# Upcoming OHDSI Community Calls

| Date | Topic |
|------|-------|
| Nov. 15 | Open Network Studies |
| Nov. 22 | 10-Minute Tutorials |
| Nov. 29 | Workgroup Updates |
| Dec. 6 | Fall Publications |
| Dec. 13 | How Did We Do In 2022? |
| Dec. 20 | Holiday-Themed Final Call of 2022 |

# Upcoming OHDSI Community Calls

| Date | Topic |
|------|-------|
| Nov. 15 | Open Network Studies |
| Nov. 22 | 10-Minute Tutorials |
| Nov. 29 | Workgroup Updates |
| Dec. 6 | Fall Publications |
| Dec. 13 | How Did We Do In 2022? |
| Dec. 20 | Holiday-Themed Final Call of 2022 |

# OHDSI Shoutouts!

👏

Congratulations to the team of **Wallis C.Y. Lau, Carmen Olga Torre, Kenneth K.C. Man, Henry Morgan Stewart, Sarah Seager, Mui Van Zandt, Christian Reich, Jing Li, Jack Brewster, Gregory Y.H. Lip, Aroon D. Hingorani, Li Wei, and Ian C.K. Wong** on the publication of **Comparative Effectiveness and Safety Between Apixaban, Dabigatran, Edoxaban, and Rivaroxaban Among Patients With Atrial Fibrillation** in Annals of Internal Medicine.

## Annals of Internal Medicine®

Search Journal

LATEST    ISSUES    IN THE CLINIC    JOURNAL CLUB    MULTIMEDIA    CME / MOC    AUTHORS / SUBMIT

Original Research

**Comparative Effectiveness and Safety Between Apixaban, Dabigatran, Edoxaban, and Rivaroxaban Among Patients With Atrial Fibrillation**

A Multinational Population-Based Cohort Study

Wallis C.Y. Lau, PhD* iD,   Carmen Olga Torre, MSc* iD,   Kenneth K.C. Man, PhD iD,   ... See More +

Author, Article, and Disclosure Information

https://doi.org/10.7326/M22-0511

Eligible for CME Point-of-Care

PDF  |  FULL  |  Tools  |  Share

**Background:**

Current guidelines recommend using direct oral anticoagulants (DOACs) over warfarin in patients with atrial fibrillation (AF), but head-to-head trial data do not exist to guide the choice of DOAC.

# OHDSI Shoutouts! 👏

Congratulations to the team of **Xiao Wang, Wenwang Rao, Xueyan Chen, Xinqiao Zhang, Zeng Wang, Xianglin Ma, Qinge Zhan** on the publication of **The sociodemographic characteristics and clinical features of the late-life depression patients: results from the Beijing Anding Hospital mental health big data platform** in BMC Psychiatry.

**BMC Psychiatry**

**RESEARCH**  **Open Access**

# The sociodemographic characteristics and clinical features of the late-life depression patients: results from the Beijing Anding Hospital mental health big data platform

Xiao Wang[1], Wenwang Rao[2], Xueyan Chen[1], Xinqiao Zhang[1], Zeng Wang[1], Xianglin Ma[1] and Qinge Zhang[1*]

**Abstract**

**Background:** The sociodemographic characteristics and clinical features of the Late-life depression (LLD) patients in psychiatric hospitals have not been thoroughly studied in China. This study aimed to explore the psychiatric outpatient attendance of LLD patients at a psychiatric hospital in China, with a subgroup analysis, such as with or without anxiety, gender differences.

**Methods:** This retrospective study examined outpatients with LLD from January 2013 to August 2019 using data in the Observational Medical Outcomes Partnership Common Data Model (OMOP-CDM) in Beijing Anding Hospital. Age, sex, number of visits, use of drugs and comorbid conditions were extracted from medical records.

**Results:** In a sample of 47,334 unipolar depression patients, 31,854 (67.30%) were women, and 15,480 (32.70%) were men. The main comorbidities of LDD are generalized anxiety disorder (GAD) (83.62%) and insomnia (74.52%).Among patients with unipolar depression, of which benzodiazepines accounted for the largest proportion (77.77%), Selective serotonin reuptake inhibitors (SSRIs) accounted for 59.00%, a noradrenergic and specific serotonergic antidepressant (NaSSAs) accounted for 36.20%. The average cost of each visit was approximately 646.27 yuan, and the cost of each visit was primarily attributed to Western medicine (22.97%) and Chinese herbal medicine (19.38%). For the cost of outpatient visits, depression comorbid anxiety group had a higher average cost than the non-anxiety group ($p < 0.05$). There are gender differences in outpatient costs, men spend more than women, for western medicine, men spend more than women, for Chinese herbal medicine, women spend more than men (all $p < 0.05$). The utilization rate of SSRIs and benzodiazepines in female patients is significantly higher than that in male patients ($p < 0.05$).

**Conclusion:** LLD patients are more commonly women than men and more commonly used SSRIs and NaSSAs. Elderly patients with depression often have comorbid generalized anxiety. LLD patients spend most of their visits on medicines, and while the examination costs are lower.

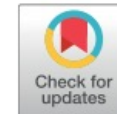**Keywords:** LLD, Outpatient, Antidepressants

# OHDSI Shoutouts!

👏

Congratulations to the team of **Tianchu Lyu, Chen Liang, Jihong Liu, Berry Campbell, Peiyin Hung, Yi-Wen Shih, Nadia Ghumman, Xiaoming Li, and members of the National COVID Cohort Collaborative Consortium** on the publication of **Temporal Events Detector for Pregnancy Care (TED-PC): A rule-based algorithm to infer gestational age and delivery date from electronic health records of pregnant women with and without COVID-19** in PLOS ONE.

## Temporal Events Detector for Pregnancy Care (TED-PC): A rule-based algorithm to infer gestational age and delivery date from electronic health records of pregnant women with and without COVID-19

Tianchu Lyu[1], Chen Liang[1*], Jihong Liu[2], Berry Campbell[3], Peiyin Hung[1], Yi-Wen Shih[1], Nadia Ghumman[1], Xiaoming Li[4], on behalf of the National COVID Cohort Collaborative Consortium[¶]

1 Department of Health Services Policy and Management, Arnold School of Public Health, University of South Carolina, Columbia, South Carolina, United States of America, 2 Department of Epidemiology & Biostatistics, Arnold School of Public Health, University of South Carolina, Columbia, South Carolina, United States of America, 3 Department of Obstetrics and Gynecology, School of Medicine, University of South Carolina, Columbia, South Carolina, United States of America, 4 Department of Health Promotion Education and Behaviors, Arnold School of Public Health, University of South Carolina, Columbia, South Carolina, United States of America

¶ Membership of the National COVID Cohort Collaborative Consortium is provided in the Acknowledgments.
* cliang@mailbox.sc.edu

## Abstract

### Objective

Identifying the time of SARS-CoV-2 viral infection relative to specific gestational weeks is critical for delineating the role of viral infection timing in adverse pregnancy outcomes. However, this task is difficult when it comes to Electronic Health Records (EHR). In combating the COVID-19 pandemic for maternal health, we sought to develop and validate a clinical information extraction algorithm to detect the time of clinical events relative to gestational weeks.

# OHDSI Shoutouts!

# OHDSI Shoutouts! 👏

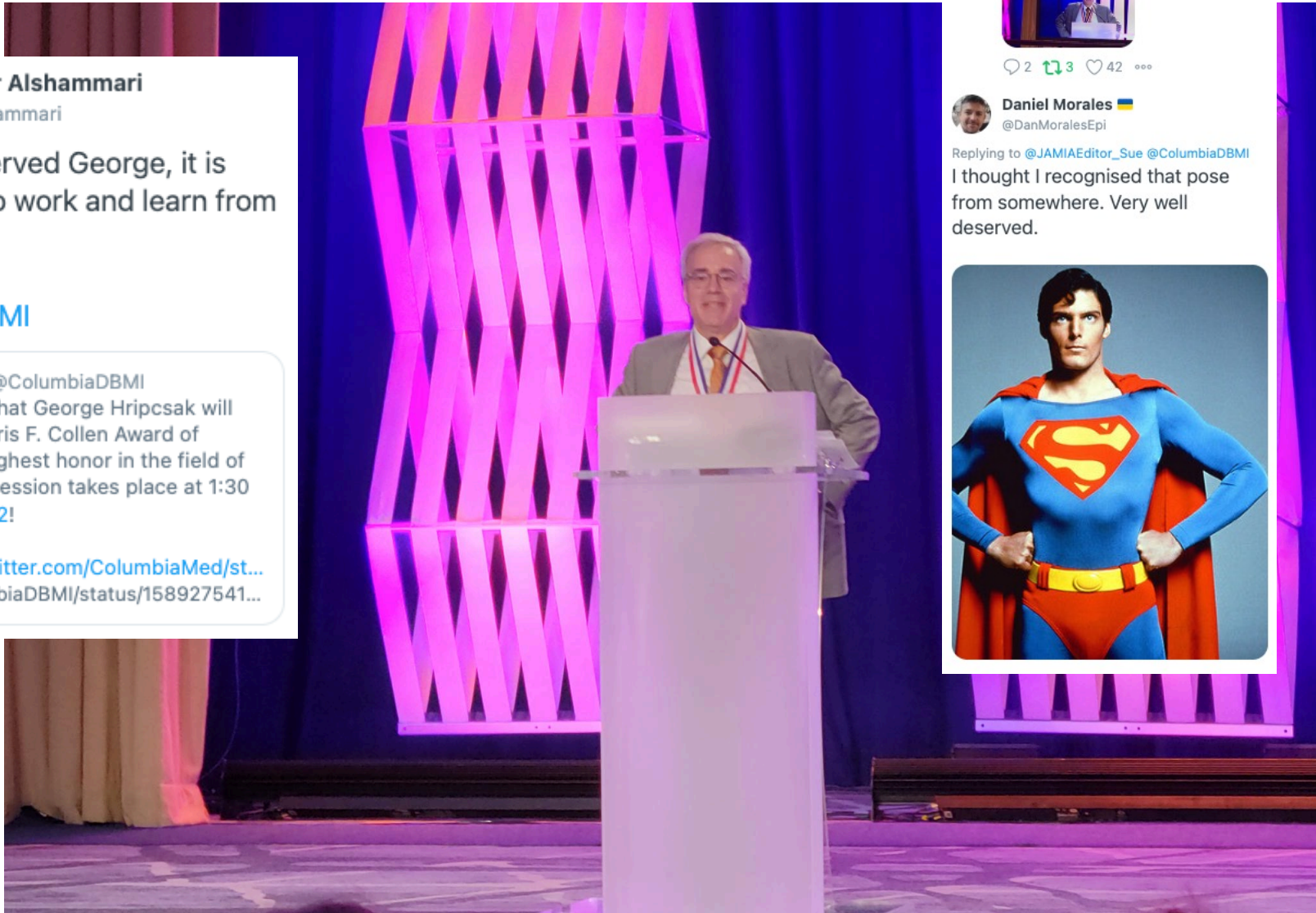**Dr. Thamir Alshammari**
@T_M_Alshammari

Well we'll deserved George, it is always great to work and learn from you.
@OHDSI
@ColumbiaDBMI

**Columbia DBMI** @ColumbiaDBMI
Today is the day that George Hripcsak will receive 2022 Morris F. Collen Award of Excellence, the highest honor in the field of informatics. The session takes place at 1:30 during #AMIA2022!

👏🎉👏🎉👏🎊 twitter.com/ColumbiaMed/st...
twitter.com/ColumbiaDBMI/status/158927541...

**Kristin Kostka** @kricketchirps    15h
The most humble leader I know! So awesome to see our fearless @OHDSI leader recognized. I know I wouldn't be where I am today without George's mentorship! He's a true star in the observational health data science and informatics field. 🌟

**Columbia DBMI** @ColumbiaDBMI
Congrats to DBMI chair George Hripcsak on receiving the 2022 Morris F. Collen Award of Excellence yesterday! #AMIA2022 @Columbia @ColumbiaPS @DataSciColumbia @AMIAinformatics

# OHDSI Shout

**Dr. Thamir Alshammari**
@T_M_Alshammari

Well we'll deserved George, it is always great to work and learn from you.
@OHDSI
@ColumbiaDBMI

**Columbia DBMI** @ColumbiaDBMI
Today is the day that George Hripcsak will receive 2022 Morris F. Collen Award of Excellence, the highest honor in the field of informatics. The session takes place at 1:30 during #AMIA2022!

👏🎉👏🎊👏🎉 twitter.com/ColumbiaMed/st...
twitter.com/ColumbiaDBMI/status/158927541...

**Suzanne Bakken** @JAMIAEditor_Sue 1d
My informatics boss wins #ACMI highest award Morris F. Collen Award @ColumbiaDBMI !

○ 2   ⟳ 3   ♡ 42   ○○○

**Daniel Morales** 🇺🇦
@DanMoralesEpi

Replying to @JAMIAEditor_Sue @ColumbiaDBMI
I thought I recognised that pose from somewhere. Very well deserved.

**Kristin Kostka** @kricketchirps    15h
The most humble leader I know! So awesome to see our fearless @OHDSI leader recognized. I know I wouldn't be where I am today without George's mentorship! He's a true star in the observational health data science and informatics field. 🌟

**Columbia DBMI** @ColumbiaDBMI
Congrats to DBMI chair George Hripcsak on receiving the 2022 Morris F. Collen Award of Excellence yesterday! #AMIA2022 @Columbia @ColumbiaPS @DataSciColumbia @AMIAinformatics

# 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 so we can share during this call and on our social channels.
Let's work together to promote the collaborative work happening in OHDSI!

# Three Stages of The Journey

## Where Have We Been?
## Where Are We Now?
## Where Are We Going?

# Upcoming Workgroup Calls

| Date | Time (ET) | Meeting |
|---|---|---|
| Tuesday | 12 pm | Common Data Model Vocabulary Subgroup |
| Tuesday | 6 pm | Eyecare & Vision Research |
| Wednesday | 9 am | Patient-Level Prediction |
| Wednesday | 10 am | FHIR and OMOP Digital Quality Measurements Subgroup (ZOOM) |
| Wednesday | 11 am | Open-Source Community |
| Wednesday | 2 pm | Natural Language Processing |
| Thursday | 12 pm | FHIR and OMOP Oncology Subgroup |
| Thursday | 7 pm | Dentistry |
| Friday | 9 am | GIS – Geographic Information System Development |
| Friday | 9 am | Phenotype Development & Evaluation |
| Friday | 10:15 am | Clinical Trials |
| Friday | 10 pm | China Chapter |
| Monday | 10 am | Africa Chapter |
| Monday | 11 am | Early-Stage Researchers |

## ohdsi.org/upcoming-working-group-calls/

# Open-Source Community WG Meeting

Please join the next Open-Source Community WG meeting, which will include a presentation from **Laurie Arp** around supporting open-source sustainability planning.

Laurie Arp's professional interests focus on the intersection of collections, technology, and people. As the Director of DuraSpace Community Supported Software, Laurie directs community supported open-source programs housed at LYRASIS including ArchivesSpace, CollectionSpace, DSpace, Fedora, and VIVO.

## Wednesday, 11 am ET

# 2022 OHDSI APAC Symposium

## Day 1 (Nov. 12) — Tutorial Workshop

8:30 – 9:00 · Registration
9:00 – 9:30 · Overview of the OHDSI Journey: where are we going
9:30 – 10:20 · OMOP Common Data Model and vocabulary
10:20 – 10:30 · Break
10:30 – 11:20 · ETL a source database into OMOP CDM
11:20 – 11:30 · Break
11:30 – 12:20 · Creating cohort definitions
12:20 – 13:30 · Lunch
13:30 – 14:20 · Phenotype evaluation
14:20 – 14:30 · Break
14:30 – 15:20 · Characterization
15:20 – 15:30 · Break
15:30 – 16:20 · Estimation
16:20 – 16:30 · Break
16:30 – 17:20 · Prediction
17:20 – 17:30 · Recap of the OHDSI Journey, where do we go from here

**Register for Day 1 Here**

### Day 1 Registration Fees (In-Person)
International Student/Trainee: $30
International Academia/Government: $70
International Industry/Corporate: $170
Local Registrants: Free

### 2022 APAC OHDSI Symposium
**Nov. 12 - 13 · Taipei Medical University**

## Day 2 (Nov. 13) — Main Conference

08:00 – 09:00 · Registration & Light Breakfast
09:00 – 09:20 · Welcome Session
09:20 – 09:40 · Group Photo

*Session 1: Envisioning of OHDSI Global & OHDSI APAC*
09:40 – 10:00 · Keynote – OHDSI Global Presentation
10:00 – 10:20 · OHDSI APAC Introduction
10:20 – 10:30 · Break

*Session 2: The Implication Experiences in OHDSI Region*
10:30 – 11:30 · Researches in OHDSI APAC
11:30 – 11:45 · Researches using Taiwan National Data
11:45 – 12:00 · Researches using TMUCRD Data
12:00 – 13:00 · Lunch & Poster Presentation

*Session 3: The Challenges of Research in OHDSI APAC*
13:00 – 14:00 · Panel – Standardization & Common Data Models
14:00 – 15:00 · Panel – APAC Regional Adaption to Standardization
15:00 – 15:15 · Break
15:15 – 16:15 · Poster & Networking Session
16:15 – 17:00 · Closing Remarks

**Register for Day 2 Here**

### Day 2 Registration Fees (In-Person)
International Student/Trainee: $50
International Academia/Government: $100
International Industry/Corporate: $200
Local Registrant: Free

### Day 2 Registration Fees (Virtual)
International Student/Trainee: $25
International Academia/Government: $50
International Industry/Corporate: $100
Local Registrant: Free

**ohdsi.org/2022apacsymposium**

# 2022 OHDSI APAC Symposium



**Open call**
Open call for data partners (October 12 – November 11)

**Harmonisation**
Initiation of the data harmonisation.

| SEP | OCT | NOV | DEC | JAN | FEB | MAR | APR |

2022 | 2023

**Evaluation**
Evaluation of all applications by our committee of both internal and external experts.

**SME linking**
Identification and linking up with the SME of choice.

**Agreement**
Grant awarding and signing of grant agreement.

**Ehden.eu**

# #OHDSISocialShowcase This Week



**MONDAY**

**PHAROS, Platform for Harmonizing and Accessing Data in Real-time on Infectious Disease Surveillance Based on OMOP-CDM in Korea** (Chungsoo Kim, Jimyung Park, Byungjin Choi, Seongwon Lee, Rae Woong Park)

# #OHDSISocialShowcase This Week

## Understanding Circe-be Logic Through Capr for Generating Complex Cohort Definitions

AUTHOR
Martin Lavallee

## 1 Introduction

### 1.1 ATLAS

Typically, we define cohort definitions for OHDSI studies using ATLAS. ATLAS has several benefits, in particular having a nice user interface to visual the cohort definition we are trying to create. However, there are times when ATLAS can be a bit tedious particularly when we must create several cohort definitions with a similar structure (template). We can deal with this situations by copying and pasting, however this can lead to errors in cohort logic and can also be quite time consuming.

**TUESDAY** — **Understanding circe-be logic through Capr for generating complex cohort definitions** (Martin Lavallee, Adam Black, Asieh Golozar)

# #OHDSISocialShowcase This Week

Characterization of first-line treatment for Breast Cancer and Multiple Myeloma using Electronic Health Record and Claims Databases

Maura Beaton[1], Matthew Spotnitz[1], Thomas Falconer[1], Melissa Accordino[2], Divaya Bhutani[2], Alison Callahan[3], Nigam Shah[3], Andrew Williams[4], Karthik Natarajan[1]

[1]Columbia University Irving Medical Center, Department of Biomedical Informatics; [2]Columbia University Irving Medical Center, Department of Medicine; [3]Stanford University Department of Biomedical Data Science; [4]Tufts University Department of Biomedical Informatics

## Background

Cancer treatment has been shown to vary over time and geography[1,2]. Numerous factors have been proposed to explain these variations such as patient-level factors which include age, comorbidities, insurance coverage[3]. Treatment variation has also been associated the clinical management of cancer and varying rates of adoption of new treatments[4]. As a first step to better understanding cancer treatment variation, the current study sought to characterize temporal and geographic variations in first-line treatments for two cancers: breast cancer and multiple myeloma.

## Methods

### Cancer Phenotypes
For each cancer, we created a cohort of adult patients who received a cancer diagnosis within 90 days following a biopsy.

Index Event
- Breast Cancer: Breast biopsy procedure code
- Multiple Myeloma: Bone marrow biopsy procedure code

### Analysis
Characterization of all intervention types patients received within 1-year following their cancer diagnosis.

### Data Sources
Analysis was run on the following 6 databases, which had been converted to OMOP version 5.3.1.

EHR Data:
- Columbia University Irving Medical Center (CUIMC)
- Stanford University Medical Center (Stanford),
- Tufts Research Data Warehouse OMOP (Tufts),

Claims Data:
- IBM MarketScan Commercial Claims and Encounters (CCAE)
- IBM MarketScan Medicare Supplemental Beneficiaries (MDCR)
- IBM MarketScan Multi-state Medicaid (MDCD).

## Results: Database Counts

Table 1. Number of patients with each cancer by database.

| Database | Breast Cancer | Multiple Myeloma |
|---|---|---|
| CUIMC | 5,165 | 1,014 |
| Stanford | 6,696 | 1,127 |
| Tufts | 392 | 134 |
| CCAE | 198,575 | 17,499 |
| MDCR | 23,986 | 9,853 |
| MDCD | 21,438 | 4,054 |

Contact: mb4023@cumc.columbia.edu

## Results: Breast Cancer



Figure 1. Percent distribution of breast cancer interventions, by year and database
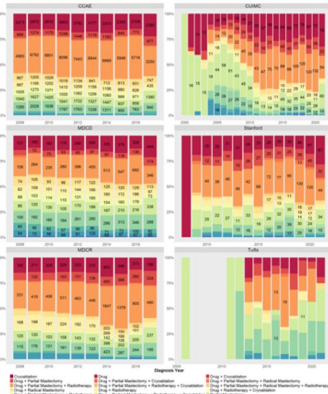
## Results: Multiple Myeloma



Figure 2. Percent distribution of multiple myeloma interventions, by year and database

## Conclusions

Across all databases, the majority of breast cancer patients received treatment that was consistent with the standard of care for HR+ breast cancer (surgical intervention + drug therapy or radiotherapy). Similarly, the majority of multiple myeloma patients received care that was consistent with National Comprehensive Cancer (NCCN) guidelines (initial systemic therapy, with some patients receiving an autologous stem cell transplant within one year following diagnosis).

Next steps for this work will be to develop an OMOP-based algorithm to detect cancer treatment regimens and compare the treatment regimens patients receive to the recommended regimens for their cancer. This work will also include linking Surveillance, Epidemiology and End Results (SEER) data to OMOP to ingest data about cancer stage and grade.
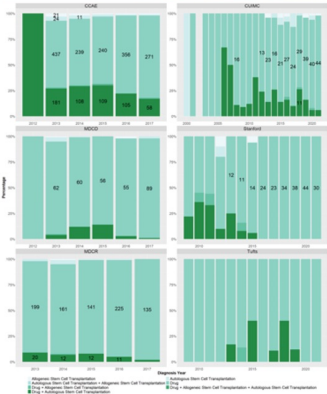
## References

1. Caram MEV, Estes JP, Griggs JJ, Lin P, Mukerjee B. Temporal and geographic variation in systemic treatment of advanced prostate cancer. BMC Cancer 2018;18:258
2. Derks MGM, Bastiaannet E, Kinderlen M, et al. Variation in treatment and survival of older patients with non-metastatic breast cancer in five European countries: A population-based cohort study from the EURECCA breast cancer group. Br J Cancer 2018;119(1):121-9
3. Tariman J, Berry D, Cochrane B, Doorenbos A, Schepp K. Physician, patient and contextual factors affecting treatment decisions in older adults with cancer: A literature review. Oncol Nurs Forum 2012; 39(1):E70-83.
4. Møller H, Coupland VH, Tataru D, et al. Geographic variations in the use of cancer treatments are associated with survival of lung cancer patients. Thorax 2018;73:530-537.

**WEDNESDAY** — Characterization of first-line treatment for Breast Cancer and Multiple Myeloma using Electronic HealthRecord and Claims Databases (Maura Beaton, Matthew Spotnitz, Thomas Falconer, Melissa Accordino,DivayaBhutani, Alison Callahan, Nigam Shah, Jake Gillberg, Andrew Williams, Karthik Natarajan)

# #OHDSISocialShowcase This Week

**The Seasonality Score: A Quantitative Complement to Qualitative Seasonality Assessment.**

PRESENTERS:
Anthony Molinaro

**Introduction:**
- Methods for seasonality classification of time series have been developed independently by researchers working in disparate fields.
- Consequently, these methods have been shown to be mutually discordant, thus limiting generalizability.
- Additionally, seasonality methods that assess qualitative aspects of a time series have difficulty yielding quantitative insight.

**Methods:**
- The OHDSI package ACHILLES is used for data retrieval and aggregation.
- The OHDSI package CASTOR is used for time series creation and metric computation.
- The seasonality score metric was implemented as part of the CASTOR R package to provide a quantitative method of characterizing seasonality.

## Quantitative Seasonality

**How do you know:**

If one time series is more or less seasonal than another?

If a time series is becoming more or less seasonal?

What the most seasonal events in a given database are?

If a time series is truly seasonal when methods disagree?

Database: IBM CCAE
Seasonality Score: 0.67
Concept_id: 40213152
Seasonal, quadrivalent, recombinant, injectable influenza vaccine, preservative free

Database: OPTUM EHR
Seasonality Score: 0.05
Concept_id: 42872722
Severe major depression

**Algorithm:**
Let u = 1/12 be a strictly non-seasonal proportion.
Let w = $(11 \times (1/12) + (1 - 1/12))$ be the normalizing value.
Let $\mathbf{M} = M_{mx12}$ be the time series.
Let $\mathbf{1}_{12}$ be a summing vector.
Let $\mathbf{1}_m$ be a summing vector.
Let $y = \mathbf{1}_m^T \mathbf{M}$ be the monthly sum over all years.
Let $g = \mathbf{1}_m^T \mathbf{M} \mathbf{1}_{12}$ be the grand sum.
Let $\mathbf{p} = y^T/g$ be the monthly proportion over all years.
Let d = $\mathbf{1}_m^T |\mathbf{p}\text{-u}|$ be the total deviation from strict non-seasonality.
Let s = d/w be the seasonality score.

Anthony Molinaro,
Frank DeFalco

**Results:**
- A quantitative seasonality score was established to be a complement to existing qualitative methods.
- The seasonality score provides a distribution-free metric that facilitates quantitative characterization and comparison.
- The seasonality score is a numeric value between 0 and 1 (inclusive), that is currently designed to quantify monthly seasonality.
- The seasonality score for all event table domains was computed for fifteen databases converted to the OMOP CDM.

**THURSDAY** — The Seasonality Score: A Quantitative Complement to Qualitative Seasonality Assessment (Anthony Molinaro, Frank DeFalco)

# #OHDSISocialShowcase This Week

## Development of Lung Cancer Survival Prediction Models Based on Real-world Data and Machine Learning

Jason C. Hsu; Phung-Anh Nguyen; Phan Thanh Phuc; Tsai-Chih Lo; Min-Huei Hsu; Chi-Tsun Cheng; Tzu-Hao Chang; Cheng-Yu Chen
Taipei Medical University, Taiwan

Jason C. Hsu     Phung-Anh Nguyen     Phan Thanh Phuc     Chi-Tsun Cheng

## Abstract

**Background**
The development of disease risk and prognosis prediction models using machine learning or deep learning algorithms with big data is a major area of academic research based on AI in the medical field. Various researchers have used machine learning or deep learning algorithms to develop lung cancer risk and prognosis prediction models.1-6

**Objectives**
The purpose of this study was to use clinical real-world data with multiple attributes and multiple machine learning algorithms to establish a prediction model for the survival of lung cancer patients and to determine the key factors that affect overall survival.

**Methods**
This study used Taipei Medical University Clinical Research Database (TMUCRD) with data from 3 hospitals as the data source, the data were mapped to OHDSI OMOP CDM. We selected non-small-cell lung cancer patients from a retrospective development dataset of TMUCRD and Taiwan Cancer Registry between January 2008 and December 2018. All patients were monitored from the index date of cancer diagnosis until the event of death or the last visit to hospitals. Variables including demographics, comorbidities, medications, laboratories, and gene tests of patients were retrieved and used to develop the machine learning models. Nine machine learning algorithms with various modes (e.g., integrating different variables) were used to develop the predicted models. The performance of the algorithms was measured by the area under the receiver operating characteristic curve (AUC), accuracy, sensitivity (Recall), specificity, positive predictive value (Precision), and F1-score.

**Results**
In total, 3,714 patients were included (2,280 for the training dataset and 1,434 for the testing dataset). The artificial neural network (ANN) AUC values of different modes were observed with the highest score of 89%. The best performance of the ANN model was achieved when integrating all variables with the AUC, accuracy, precision, recall, and F1-score of 0.89, 0.82, 0.91, 0.75, and 0.65, respectively. The most important features were the cancer stage, cancer size, diagnosed age, smoking, drinking status, EGFR gene, and body mass index.

**Conclusion**
In this evaluation of lung cancer survival, the ANN model led to a better predictive performance with high AUC, precision, and recall when integrating different data types. Further research is necessary to determine the feasibility of applying the algorithm in the clinical setting and explore whether using this tool could improve care and outcomes. This study is expected to be developed into a multinational cooperative research using OHDSI tools and OMOP CDM in the future.

## Methods

**Feature Selection**
Based on a literature review and consultation with clinicians, we selected features that may lead to the mortality of NSCLC patients to build prediction models. Those features consisted of: (1) Demographic information, (2) Cancer conditions, (3) Comorbidities, (4) Medications, (5) Laboratory tests, (6) Genomic tests.

**Development of the Algorithms**
This study established prediction models based on four modes and different algorithms.
(1) The primary mode (e.g., mode 1) included demographic information, cancer conditions, comorbidities, and medications.
(2) The second mode (mode 2) included the data of mode one and the laboratory tests.
(3) The third mode (mode 3) included the data of mode one and genomic tests.
(4) The fourth mode (mode 4) considered all the above features.
The study aims to predict the survival of lung cancer patients; therefore, the problem can be formulated as a classification model as it could occur in the same patients. We used those possible machine learning techniques such as logistic regression (LR), linear discriminant analysis (LDA), light gradient boosting machine (LGBM), gradient boosting machine (GBM), extreme gradient boosting (XGBoost), random forest (RF), AdaBoost, support vector machine (SVC), and artificial neural network (ANN). These methods were briefly introduced as follows.

**Evaluating the Algorithms**
The training dataset contained the data of patients from TMUH and WFH. The stratified 5-fold cross-validation was applied in the training set to assess the different machine learning models' performance and general errors. In words, patients in the training set were divided into five groups, each used repeatedly as the internal validation set. We recruited data from SHH and used it for external testing set for generalizing the model.
The performance of the algorithms was measured by the area under the receiver operating characteristic curve (AUC), accuracy, sensitivity (Recall), specificity, positive predictive value (PPV, Precision), negative predictive value (NPV), and F1-score. We defined the best model using the highest AUC by comparing various models based on the external testing set. We, furthermore, analyzed the feature's contribution (i.e., features importance) of the best model using SHAP values (SHapley Additive exPlanations).
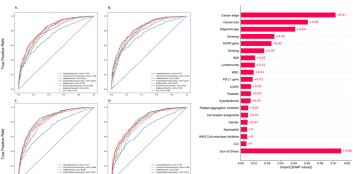
## Results



**Figure 2.** The Performance of the Prediction Models in the Testing dataset by different Modes
**Note:** A, Mode 1; B, Mode 2 ; C, Mode 3 ; D, Mode 4

**Figure 3.** Feature Importance of the ANN Prediction Model in Mode 4
**Note:** BMI, Body mass index; EGFR, Epidermal growth factor receptor; WBC, White blood cell; PD-L1, Programmed death-ligand 1; COPD, Chronic obstructive pulmonary disease; CCI, Charlson comorbidity index

## Methods

**Study Design and Data Source**
We conducted a retrospective study in which we obtained the data from the Taipei Medical University Clinical Research Database (TMUCRD), which were mapped to OHDSI OMOP CDM. The TMUCRD retrieved data from various electronic medical records (EHR) of three hospitals, Taipei Medical University Hospital (TMUH), Wan-Fang Hospital (WFH), and Shuang-Ho Hospital (SHH). The database contains the electronic medical record data of 3.8 million people accumulated from 1998 to 2020. This study has been approved by the Joint Institute Review Board of Taipei Medical University (TMU-JIRB), Taipei, Taiwan (approved number, N202101080).

**Cohort Selection**
This study selected patients with lung cancer (ICD-O-3 code: C33, C34) from 2008 to 2018 in the TCR database. Exclusion criteria included individuals with ages under 20, SCLC patients, and patients who did not have any medical history in the three hospitals (TMUH, WFH, SHH). These 3,714 patients were included in this study, including 960 patients from TMUH, 1,320 from WFH, and 1,434 from SHH.



Figure 1. Cohort Selection Process

**Outcome Measurement**
We ascertained the study outcomes using TMUCRD EHR and vital status data from the Taiwan Death Registry (TDR). We used the diagnosis date of NSCLC as the index date, and 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, 2018.

## Results

**Table 1. Performance of various Prediction Models by Modes**

| Modes | Models | AUC Training | AUC Testing | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|---|---|
| Mode 1 | LR | 0.70 | 0.72 | 0.65 | 0.88 | 0.64 | 0.75 |
| | LDA | 0.78 | 0.78 | 0.71 | 0.90 | 0.70 | 0.80 |
| | LGBM | 0.98 | 0.81 | 0.73 | 0.92 | 0.72 | 0.81 |
| | GBM | 0.82 | 0.83 | 0.75 | 0.91 | 0.76 | 0.84 |
| | XGBoost | 0.99 | 0.80 | 0.75 | 0.90 | 0.77 | 0.84 |
| | RF | 0.90 | 0.82 | 0.72 | 0.92 | 0.70 | 0.80 |
| | AdaBoost | 0.94 | 0.81 | 0.73 | 0.91 | 0.72 | 0.81 |
| | SVC | 0.78 | 0.78 | 0.71 | 0.89 | 0.72 | 0.79 |
| | ANN* | 0.89 | 0.88 | 0.82 | 0.90 | 0.75 | 0.64 |
| Mode 2 | LR | 0.74 | 0.75 | 0.60 | 0.93 | 0.53 | 0.67 |
| | LDA | 0.81 | 0.79 | 0.71 | 0.90 | 0.70 | 0.80 |
| | LGBM | 0.99 | 0.83 | 0.78 | 0.91 | 0.79 | 0.86 |
| | GBM | 0.96 | 0.84 | 0.78 | 0.91 | 0.80 | 0.87 |
| | XGBoost | 1.00 | 0.81 | 0.76 | 0.90 | 0.81 | 0.86 |
| | RF | 0.92 | 0.83 | 0.69 | 0.94 | 0.64 | 0.76 |
| | AdaBoost | 0.95 | 0.80 | 0.74 | 0.90 | 0.76 | 0.83 |
| | SVC | 0.81 | 0.79 | 0.70 | 0.91 | 0.68 | 0.78 |
| | ANN* | 0.89 | 0.89 | 0.80 | 0.91 | 0.75 | 0.64 |
| Mode 3 | LR | 0.70 | 0.73 | 0.65 | 0.88 | 0.63 | 0.74 |
| | LDA | 0.80 | 0.81 | 0.75 | 0.91 | 0.76 | 0.83 |
| | LGBM | 0.98 | 0.85 | 0.80 | 0.92 | 0.81 | 0.87 |
| | GBM | 0.96 | 0.85 | 0.79 | 0.92 | 0.79 | 0.86 |
| | XGBoost | 1.00 | 0.83 | 0.79 | 0.91 | 0.80 | 0.86 |
| | RF | 0.91 | 0.84 | 0.72 | 0.93 | 0.69 | 0.80 |
| | AdaBoost | 0.95 | 0.83 | 0.79 | 0.91 | 0.80 | 0.86 |
| | SVC | 0.80 | 0.81 | 0.75 | 0.90 | 0.75 | 0.83 |
| | ANN* | 0.89 | 0.89 | 0.83 | 0.89 | 0.81 | 0.64 |
| Mode 4 | LR | 0.74 | 0.75 | 0.61 | 0.93 | 0.53 | 0.67 |
| | LDA | 0.83 | 0.82 | 0.76 | 0.90 | 0.77 | 0.84 |
| | LGBM | 0.99 | 0.86 | 0.81 | 0.92 | 0.83 | 0.88 |
| | GBM | 0.97 | 0.85 | 0.79 | 0.92 | 0.81 | 0.87 |
| | XGBoost | 1.00 | 0.84 | 0.77 | 0.92 | 0.77 | 0.85 |
| | RF | 0.93 | 0.85 | 0.75 | 0.93 | 0.73 | 0.82 |
| | AdaBoost | 0.95 | 0.85 | 0.76 | 0.92 | 0.75 | 0.83 |
| | SVC | 0.83 | 0.81 | 0.75 | 0.90 | 0.76 | 0.84 |
| | ANN* | 0.89 | 0.89 | 0.82 | 0.91 | 0.75 | 0.65 |

## Conclusions

In summary, to observe the expected survival of NSCLC patients during two years period, an artificial neural network model had high AUC, precision, and recall. Thus, integrating different data types (especially laboratory and genomic data) led to better predictive performance. Further research is necessary to determine the feasibility of applying the algorithm in the clinical setting and explore whether using this tool could improve care and outcomes.

**Note**
This study used TMUCRD with data from 3 hospitals as the data source, the data were mapped to OHDSI OMOP CDM. It is expected to be developed into a multinational cooperative research using OHDSI tools and OMOP CDM in the future as well.

**References**
1. Lynch CM, Abdollahi B, Fuqua JD, et al. Prediction of lung cancer patient survival via supervised machine learning classification techniques. Int J Med Inform. 2017;108:1-8.
2. Bartholomai JA, Frieboes HB. Lung Cancer Survival Prediction via Machine Learning Regression, Classification, and Statistical Techniques. Proc IEEE Int Symp Signal Proc Inf Tech. 2018;2018:632-637.
3. Stah KW, Khezin S, Wong CH, Lo AW. Machine Learning and Stochastic Tumor Growth Models for Predicting Outcomes in Patients With Advanced Non-Small-Cell Lung Cancer. JCO Clin Cancer Inform. 2019;3:1-11.
4. Cui L, Li H, Hui W, et al. A deep learning-based framework for lung cancer survival analysis with biomarker interpretation. BMC Bioinformatics. 2020;21(1):112.
5. He J, Zhang JX, Chen CT, et al. The Relative Importance of Clinical and Socio-demographic Variables in Prognostic Prediction in Non-Small Cell Lung Cancer: A Variable Importance Approach. Med Care. 2020;58(5):461-467.
6. She Y, Jin Z, Wu J, et al. Development and Validation of a Deep Learning Model for Non-Small Cell Lung Cancer Survival. JAMA Netw Open. 2020;3(6):e205842.
7. Dreiseitl S, Ohno-Machado L. Logistic regression and artificial neural network classification models: a methodology review. J Biomed Inform. 2002;35(5):352-359.
8. Izenman AJ. Linear discriminant analysis. In: Modern multivariate statistical techniques. Springer 2013:237-280.
9. Ke G, Meng Q, Finley T, et al. Lightgbm: A highly efficient gradient boosting decision tree. Adv Neural Inf Process Syst. 2017;30.
10. Friedman JH. Greedy function approximation: a gradient boosting machine. Annals of statistics. 2001:1189-1232.
11. Chen T, He T, Benesty M, et al. Xgboost: extreme gradient boosting. R package version 04-2. 2015;1(4):1-4.
12. Ho TK. Random decision forests. Paper presented at: Proceedings of 3rd international conference on document analysis and recognition1995.
13. Hastie T, Rosset S, Zhu J, Zou H. Multi-class adaboost. Statistics and its interface. 2009;2(3):349-360.
14. Gunn SR. Support vector machines for classification and regression. ISIS technical report. 1998;14(1):5-16.
15. Agatonovic-Kustrin S, Beresford R. Basic concepts of artificial neural network (ANN) modeling and its application in pharmaceutical research. Journal of pharmaceutical and biomedical analysis. 2000;22(5):717-727.

---

## FRIDAY

**Development of Lung Cancer Survival Prediction Models Based on Real-world Data and Machine Learning** (Jason C. Hsu, Phung-Anh Nguyen, Phan Thanh Phuc, Tsai-Chih Lo, Min-Huei Hsu, Chi-Tsun Cheng, Tzu-Hao Chang, Cheng-Yu Chen)

# Openings

FDA/CDER's Division of Hepatology and Nutrition is seeking a clinician with bioinformatics or biostatistics training to work with the Drug-Induced Liver Injury (DILI) Team to evaluate large datasets of liver-related data, collaborate on the Team's review of drugs with hepatotoxicity signals, and help develop informatics-based processes in DILI evaluation across the Agency.

Contact **Judy Racoosin** at **judith.racoosin@fda.hhs.gov** for information about the application process (that will be through USAJOBS).

# Openings

**Andrew Williams** recently announced two exciting new openings at Tufts Medicine.

1) Senior Project Manager for a multisite multiyear grant standardizing critical care EHR and waveform data. (CHoRUS Bridge2AI)

2) Lead software developer and research data warehouse manager for Tufts Medicine's OMOP instance and related services.

Remote work is possible for both positions.

1.  Link for Senior Project Manager position: https://smrtr.io/bBVzh

2.  Link for Lead Software Developer and Research Data Warehouse Manager position: https://jobs.smartrecruiters.com/TuftsMedicalCenter1/743999857980631-software-development-lead-res-g-c-ctsi

Andrew's email: awilliams15@tuftsmedicalcenter.org

# Openings

## Research Associate (Data Scientist/Statistical Engineer), Johns Hopkins inHealth and Biostatistics Center

- Execute OHDSI studies (e.g. for cohort characterizations and comparative effectiveness) on Johns Hopkins's EHR data to support clinicians;

- Collaborate with statisticians and clinicians to continuously integrate state-of-the-art statistical tools to the inHealth/OHDSI tool stack for deployment;

- Mentor trainees on data science and software development skills;

- Co-teach courses on observational health data analytics and data science skills at School of Medicine and Public Health;

- Facilitate adoption of the inHealth tools among the broader OHDSI community by contributing to OHDSI's Health Analytics Data-to-Evidence Suite.

- https://apply.interfolio.com/114436

# Where Are We Going?

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

the journey

# Three Stages of The Journey

**Where Have We Been?**

**Where Are We Now?**

**Where Are We Going?**

# Best Community Contribution Awards

## Methods Research



**Assessing Racial Fairness of Dialysis Allocation in End-Stage Renal Disease (Linying Zhang, Lauren R. Richter, David M. Blei, Yixin Wang, Anna Ostropolets, Noemie Elhadad, George Hripcsak)**

# Best Community Contribution Awards

## Data Standards



Analyzing the Effect of Hypertension on Retinal Thickness Using Radiology Common Data Model (R-CDM) (**Chul Hyoung Park**, Rae Woong Park, Sang Jun Park, Da Yun Lee, Seng Chan You, Ki Hwang Lee)

# Best Community Contribution Awards

**Open-Source Analytics**

## Cohort Definition Validation in Atlas

Charity Hilton MS[1], Saul Crumpton MS[1], Jon Duke MD, MS[1,2]

[1]Georgia Tech Research Institute, [2]Georgia Institute of Technology

### Background

OHDSI Atlas has long been an effective tool for developing rule-based cohort definitions in observational data. In the public version of Atlas, thousands of cohort definitions have been created. While patient record verification is a common method of cohort definition validation, it is not without difficulties, including but not limited to the need for clinical experts to access data, a tool to review all in-cohort patients, a method to gather review data, and a system of tabulation to determine in-cohort (case/no-case) participation or not[1].

Until now, there has not been an Atlas-based system for clinical expert review. For this effort, we introduce the Atlas Cohort Definition Validation tool (ACDV). This tool aims to solve some of the primary concerns around cohort definition validation, while having the chief benefit of being cohesively integrated into the OHDSI Atlas stack. Additionally, the tool allows for creation of more complex validation question sets, beyond the standard case/no-case assessment.



Figure 1: Question Set Creation

### Methods

We designed and developed two modules around cohort definition validation. The first (1) allows for validation study creation and management, and the second (2) allows for validation of study questions for clinical reviewers in the Atlas Patient Profile tool.

The ACDV tool introduces a 'Validation' section to Atlas cohort definition creation, which allows for cohort managers to complete a cohort definition validation workflow. This workflow begins by the creation of question set. Question sets in the ACDV tool, shown in Figure 1, allow for common types of questions (including text, radio, checkbox, numbers, and dates). Multiple questions in a question set can be created and a case/no-case distinction can be selected at the question level. After a question set has been created, it can be linked to a cohort definition sample, this creates the validation study.

After a validation study is created, cohort managers can assign patients for review in the Atlas Patient Profile tool to clinical reviewers. Study questions are displayed to clinical reviewers at the patient level in a collapsible sidebar (see Figure 3). The study question set at the patient profile-level can be accessed via the Cohort Definition tool, the Patient Profile tool, or via a customized link. Once reviewers have viewed patient profiles and answered study questions, study results can be viewed by cohort managers in Atlas or exported to CSV (Figure 4).

### Results

Primary development efforts of the ACDV tool are complete, and final modifications and integrations to the tool are being prepared for inclusion in an upcoming OHDSI release. We have validated the tool internally with a clinician-informaticist.



Figure 2: Annotation Study Manager View



Figure 3: Profile Level Validation

### Conclusions

The Atlas Cohort Definition Validation tool will provide an integrated way for clinical chart reviewers to validate cohorts well beyond the question of cohort inclusion or not.

This tool will support research in the OHDSI community by living firmly within the active OHDSI Atlas ecosystem of tools. Additionally, this tool will continue the OHDSI legacy of open and community-driven tools to advance research in observational health data.



Figure 4: Study Results

### Bibliography

1. Observational Health Data Sciences and Informatics. The Book of OHDSI; 2020. Available from: https://ohdsi.github.io/TheBookOfOhdsi/

**Cohort Definition Validation in Atlas (Charity Hilton, Saul Crumpton, Jon Duke)**

# Best Community Contribution Awards

**Clinical Applications**

A Pilot Characterization Study Assessing Health Equity in Mental Healthcare Delivery within the State of Georgia (**Jacob Zelko**, Malina Hy, Varshini Chinta, Emily Liau, Morgan Knowlton, Jon Duke)