team of **Raquel Paradinha, Vicente Barros, João Rafael Almeida, and José Luís Oliveira** on the publication of **A Semantic-Driven for Cohort Data Harmonisation into OMOP CDM Schema** in *Volume 332 of Studies in Health Technology and Informatics: Good Evaluation - Better Digital Health.*

## A Semantic-Driven for Cohort Data Harmonisation into OMOP CDM Schema

Raquel PARADINHA[a], Vicente BARROS[a], João Rafael ALMEIDA[a] and José Luís OLIVEIRA[a]

[a] *IEETA / DETI, LASI, University of Aveiro, Portugal*
ORCiD ID: RP 0009-0006-3983-7926; VB 0009-0007-4526-2249; JRA 0000-0003-0729-2264; JLO 0000-0002-6672-6176

**Abstract.** Clinical research often requires integrating data from diverse sources, which differ not only in structure but also in semantics and language. Traditional extract-transform-load (ETL) pipelines struggle to handle semantic variability and lack built-in support for multilingual or ontology-driven harmonisation. This fragmentation limits the interoperability and reuse of clinical datasets in large-scale analyses. In this paper, we propose an integrated framework that combines an embedding-based concept mapping engine with an automated ETL pipeline using Apache Airflow. The mapping engine uses transformer-based embeddings to align clinical terms with standard concepts, producing outputs in White Rabbit and Usagi-compatible formats to ensure backward interoperability. We validated the system using multilingual real-world datasets demonstrating its ability to handle heterogeneous inputs and maintain end-to-end reproducibility.

**Keywords.** OMOP CDM, Concept mapping, ETL, Clinical data harmonisation

# OHDSI Shoutouts! 👏

Congratulations to the team of **Somayeh Abedian, Eugene Yesakov, Stanislav Ostrovskiy, and Rada Hussein** on the publication of **Integrating Garmin Wearable Data into FHIR-Based Health Systems for Improved Interoperability** in *Volume 332 of Studies in Health Technology and Informatics: Good Evaluation - Better Digital Health.*

## Integrating Garmin Wearable Data into FHIR-Based Health Systems for Improved Interoperability

Somayeh ABEDIAN [a, b,1], Eugene YESAKOV [c], Stanislav OSTROVSKIY [c] and Rada HUSSEIN [a]

[a] The Ludwig Boltzmann Institute for Digital Health and Prevention, Salzburg, Austria
[b] Institute of Health Policy, Management and Evaluation, Dalla Lana School of Public Health, University of Toronto, Canada
[c] Edenlab, Innovative Digital Health Solutions, Estonia

ORCiD ID: Somayeh Abedian https://orcid.org/0000-0002-6427-6941, Rada Hussein https://orcid.org/0000-0003-1257-4848

**Abstract.** As wearable technologies become more common in everyday life, integrating Patient-Generated Health Data (PGHD) into clinical systems has emerged as a critical area in digital health. This study explores how data such as heart rate, step count, sleep patterns, and activity levels (captured in this study via the Garmin Vívoactive 4 smartwatch) can be brought into FHIR-based healthcare systems through the Fitrockr platform. We explore how these data align with key Fast Healthcare Interoperability Resources (FHIR), such as Observation, Device, and Patient. Additionally, we evaluate the compatibility of collected datasets by the Modular Open Research Environment (MORE) platform with FHIR and examine the feasibility of transferring these records to FHIR servers. This level of semantic interoperability could simplify the integration of PGHD into hospital information systems or other healthcare information systems and especially EHRs, thus enhancing their contribution to care delivery, especially in medical decision making and as a source for Clinical Decision Support Systems (CDSS). The paper also discusses how standards like FHIR, openEHR, and Observational Medical Outcomes Partnership (OMOP) can work together to ensure consistent, meaningful integration of wearable data for both clinical practice and secondary analysis. In summary, we reflect on the importance of real-time wearable data availability, reliability, and privacy in supporting a more personalized, data-driven healthcare experience.

**Keywords.** Fast Healthcare Interoperability Resources (FHIR), Healthcare Interoperability, Patient-Generated Health Data (PGHD), Wearables Data

# OHDSI Shoutouts! 👏

Congratulations to the team of **Pawel Rajwa, Angelika Borkowetz, Thomas Abbott, Andrea Alberti, Katharina Beyer, Anders Bjartell, James T Brash, Andrew Chilelli, Eleanor Davies, Bertrand De Meulder, Tamas Fazekas, Asieh Golozar, Ayman Hijazy, Andreas Josefsson, Veeru Kasivisvanathan, Raivo Kolde, Daniel Kotik, Michael S Leapman, Marcin Miszczyk, Rossella Nicoletti, Peter Prinsen, Sebastiaan Remmers, Maria J Ribal, Juan Gómez Rivas, Lara Rodriguez-Sanchez, Monique J Roobol, Emma Smith, Robert Snijder, Carl Steinbeisser, Hein V Stroomberg, Giorgio Gandaglia, Philip Cornford, Susan Evans-Axelsson, James N'Dow, Peter-Paul M Willemse and the PIONEER Consortium** on the publication of **Observational Health Data Analysis of the Cardiovascular Adverse Events of Systemic Treatment in Patients with Metastatic Hormone-sensitive Prostate Cancer: Big Data Analytics Using the PIONEER Platform** in *European Urology Focus.*

**Observational Health Data Analysis of the Cardiovascular Adverse Events of Systemic Treatment in Patients with Metastatic Hormone-sensitive Prostate Cancer: Big Data Analytics Using the PIONEER Platform**

Pawel Rajwa [a,b,c,*], Angelika Borkowetz [d,e], Thomas Abbott [f], Andrea Alberti [g], Katharina Beyer [h], Anders Bjartell [i], James T. Brash [j], Andrew Chilelli [k], Eleanor Davies [j], Bertrand De Meulder [f,l], Tamas Fazekas [c,m], Asieh Golozar [n,o], Ayman Hijazy [l], Andreas Josefsson [p], Veeru Kasivisvanathan [a], Raivo Kolde [q], Daniel Kotik [r,s], Michael S. Leapman [t], Marcin Miszczyk [c,u], Rossella Nicoletti [g], Peter Prinsen [v], Sebastiaan Remmers [h], Maria J. Ribal [w], Juan Gómez Rivas [x], Lara Rodriguez-Sanchez [y], Monique J. Roobol [h], Emma Smith [z], Robert Snijder [k], Carl Steinbeisser [aa], Hein V. Stroomberg [bb,cc], Giorgio Gandaglia [dd], Philip Cornford [ee], Susan Evans-Axelsson [f,ff], James N'Dow [gg], Peter-Paul M. Willemse [hh], on behalf of the PIONEER Consortium

a Division of Surgery and Interventional Sciences, University College London and University College London Hospital, London, UK; b Second Department of Urology, Centre of Postgraduate Medical Education, Warsaw, Poland; c Department of Urology, Comprehensive Cancer Center, Medical University of Vienna, Vienna, Austria; d Department of Urology, University Hospital, University of Rostock, Rostock, Germany; e Department of Urology, University Hospital Carl Gustav Carus, TU Dresden, Dresden, Germany; f European Association of Urology, Nijmegen, The Netherlands; g Unit of Urological Robotic Surgery and Renal Transplantation, University of Florence, Careggi Hospital, Florence, Italy; h Department of Urology, Erasmus MC Cancer Institute, Erasmus University Medical Center, Rotterdam, The Netherlands; i Department of Translational Medicine, Lund University, Lund, Sweden; j IQVIA, Real World Solutions, Brighton, UK; k Astellas Pharma Europe Ltd., Surrey, UK; l Association EISBM, Vourles, France; m Department of Urology, Semmelweis University, Budapest, Hungary; n Odysseus Data Services, New York, NY, USA; o OHDSI Center, Northeastern University, Boston, MA, USA; p Department of Urology, Institute of Clinical Sciences, Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden; q Institute of Computer Science, University of Tartu, Tartu, Estonia; r Center for Advanced Systems Understanding, Görlitz, Germany; s Helmholtz-Zentrum Dresden-Rossendorf, Dresden, Germany; t Department of Urology, Yale School of Medicine, New Haven, CT, USA; u Collegium Medicum - Faculty of Medicine, WSB University, Dąbrowa Górnicza, Poland; v Netherlands Comprehensive Cancer Organisation (IKNL), Utrecht, The Netherlands; w Uro-Oncology Unit, Hospital Clinic, University of Barcelona, Barcelona, Spain; x Department of Urology, Hospital Clinico San Carlos, Madrid, Spain; y Department of Urology, Institut Mutualiste Montsouris, Paris, France; z Guidelines Office, European Association of Urology, Arnhem, The Netherlands; aa Collaborate Project Management, Munich, Germany; bb Copenhagen Prostate Cancer Center, Department of Urology, Copenhagen University Hospital – Rigshospitalet, Copenhagen, Denmark; cc Biotech Research & Innovation Centre (BRIC), University of Copenhagen, Copenhagen, Denmark; dd Division of Oncology/Unit of Urology, Soldera Prostate Cancer Lab, URI, IRCCS San Raffaele Scientific Institute, Milan, Italy; ee Liverpool University Hospitals NHS Trust, Liverpool, UK; ff Bayer AG, Berlin, Germany; gg Academic Urology Unit, University of Aberdeen, Aberdeen, UK; hh Department of Urology, Cancer Center, University Medical Center Utrecht, Utrecht, The Netherlands

# OHDSI Shoutouts! 👏

Congratulations to the team of **Parvaneh Badri, Ivonne Hernández, Justin Long, Maryam Amin, and Reid Friesen** on the publication of **Chronic orofacial pain and psychological distress: findings from a multidisciplinary university clinic** in the *Journal of Oral & Facial Pain and Headache.*

ORIGINAL RESEARCH

## Chronic orofacial pain and psychological distress: findings from a multidisciplinary university clinic

Parvaneh Badri[1]⊙, Ivonne Hernández[1]⊙, Justin Long[1], Maryam Amin[2],†, Reid Friesen[1],*,†⊙

[1]Oral Medicine, Mike Petryk School of Dentistry, Faculty of Medicine & Dentistry, University of Alberta, Edmonton, AB T6G 1C9, Canada
[2]Mike Petryk School of Dentistry, Faculty of Medicine & Dentistry, University of Alberta, Edmonton, AB T6G 1C9, Canada

*Correspondence
rtfriese@ualberta.ca
(Reid Friesen)

† These authors contributed equally.

## Abstract

**Background**: Chronic orofacial pain (COFP) is a complex condition that requires multidisciplinary management grounded in the biopsychosocial model. This study examined the associations between temporomandibular disorders (TMD) and headache symptoms and psychological factors within a university-based multidisciplinary care setting, providing insight into the integration of mental health in COFP management. **Methods**: A retrospective review of 162 patient records from the University of Alberta Multidisciplinary Orofacial Pain Clinic (2020–2023) was conducted. Psychological assessments included the Adverse Childhood Experiences (ACE) scale, Pain Catastrophizing Scale (PCS) and Injustice Experience Questionnaire (IEQ). Logistic regression was used to evaluate associations between psychological factors and pain severity. **Results**: The cohort (aged 13–93) was predominantly female (84.0%). Fifteen percent declined psychological measures. Significant associations were observed between PCS ($p = 0.036$) and IEQ ($p = 0.005$) scores and reported pain severity. Moderate-to-high PCS scores were associated with a 3.67-fold increase in the odds of moderate to severe TMD symptoms (Odds Ratio (OR): 3.67, 95% Confidence Interval (CI): 1.09–12.35), while high PCS scores predicted severe headaches (OR: 3.91, 95% CI: 1.50–10.17, $p = 0.005$). Elevated IEQ scores were similarly associated with increased odds of severe headaches (OR: 2.76, 95% CI: 1.08–7.05, $p = 0.034$). **Conclusions**: Psychological factors such as pain catastrophizing and perceived injustice are strongly associated with symptom severity of TMD and headache symptoms in COFP. These findings underscore the importance of integrating targeted psychological assessments into multidisciplinary care. Further research should explore barriers to implementation and advance biopsychosocial approaches to improve outcomes for patients with COFP.

## Keywords

Orofacial pain; Multidisciplinary clinic; Psychological distress; Temporomandibular joint (TMJ) pain; Chronic pain management

# OHDSI Shoutouts! 👏

Congratulations to the team of **Justin Bohn, James P. Gilbert, Christopher Knoll, David M. Kern and Patrick B. Ryan** on the publication of **Large-scale Empirical Identification of Candidate Comparators for Pharmacoepidemiological Studies** in *Drug Safety.*

## Large-scale Empirical Identification of Candidate Comparators for Pharmacoepidemiological Studies

Justin Bohn[1] · James P. Gilbert[1] · Christopher Knoll[1] · David M. Kern[1] · Patrick B. Ryan[1]

## Abstract

**Background and Objective** The new user cohort design has emerged as a best practice for the estimation of drug effects from observational data. However, despite its advantages, this design requires the selection and evaluation of comparators for appropriateness, a process that can be challenging. The objective of this work was to introduce an empirical approach to rank candidate comparators in terms of their similarity to a target drug in high-dimensional covariate space.

**Methods** We generated new user cohorts for each RxNorm ingredient and Anatomic Therapeutic Chemical level 4 class in five administrative claims databases then extracted aggregated pre-treatment covariate data for each cohort across five clinically oriented domains. We formed all pairs of cohorts with ≥ 1000 patients and computed a scalar similarity score, defined as the average of cosine similarities computed within each domain, for each pair. We then generated ranked lists of candidate comparators for each cohort.

**Results** Across up to 1350 cohorts forming 922,761 comparisons, drugs that were more similar in the Anatomic Therapeutic Chemical hierarchy had higher cohort similarity scores. The most similar candidate comparators for each of six example drugs corresponded to alternative treatments used in the target drug's indication(s), and choosing the top-ranked comparator for randomly selected drugs tended to produce balance on most covariates. This approach also ranked highly those comparators chosen in high-quality published new user cohort design studies.

**Conclusion** Empirical comparator recommendations may serve as a useful aid to investigators and could ultimately enable the automated generation of new user cohort design-derived evidence, a process that has previously been limited to self-controlled designs.

# 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 | ATLAS/WebAPI |
| Tuesday | 12 pm | Generative AI and Analytics |
| Wednesday | 8 am | Psychiatry |
| Wednesday | 11 am | Common Data Model |
| Wednesday | 1 pm | Perinatal & Reproductive Health |
| Wednesday | 7 pm | Medical Imaging |
| Thursday | 8 am | India Community Call |
| Thursday | 11 am | Themis |
| Thursday | 12 pm | Medical Devices |
| Thursday | 12 pm | HADES |
| Thursday | 7 pm | Dentistry |
| Friday | 10 am | Transplant |
| Friday | 10 am | GIS-Geographic Information System |
| Friday | 11:30 am | Steering |
| Monday | 10 am | Healthcare Systems |
| Monday | 11 am | Data Bricks User Group |
| Monday | 2 pm | Electronic Animal Health Records |
| Tuesday | 9 am | Data2Evidence |

# Upcoming Workgroup Calls

| Date | Time (ET) | Meeting |
|---|---|---|
| Tuesday | 12 pm | ATLAS/WebAPI |
| Tuesday | 12 pm | Generative AI and Analytics |
| Wednesday | 8 am | Psychiatry |
| Wednesday | 11 am | Common Data Model |
| Wednesday | 1 pm | Perinatal & Reproductive Health |
| Wednesday | 7 pm | Medical Imaging |
| Thursday | 8 am | India Community Call |
| Thursday | 11 am | Themis |
| Thursday | 12 pm | Medical Devices |
| Thursday | 12 pm | HADES |
| Thursday | 7 pm | Dentistry |
| Friday | 10 am | Transplant |
| Friday | 10 am | GIS-Geographic Information System |
| Friday | 11:30 am | Steering |
| Monday | 10 am | Healthcare Systems |
| Monday | 11 am | Data Bricks User Group |
| Monday | 2 pm | Electronic Animal Health Records |
| Tuesday | 9 am | Data2Evidence |

# OHDSI 2025



ohdsi.org/ohdsi2025

# OHDSI 2025



ohdsi.org/ohdsi2025

If you have a great photo from OHDSI2025, please share it! Use the link in the chat (same link for sharing Showcase posters)

ohdsi.org/ohdsi2025

2025 Titan Awards

Titan Awards

...ina Talapova

OHDSI

Building Community
•
Advancing Data

Data Standards

**Open-Source Development**

Clinical Applications

#JoinTheJourney

Community Collaboration

Community Support

2025 Titan Awards

**Community Leadership**

# Best Community Contribution Winners

## Data Standards

Jared Houghtaling, Polina Talapova, Brian Gow, Manlik Kwong, Andrew J King, Benjamin Moody, Mike Kriley, Tom Pollard, Andrew E Williams



ohdsi.org/2025-global-collaborator-showcase

# Best Community Contribution Winners

## Methods Research

Lu Li, Qiong Wu, Yiwen Lu, Kyra S. O'Brien, Bingyu Zhang, Ting Zhou, Jiayi Tong, Dazheng Zhang, Yuqing Lei, Huilin Tang, Yun Lu, David Asch, Yong Chen



ohdsi.org/2025-global-collaborator-showcase

# Best Community Contribution Winners

**Open-Source Community**

Adil Ahmed, Selvin Soby, Boudewijn Aasman, Parsa Mirhaji



## Summary

A LLM-workflow that maps clinical terminologies to standard OMOP concepts. The pipeline consists of 4 stages:

Uses the harmonized concept name to retrieve similar concepts and traverses OMOP Knowledge Graph to extend context

Final concept selection along with alternative IDs

**Concept Expansion** → **Semantic Retrieval** → **Concept Filteration** → **Concept Selection**

Expands on the input term and generates a harmonized version of it

Filters out irrelevant concepts; minimize signal to noise ratio

Pipeline Workflow



**ohdsi.org/2025-global-collaborator-showcase**

# Best Community Contribution Winners

**Clinical Applications**

Hsin Yi "Cindy" Chen, Thomas Falconer, Anna Ostropolets, Tara V. Anand, Xinzhuo Jiang, David Dávila-García, Linying Zhang, Ruochong Fan, Hannah Morgan-Cooper, George Hripcsak



ohdsi.org/2025-global-collaborator-showcase

# Best Community Contribution Winners

## Community

Clair Blacketer, Haeun Lee, Benjamin Martijn, Evanette Burrows, Patricia Mabry, Deran McKeen, Sam Patnoe, Ben Gerber, Pantelis Natsiavas, Aamirah Vadsariya, Hanieh Razzaghi, Paul Nagy



**Building the OHDSI Evidence Network:** A Global, Open, Federated Collaboration

PRESENTER: **Clair Blacketer**

### INTRODUCTION
- **Real-world data is plentiful** and reflects natural conditions, **but is siloed**, due to **privacy concerns**, preventing the benefits of dataset integration from being realized, e.g. study of rare events, generalizability, use of data-hungry AI tools to reveal new insights
- **Federated networks address this problem** by sharing only aggregated results (not record level data) to preserve data privacy
- The OHDSI Evidence Network was **launched in 2024** inspired by the success of other federated networks, e.g., European Health Data and Evidence Network (EHDEN) and the Data Analysis and Real World Interrogation Network (DARWIN EU).

### METHODS
- **The Evidence Network (EN) is composed of "Data Partner Organizations" (DPOs)** who volunteer to run analytic code on their organization's data.
- **Membership in EN is voluntary** - no contracts or centralized data sharing!
- **Governance is decentralized**; each DPO adheres to its local IRB requirements.
- To catalog data available in the EN, each DPO is sent a **Database Diagnostics** software package which they run **locally** to produce a standardized DbProfiles - aggregated metadata describing the DPO's database(s)
- **All EN activities are opt-in** and include EN workgroup meetings, steering committee representation, monthly data partner calls, and EN study co-development
- A pilot study, "Save Our Sisyphus", measured partner engagement. Results led to the **adoption of best practices by the EN** (learning, clear protocols, transparent communication).

The **OHDSI Evidence Network** demonstrates that **open, federated, community-led research** is **inclusive and effective** on a **global** scale

| JOIN | DIAGNOSE | SHARE | RESEARCH |
| --- | --- | --- | --- |
| Voluntary No Contracts | Run DB Package Locally | Db Profiles Metadata Only | Co-develop Studies |

| 28 Partner Organizations | 48 Databases Connected | 4 Continents Represented | 20+ Studies Supported |

**Scale Without Contracts** — 28 partners joined with zero legal agreements

**Open Tools (R pkg)** — Partners can test locally before committing

**Global Community** — Decentralized governance scales

**Rapid Study Assessment** — From months to weeks - shared Db Profiles accelerate feasibility checks

### RESULTS
- **28 DPOs onboarded since inception**, contributing access to 48 databases across 4 continents (see map)
- The EN supported 20 rapid fit-for-purpose assessments and study co-developments in 2025

*Figure 1: Global map of current OHDSI Evidence Network data partner organizations.*

### KEY LESSONS:
- **Decentralized**, federated, community-led **governance is feasible** and effective at a global scale
- **Trust and transparency** drive collaboration
- Low-burden participation lowers barriers
- Shared tools enable shared learning

### FUTURE GOALS:
- Address funding/sustainability challenges
- Develop and test a process for study development and support
- Refine DPO–study matching
- Expand DPO membership

Clair Blacketer[1,2,4], Haeun Lee[1,8], Benjamin Martijn[1,8], Evanette Burrows[1,4], Patricia Mabry[1,10], Deran McKeen[1,10], Sam Patnoe[1,10], Elizabeth Grossman[1,10], Ben Gerber[1,5], Pantelis Natsiavas[1,6], Aamirah Vadsariya[1,7], Hanieh Razzaghi[1,9], Paul Nagy[1,8]

1. OHDSI Collaborators, Observational Health Data Sciences and Informatics (OHDSI), New York, NY, USA 2. Department of Medical Informatics, Erasmus University Medical Center, Rotterdam, NL 3. Department of Biomedical Informatics, Columbia University, New York, NY, USA 4. Johnson & Johnson, Raritan, NJ, USA 5. Department of Population and Quantitative Health Sciences, UMass Chan Medical School, Boston, MA, USA 6. Institute of Applied Biosciences, Centre for Research & Technology Hellas, Thessaloniki, GR 7. Clinical Informatics Center, University of Texas Southwestern Medical Center, Dallas, TX, USA 8. Johns Hopkins University, Baltimore, MD, USA 9. Applied Clinical Research Center, Children's Hospital of Philadelphia, Philadelphia, PA 10.Health Partners Institute, Bloomington, MN, USA

Take a picture to learn more

Start making steps to join today!

**OHDSI**

**ohdsi.org/2025-global-collaborator-showcase**

# Africa Symposium: Nov. 10-12

The first-ever OHDSI Africa Symposium will be held Nov. 10-12 in Kampala, Uganda, at the Joint Clinical Research Centre (JCRC) and Mestil Hotel. The event will begin with a dedicated one-day training course at JCRC, followed by a two-day main conference at the Mestil Hotel.



**ohdsi.org/africa2025**

# APAC Symposium: Dec. 6-7

The 2025 OHDSI APAC Symposium will be held Dec. 6-7 in Shanghai, China at the Shanghai Jiao Tong University. It will feature a 1-day tutorial and a 1-day main conference.





ohdsi.org/apac2025

# #OHDSISocialShowcase This Week

## Monday

**Taxonomy development as an approach to harmonize source-level data**

(**Maryia Rahozhkina**, Vlad Korsik, Aliaksei Katyshou, Oleg Zhuk, Imelda Henrikson, Matthew Littman)

## Tuesday

**Coordinating center-based, rather than self-deployed, data readiness assessment and improvement for oncology RWE**

(Asieh Golozar, Henry Morgan Stewart, Patrick Alba, Stelios Theophanous, Eric Fey, Benjamin Martin, Jared Houghtaling, Roshanthi Weerasinghe, Thomas Falconer, Benjamin May, Espen Enerly, Shantha Bethusamy, Priya Desai, Annelies Verbiest, Patricia Mabry, Qi Yang, Jonas Minne, Maryna Borshchivska, John Methot, Alvaro Andres Alvarez Peralta, Katja Hoffmann, Michael Franz, Jasmin Carus, Andreas Bjerrum, Elin Hallan Naderi, Ayman Hijazy, Daniel Smith, Petr Domecký, Talita Duarte Salles, Clara L. Oeste Aiara Lobo Gomes, Georgina Kennedy, Thomas Stone, Vagelis Chandakas, Dmytro Dymshyts, Kukkurainen Sampo, Pia Tajanen-Doumbouya, Kimmo Porkka, Ben Gerber, Christian Reich)



Distributed data quality check and study feasibility are notoriously difficult - a central approach can help

Coordinating center-based, rather than self-deployed, data readiness assessment and improvement for oncology RWE

# #OHDSISocialShowcase This Week

## Friday

**Enhancing Data Quality Assessment in Healthcare Research: A Comprehensive Evaluation Framework Using OMOP CDM**

(Júlia Moita, Jorge Cerejo, Inês Mota, Simão Gonçalves, Bernardo Neves, Nuno André da Silva, Francisca Leite, Maria Rosário Oliveira, José Maria Moreira)

# 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?

# Mad Minutes

**Dmytro Dymshyts (148):** Evaluating the OHDSI Phenotype library concept sets using Large Language Models

**Qingrui (Carrie) Wang (115):** Automated Anatomical Identification and Standardization for Medical Images

**Gabriel Salvador (403):** Replicating Alzheimer's Research using standardized phenotyping with the OMOP common data model imaging extension

**Melanie Philofsky (141):** Maximizing EHR Semantic Meaning for Rare Diseases Utilizing a Direct Mapping Strategy

**Erik Benton (507):** OMOP Annotator: A Database agnostic tool for reviewing and augmenting the patient record

**Niko Möller-Grell (310):** Agentic conversation on OMOP CDM: the OMCP-A2A foundation library

**Jared Houghtaling (602):** OMOP Waveform Extension: A Schema for Integrating Physiological Signals and Derived Features into the OMOP CDM

**Jen Park (113):** Real-World Implementation of the Medical Imaging CDM: An Alzheimer's Disease Use Case

**Robert Barrett (603):** Improving VSAC to OMOP Mapping Using LLM Assisted Curation

**Christelle Xiong (205):** AgentDose: Towards Accurate and Scalable Steroid Dose Extraction in OMOP Using NLP Parsers and LLM Agents

**The weekly OHDSI community call is held every Tuesday at 11 am ET.**

**Everybody is invited!**

**Links are sent out weekly and available at:**
**ohdsi.org/community-calls-2025**