



Spotlighting OHDSI's Early-Stage Researchers

OHDSI Community Call
Nov. 25, 2025 • 11 am ET



Upcoming Community Calls

Date	Topic
Nov. 25	Early-Stage Researcher Presentations
Dec. 2	OHDSI/OMOP Research Spotlight
Dec. 9	How Did OHDSI Do This Year?
Dec. 16	Holiday Farewell To 2025
Dec. 23	No Meeting
Dec. 30	No Meeting
Jan. 6	No Meeting
Jan. 13	Where Can We Go Together in 2026?



Dec. 2: OHDSI/OMOP Research Spotlight



Ágota Mészáros

PhD Student, Semmelweis University

Semiautomatic mapping of a national drug terminology to standardised OMOP drug concepts using publicly available supplementary information (*BMC Medical Research Methodology*)



Marta Pineda Moncusí

Postdoctoral Researcher In Health Data, University of Oxford

Changes in use and utilisation patterns of drugs with reported shortages between 2010 and 2024 in Europe and North America: a network cohort study (*The Lancet Public Health*)



Hanieh Razzaghi

Associate Director, PEDSnet Data Coordinating Center

A multifaceted approach to advancing data quality and fitness standards in multi-institutional networks (*JAMIA*)

Lucía Bellas

Real World Evidence Epidemiologist, IMI_EHDEN

Secular Trends in the Use of Valproate-Containing Medicines in Women of Childbearing Age in Europe: A Multinational DARWIN EU Network Study (*Pharmacoepidemiology & Drug Safety*)



Three Stages of The Journey

Where Have We Been?

Where Are We Now?

Where Are We Going?



OHDSI Shoutouts!



Congratulations to the team of **Pablo Guerrero**, **Morten Ernebjerg**, **Thomas Holst**, **David Weese**, **Herve DiBello**, **Susanne Ibing**, **Linea Schmidt**, **Ryan Ungaro**, **Bernhard Renard**, **Christoph Lippert**, **Eugenia Alleva**, **Timothy David Quinn**, **Patricia Kovatch**, **Esther-Maria Antao**, **Elmien Heyneke**, **Aadil Rasheed**, **Stefan Kalabakov**, **Bert Arnrich**, **Alexander Charney**, **Lothar H Wieler**, and **Girish Nadkarni** on the publication of **The AIR·MS data platform for artificial intelligence in healthcare** in *JAMIA Open*.

JAMIA Open, 2025, 8(6), ooaf145
<https://doi.org/10.1093/jamiaopen/ooaf145>
Research and Applications



Research and Applications

The AIR·MS data platform for artificial intelligence in healthcare

Pablo Guerrero , PhD^{*,1,2}, **Morten Ernebjerg** , PhD¹, **Thomas Holst** , MSc¹, **David Weese** , PhD¹, **Herve DiBello** , PhD^{2,3,4}, **Susanne Ibing** , PhD^{2,3,4}, **Linea Schmidt** , MSc², **Ryan Ungaro** , MD⁵, **Bernhard Renard** , PhD^{2,3,4}, **Christoph Lippert** , PhD^{2,3,4}, **Eugenia Alleva** , MD^{3,4}, **Timothy David Quinn** , PhD⁶, **Patricia Kovatch** , BS^{3,4}, **Esther-Maria Antao** , PhD², **Elmien Heyneke** , PhD², **Aadil Rasheed** , MSc², **Stefan Kalabakov** , MSc^{2,3,4}, **Bert Arnrich** , PhD^{2,3,4}, **Alexander Charney** , MD, PhD^{3,4,7}, **Lothar H. Wieler** , PhD^{2,3,4}, **Girish Nadkarni** , MD, MPH^{3,4,7}

¹D4L Data4Life gGmbH, Potsdam, Brandenburg 14482, Germany, ²Hasso Plattner Institute, Digital Engineering Faculty, University of Potsdam, Potsdam, Brandenburg 14482, Germany, ³Hasso Plattner Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai, New York, NY 10029-5674, United States, ⁴Windreich Department of Artificial Intelligence and Human Health, Icahn School of Medicine at Mount Sinai, New York, NY 10029-5674, United States, ⁵Division of Gastroenterology, Icahn School of Medicine at Mount Sinai, New York, NY 10029-5674, United States, ⁶Icahn School of Medicine at Mount Sinai, New York, NY 10029-5674, United States, ⁷The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY 10029, United States

Present address: Morten Ernebjerg, Doctolib GmbH, Berlin, Germany.

Drs P. Guerrero and M. Ernebjerg contributed equally to this work.

*Corresponding author: Pablo Guerrero, PhD, D4L Data4Life gGmbH, Rudolf-Breitscheid-Straße 187, Potsdam, Brandenburg 14482, Germany (pablo.guerrero@data4life.care)

Abstract

Objective: To present the Artificial Intelligence-Ready Mount Sinai (AIR·MS) platform-unified access to diverse clinical datasets from the Mount Sinai Health System (MSHS), along with computational infrastructure for AI-driven research and demonstrate its utility with 3 research projects.

Materials and Methods: AIR·MS integrates structured and unstructured data from multiple MSHS sources via the OMOP Common Data Model on an in-memory columnar database. Unstructured pathology and radiology data are integrated through metadata extracted from and linking the raw source data. Data access and analytics are supported from the HIPAA-compliant Azure cloud and the on-premises Minerva High-Performance Computing (HPC) environment.

Results: AIR·MS provides access to structured electronic health records, clinical notes, and metadata for pathology and radiology images, covering over 12M patients. The platform enables interactive cohort building and AI model training. Experimentation with complex cohort queries confirm a high system performance. Three use cases demonstrate, risk-factor discovery, and federated cardiovascular risk modeling.

Discussion: AIR·MS demonstrates how clinical data and infrastructure can be integrated to support large-scale AI-based research. The platform's performance, scale, and cross-institutional design position it as a model for similar initiatives.

Conclusion: AIR·MS provides a scalable, secure, and collaborative platform for AI-enabled healthcare research on multimodal clinical data.



OHDSI Shoutouts!



Congratulations to the team of **Meredith Adams, Robert Hurley, Karsten Bartels, Matthew Perkins, Cody Hudson, Umit Topaloglu, J Perren Cobb, Karin Reuter-Rice, Jacqueline Stocking, and Ashish Khanna** on the publication of **Extending the Observational Medical Outcomes Partnership (OMOP) Common Data Model for Critical Care Medicine: A Framework for Standardizing Complex ICU Data Using the Society of Critical Care Medicine's Critical Care Data Dictionary (C2D2)** in *Critical Care Medicine*.

Critical Care Medicine

XXX 2025 • Volume 54 • Number 00 • Pages XX-XX

CLINICAL INVESTIGATION

Extending the Observational Medical Outcomes Partnership (OMOP) Common Data Model for Critical Care Medicine: A Framework for Standardizing Complex ICU Data Using the Society of Critical Care Medicine's Critical Care Data Dictionary (C2D2)

OBJECTIVES: To evaluate the compatibility of the Society of Critical Care Medicine's (SCCM) Critical Care Data Dictionary (C2D2) with the Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM) and initiate a set of steps extending OMOP to accommodate specialized critical care data elements.

DESIGN: Systematic analysis and mapping study using a three-tiered semantic matching approach to demonstrate technical feasibility and identify fundamental challenges in critical care data standardization.

SETTING: Critical care medicine informatics research environment.

SUBJECTS: The SCCM's C2D2 elements.

INTERVENTIONS: None.

MEASUREMENTS AND MAIN RESULTS: We evaluated the compatibility of C2D2 clinical variables with the OMOP CDM using a three-tier classification system (full match, partial match, and no match). Our analysis of 226 C2D2 elements revealed that 49.6% of concepts had full OMOP equivalents, 46.4% required modification, and 4.0% had no suitable mapping. Key incompatibilities were identified in ventilator parameters, composite scoring systems, and advanced organ support documentation. A large language model-based semantic matching system yielded a precision of 59.5%, recall of 87.0%, and F1 score of 70.7% at an optimized similarity threshold of 0.90. These findings highlight the need to harmonize data standardization approaches within the field of critical care, including how to handle concept stacking within single variables, age-specific criteria, and specialized constructs that were curated through the SCCM Delphi process, but reveal an OMOP mapping incompatibility or missing variables.

CONCLUSIONS: Extending the OMOP CDM for critical care is technically feasible and requires targeted modifications to accommodate composite scores, temporal precision, and specialized critical care concepts as well as the resources needed to support this build. The community acutely faces crucial decisions about whether to pursue OMOP integration, adapt the C2D2 for version 2.0 compatibility, work toward OMOP vocabulary inclusion through Observational Health Data Sciences and Informatics pro-

Meredith C. B. Adams¹, MD, MS¹
Robert W. Hurley², MD, PhD^{2,3}
Karsten Bartels, MD, PhD, MBA⁴
Matthew L. Perkins, BS⁵
Cody Hudson, MS⁶
Umit Topaloglu, PhD^{7,8}
J. Perren Cobb, MD, FCCM⁹
Karin Reuter-Rice, PhD, NP¹⁰
Jacqueline C. Stocking, PhD,
MBA, MSN¹¹
Ashish K. Khanna, MD, MS^{12,13}

Written work prepared by employees of the Federal Government as part of their official duties is, under the U.S. Copyright Act, a "work of the United States Government" for which copy-



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
Wednesday	8 am	Psychiatry
Wednesday	10 am	Oncology Outreach/Research Subgroup
Wednesday	10 am	Women of OHDSI
Wednesday	12 pm	Latin America
Monday	9 am	Vaccine Vocabulary
Monday	10 am	Healthcare Systems Interest Group
Tuesday	9 am	Data2Evidence



India Symposium: Dec. 2

9:00am - 9:30am	Registration and Coffee
9:30am - 9:45am	Lamp Lighting Ceremony
9:45am - 10:00am	Opening & Welcome Address Dr. Vikram Patil, Deputy Dean – Research (Clinical & Translational), JSS AHER
10:00am - 10:30am	OHDSI India: Progress, Scope & Vision Parthiban Sulur, VP Innovation & Growth, GVV
10:40am - 11:10am	Standardizing Mental Health Data: Africa's Contribution to the Global OHDSI Journey Dr. Tathagata Bhattacharjee, Population Health Data Scientist
11:10am - 11:30am	Coffee Break
11:30am - 12:00pm	Where Does India Stand in the Global Health Data Race? Dr. Rintu Kutum, Group Lead, Augmented Health Systems Laboratory, Faculty Fellow of Computer Science, (KCDH-A); Ashoka University
12:10pm - 12:40pm	India's Healthcare Landscape & National Priorities Dr. Thanga Prabhu, Clinical Informatics Expert
12:40pm - 1:10pm	OHDSI & FHIR: Building Interoperability Together Kumar Sathyam, HL7 Chair, Technical Committee
1:15pm - 2:15pm	Lunch Break
2:15pm - 2:30pm	Inauguration of the OHDSI India Training Program Dr. Kavitha Lamror, Partner, RWE & Digital Transformation

REGISTER NOW		
2:30pm - 3:15pm	Panel Discussion - Shaping India's Real-World Data Journey – From Current Practices to Collaborative Innovation Moderator: Dr. Kavitha Lamror, Partner, RWE & Digital Transformation Panel Members: Manish Sharma, yajur.ai, Founder & Director Dr. Prasan Shankar, Medical director, IAIM Healthcare center, Bangalore Dr. Chandil Kumar, Co-Founder SVM Hospital Dr. Mayank Agarwal, Physician Executive, Clinical Informaticist at BG Tek Personalized Medicine (OPC) Private Limited	
3.25 pm - 4.00 pm	Poster Showcase & Break	
4.00 pm - 5:15pm	Lightning Talks	
5:15pm - 5:30pm	Closing Remarks & OHDSI Team Recognition	
5:30pm - 6:30pm	Networking Session & Closing	

www.ohdsi-india.org/events



APAC Symposium: Dec. 6-7

Day 1 (December 6) – Tutorial at Room 102, Dongxia Yuan Building (Zheng-Cai Cuiju Teaching Building)

Morning Session

- 09:00-09:20 Introduction of OHDSI/OMOP
- 09:20-10:00 OMOP CDM and Vocabulary
- 10:00-10:30 OMOP Conversion Process
- 10:40-12:00 ETL Exercises

Afternoon Session

- 13:30-14:50 OHDSI Analyses: Building Cohorts & Hands-on
- 14:50-15:30 CohortDiagnostics and Population-Level Estimation
- 15:50-16:30 Interpreting Results

Day 2 (December 7) – Main conference at Room A100, 1F, Student Center

Session 1 – From Global to Regional Impact: OHDSI across APAC & Africa

- 09:00 – 09:15 Opening Speech
- 09:15 – 09:45 Keynote Speech from OHDSI Global
- 09:45 – 10:45 APAC Regional Chapter Updates
- 10:45 – 11:00 OHDSI Africa

Day 2 (December 7) – Main conference at Room A100, 1F, Student Center

Session 1 – From Global to Regional Impact: OHDSI across APAC & Africa

- 09:00 – 09:15 Opening Speech
- 09:15 – 09:45 Keynote Speech from OHDSI Global
- 09:45 – 10:45 APAC Regional Chapter Updates
- 10:45 – 11:00 OHDSI Africa

Session 2 – From Research to Reflection: 2025 APAC Studies and Lessons Learned

- 11:15 – 11:30 2025 APAC Study 1 by Fudan University
- 11:30 – 11:45 2025 APAC Study 2 by Peking University
- 11:45 – 12:00 2025 APAC Study 3 by University of Science and Technology of China (USTC)
- 12:00 – 12:10 Journal's Perspectives
- 12:10 – 12:30 Panel Discussion

Session 3 – From Regional Insights to Local Challenges: Real-World Evidence and OHDSI/OMOP in China

- 13:30 – 14:30 Collaborator Showcase: Lightning Talks
- 14:30 – 14:45 Real-World Evidence Talk 1
- 14:45 – 15:00 Real-World Evidence Talk 2
- 15:00 – 15:15 Real-World Evidence Talk 3
- 15:30 – 15:50 Real-World Evidence Using OHDSI/OMOP
- 15:50 – 16:10 Panel Discussion: Opportunities and Challenges Using OHDSI/OMOP for Real-World Evidence in China
- 16:10 – 16:50 Closing & Networking



ohdsi.org/apac2025



2026 Global Symposium

The 2026 OHDSI Global Symposium will return to the Hyatt Regency Hotel in New Brunswick, N.J., on **Oct. 20-22.**





2026 Global Symposium

2026 OHDSI Global Symposium Call for Plenary Sessions

Symposium plenaries provide opportunities to share innovative, community-developed content to empower researchers to generate reliable real-world evidence. The community is currently seeking proposals for our #OHDSI2026 plenaries. These sessions will be 60 minutes in duration and must touch on at least two of following pillars of our community:

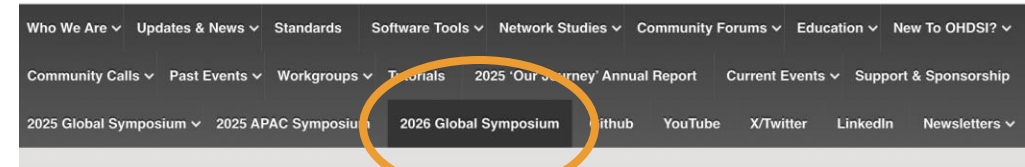
- Open community data standards
- Methodological research
- Open-source development
- Clinical applications

Plenary sessions must also involve three or more on-stage participants across at least two organizations. Sessions may include a combination of keynote talks, panel discussions, interactive activities, and more. We strongly encourage using multiple formats and synthesizing completed research, current perspectives and future calls-to-action to maximize community engagement.

The deadline for proposal submissions is January 30, 2026. Please use the link below to submit your proposal by answering the following questions:

- Name(s) of plenary session organizers:
- Your email address(es):
- Short (2,500 character max) description / abstract of your proposed session:
- Which pillars are you targeting:
- One sentence "pitch" of your session to excite the community:
- Names and roles of individuals who have tentatively agreed to participate in your session:

**Deadline to submit
proposals for #OHDSI2026
plenaries or tutorials is
Jan. 30, 2026!**



2026 OHDSI Global Symposium

Oct. 20-22 • New Brunswick, N.J. • Hyatt Regency Hotel

2026 OHDSI Global Symposium Call for Tutorials

Tutorial sessions aim to deliver educational content, led by community members who wish to train our global collaborators on scientific, technical, and other skills that can support advancing OHDSI's mission and the effective use of real-world data and the generation and dissemination of reliable real-world evidence. Examples of prior tutorials offered are provided here: <https://www.ohdsi.org/tutorials>.

Tutorial sessions are 4 hours in duration. Registrants for your tutorial will be requested to pay a registration fee. The fees will be used to offset the costs of the symposium and other OHDSI expenses. Sessions may include a combination of talks, interactive activities, and more. We strongly encourage using multiple formats to maximize community engagement. Your session must include at least three people from at least two different organizations.

The deadline for tutorial proposal submissions is January 30, 2026. Please use the link below to submit your proposal by answering the following questions:

- Name(s) of tutorial session organizers:
- Your email address(es):
- Short (2,500 character) description / abstract of your proposed session:
- Names and roles of individuals who have tentatively agreed to participate in your session:



2026 Europe Symposium

The 2026 OHDSI Europe Symposium returns to Rotterdam next year and will be held **April 18-20**.

The deadline for abstract submissions will be Feb. 6, 2026.





Tutorials Homepage

OHDSI Tutorials

Education is at the heart of OHDSI's mission, and these tutorials showcase the community's commitment to sharing knowledge. Developed and taught by OHDSI faculty, they highlight tools, standards, and best practices that empower collaborators at every level to engage in open science and generate reliable evidence.

2025 Global Symposium

An Introduction to the Journey from Data to Evidence Using OHDSI

Introduction to OHDSI

OHDSI2025 Tutorial: An Introduction to the Journey from Data to Evidence Using OHDSI

Observational Health Data Sciences and Informatics (OHDSI, pronounced "Odyssey") was founded in 2014.

- Central coordinating center housed at Columbia University.
- A multi-stakeholder, interdisciplinary evidence collaborative to bring out the value of health data for large-scale analytics.

Faculty: Erica Voss, Yong Chen, Katy Sadowski, Nicole Pratt, Roger Carlson, Chongliang (Jason) Luo

Developing and Evaluating Your Extract, Transform, Load (ETL) Process to the OMOP Common Data Model

OMOP Common Vocabulary Model

reference data for the OMOP CDM

What it is

- Compiled standards from disparate public and private sources and some OMOP-grown concepts
- Standardized structure to house existing vocabularies used in the public domain

What it's not

- Static dataset - the vocabulary updates regularly to keep up with the continual evolution of the sources
- Product - vocabulary and improvement is ongoing activity that requires community participation and support

Faculty: Clair Blacketer, Karthik Natarajan, Evanette Burrows, Max Adulyanuksol, Maxim Molnat

Using the OHDSI Standardized Vocabularies for Research

OHDSI Standardized Vocabularies

OHDSI2025 Tutorial: Using the OHDSI Standardized Vocabularies for Research

Data in US after 2005 ICD-10

Data in US after 2005 ICD-9

Data in UK after 2005 ICD-10

Data in UK after 2005 ICD-9

Common reference standard: SNOMED

Faculty: Anna Ostropelets, Vlad Korsik, Polina Talapova, Masha Khitrin

Population-Level Effect Estimation Applications to Generate Reliable Real-World Evidence

OHDSI2025 Tutorial: Population-Level Effect Estimation Applications to Generate Reliable Real-World Evidence

Population-Level Effect Estimation Applications to Generate Reliable Real-World Evidence

George Hripcsak
Martijn Schuemie
Linying Zhang
Tara Anand

Faculty: George Hripcsak, Martijn Schuemie, Linying Zhang, Tara Anand

Clinical Characterization Applications to Generate Reliable Real-World Evidence

OHDSI2025 Tutorial: Clinical Characterization Applications to Generate Reliable Real-World Evidence

Complementary evidence to inform the patient

Clinical characterization: What happened to me?

Patient profile: What is it like to be me?

Population-level effect estimation: What is the causal effects?

Faculty: Patrick Ryan, Aniek Markus, Hsin Yi "Cindy" Chen, Azza Shoaibi

Patient-Level Prediction Applications to Generate Reliable Real-World Evidence

OHDSI2025 Tutorial: Patient-Level Prediction Applications to Generate Reliable Real-World Evidence

Prediction Problem Definition

Observation Window

Time-at-risk

outcome

Faculty: Jenna Reys, Egill Fridgerisson, Ross Williams

ohdsi.org/tutorials

www.ohdsi.org

#JoinTheJourney





#OHDSISocialShowcase This Week

Monday

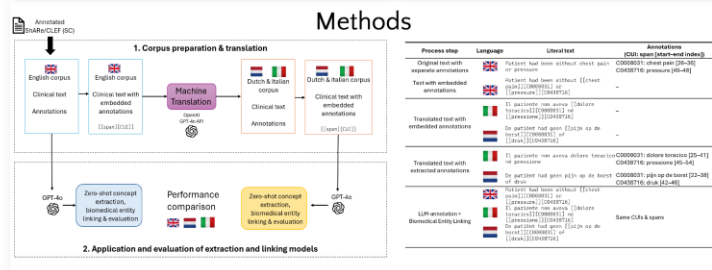
The potential of LLMs for multilingual NER in clinical notes

(Sara Mazzucato, Erik M. van Mulligen, Tom M. Seinen)

LLMs enable unsupervised multilingual clinical data extraction supporting healthcare data standardization across languages

The potential of LLMs for multilingual Name Entity Recognition and Biomedical Entity Linking in clinical notes

Background: extracting biomedical concepts exist for English but are scarce for other languages like Dutch and Italian. Translating corpora while preserving annotation integrity is challenging. This study explores how effectively a large language model can extract and link clinical entities from translated notes while maintaining structured information across languages.



Limitations: GPT-4o performs well on translated corpora, but real-world EHRs are noisier. Precision is high but recall—especially in English—needs improvement. Future work will address robustness, hallucinations, integrate outputs into OMOP CDM, test other LLMs and real-data validation.

Conclusions: LLMs can effectively extract biomedical concepts across languages. Despite recall gaps, results support their use in multilingual clinical NLP and integration into standardized data models.



Sara Mazzucato, Erik M. van Mulligen, Tom M. Seinen
Department of Medical Informatics, Erasmus MC





#OHDSISocialShowcase This Week

Tuesday

An Iterative LLM Based Pipeline for Extracting Clinical Data from Medical Records

(**Erik Calcina**, Tinkara Meterc, Nadya Shpanko, Lahin Spindari, Eva Lindner, Lisa Hoogendam, Harm Slijper, Ruud Selles, Erik Novak)

Turn free-text medical notes into structured data, saving time and maintaining data privacy with large language models.

An Iterative LLM Based Pipeline for Extracting Clinical Data from Medical Records

Background: Clinical notes in electronic health records (EHRs) are **unstructured** and often contain **abbreviations** and **inconsistent terminology**, making it difficult to extract **accurate information** and transform it into the **OMOP CDM** for research and decision making.

1. Zero-Shot Prompting

- **General LLM** for initial extraction from clinical notes.
- Tailored prompts to minimize **hallucinations** and **irrelevant outputs**.
- Created a **baseline** to support **human annotation**.

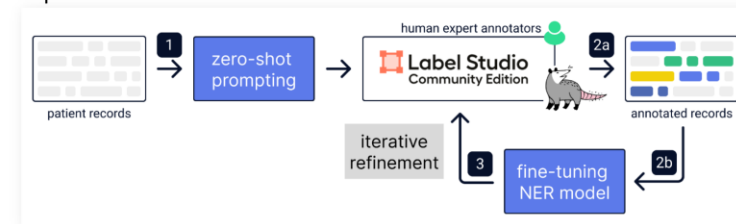
2. Data Labeling & Fine-Tuning

- **2a: Human experts** annotated clinical records using **Label Studio**, assisted by **zero-shot outputs**.
- **2b: Fine-tuned the NER model** to improve extraction accuracy.

3. Iterative Refinement

- After **fine-tuning** the model is used on **broader data**.
- Outputs go back into annotation (**iterative loop**)
- **Enhances model accuracy and reliability**.

Pipeline



Limitation: The model still struggles with **abbreviations** and **inconsistent clinical terminology**. Its performance depends on the **quality** and **consistency** of human-annotated training data. It may also miss **rare** or **subtle clinical details** that are not well represented in the training data.





#OHDSISocialShowcase This Week

Wednesday

Facilitating OHDSI ATLAS Cohort creation via a Custom ChatGPT: a preliminary evaluation based on system instructions and file uploads model

(Miguel-Angel Mayer, Angela Leis, Juan Manuel Ramírez-Angueta)

- ✳ This initial evaluation demonstrates the **potential of LLMs to assist in ATLAS cohort building** and lays the groundwork for more advanced integrations (more refinements are necessary)
- ✳ Training methods used in similar ChatGPT solutions, such as providing additional structured examples, or performing domain-specific fine-tuning, have the potential to boost the robustness of these **custom assistants for supporting OHDSI tools**

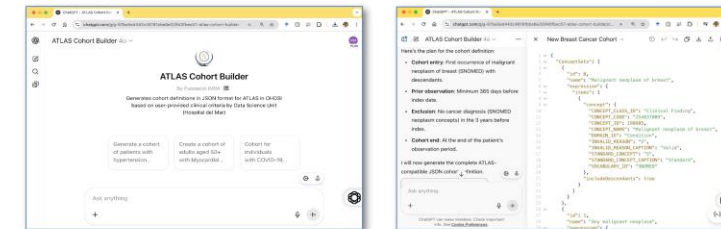
Facilitating OHDSI ATLAS Cohort creation via a Custom ChatGPT: a preliminary evaluation of a knowledge-based model using system instructions and file uploads

Background:

Despite ATLAS's user-friendly intent, creating complex cohort definitions can be challenging for non-expert users. Recent advances in large language models (LLMs) such as OpenAI's ChatGPT offer a potential solution to lower this barrier. LLMs can understand natural language descriptions and generate structured outputs, which suggests they could translate a researcher's description of a cohort into the formal definition needed by ATLAS. Three main methods exist for customising ChatGPT: First, using System Instructions and File Uploads in ChatGPT's interface allows easy, rapid prototyping by directly providing structured examples, requiring minimal effort and cost but offering less control, and it was the method used in this evaluation. Second, API fine-tuning and third, API integration (Retrieval-Augmented Generation).

Results

Figure 1 shows the customChatGPT created. Figure 2 shows the results obtained when creating a cohort in a JSON file format



Methods

- Assistant Development:** We constructed a prototype ChatGPT-based assistant (using GPT-4) tailored to cohort definition tasks through prompt engineering and provided context. Rather than training a new model, we leveraged ChatGPT's system message and file-upload features to supply domain-specific knowledge:
 - First, we uploaded the cohort definition chapter of the OHDSI Book and several example cohort definition JSON files (exported from ATLAS) and documentation snippets to serve as a format template and reference for the model
 - Second, we also crafted detailed system instructions describing the ATLAS cohort JSON syntax and the mapping from common clinical criteria to the JSON fields
- Evaluation:** We performed a preliminary evaluation using multiple test scenarios to gauge the assistant's feasibility and usefulness. We designed a set of natural language cohort requests reflecting use cases a clinical researcher might have
 - To validate correctness, we imported each JSON file into the ATLAS platform to see if it was accepted and interpreted as intended
 - Several iterations were necessary

Research Programme on Biomedical Informatics (GRIB), Hospital del Mar Research Institute and Data Science Unit, Hospital del Mar Barcelona (Spain)

Hospital del Mar

Hospital del Mar Research Institute



Miguel-Angel Mayer, Angela Leis, Juan Manuel Ramírez-Angueta

OHDSI EUROPE'25 Symposium 5-7 July 2025
Old Prison - Hasselt University, Belgium



#OHDSISocialShowcase This Week

Thursday

Using OHDSI Usagi as a validation tool after LLM translation and auto-mapping of standard procedure classifications from different health systems

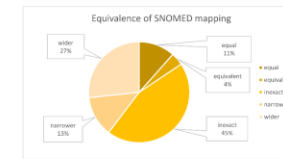
(Karen Triep, Hugo Guillen-Ramirez, Florian Duss, Stefanie Marti, Guido Beldi, Olga Endrich)

Usagi tool: supports automapping and efficient **validation** of non-English standard procedure **classifications**

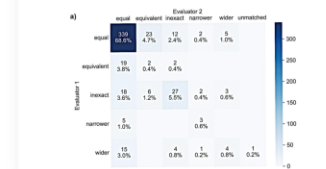
Using OHDSI USAGI as a validation tool after LLM translation and auto-mapping of standard procedure classifications from different health systems.

Background: Procedure classifications like the Swiss procedure catalogue CHOP or the German OPS often contain similar content, but differ on granular level and in hierarchical structure and format. Therefore, translation and/or mapping done either manually or automated is time-consuming and prone to errors.

Result 1: Validation results of LLM translated terms



Result 2: confusion matrix: distribution of labels across evaluators.

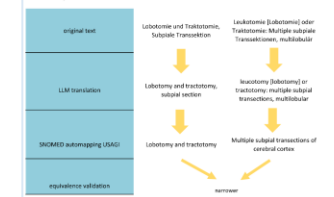


Methods

- 1 Terms of 2 different catalogues translated by LLM
- 2 Translated terms loaded into USAGI
- 3 Automapping of both catalogues by USAGI to SNOMED CT
- 4 Assessment of equivalence for validation, 2 evaluators
- 5 Distribution of labels across evaluators

We tested the Usagi term similarity approach as a validation method for datasets of automatically translated and mapped catalogues from different coding systems.

Example



Limitation: Only 55% acceptable matches (equal, equivalent, narrower, wider) for automatically translated terms from different national health systems when translated automatically by the LLM. Manual validation and correction necessary.



Karen Triep, Hugo Guillen-Ramirez, Florian Duss, Stefanie Marti, Guido Beldi, Olga Endrich



#OHDSISocialShowcase This Week

Friday

Expanding the Semmelweis Clinical Database with parameters extracted from free text documents using large language models

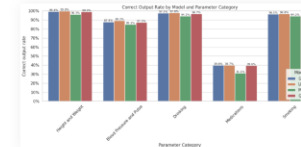
(**Zsolt Bagyura**, Ágota Mészáros, Hujter Mónika, Práger Zsófia, Attila Kovács, Eszter Kővári, Alexandra Assabiny, Ádám Tabák, Tibor Héja, Loretta Kiss)

High-quality **targeted parameter extraction** from **free text documents** using **LLMs** is feasible due to **very high overall accuracy**. After post-processing and mapping to the OMOP standard, the acquired data can be integrated into our data warehouse.

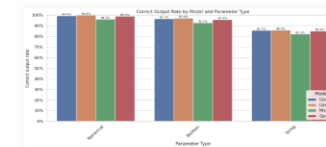
Expanding the Semmelweis Clinical Database with parameters extracted from free text documents using large language models

Background: The Semmelweis Clinical Database stores OMOP-standardized clinical data from approximately 2 million patient visits since 2011, including diagnoses, procedures, drugs, and laboratory observations. However, significant information, such as family history, health behaviors, physical status, medication, and diet, remains in unstructured free-text format, hindering research. Applying Large Language Models (LLMs) for medical data extraction lacks out-of-the-box solutions. Our aim was to identify techniques and test the usability of different LLMs for extracting and structuring selected parameters.

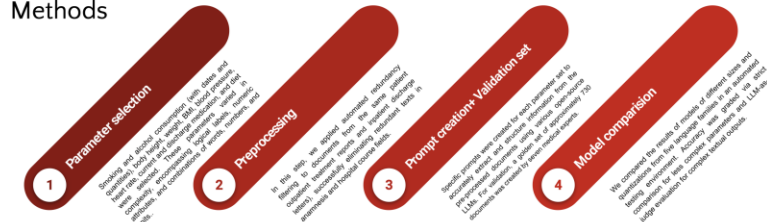
Result 1: For efficient data extraction different prompting techniques needed to be utilized. Finding these category specific prompts was a crucial step during the research, and seemingly successful, given the overall good results.



Result 2: While the numerical and boolean output types are grasped easier by the models, the string variables might pose greater complexity. This is due to the higher variance of the potential outputs and the greater importance of understanding complex context.



Methods



LLMs successfully identified parameter instances not found through manual labeling. This highlights the superior efficiency and comprehensiveness of our LLM-based method compared to traditional manual extraction, underscoring its power in uncovering valuable insights that might otherwise be overlooked.



Zsolt Bagyura¹, Ágota Mészáros¹, Mónika Hujter², Zsófia Práger², Attila Kovács¹, Eszter Kővári¹, Alexandra Assabiny¹, Ádám Tabák¹, Tibor Héja¹, Loretta Kiss¹,
1 Institute for Clinical Data Management, Semmelweis University, 2 Hiflylabs Zrt



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?



Nov. 25: Spotlighting OHDSI's Early-Stage Researchers



Sumin Lee and Kyulee Jeon, Yonsei College of Medicine

ARKE: An Ontology-Driven Framework for Standardizing Radiology Procedure Terminology Using LLMs and RAG



Markian Hromiak and Jacob Zelko, George Institute of Technology

**Exploring Efficient and Scalable OMOP
CDM Workflows by Leveraging dbt-synthea**



Bingyu Zhang, Univ. of Pennsylvania

**The Fine Art of Tolerance: Robustify p-value
Calibration in Observational Studies with Partially Valid Negative Control Outcomes**





**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](https://www.ohdsi.org/community-calls-2025)