



# Collaborations for Strategic Priorities

**OHDSI Community Call**  
**Jan. 24, 2023 • 11 am ET**



# Upcoming OHDSI Community Calls

Date	Topic
Jan. 31	Introduction to Phenotype Phebruary
Feb. 7	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 14	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 21	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 28	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023



# Upcoming OHDSI Community Calls

Date	Topic
Jan. 31	Introduction to Phenotype Phebruary
Feb. 7	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 14	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 21	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023
Feb. 28	Phenotype Phebruary Weekly Update + Workgroup Plans for 2023



# Jan. 31: Introduction to Phenotype Phebruary



**Patrick Ryan**

Vice President, Observational Health Data Analytics, Janssen Research and Development, Inc.; Adjunct Assistant Professor, Columbia University



**Gowtham Rao**

Senior Director, Observational Health Data Analytics, Janssen Research and Development, Inc.; Phenotype Development & Evaluation Workgroup Lead



**Azza Shoaibi**

Associate Director, Observational Health Data Analytics, Janssen Research and Development, Inc.; OHDSI2022 presenter on “OHDSI Phenotype Phebruary: lessons learned”





# Three Stages of The Journey

**Where Have We Been?**

**Where Are We Now?**

**Where Are We Going?**





# OHDSI HADES releases: Characterization v0.0.5

Characterization 0.0.4



Get started

Reference

Articles ▾

Changelog

 HADES



## Characterization

## Introduction

Characterization is an R package for performing characterization of a target and a comparator cohort.

## Features

- Compute time to event
- Compute dechallenge and rechallenge
- Computer characterization of target cohort with and without occurring in an outcome cohort during some time at risk
- Run multiple characterization analyses efficiently
- upload results to database
- export results as csv files

### Links

[Browse source code](#)

[Report a bug](#)

[Ask a question](#)

### License

Apache License 2.0

### Citation

[Citing Characterization](#)

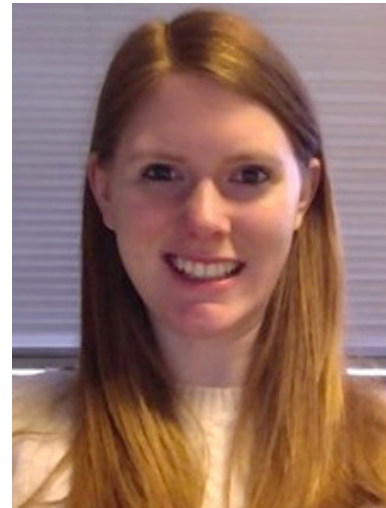
### Developers

Jenna Reps

Author, maintainer

Patrick Ryan

Author





# New Demos: PatientLevelPrediction v6/Strategus

**Jenna Reys**, co-lead of the PLP workgroup, recently shared several video tutorials of version 6 of the PatientLevelPrediction tool. The demos are available on both our website and our YouTube page.

## Videos

- how to extract data and develop single model using PLP v6
- how to design prediction models and develop multiple models using PLP v6
- demonstrating the PLP v6 shiny app that enables users to interactively explore prediction model results
- how to use the new OHDSI R package Strategus and OHDSI modules to develop an OHDSI prediction development network study
- how to run an OHDSI prediction network study using the new Strategus approach

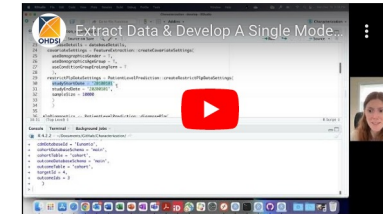
[ohdsi.org/plp-v6-demos/](https://ohdsi.org/plp-v6-demos/)

## Learn More About Version 6 Of The PatientLevelPrediction Package

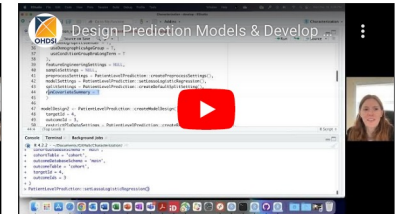
PatientLevelPrediction, a part of the [HADES open-source tool library](#), is an R package for building and validating patient-level predictive models using data in the OMOP Common Data Model format. Check out the [PatientLevelPrediction \(PLP\) github page](#) for more information.

PLP workgroup co-lead and package maintainer Jenna Reys created a series of demo videos to provide assistance with using v6 of the package. You can check out the descriptions and videos here, or on [our OHDSI YouTube page](#) (check out the tutorials playlist).

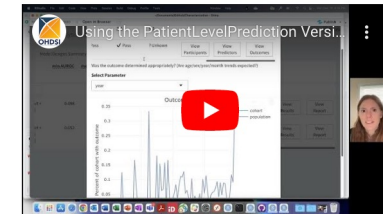
This video demonstrates how to extract data and develop single model using PatientLevelPrediction version 6.



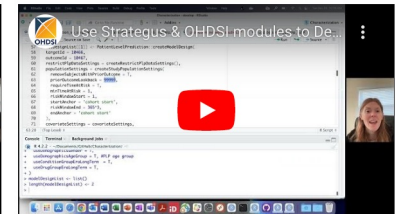
This video demonstrates how to design prediction models and develop multiple models using PatientLevelPrediction version 6.



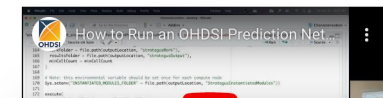
This video demonstrates the PatientLevelPrediction version 6 shiny app that enables users to interactively explore prediction model results.



This video explains how to use the new OHDSI R package Strategus and OHDSI modules to develop an OHDSI prediction development network study. [Text instructions are available here.](#)



This video explains how to run an OHDSI prediction network study using the new Strategus approach. [Text instructions are available here.](#)





# 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
Wednesday	7 am	Medical Imaging
Wednesday	11 am	Latin America
Thursday	9:30 am	Data Quality Dashboard
Thursday	7 pm	Dentistry
Friday	9 am	GIS – Geographic Information System General
Friday	9 am	Phenotype Development and Evaluation
Friday	10 am	Education
Friday	11 am	Clinical Trials
Monday	10 am	Healthcare Special Interest Group
Tuesday	9 am	OMOP CDM Oncology Genomic Subgroup

[ohdsi.org/workgroups](https://ohdsi.org/workgroups)





# Save the Date!

## October 20-22, OHDSI Global Symposium



*Location and more details coming soon*





# Next CBER Best Seminar

The CBER BEST Seminar Series returns Wed., Feb. 8, at 11 am ET, as 2022 Titan Award recipient **Fan Bu** will provide a presentation on Bayesian Safety Surveillance with Adaptive Bias Correction.



Speaker: Dr. Fan Bu (UCLA)

Description: In this presentation, we will discuss a collaborative project with the FDA CBER BEST Initiative to improve on post-market vaccine safety surveillance procedures through Bayesian sequential analysis. Post-market surveillance on approved vaccine products is essential for addressing safety concerns. The goal is to detect rare or high-risk adverse events that often go undetected in clinical trials due to limited sample sizes. Collaborating with FDA CBER, we have developed a Bayesian alternative surveillance procedure that tackles these challenges in sequential analysis of observational data. The standard statistical approach for surveillance is Maximum Sequential Probability Ratio Test (MaxSPRT). Through comprehensive empirical evaluations on large-scale observational healthcare databases, we show that, compared to MaxSPRT, our Bayesian method offers more flexibility on the surveillance schedule, more transparency and interpretability in decision-making, and better error control through statistical correction of bias in observational data.





# Collaboration Opportunity Spotlight: Joint Statistical Meeting (JSM) – Feb. 1 Deadline!



Conference Information

Program

EXPO

Career Service

Professional Development

[Be on the Program](#)

Sponsors

Home

Follow: [f](#) [t](#)

Overview

Participant Guidelines

Sessions

**Submissions**

Professional Development

## Submissions

Invited Session Proposals – Deadline September 8, 2022



Topic-Contributed Session Proposals – Deadline December 8, 2022



Abstract Submissions – Deadline February 1, 2023 **NOW OPEN**



Presentations may be given on any topic of statistical interest; however, authors are encouraged to submit papers on the theme set by 2023 ASA President Dionne Price, *"One Community: Informing Decisions and Driving Discovery."* Additionally, abstracts with a primary focus on statistical applications are encouraged.

Opens: December 1, 2022

Closes: February 1, 2023

Decision Due: March 31, 2023

## Key Dates

November 15, 2022 - December 8, 2022  
Online submission of topic-contributed session proposals

**December 1, 2022 - February 1, 2023**  
Online submission of abstracts (all except invited papers and panels)

January 15, 2023  
Computer Technology Workshop (CTW) proposal deadline

**January 25, 2023 - April 5, 2023**  
Online submission of JSM Meeting & Event Requests

**January 25, 2023**  
Deadline to request registration extension for government agencies

**May 1, 2023**  
Registration and housing open



# ICPE 2023 Abstract Deadline: Feb. 13



**ICPE 2023**

**August 23 - 27**

HALIFAX, NOVA SCOTIA, CANADA  
HALIFAX CONVENTION CENTRE

ispe  
pharmacoepi.org  
#ICPE23 | @IntPharmacoEpi

---

**ICPE 2023 Call for Abstracts**  
**Submission Deadline: February 13, 2023**

---

**Abstract submissions for the 39th International Conference on Pharmacoepidemiology and Therapeutic Risk Management (ICPE 2023) are now being accepted online**

**Call for Abstracts**  
ICPE 2023 will be a live event held at the Halifax Convention Centre, Halifax, Nova Scotia, Canada, August 23-27, 2023. Virtual presentations are not permitted for the event; all presentations must be delivered in person. If you submit an abstract, it is with the intention that you will physically attend the conference to present it.

The ICPE 2023 is a unique forum for the exchange of scientific information from the fields of pharmacoepidemiology and therapeutic risk management among those in the pharmaceutical industry, government, academia, service

[pharmacoepi.org/meetings/annual-conference/](https://pharmacoepi.org/meetings/annual-conference/)



# Grant Opportunity

## VIEW GRANT OPPORTUNITY



PAR-23-034

NLM Research Grants in Biomedical Informatics and Data Science (R01 Clinical Trial Optional)

Department of Health and Human Services

National Institutes of Health

[« Back](#) | [Link](#)

Apply

Subscribe

SYNOPSIS

VERSION HISTORY

RELATED DOCUMENTS

PACKAGE

[Print Synopsis Details](#)



### General Information

<b>Document Type:</b> Grants Notice	<b>Version:</b> Synopsis 1
<b>Funding Opportunity Number:</b> PAR-23-034	<b>Posted Date:</b> Oct 06, 2022
<b>Funding Opportunity Title:</b> NLM Research Grants in Biomedical Informatics and Data Science (R01 Clinical Trial Optional)	<b>Last Updated Date:</b> Oct 06, 2022
<b>Opportunity Category:</b> Discretionary	<b>Original Closing Date for Applications:</b> Jan 07, 2026
<b>Opportunity Category Explanation:</b>	<b>Current Closing Date for Applications:</b> Jan 07, 2026
<b>Funding Instrument Type:</b> Grant	<b>Archive Date:</b> Feb 12, 2026
<b>Category of Funding Activity:</b> Education Health	<b>Estimated Total Program Funding:</b>
<b>Category Explanation:</b>	<b>Award Ceiling:</b> \$250,000
<b>Expected Number of Awards:</b>	<b>Award Floor:</b>
<b>CFDA Number(s):</b> 93.310 -- Trans-NIH Research Support 93.879 -- Medical Library Assistance	
<b>Cost Sharing or Matching Requirement:</b> No	



# Upcoming OHDSI APAC Community Calls

Date	Topic
Feb. 16	Training Session #1
Mar. 16	Training Session #2
Apr. 20	Training Session #3
May 18	Training Session #4
June 15	Regional Chapter Mid-Year Updates



# Oxford Real World Evidence Summer School

## Oxford Summer School 2023: Real World Evidence using the OMOP Common Data Model

### COURSE DIRECTORS

**Daniel Prieto-Alhambra**

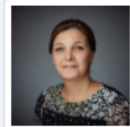
Professor of Pharmaco- and Device Epidemiology



### COURSE ADMINISTRATOR

**Mahkameh Mafi**

Personal Assistant to Professor Prieto-Alhambra



### OTHER COURSES

**Statistics: Designing clinical research and biostatistics**

### Brief Description:

Our Real World Evidence Summer School will provide participants with the tools and concepts necessary to plan and execute Real World Evidence studies, with a focus on the use of the OMOP common data model. The course will have morning lectures followed by afternoon practicals where concepts discussed in the morning will be put in practice with hands-on sessions. Practical sessions will have two tracks: a) for those interested in the design of studies and use of existing analytical and data curation tools; and b) for more advanced data scientists and programmers interested in the development or modification of analytical code using R.

**Registration:** It is now open

**Venue:** Lady Margaret Hall Talbot Hall Theatre, Norham Gardens, Oxford OX2 6QA

**Date:** 19th- 23rd June 2023

For booking please use **Booking information**

**Please see the Preliminary Programme here**

### AUDIENCE:

Pharmacists, clinicians, academics (including statisticians, epidemiologists, and related MSc/PhD students); Industry (pharmacy or device) or Regulatory staff with an interest in the use of routinely collected data for research.

### LEARNING GOALS:



# #OHDSISocialShowcase This Week



## Development of Machine Learning models for Cancer Survival among Lung cancer patients with Tyrosine Kinase Inhibitors (TKIs) treatment

Alex PA. Nguyen<sup>1</sup>, Phuc T. Phan<sup>2</sup>, Min-Huei Hsu<sup>1</sup>, Jason C. Hsu<sup>1,2\*</sup>  
<sup>1</sup> Clinical Data Center, Office of Data Science, Taipei Medical University, Taipei, Taiwan;  
<sup>2</sup> International PhD Program in Biotech and Healthcare Management, College of Management, Taipei Medical University, Taipei, Taiwan;



### Background

Lung cancer is the most common cause of cancer death worldwide, including in Taiwan. Mutation in the EGFR gene is a driver in lung adenocarcinoma, as this gene is overexpressed in more than 50% of non-small cell lung cancer (NSCLC) in Asia. Most patients benefited from TKI therapies, but 5%-10% of patients did not achieve disease control when administered EGFR-TKIs and therefore acquired drug resistance within 10-12 months.

In this study, we aimed to develop prediction models for lung cancer survival among patients with TKI treatment using a larger number of samples, different data types, and various machine learning algorithms.

### Methods

#### Study Design and Data Source

We conducted a retrospective study in which we obtained the data from the Taiwan Cancer Registry (TCR) database and the Taipei Medical University Clinical Research Database (TMUCRD).

#### Cohort Selection

Patients with lung cancer (ICD-O-3 code: C33, C34.1) from 2008 to 2018 in the TCR database. Exclusion criteria included individuals under 20, small cell lung cancer (SCLC) patients, and patients who did not receive lung cancer treatment in the three hospitals. Following that, only cancer patients who were undergoing TKIs (i.e., patients using EGFR-TKIs, ATC codes L01EB) were included in our study cohorts.

#### Outcome Measurement

The outcome of this study was death within two years following diagnosis. Data were censored at the date of death or loss to follow-up, insurance termination, or the study's end on December 31, 2020.

#### Feature Selection

The selected features were as follows: (1) Demographic information; (2) Cancer condition; (3) Comorbidities; (4) Current medications use; and (5) Laboratory test results. All the features were defined before the time patients were prescribed TKI drugs.

#### Developing the Machine Learning models

Six machine learning algorithms were used including Logistic Regression (LR), bootstrap aggregation (bagging), gradient boosting machine (GBM), AdaBoost, random forest (RF), and extreme gradient boosting (XGBoost), to develop the prediction models.

The training set, containing the data of Taipei Medical University Hospital and Wang Fang Hospital. The testing set, including the data of Shuang Ho Hospital, was used to validate the models. The 5-fold cross-validate was applied.

The area under the receiver operating characteristic curve (AUC), accuracy, sensitivity, specificity, and F1-score were computed to evaluate and compare the performance of all prediction models.

### Results

Table 1. Baseline demographic of cohort patients in the study

Feature	Training cohort (n=733)	Testing cohort (n=454)	Feature	Training cohort (n=733)	Testing cohort (n=454)
<b>Demographic</b>			<b>Cancer Condition</b>		
Gender, No. (%)			Tumor size, No. (%)		
Female	327 (44.7%)	206 (45.4%)	T<3cm	212 (29.0%)	103 (22.7%)
Male	404 (55.3%)	248 (54.6%)	3<=T<=7cm	321 (43.9%)	210 (46.3%)
Age, Mean (SD), y	67.8 (13.2)	67.3 (12.7)	T>7cm	63 (8.6%)	48 (10.6%)
BMI			Missing	135 (18.5%)	93 (20.5%)
Mean (SD)	23.4 (3.85)	23.2 (4.00)	Cancer stage, No. (%)		
Median (Min, Max)	23.1 (13.0, 61.3)	22.9 (13.2, 38.1)	stage = 0	52 (7.1%)	28 (6.2%)
Missing	238 (32.6%)	94 (20.7%)	stage = 1	17 (2.3%)	9 (2.0%)
Smoking, No. (%)			stage = 2	81 (11.1%)	33 (7.3%)
No	356 (48.7%)	222 (48.9%)	stage = 3	547 (74.8%)	368 (81.1%)
Yes	156 (21.3%)	139 (30.6%)	stage = 4	34 (4.7%)	16 (3.5%)
Unknown	219 (30.0%)	93 (20.5%)	Unknown	52 (7.1%)	28 (6.2%)
Drinking, No. (%)			Mortality, No. (%)	609 (83.3%)	368 (81.1%)
No	425 (58.1%)	316 (69.6%)			
Yes	85 (11.6%)	45 (9.9%)			

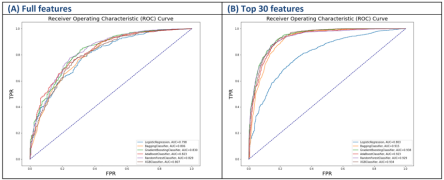


Figure 1. Receiver Operating Characteristic (ROC) Curve of various models

### Conclusions

Random forest was observed as the best model when using all features. Moreover, while choosing the top 30 features, Gradient Boosting Classifier was found with the highest AUC of 0.94.

In summary, the model developed using the Gradient Boosting Classifier algorithm had the highest AUC regardless of the mode and was the most suitable tool for NSCLC survival prediction among patients who underwent TKI treatment. In addition, using more types of data (especially laboratory and genomic test results) led to better predictive performance. Cancer stage, cancer size, gender, diagnosis age, and body mass index were the essential features for NSCLC survival prediction.

Contact: Jason C. Hsu, International Ph.D. program in Biotech and Health Management, College of Management, Taipei Medical University, Taipei, Taiwan; 11F, No.172-1, Sec. 2, Keelung Rd., Daan Dist., Taipei City 106, Taiwan (R.O.C.); E-mail: jasonhsu@tmu.edu.tw

**MONDAY** Development of Machine Learning models for Cancer Survival among Lung cancer patients with Tyrosine Kinase Inhibitors (TKIs) treatment (Alex PA. Nguyen, Phuc T. Nguyen, Min-Huei Hsu, Jason C. Hsu)





# #OHDSISocialShowcase This Week

## Extending the OMOP Standard Vocabulary to Include Botanical Natural Products

Sanya B. Taneja, Mary F. Paine, Sandra L. Kane-Gill, Richard D. Boyce

### INTRODUCTION

**OBJECTIVE:** extend the OMOP vocabulary to include natural products, their synonyms, phytoconstituents, and name variations to standardize the natural product reports in spontaneous reporting systems.

- Increase in consumption of natural products and/or dietary supplements has led to adverse event concerns.
- Spontaneous reporting systems (e.g., FAERS) can be used for natural product pharmacovigilance by identifying reports with natural products.
- Lack of interoperability in natural product data sources, coverage of synonyms, scientific names and common names, and ambiguity in natural product names are major challenges.

### METHODS



SQL Queries for custom vocabulary - concept, concept\_relationship tables

RxNorm mappings

### RESULTS

- 303 unique natural product Latin binomials
- 2,289 unique concepts in concept table
- 2,772 manually curated name variations for 65 natural products from FAERS
- Relationships: *napdi\_pt*, *napdi\_is\_pt\_of*, *napdi\_has\_const*, *napdi\_is\_const\_of*, *napdi\_spell\_vr*, *napdi\_is\_spell\_vr\_of*, *napdi\_np\_maps\_to*, *napdi\_const\_maps\_to*
- 47,601 reports matched to natural product names, 60,223 reports matched to natural product names & name variations, & 100,522 reports matched to natural product constituents.

303 botanical natural products, 2,289 concepts, and 2,772 name variations added to extended OHDSI vocabulary.

160,745 adverse event reports identified using terms for 65 natural products from the extended vocabulary.

Includes relationships to natural product constituents and RxNorm concepts.



Take a picture to download the full paper

### EXTENDED VOCABULARY

Table 1: concept table with green tea concepts.

concept_id	concept_name	vocab_id	concept_class_id
-7000189	Black tea [Camellia sinensis]	NAPDI	Green tea
-7000190	Green tea [Camellia sinensis]	NAPDI	Green tea
-7000191	Oolong tea [Camellia sinensis]	NAPDI	Green tea
-7000192	Tea [Camellia sinensis]	NAPDI	Green tea
-7000193	White Tea [Camellia sinensis]	NAPDI	Green tea
-7000293	Camellia sinensis [Camellia sinensis]	NAPDI	Green tea

Table 2: concept table with green tea constituents.

concept_name	constituent_name	concept_id
Green tea	EPICATECHIN	-7001895
Green tea	EPICATECHIN GALLATE	-7002175
Green tea	EPIGALLOCATECHIN	-7001785
Green tea	EPIGALLOCATECHIN GALLATE	-7002248
Green tea	GALLOCATECHIN	-7002061
Green tea	GALLOCATECHIN GALLATE	-7001793

Table 3: concept table with green tea name variations.

concept_name	name_variation	concept_id
Green tea	GUARANA GREEN TEA	-7004112
Green tea	CAMELLIA SINENSIS/PANAX GINSENG EXTRACT	-7004069
Green tea	APPLE CIDER VINEGAR + GREEN TEA SUPPLEMENT	-7003800
Green tea	UNSPECIFIED GREEN TEA EXTRACT SUPPLEMENT	-7002714
Green tea	TEA, GREEN (TEA, GREEN)	-7002713

Table 4: Green tea concepts mapped to RxNorm terms.

rxnorm_id	napdi_concept_id	rxnorm_concept	rxnorm_class
19121499	-7001008	GREEN TEA PREPARATION 25 MG	Clinical Drug Comp
1304239	-7001008	GREEN TEA LEAF EXTRACT	Ingredient
1304273	-7001008	GREEN TEA LEAF EXTRACT 1000 MG ORAL TABLET	Clinical Drug
1396861	-7001008	GREEN TEA EXTRACT 315 MG ORAL CAPSULE	Clinical Drug

Details: [github.com/dbmi-pitt/np-terminology-imports](https://github.com/dbmi-pitt/np-terminology-imports)

FAERS: FDA Adverse Event Reporting System

University of Pittsburgh

NAPDI Center of Excellence for Natural Product Drug Interaction Research



## Analyzing the Effect of Hypertension on Retinal Thickness Using Radiology

TUESDAY

Common Data Model (R-CDM) (Chul Hyoung Park, Rae Woong Park, Sang Jun Park, Da Yun Lee, Seng Chan You, Ki Hwang Lee)



# #OHDSISocialShowercase This Week

## Comparing the impact of clean windows across cohorts and databases

▲ PRESENTER: Rupa Makadia

### INTRO:

- Clean periods of observed person-time allow for the removal of prior exposures (conditions, drugs or procedures) within a cohort.
- The selection of this time can vary based on design of the study, prior knowledge, or random assignment by researchers.
- In this study we examine the trade-offs between various time-windows for a clean window to identify new events within a phenotype across a variety of databases.

### METHODS:

- Database
  - IBM MarketScan® Databases
  - Commercial Claims (CCAE)
- Phenotypes (10)
  - Acute myocardial infarction, myocarditis/pericarditis, deep vein thrombosis, pulmonary embolism, disseminated intravascular coagulation, non-hemorrhagic stroke, hemorrhagic stroke, cerebral venous thrombosis, peripheral arterial thrombosis, and thrombocytopenia
  - Clean window times (0, 14, 28, 90, 180, 270, 365)
  - Clean windows applied during cohort creation in ATLAS, 70 cohorts were created, 10 for each phenotype
  - Custom code was created to summarize the number of persons and events across each time window.

### RESULTS:

- Table 1 presents the number of patients identified (next to phenotype name), with the counts of persons by events (max=10). The total persons represents the total patients identified with 10 events.
- Phenotypes restricted by inpatient are banded in black.

## Clean windows should be empirically derived in cohort studies to eliminate prevalent cases for identification of new events

Table 1. Phenotypes by clean window and number of persons by event with a maximum of 10.

	Days	1	2	3	4	5	6	7	8	9	10	Total events	Total persons
Acute myocardial infarction (n=527,411)	0	371,577	62,467	24,327	13,554	6,200	4,576	2,334	1,294	773	513	81,813	423,993
	14	447,805	53,971	14,315	6,462	3,014	1,046	296	166	86	39	452,772	527,300
	28	462,567	48,307	11,280	5,466	1,075	335	141	53	32	23	463,518	527,358
	90	488,641	32,747	4,648	927	278	93	33	12	4	3	574,207	527,406
	180	496,228	22,536	3,275	567	141	52	17	7	2	3	545,623	527,411
	270	500,466	23,946	2,550	349	82	14	3	1	-	-	557,928	527,411
	365	503,873	21,280	1,904	222	45	4	-	-	-	-	553,543	527,411
Cerebral venous thrombosis (n=3,329)	0	2,212	732	215	79	29	22	17	8	3	2	3,334	3,319
	14	3,060	45	18	6	1	1	1	1	-	-	3,784	3,329
	28	3,081	199	43	3	3	-	-	-	-	-	3,433	3,329
	90	3,285	137	13	1	-	-	-	-	-	-	3,496	3,329
	180	3,259	63	7	-	-	-	-	-	-	-	3,496	3,329
	270	3,281	46	3	-	-	-	-	-	-	-	3,379	3,329
	365	3,299	30	-	-	-	-	-	-	-	-	3,339	3,329
Disseminated intravascular coagulation (n=22,667)	0	14,900	4,908	1,426	585	307	159	86	45	34	33	17,424	22,665
	14	20,751	1,577	243	56	18	9	7	2	2	2	25,140	22,665
	28	21,667	964	303	22	4	4	2	1	1	-	23,877	22,665
	90	22,284	364	32	3	1	1	-	-	-	-	23,095	22,667
	180	22,654	296	17	-	-	-	-	-	-	-	22,963	22,667
	270	22,517	139	10	1	-	1	-	-	-	-	22,829	22,667
	365	22,327	99	-	-	-	-	-	-	-	-	22,781	22,667
Non-hemorrhagic stroke (n=379, 953)	0	251,158	64,854	24,108	12,649	7,734	4,998	3,001	2,453	1,865	1,330	718,190	374,600
	14	314,736	42,913	12,334	4,665	2,246	1,174	607	297	238	165	491,947	379,243
	28	324,811	37,891	9,231	3,145	1,298	644	302	178	118	57	466,237	379,228
	90	344,836	28,114	5,076	1,283	369	154	54	23	6	2	420,033	379,949
	180	353,270	21,441	3,246	643	169	49	18	3	1	1	411,952	379,953
	270	358,314	18,905	2,351	377	-	-	-	-	-	-	404,951	379,953
	365	361,654	14,072	1,657	223	41	2	2	2	1	-	400,208	379,953
Hemorrhagic stroke (n=123,178)	0	73,152	25,413	9,750	5,067	2,978	1,681	1,215	887	598	454	255,200	121,395
	14	104,123	12,451	2,799	914	397	181	102	64	40	23	149,112	123,154
	28	111,491	8,977	1,677	456	181	81	42	19	21	11	138,957	123,156
	90	117,829	4,747	492	73	20	6	3	3	-	-	129,287	123,178
	180	119,979	2,975	187	53	3	1	-	-	-	-	126,443	123,178
	270	120,886	1,373	104	31	2	-	-	-	-	-	125,668	123,178
	365	121,429	660	373	-	-	-	-	-	-	-	123,004	123,178
Thrombocytopenia (n=1,157,870)	0	736,347	208,758	99,525	59,930	40,300	29,304	21,567	14,741	11,279	10,814	1,016,615	1,157,870
	14	796,446	186,179	86,318	50,807	32,549	21,959	15,776	11,753	8,687	6,727	2,768,501	1,157,870
	28	798,373	179,992	86,245	46,475	26,778	17,424	12,185	8,475	6,223	4,998	2,332,175	1,157,870
	90	861,842	141,852	61,145	29,488	16,197	9,512	5,854	3,745	2,511	1,719	1,789,701	1,157,870
	180	899,171	141,314	43,945	17,602	7,896	3,747	1,881	1,053	572	277	1,520,609	1,157,870
	270	990,580	122,008	30,294	9,422	3,371	1,382	508	242	76	51	1,394,838	1,157,870
	365	1,022,997	105,495	21,311	5,522	1,576	448	137	63	15	4	1,352,756	1,157,870
Deep vein thrombosis (n=785, 378)	0	346,081	119,529	67,203	46,311	31,643	21,564	20,049	14,365	13,588	11,351	2,880,484	785,378
	14	432,790	122,825	68,464	42,994	28,029	19,457	14,002	10,008	7,697	5,978	2,462,450	785,378
	28	478,515	130,111	64,803	35,560	21,333	13,923	9,297	4,724	5,100	3,925	1,838,306	785,378
	90	506,705	117,555	41,452	18,604	9,949	5,414	2,336	1,444	1,266	824	1,180,265	785,378
	180	652,404	93,492	25,795	9,695	4,044	1,825	847	409	192	91	995,882	785,378
	270	696,425	78,653	17,223	5,097	1,673	365	191	40	25	4	921,958	785,378
	365	703,876	63,716	11,993	2,834	729	149	46	8	2	-	885,701	785,378
Myocarditis/pericarditis (n=139,677)	0	87,723	28,214	11,173	5,896	3,246	2,019	1,402	864	706	510	297,728	139,677
	14	104,185	18,149	6,742	3,174	1,779	1,048	718	483	359	261	222,873	139,677
	28	112,463	15,951	5,240	2,344	1,317	776	498	336	271	168	200,007	139,677
	90	122,477	11,337	3,175	1,314	638	376	245	135	69	40	170,880	139,677
	180	127,880	8,707	2,017	728	295	136	51	32	15	12	157,412	139,677
	270	131,075	6,862	1,272	362	94	50	14	4	-	2	151,008	139,677
	365	133,007	5,568	861	187	56	13	4	-	-	-	145,811	139,677
Pulmonary embolism (n=437,437)	0	142,867	59,389	36,320	27,398	21,649	17,222	13,921	11,448	9,584	8,426	2,801,205	437,437
	14	214,479	63,785	40,001	28,961	20,019	14,590	10,391	7,814	5,892	4,740	1,634,995	437,437
	28	245,244	70,931	40,841	24,542	15,286	10,075	6,934	5,042	3,865	2,810	1,175,729	437,437
	90	237,798	49,921	25,781	11,824	6,212	3,160	1,797	1,247	791	451	480,233	437,437
	180	361,473	51,753	13,521	5,600	2,287	1,102	464	213	101	50	558,208	437,437
	270	382,542	43,119	9,322	2,437	742	382	202	98	4	4	437,250	437,437
	365	395,488	33,736	6,210	1,391	306	82	21	2	1	-	489,548	437,437
Peripheral arterial thrombosis (n=49,936)	0	14,287	5,547	3,752	1,761	970	520	322	198	159	111	111,992	49,936
	14	14,627	5,579	3,771	1,823	994	526	294	206	159	111	111,992	49,936
	28	14,643	7,743	2,149	829	377	141	92	31	20	11	88,141	49,936
	90	14,873	6,000	1,373	375	115	52	8	1	-	-	88,141	49,936
	180	14,873	6,000	826	158	34	3	-	-	-	-	88,141	49,936
	270	14,873	6,000	826	158	34	3	-	-	-	-	88,141	49,936
	365	14,873	6,000	826	158	34	3	-	-	-	-	88,141	49,936

### RESULTS:

- 94% (496,228) of people have a single acute myocardial infarction diagnosis within 180 days whereas only 70% (371,577) people have a single acute myocardial infarction diagnosis with no clean window.
- Phenotypes that occur in any place show the most variability of events compared to those phenotypes with an inpatient restriction.
- Rarer events (cerebral venous thrombosis, disseminated intravascular coagulation) encapsulate most events within 90 days.

### CONCLUSIONS:

- This analysis highlights the need to consider this diagnostic when building a cohort. The results from these phenotypes show wide variability among them and highlight the need for careful consideration in a study design.
- There is a possibility that new events may be missed especially in cases such as deep vein thrombosis where the clean window is 365 days.
- Plausibility of events occurring in observational data and identification of new events or continuation of care from a prior event should be carefully evaluated for conditions that could have recurrence.

Rupa Makadia<sup>1</sup>, Kevin Haynes<sup>1</sup>, Patrick Ryan<sup>1</sup>  
Janssen Research & Development, Titusville, NJ



WEDNESDAY Comparing the impact of clean windows across cohorts and databases (Rupa Makadia, Kevin Haynes, Patrick Ryan)





# #OHDSISocialShowcase This Week



## Assessing and Benchmarking Data Quality and Diversity in the *All of Us* Research Program

Lina Suleiman, PhD;<sup>1</sup> Karthik Natarajan<sup>2</sup>, PhD; Kayla Marginean, MS;<sup>1</sup> Robert Carroll, PhD;<sup>1</sup> Paul Harris, PhD<sup>1</sup>

<sup>1</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center; <sup>2</sup>Department of Biomedical Informatics, Columbia University, New York, NY

### Objective

Assessing the data quality of Electronic Health Records (EHR) in the *All of Us* Research Program using OHDSI tools and libraries

### Background/Introduction

- Quantifying and benchmarking data quality in research clinical repositories is crucial
  - Ensuring the utility of the data
  - Improving the reproducibility and research credibility of research
- Quality assessment methods:
  - Calculating clinical data quality metrics: completeness, plausibility, conformance
  - Replication: replicating existing clinical studies (e.g., phenotype algorithms)
  - Benchmarking against expected published numbers
- The *All of Us* Research Program is a national initiative collecting Electronic Health Records, surveys, and genetic data
- The *All of Us* Research Program collects data from over 50 EHR sites
- Aim: Assessing the *All of Us* EHR data quality by quantifying: completeness, plausibility, conformance, and the prevalence of phenotypes

### Methods

- Dataset: *All of Us* controlled tier launched in March 2022, we extracted participants who have EHR
- Quality assessment methods:
  - Phenotype replications:
    - Assessing the prevalence of phenotypes
    - Implemented OHDSI phenotype library to extract 212 phenotypes using 763 algorithms
  - Benchmarking: Comparing the prevalence of phenotypes in *All of Us* to the prevalence reported by the Center for Disease and Control (CDC)
- Quality metrics: Running OHDSI data quality dashboard package and other ad-hoc metrics to calculate the following:
  - Plausibility: the extent to which the values agree with internal and external knowledge
  - Conformance: the percentage of the dataset that complies with standards and constraints
  - Completeness: the percentage of data that is expected to be present
  - Additional completeness metrics: the percentage of participants who have core measurements: height, weight, Body Mass Index (BMI), cholesterol, and heart rate per EHR site

### Results

- Dataset: *All of Us* Research Program controlled tier included 331,382 participants
  - 76% were considered underrepresented in biomedical research
  - White participants: 55.81% compared to 76.30% white in the US general population
  - Female participants: 50%
  - Have any EHR data: 224,507 (67.7%) participants

### Results

- Replication (Figure 1):
  - 223,018 participants with at least one of the 212 EHR-based phenotypes
  - Racial distributions varied: 413 (55%) of the cohorts had 60% or lower white participants
- Benchmarking:
  - Mean prevalence of the US leading causes of death phenotypes: close to or slightly higher than the prevalence reported by CDC (except for Alzheimer's disease, suicide, influenza) as Figure 2 shows
- Quality metrics:
  - OHDSI quality dashboard (Figure 3)
  - Common quality problems:
    - Plausibility: out-of-range values, gender-specific conditions and procedures concepts
    - Conformance: Non-standard/non-existing
    - Completeness: Null or zero values
  - Core measurement completeness per site measured by the percentage of participants:
    - Height: 72.0% (3%-100%)
    - Weight: 76.5% (3%-100%)
    - BMI: 76.8% (11%-100)
    - Heart rate: 80.5% (11%-100)

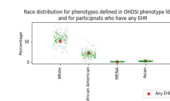


Figure 1. Race distribution in phenotype algorithms applied on the *All of Us* Research Program dataset compared to the race distribution for participants with any EHR

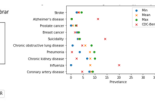


Figure 2. The prevalence of phenotypes that are the leading causes of death as reported by the CDC. We reported the maximum, minimum, and mean prevalence for all algorithms per phenotype

	Min	Max	Mean
Plausibility	0.00	0.00	0.00
Conformance	0.00	0.00	0.00
Completeness	0.00	0.00	0.00
Height	0.00	1.00	0.72
Weight	0.00	1.00	0.77
BMI	0.00	1.00	0.77
Heart rate	0.00	1.00	0.81

Figure 3. The plausibility, conformance, and completeness metrics in the *All of Us* Research Program

### Conclusions

- Data quality and diversity are essential factors that can improve clinical research reproducibility in major clinical research repositories such as *All of Us* Research Program
- We used OHDSI tools to assess data quality using: replication, benchmarking, and quality dimension metrics (plausibility, conformance, completeness)
- Our analysis demonstrated the diversity of the *All of Us* EHR
  - Non-white participants' representation is high compared to other research repositories
- Quality aspects to further investigate:
  - Identifying possible reasons for having different prevalence values compared to CDC
  - Simpler versions of phenotype algorithms
  - Site recruitment: might influence the disease prevalence within *All of Us* Research Program (recruiting from breast cancer clinic)
- Low conformance and core measurement completeness: mapping differences in EHR sites



Email: [lina.suleiman@vumc.org](mailto:lina.suleiman@vumc.org)  
Twitter: [@LinaSuleiman](https://twitter.com/LinaSuleiman)

THURSDAY

Assessing and Benchmarking Data Quality and Diversity in the *All of Us* (Lina Suleiman, Karthik Natarajan, Kayla Marginean, Robert Carroll, Paul Harris)



# #OHDSISocialShowcase This Week

## Moving OMOP to the Cloud With DBT and Snowflake

Roger Carlson - Spectrum Health  
Matthew Phad - Spectrum Health  
Sam Martin - Spectrum Health  
Grand Rapids, Michigan



Poster and References Available Online



1. The Snowflake Platform
2. Cloud computing with AWS
3. What is dbt?
4. Snowflake Security



Roger Carlson



Matt Phad



Sam Martin

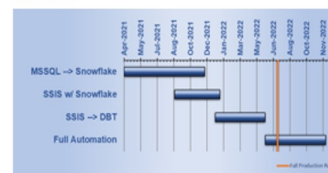
### Background

- Spectrum Health has been evaluating the use of modern cloud-computing for its analytical processing needs.
- Cloud-computing Advantages
  - Improved speed, scalability, security, data-sharing capabilities.
  - Unique opportunity to invest in modern, open-source tools and methodologies through the adoption of a unified tool set.
- Envisioned Platform (Proof-of-concept)
  - Amazon Web Services (AWS)
  - All of Us® Research Program (NIH), which transforms EPIC® Clarity EHR data into the OMOP format.
  - Relatively small-scale project (~12,000 patients)
  - Complex and robust ETL process,

### Methods

- Legacy System
  - Microsoft SQL Server database (on-prem)
  - Tools used: SQL Server Management Studio, SQL Server Integration Services, Visual Studio, Redcap, Oracle SQL Developer, R-Studio, Tortoise SVN, and Microsoft Access.
- Proposed system
  - Snowflake database on AWS platform
  - Reduced toolset: Snowflake, DBT, GitHub, REDCap, VSCode, DBeaver, and R Studio.

### Timeline



### Run Times

	MSSQL	SNOWFLAKE
SSIS	1h:50m	20m
DBT	n/a	20m

### Workflow Comparison

Workflow Process	SSIS	DBT
API download from REDCap*	❌	✅
Extract data from EPIC® Clarity	✅	✅
Transform Clarity data into OMOP	✅	✅
Built-in and custom testing features	❌	✅
Referential Integrity	✅	✅
Curation Reporting	❌	✅
Automated data export	✅	✅
Automated SMTP transfer	❌	❌
Transfer from S3 bucket to Google*	❌	✅

### Development

- **Phase 1:** Move database to cloud-based database (MSSQL → Snowflake)
  - Timeframe: Apr 2021 – Dec 2021
  - Scope: 60 tables, 35 views, 262 queries
  - OMOP v5.2 to v5.3.1 upgrade
- **Phase 2:** Convert SSIS project to work with Snowflake
  - Timeframe: Sept 2021 – Dec 2021
  - Scope: 19 packages, 171 tasks, 262 queries
- **Phase 3:** Move workflow process to open-source tool (SSIS → DBT)
  - Timeframe: Jan 2022 – May 2022
  - Scope: 578 steps, 416 models (347 views, 29 tables), 202 tests, 433 macros, 95 sources
- **Full Production Run:**
  - July 7, 2022 (21 person months total)
- **Phase 4:** Full Integration and automation, i.e., delivery of OMOP files from AWS to Google Bucket.
  - Timeframe: Jun 2022 – Nov 2022\* (envisioned)

### Conclusion

- Conversion from SQL Server using SSIS to Snowflake using DBT was timely and effective. Our result is a more robust platform featuring a collaborative workflow built on modern toolsets.
- Snowflake/DBT significantly outperforms MSSQL/SSIS.
- Snowflake is effectively unlimited in terms of scalability and complies with a wide range of compliance standards including HIPAA/HITRUST, SOC 1 Type II and SOC 2 Type II.

FRIDAY

Moving OMOP to the cloud with DBT and Snowflake (Roger Carlson, Matthew Phad, Samuel Martin)



# 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?**





# Collaborations for Strategic Opportunities



**Anna Ostropolets**

Data Scientist, Odysseus Data Services, Inc.  
PhD Graduate, Columbia University



**Clair Blacketer**

Director, Janssen Research and Development, Inc.



**Patrick Ryan**

Vice President, Observational Health Data Analytics, Janssen Research and Development, Inc.; Adjunct Assistant Professor, Columbia University