



# The missing link: Cross-species EHR data linkage offers new opportunities for improving One Health



THE UNIVERSITY  
of NORTH CAROLINA  
at CHAPEL HILL

**Kathleen Mullen**

University of North Carolina at Chapel Hill  
Translational and Integrative Sciences Lab

2024 OHDSI Global Symposium  
23 October 2024



# One Health

is an integrative multidisciplinary effort focused on achieving optimal health for people, animals, and their shared environments.



There are significant opportunities for **learning across species living in a shared environment.**

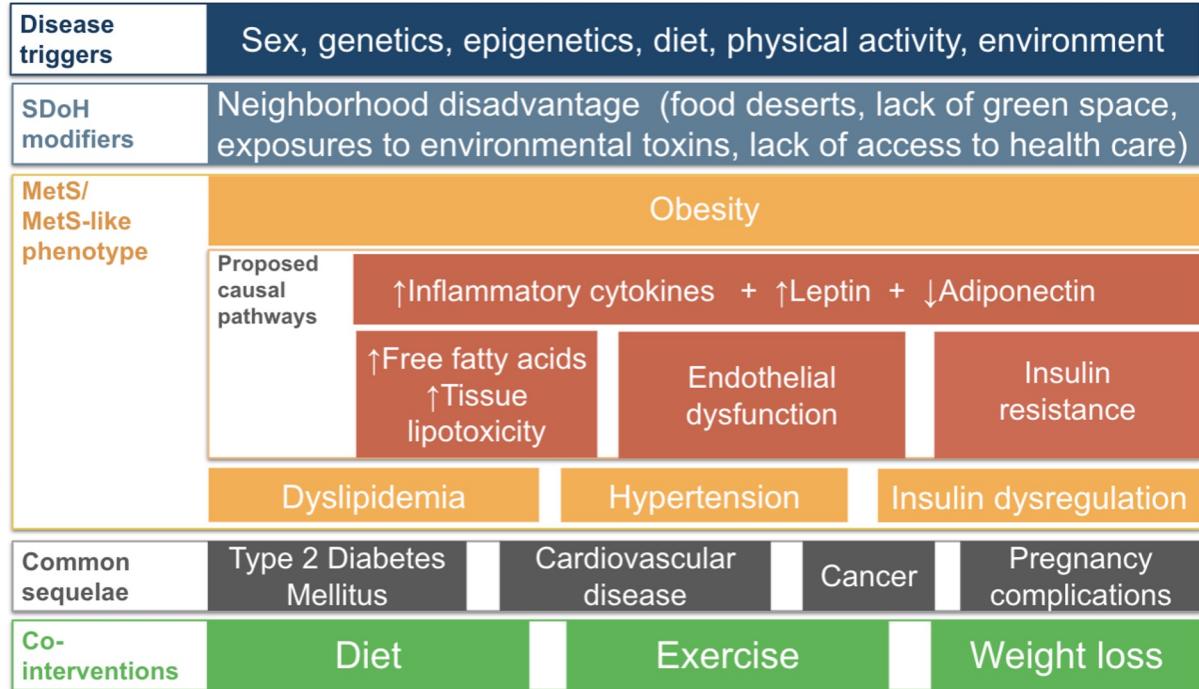
A mechanism to study the household co-occurrence of disease is needed.

# Metabolic syndrome is similar in people and animals!



Image credit: Sarah M. Reuss

**Miniature donkeys with equine metabolic syndrome phenotype.**

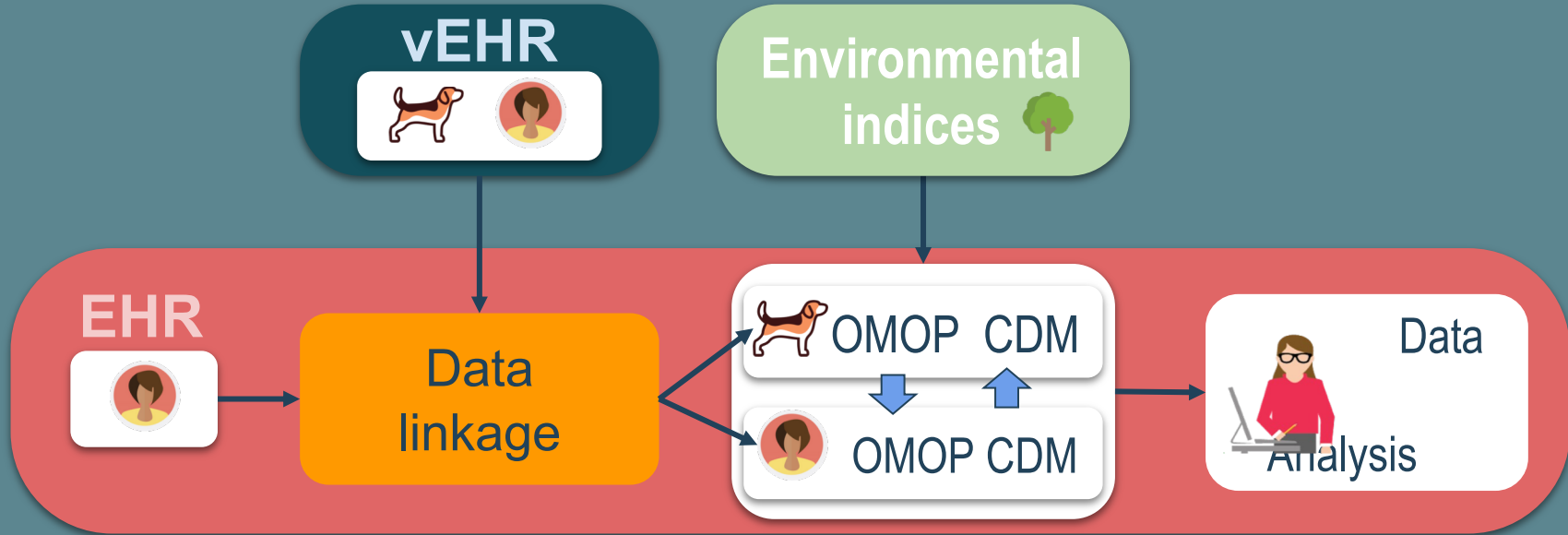


**Commonalities in MetS in people and MetS-like in companion animals.**

**How can we explore the causes of metabolic syndrome in a single household?**

## Project goal

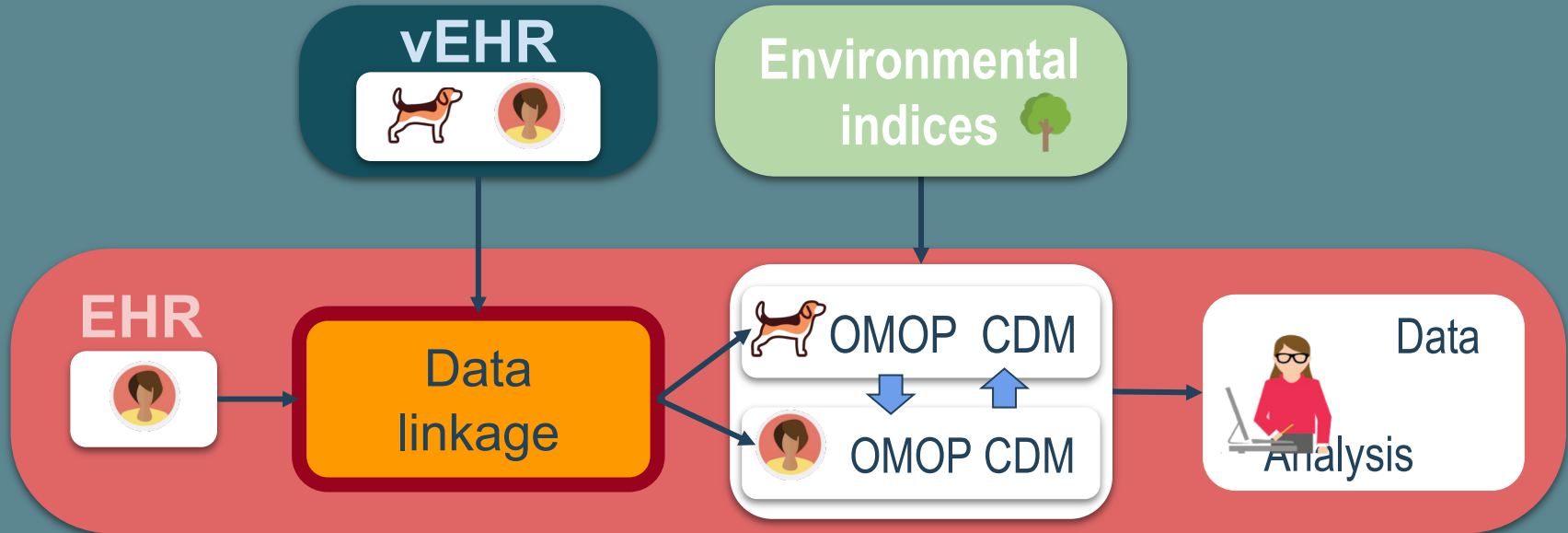
Create a secure, pet-patient registry linking people, their animals & the environment.





## Project goal

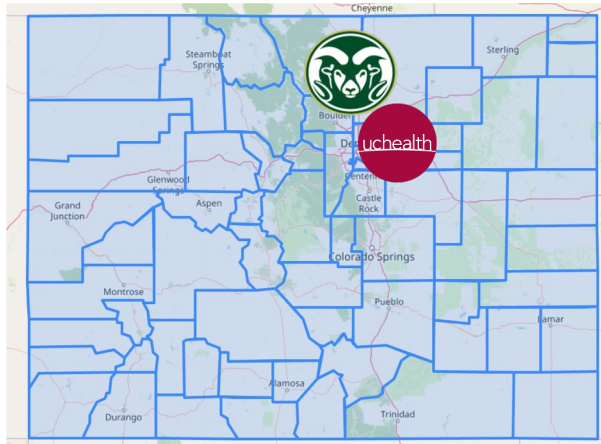
Putting the household back together with EHR data linkage.



# How many animals and owners can be linked via their EHRs?

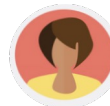
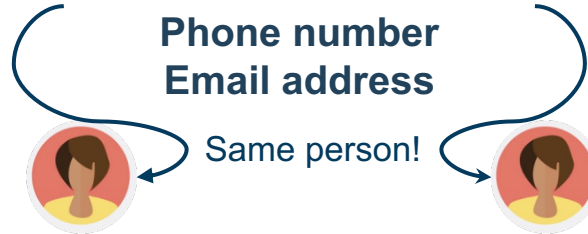


## Linkage Honest Broker



Colorado map with locations of CSU-VTH and CU Medical Campus (UCHealth).

**First name, Last name**  
**Street address**  
**Phone number**  
**Email address**



# How many animals and owners can be linked via their EHRs?



**Vet. EHR (vEHR) time span:**  
2019-2024

**# animal owners:** 41,081

**# animals patients:** 76,282

**Cats:** 13.0%

**Dogs:** 55.5%

**Horses:** 15.3%

**Other:** 16.2%

**Female animals patients:** 47.9%

## Linkage Honest Broker

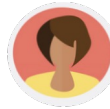
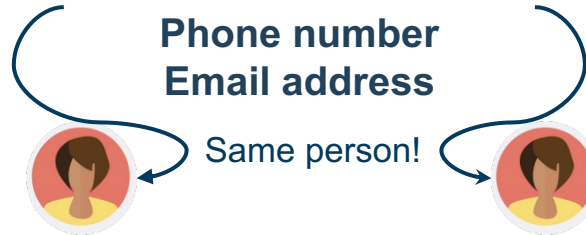


**First name, Last name**

**Street address**

**Phone number**

**Email address**



**EHR time span:**  
2015-2024

**# human patients:** 3,282,860

**Female human patients:** 53.4%

# How many animals and owners can be linked via their EHRs?



**Vet. EHR (vEHR) time span:**  
2019-2024

**# animal owners:** 41,081

**# animals patients:** 76,282

**Cats:** 13.0%

**Dogs:** 55.5%

**Horses:** 15.3%

**Other:** 16.2%

**Female animals patients:** 47.9%

## Linkage Honest Broker

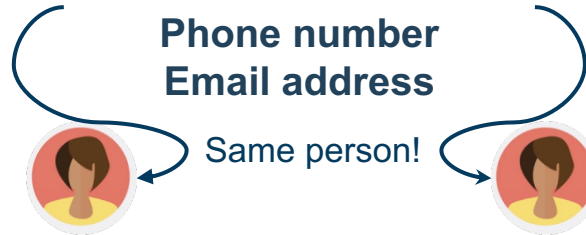


**First name, Last name**

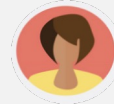
**Street address**

**Phone number**

**Email address**



**12,115 vEHR-EHR pairs!**



**uhealth**

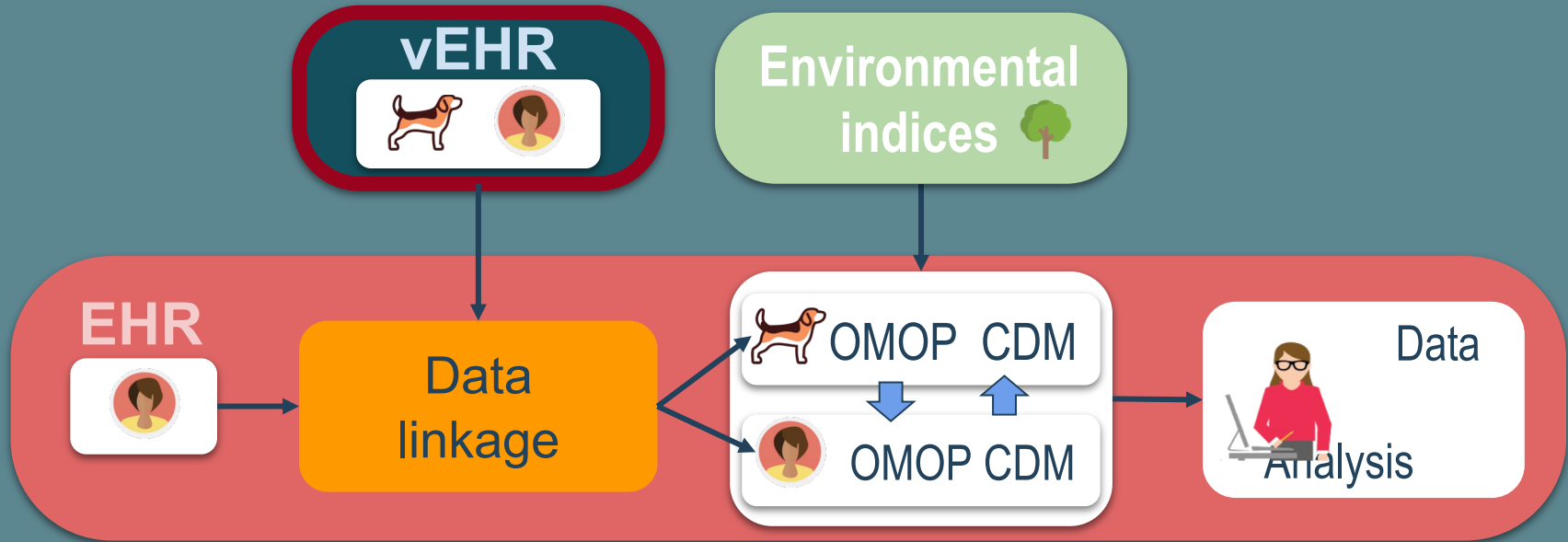
**EHR time span:**  
2015-2024

**# human patients:** 3,282,860

**Female human patients:** 53.4%

## Project goal

Interrogating the vEHR for metabolic-like syndrome features.



# Do key indicators of MetS-like exist in the vEHR?

## Key indicators

Elevated BCS

Overweight

Diabetes

Obesity

Over-conditioned\*

Equine metabolic syndrome

Cresty neck\*

## Prevalence of MetS-like key indicators in the CSU-VTH vEHR for companion animals.

Species	Animal patients (N)	Prevalence (%)
Cats*	3,037	51.0%
Dogs*	13,672	43.9%
Horses*	1,027	11.8%

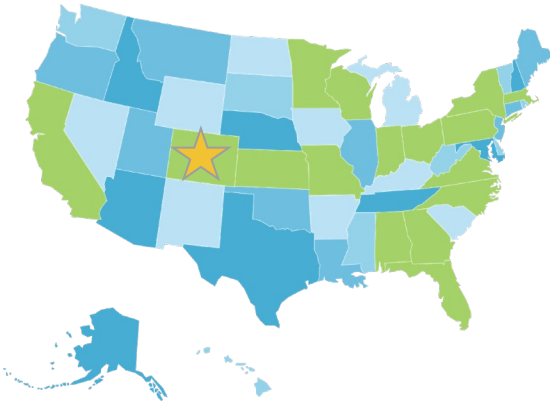
\*Significant differences in the prevalence of MetS-like key indicators by species ( $p < 0.001$ ).


**That's a lot of fat cats!**

# The CSU-CU pet-patient data registry is a blueprint for **One Health** study



Clinical and Translational Science Award  
One Health Alliance



 States with paired veterinary colleges  
and academic medical centers

**One Health** issues for the pet-patient data  
registry:

- Antimicrobial resistance
- Chronic diseases
- Environmental health
- Mental health
- Trauma
- Vector-borne diseases
- Zoonotic diseases

. . . And more!

**We need you and your use cases! Visit poster #114 for more details!**

# Acknowledgements



Nadia Saklou



Adam Kiehl



Joe Strecker



Tracy Webb



Sue VandeWoude



Ian Brooks



Toan Ong



Sabrina Toro



Melissa Haendel



OHDSI Community

Support: NIH/NIAMS K12AR084226, NIH/NCATS Colorado CTSA UM1 TR004399 & UNC Department of Genetics



THE UNIVERSITY  
of NORTH CAROLINA  
at CHAPEL HILL





# Comparing probabilistic and rule-based phenotype algorithms for hypotension and angioedema to the experience observed in randomized clinical trials.

**Joel Swerdel, PhD MS MPH<sup>1,2</sup>, Martijn Schuemie, PhD<sup>1,2</sup>,  
Judith A. Racoosin, MD, MPH<sup>2,3</sup>, and Patrick Ryan, PhD<sup>1,2</sup>**

<sup>1</sup>Janssen R&D LLC, Titusville, NJ USA; <sup>2</sup>Observational Health Data Sciences and Informatics (OHDSI), New York, NY, USA <sup>3</sup>United States Food and Drug Administration, Silver Spring, MD, USA

Presented at OHDSI Symposium; October 23, 2024; New Brunswick, NJ USA

The QR code is intended to provide scientific information for individual reference, and the information should not be altered or reproduced in any way.



## Disclosures

Joel Swerdel, Martijn Schuemie, and Patrick Ryan are employees of Janssen Research and Development and shareholders of Johnson & Johnson. Judith Racoosin has no disclosures and notes that the views and opinions expressed in this presentation are those of the authors and should not be construed to represent the views or policies of the FDA.

## Background

- Rule-based phenotype algorithms (PAs) are the standard for identifying outcomes in observational data.
- However, the performance characteristics of the PAs, such as sensitivity and positive predictive value, are estimated to be low for many phenotypes.
- Probabilistic PAs, e.g., PAs based on logistic regression models, offer an alternative to the rule-based method. Prior efforts have demonstrated the potential for probabilistic phenotyping as an alternative to rule-based PAs.<sup>[1-3]</sup>

## Objective

The objective of this study was to develop a methodology for creating probabilistic phenotypes and to show examples of its use in hypotension, i.e., low blood pressure, and angioedema, a subcutaneous tissue swelling triggered by an allergic reaction.

## **Methods – Building rule-based and probabilistic phenotypes**

# Developing Rule-based and Probabilistic Phenotypes

## Example - Rule-based phenotype

### Cohort Entry Events:

People enter the cohort when observing any of the following:

condition occurrences of 'Angioedema'.

Code	Name	Vocabulary
T78.3	<a href="#">Angioneurotic oedema</a>	ICD10
995.1	<a href="#">Angioneurotic edema, not elsewhere classified</a>	ICD9CM

### Cohort Exit:

The cohort end date will be offset from index event's end date plus 7 days.

# Developing Rule-based and Probabilistic Phenotypes

## Example - Rule-based phenotype

### Cohort Entry Events:

People enter the cohort when observing any of the following:

condition occurrences of 'Angioedema'.

Code	Name	Vocabulary
T78.3	Angioneurotic oedema	ICD10
995.1	Angioneurotic edema, not elsewhere classified	ICD9CM

### Cohort Exit:

The cohort end date will be offset from index event's end date plus 7 days.

## Example - Probabilistic phenotype

Beta	
Coefficient	Covariate Name
4.82	condition_era group during day 0 through 10 days relative to index: <b>Angioedema</b>
3.88	condition_era group during day 0 through 10 days relative to index: <b>Allergic disposition</b>
	visit_occurrence concept count during day 0 through 10 concept_count relative to index:
2.35	<b>Emergency Room Visit</b>
1.88	drug_era group during day 0 through 10 days relative to index: <b>prednisone</b>
1.87	condition_era group during day 0 through 10 days relative to index: <b>Anaphylaxis</b>
1.78	drug_era group during day 0 through 10 days relative to index: <b>ACE INHIBITORS, PLAIN</b>
1.58	drug_era group during day 0 through 10 days relative to index: <b>H2-receptor antagonists</b>
1.57	observation during day 0 through 10 days relative to index: <b>Adverse reaction to substance</b>
	condition_era group during day 0 through 10 days relative to index: <b>Angioedema and/or</b>
1.52	<b>urticaria</b>
1.38	condition_era group during day 0 through 10 days relative to index: <b>Edema</b>
	drug_era group during day 0 through 10 days relative to index: <b>Sympathomimetics in</b>
1.34	<b>glaucoma therapy</b>
	drug_era group during day 0 through 10 days relative to index: <b>CORTICOSTEROIDS FOR</b>
1.28	<b>SYSTEMIC USE, PLAIN</b>
	condition_era group during day 11 through 20 days relative to index: <b>Angioedema and/or</b>
1.26	<b>urticaria</b>
1.20	condition_era group during day 0 through 10 days relative to index: <b>Acute allergic reaction</b>
...	

# Developing Rule-based and Probabilistic Phenotypes

## Example - Rule-based phenotype

### Cohort Entry Events:

People enter the cohort when observing any of the following:

condition occurrences of 'Angioedema'.

## Example - Probabilistic phenotype

Beta	
Coefficient	Covariate Name
4.82	condition_era group during day 0 through 10 days relative to index: <b>Angioedema</b>
3.88	condition_era group during day 0 through 10 days relative to index: <b>Allergic disposition</b>
2.35	visit_occurrence concept count during day 0 through 10 concept_count relative to index: <b>Emergency Room Visit</b>
1.88	drug_era group during day 0 through 10 days relative to index: <b>prednisone</b>
1.87	condition_era group during day 0 through 10 days relative to index: <b>Anaphylaxis</b>
1.78	drug_era group during day 0 through 10 days relative to index: <b>ACE INHIBITORS, PLAIN</b>

Beta	
Coefficient	Covariate Name
4.82	condition during day 0 through 10 days: <b>Angioedema</b>
3.88	condition during day 0 through 10 days: <b>Allergic disposition</b>
2.35	Visit occurrence during day 0 through 10 days: <b>Emergency Room Visit</b>



## Developing Rule-based and Probabilistic Phenotypes (cont.)

### Rule-based phenotype

1. Create a rule-based phenotype algorithm
2. Find subjects during the time-at-risk in the cohort of interest satisfying algorithm

## Developing Rule-based and Probabilistic Phenotypes (cont.)

### Rule-based phenotype

1. Create a rule-based phenotype algorithm
2. Find subjects during the time-at-risk in the cohort of interest satisfying algorithm

### Probabilistic phenotype

1. Use noisy labeled positive and negative controls to develop a supervised learning probabilistic model using LASSO regularized regression
2. Apply model at each appropriate time point during the time-at-risk for each subject in the cohort of interest
3. Select the highest probability among the different time points within the time-at-risk for each subject
4. Use a designated probability cut-point, e.g., 70%, to determine those with the outcome

## Evaluating the model

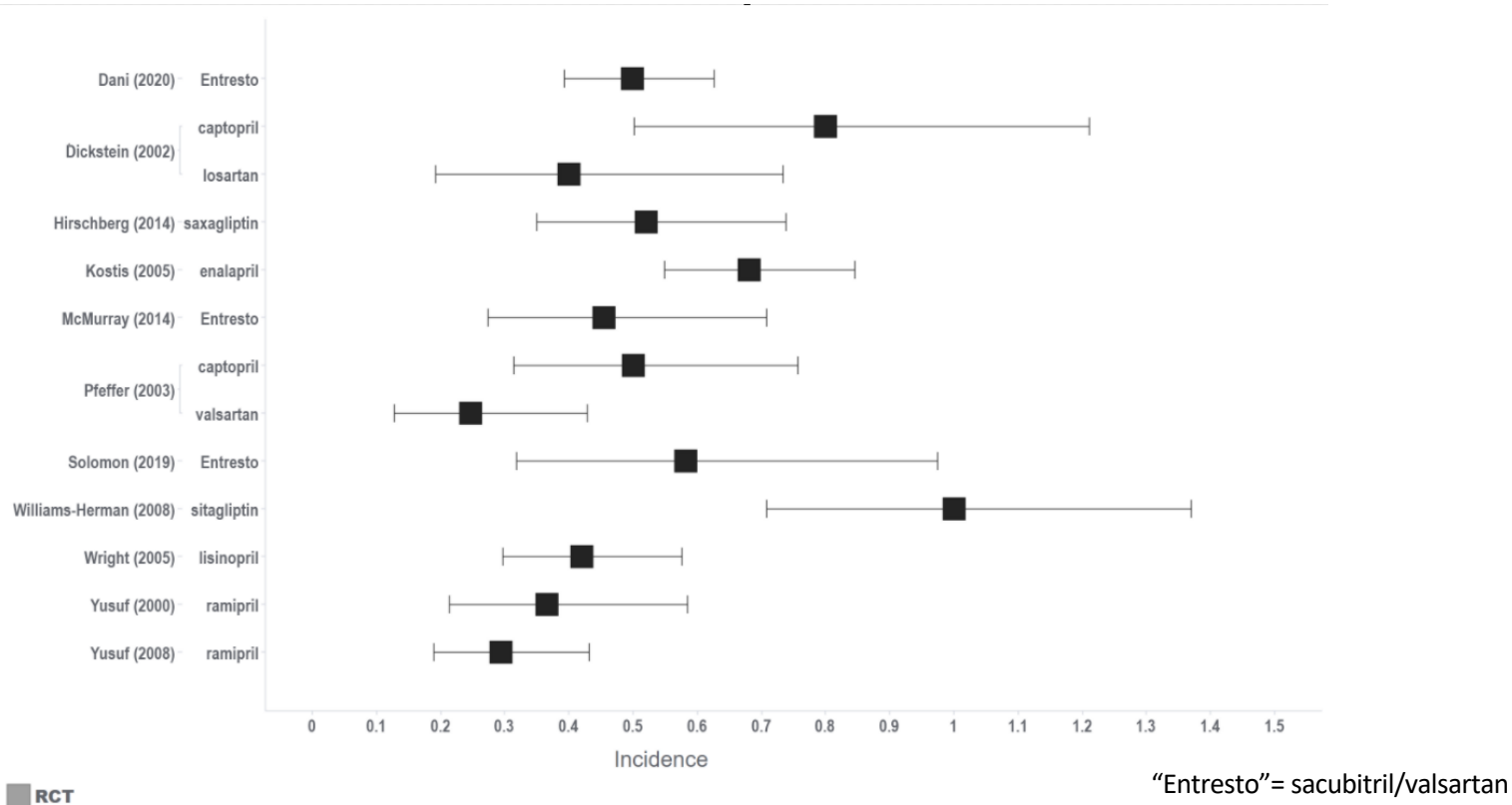
- Analysis conducted in 5 administrative claims datasets
- Rule-based algorithm used an occurrence of a diagnosis code for hypotension or angioedema
- Developed probabilistic phenotypes and examined the results using probability cut-points of 0.50, 0.60, 0.67, 0.70, 0.75, 0.80, and 0.90.
- Estimated incidence of angioedema, while on-treatment, for 7 anti-hypertensive and 2 anti-diabetic (DPP-4 inhibitors) drugs
- Using both the rule-based and probabilistic phenotypes, we performed the analysis on 9 new user drug cohorts from 2010 to 2023

## Evaluating the model - Metrics

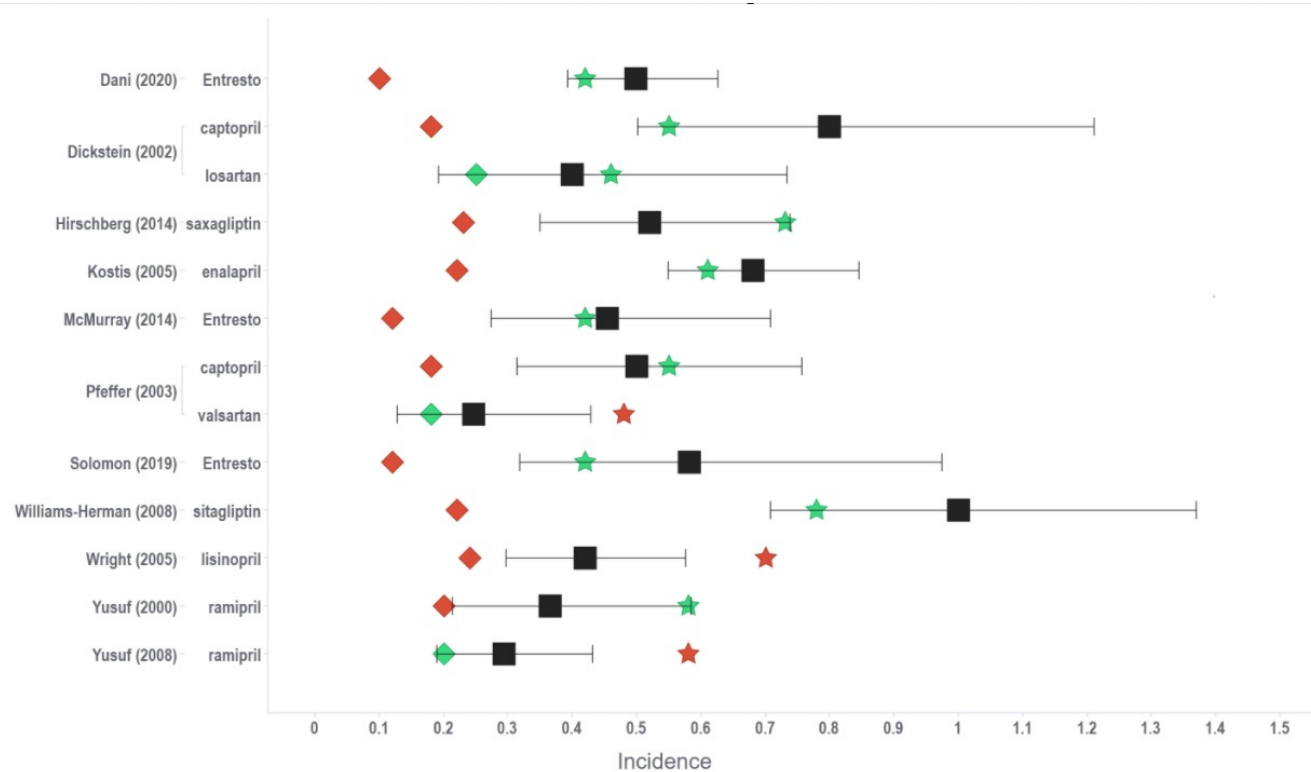
- Extracted incidence estimates and 95% confidence intervals from randomized clinical trials (RCTs) as a basis for comparison.
- Computed the proportion of incidence estimates for the rule-based and probabilistic algorithms that fell within the 95% confidence intervals (CI) of the incidence estimates from the clinical trials.
- Assessed the performance characteristics, e.g., positive predictive value (PPV) and sensitivity, of the rule-based and probabilistic phenotypes using the OHDSI tool PheValuator.

# Results

# Results: Angioedema - Probabilistic (Cut-point 0.67) v. Rule-based

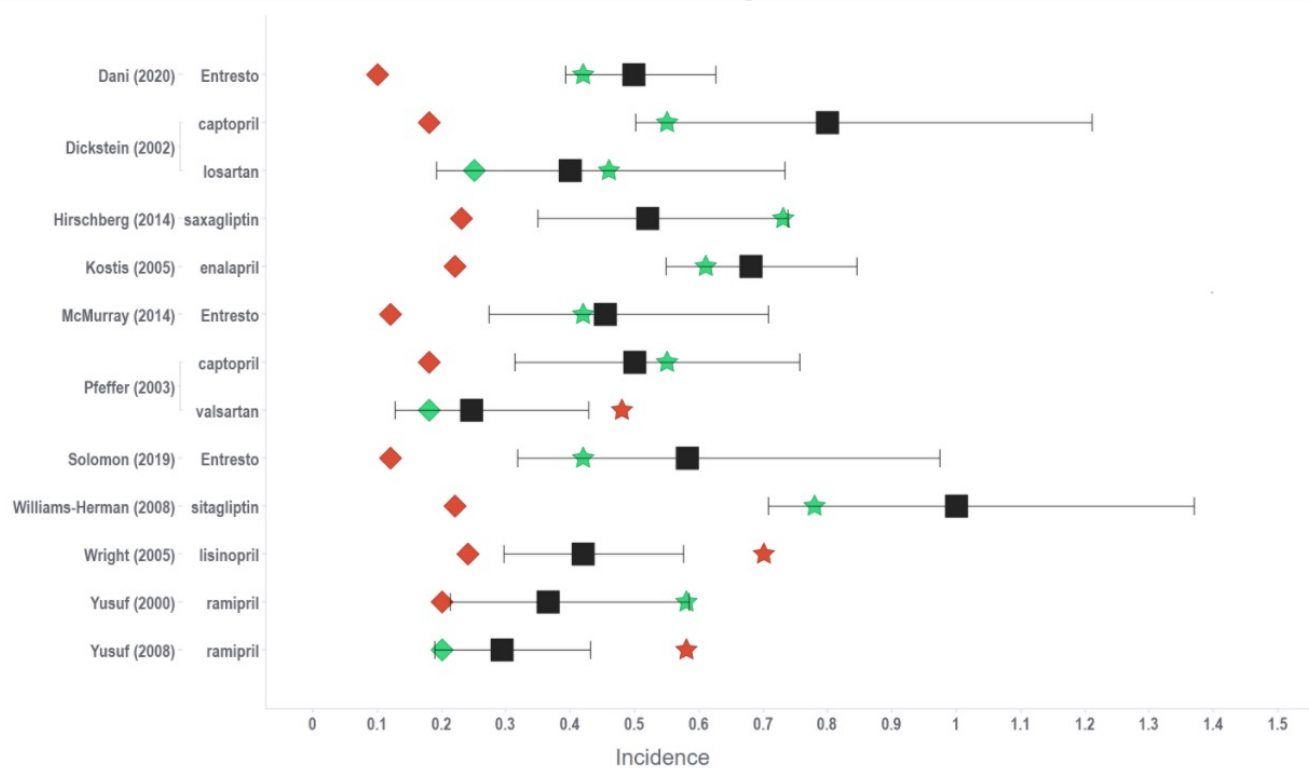


# Results: Angioedema - Probabilistic (Cut-point 0.67) v. Rule-based



“Entresto” = sacubitril/valsartan

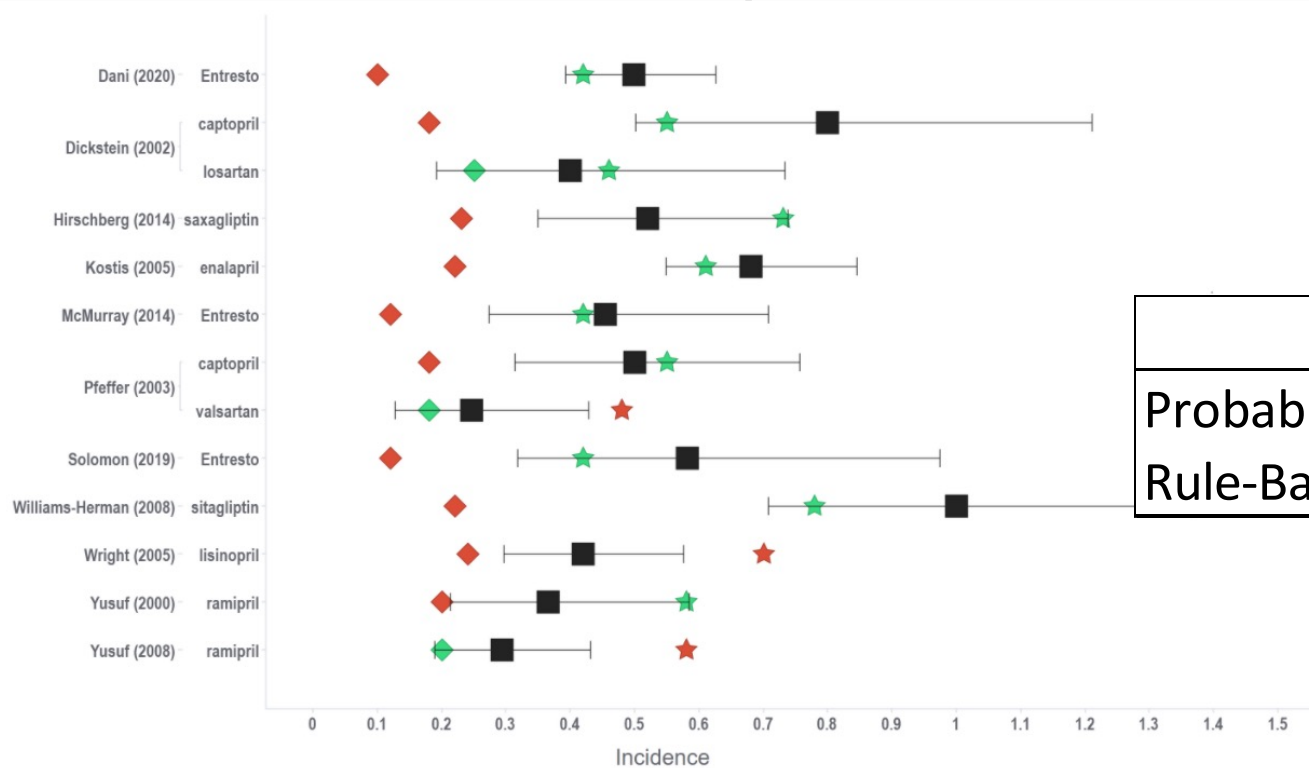
# Results: Angioedema - Probabilistic (Cut-point 0.67) v. Rule-based



**Within 95% CI:**  
Probabilistic: 77%  
Rule-based: 23%



# Results: Angioedema - Probabilistic (Cut-point 0.67) v. Rule-based

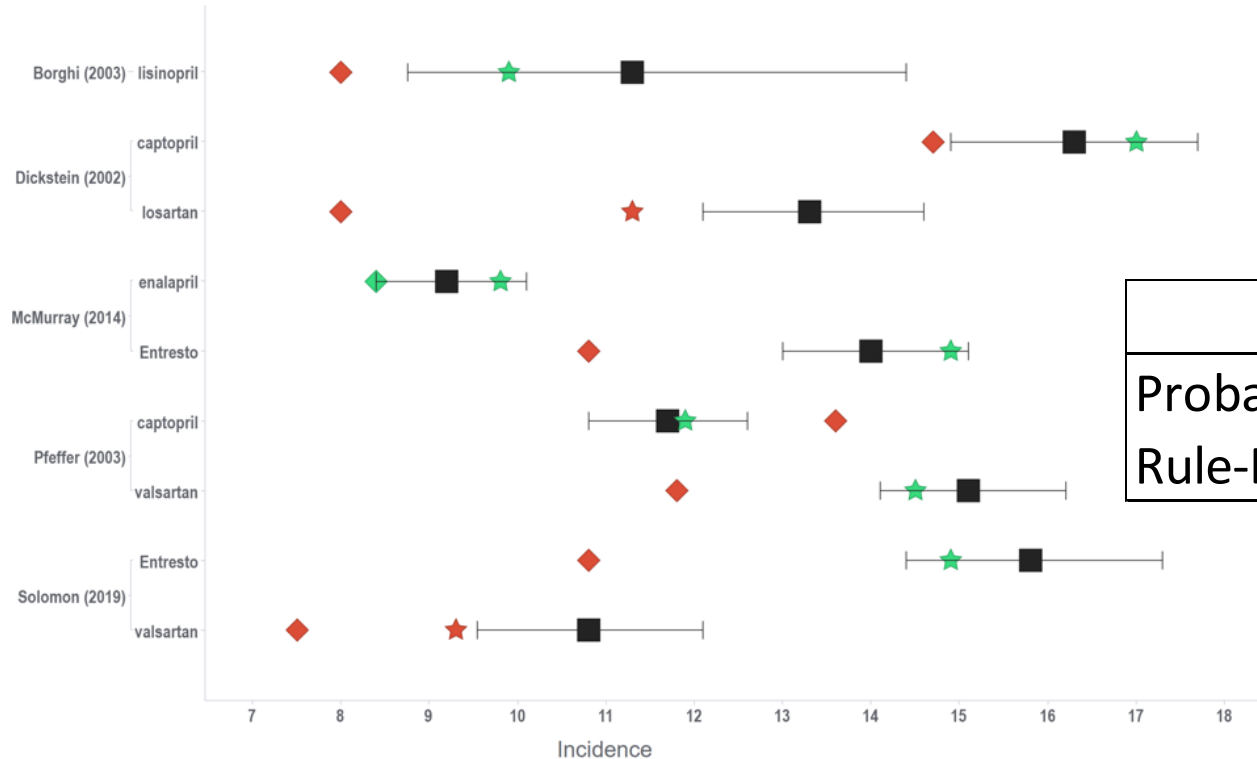


**Within 95% CI:**  
Probabilistic: 77%  
Rule-based: 23%

	Sensitivity	PPV
Probabilistic	0.71	0.87
Rule-Based	0.23	0.78

“Entresto” = sacubitril/valsartan

# Results: Hypotension - Probabilistic (Cut-point 0.67) v. Rule-based



**Within 95% CI:**  
Probabilistic: 78%  
Rule-based: 11%

	Sensitivity	PPV
Probabilistic	0.82	0.89
Rule-Based	0.56	0.72

★ Probabilistic ■ RCT ◆ Rule-based Green – within 95% CI; Red – outside 95% CI

“Entresto” = sacubitril/valsartan

## Results: Probabilistic v. Rule-based

For both angioedema and hypotension:

- The 0.67 cut-point provided the closest match for the RCT results.
- Lower cut-points – produced higher average incidence estimates compared to RCT results.
- Higher cut-points – produced lower average incidence estimates compared to RCT results.

## Conclusions

- Probabilistic phenotype algorithms (PA) for angioedema and hypotension estimated incidence closer to the results from RCTs than rule-based PAs.
- The performance of probabilistic PAs was superior to rule-based PAs on PPV and sensitivity.
- Future research is needed to evaluate the performance of probabilistic PAs in postmarket safety settings and to determine how they could potentially be used to estimate the incidence of drug adverse effects.

# References

## References:

1. Agarwal V, Podchiyska T, Banda JM, Goel V, Leung TI, Minty EP, et al. Learning statistical models of phenotypes using noisy labeled training data. *J Am Med Inform Assoc.* 2016;23(6):1166-73.
2. Banda JM, Halpern Y, Sontag D, Shah NH. Electronic phenotyping with APHRODITE and the Observational Health Sciences and Informatics (OHDSI) data network. *AMIA Jt Summits Transl Sci Proc.* 2017;2017:48-57.
3. Banda JM, Seneviratne M, Hernandez-Boussard T, Shah NH. Advances in Electronic Phenotyping: From Rule-Based Definitions to Machine Learning Models. *Annu Rev Biomed Data Sci.* 2018;1:53-68.
4. Reps JM, Schuemie MJ, Suchard MA, Ryan PB, Rijnbeek PR. Design and implementation of a standardized framework to generate and evaluate patient-level prediction models using observational healthcare data. *J Am Med Inform Assoc.* 2018;25(8):969-75.
5. Swerdel JN, Schuemie M, Murray G, Ryan PB. PheValuator 2.0: Methodological improvements for the PheValuator approach to semi-automated phenotype algorithm evaluation. *J Biomed Inform.* 2022:104177.
6. Dickstein K, Kjekshus J. Comparison of the effects of losartan and captopril on mortality in patients after acute myocardial infarction: the OPTIMAAL trial design. *Optimal Therapy in Myocardial Infarction with the Angiotensin II Antagonist Losartan.* *Am J Cardiol.* 1999;83(4):477-81.
7. Pfeffer MA, McMurray JJ, Velazquez EJ, Rouleau JL, Køber L, Maggioni AP, et al. Valsartan, captopril, or both in myocardial infarction complicated by heart failure, left ventricular dysfunction, or both. *N Engl J Med.* 2003;349(20):1893-906.
8. Kostis JB, Kim HJ, Rusnak J, Casale T, Kaplan A, Corren J, et al. Incidence and characteristics of angioedema associated with enalapril. *Arch Intern Med.* 2005;165(14):1637-42.
9. Dani SS, Ganatra S, Vaduganathan M. Angioedema with sacubitril/valsartan: Trial-level meta-analysis of over 14,000 patients and real-world evidence to date. *Int J Cardiol.* 2021;323:188-91.
10. McMurray JJ, Packer M, Desai AS, Gong J, Lefkowitz MP, Rizkala AR, et al. Angiotensin-neprilysin inhibition versus enalapril in heart failure. *N Engl J Med.* 2014;371(11):993-1004.
11. Solomon SD, McMurray JJV, Anand IS, Ge J, Lam CSP, Maggioni AP, et al. Angiotensin-Neprilysin Inhibition in Heart Failure with Preserved Ejection Fraction. *N Engl J Med.* 2019;381(17):1609-20.
12. Yusuf S, Sleight P, Pogue J, Bosch J, Davies R, Dagenais G. Effects of an angiotensin-converting-enzyme inhibitor, ramipril, on cardiovascular events in high-risk patients. *N Engl J Med.* 2000;342(3):145-53.
13. Yusuf S, Teo KK, Pogue J, Dyal L, Copland I, Schumacher H, et al. Telmisartan, ramipril, or both in patients at high risk for vascular events. *N Engl J Med.* 2008;358(15):1547-59.
14. Hirshberg B, Parker A, Edelberg H, Donovan M, Iqbal N. Safety of saxagliptin: events of special interest in 9156 patients with type 2 diabetes mellitus. *Diabetes Metab Res Rev.* 2014;30(7):556-69.
15. Williams-Herman D, Round E, Swern AS, Musser B, Davies MJ, Stein PP, et al. Safety and tolerability of sitagliptin in patients with type 2 diabetes: a pooled analysis. *BMC Endocr Disord.* 2008;8:14.
16. Borghi C, Ambrosioni E. Double-blind comparison between zofenopril and lisinopril in patients with acute myocardial infarction: results of the Survival of Myocardial Infarction Long-term Evaluation-2 (SMILE-2) study. *Am Heart J.* 2003;145(1):80-7.
17. Wright JT, Dunn JK, Cutler JA, Davis BR, Cushman WC, Ford CE, et al. Outcomes in Hypertensive Black and Nonblack Patients Treated With Chlorthalidone, Amlodipine, and Lisinopril. *JAMA.* 2005;293(13):1595-608.

# Thank you!



Questions?

Come to Poster

115

2024 Global OHDSI lightning talk

# Exploring the interplay between metabolic syndrome and brain volume in depression : Basis for Phenotype-Based Classification

**Department of Biomedical Science and Biomedical Informatics**

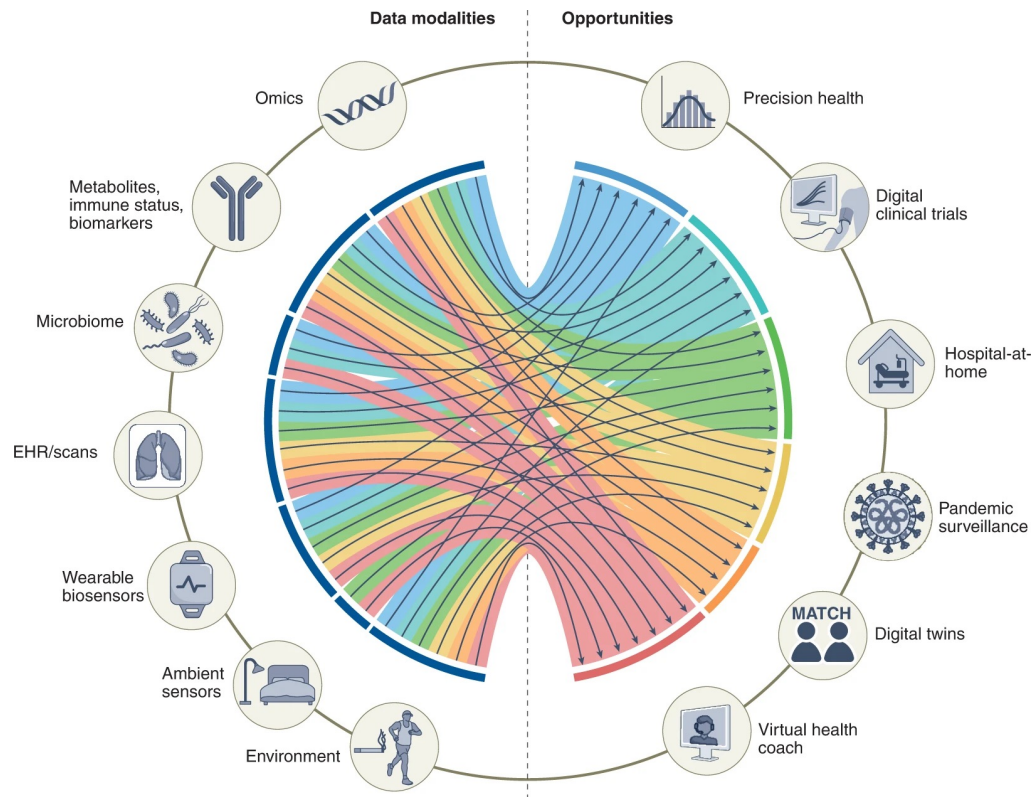
**Sujin Gan**

PhD student, Ajou University School of Medicine

Advisor: Professor Rae Woong Park

# Latest Research Trends

## Growing utilization of multimodal approaches in medical research

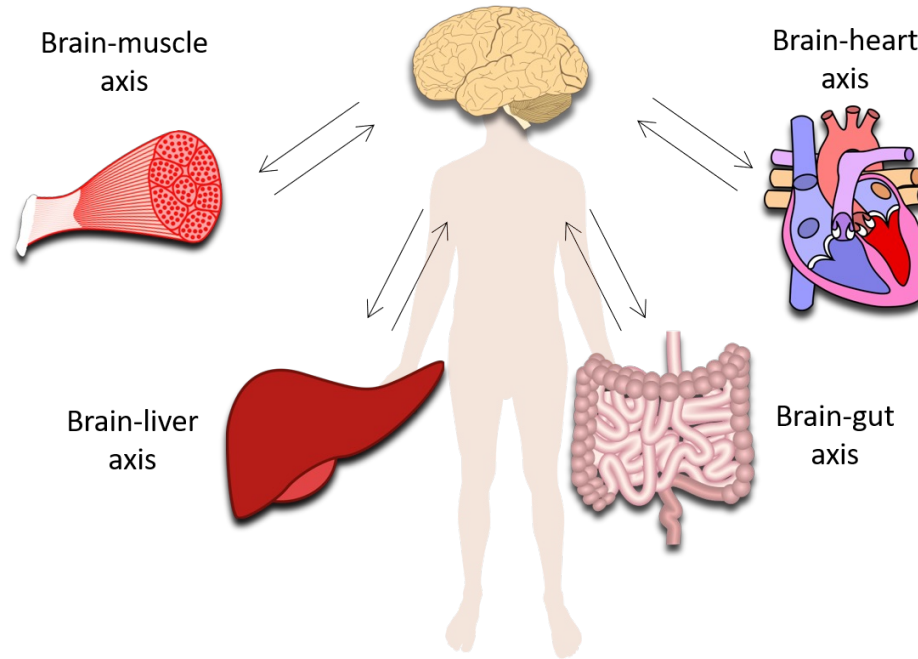




## Latest Research Trends

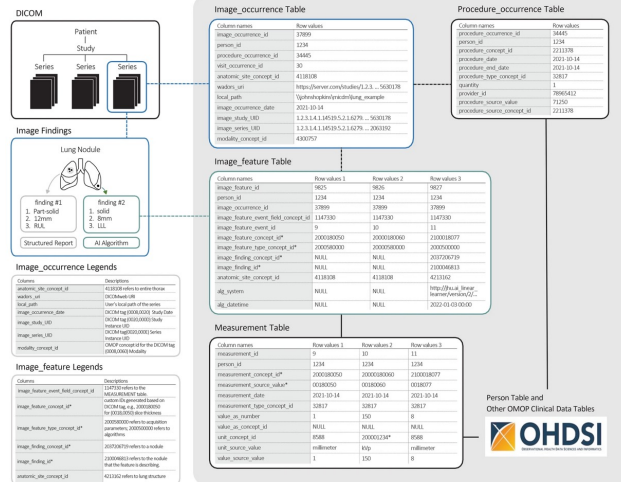
### Multi-organ Interaction

How do these interconnected systems between the various organs contribute to overall health and disease mechanisms?



Understanding the interaction axes between the brain and various

## OMOP CDM Medical Imaging Extension



Park, W.Y., Jeon, K., et al. *J Digit Imaging. Inform. med.* 2024

```

graph TD
    Person[Person  
Personal ID, Sex, Age, Race] --> RadiologyCDM
    subgraph RadiologyCDM [Radiology-CDM]
        RadiologyOccurrence[Radiology Occurrence]
        RadiologyImage[Radiology Image]
        RadiologyOccurrence -- "Occurrence date/time, Protocol concept ID, Total image count" --> RadiologyImage
    end
    RadiologyImage -- "Image resolution, Contrast administration status, Image photographing direction" --> ProcedureOccurrence
    subgraph ProcedureOccurrence [Procedure Occurrence]
        DateType[Date, Type]
    end

```

Park C, You SC, et al. Yonsei Med J. 2022

```

graph TD
    CareSite[Care Site] --> GenomicTest[Genomic Test]
    Person[Person] --> VarOcc[Variant Occurrence]
    CondOcc[Condition Occurrence] --> VarOcc
    ProcOcc[Procedure Occurrence] --> VarOcc
    Specimen[Specimen] --> VarOcc
    GenomicTest --> TargetGene[Target Gene]
    TargetGene --> VarOcc
    VarOcc --> VarAnnotation[Variant Annotation]
  
```

**Genomic Test**  
Genomic test name/version, Sequencing device, Reference genome, Analysis tools

**Target Gene**  
Target gene list (HGNC ID and Symbol)

**Variant Occurrence**  
Reference sequence, HGVS nomenclature, Read depth, Exon number, Variant type

**Variant Annotation**  
Annotation database, Variant origin, Pathogeny, Allele frequency

**Care Site**  
Care site name, Location, Place of service

**Person**  
Personal ID, Sex, Birth, Race

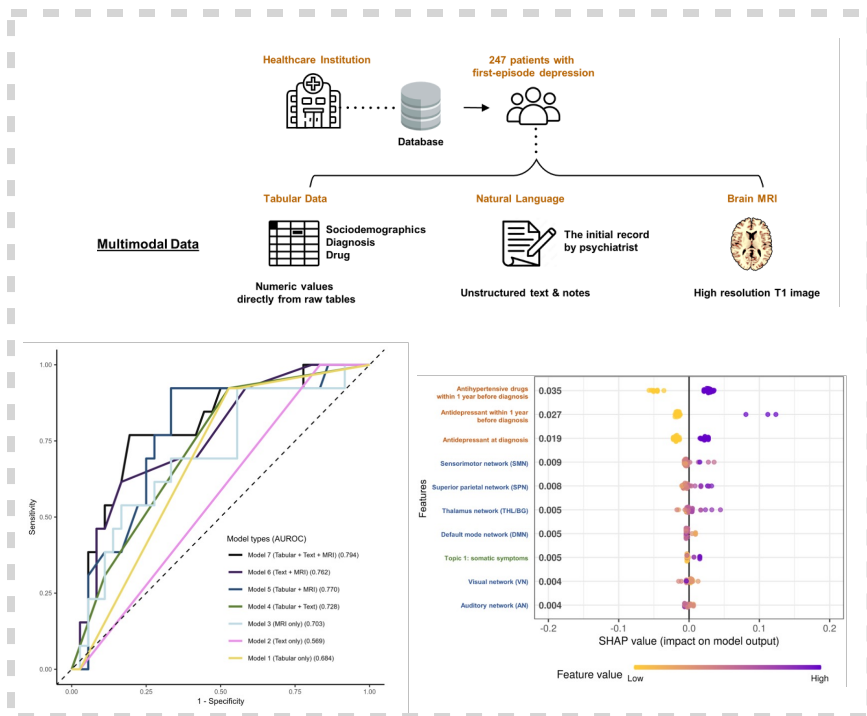
**Condition Occurrence**  
Condition, Start/end date, Type, Stop reason

**Procedure Occurrence**  
Procedure, Date, Type

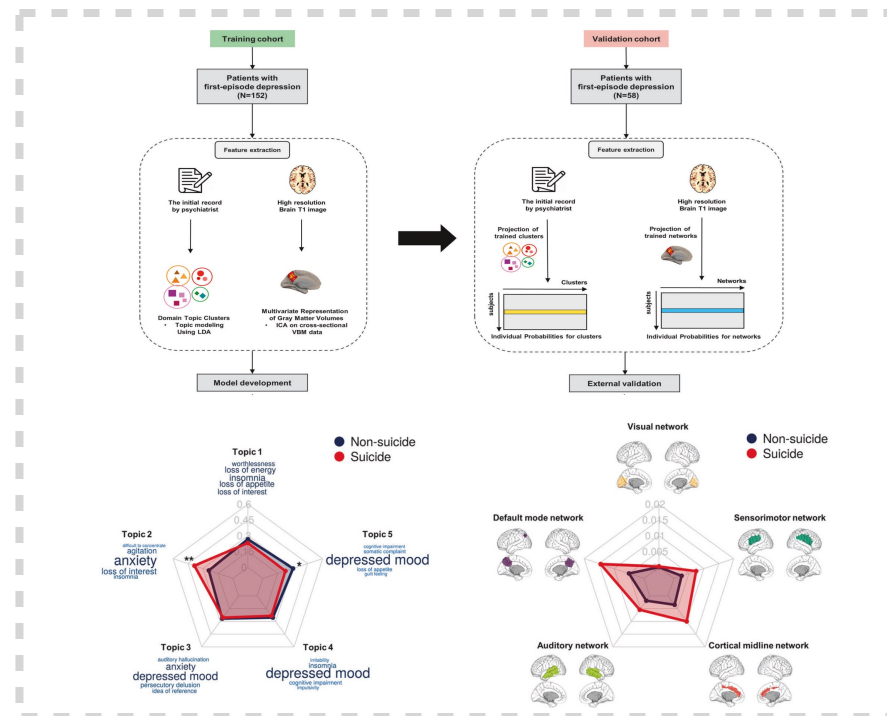
**Specimen**  
Specimen, Anatomic site, Date, Disease status

Shin SJ, You SC, Park RW. et al. J Med Internet Res. 2019

# Multimodal Integration in OMOP CDM

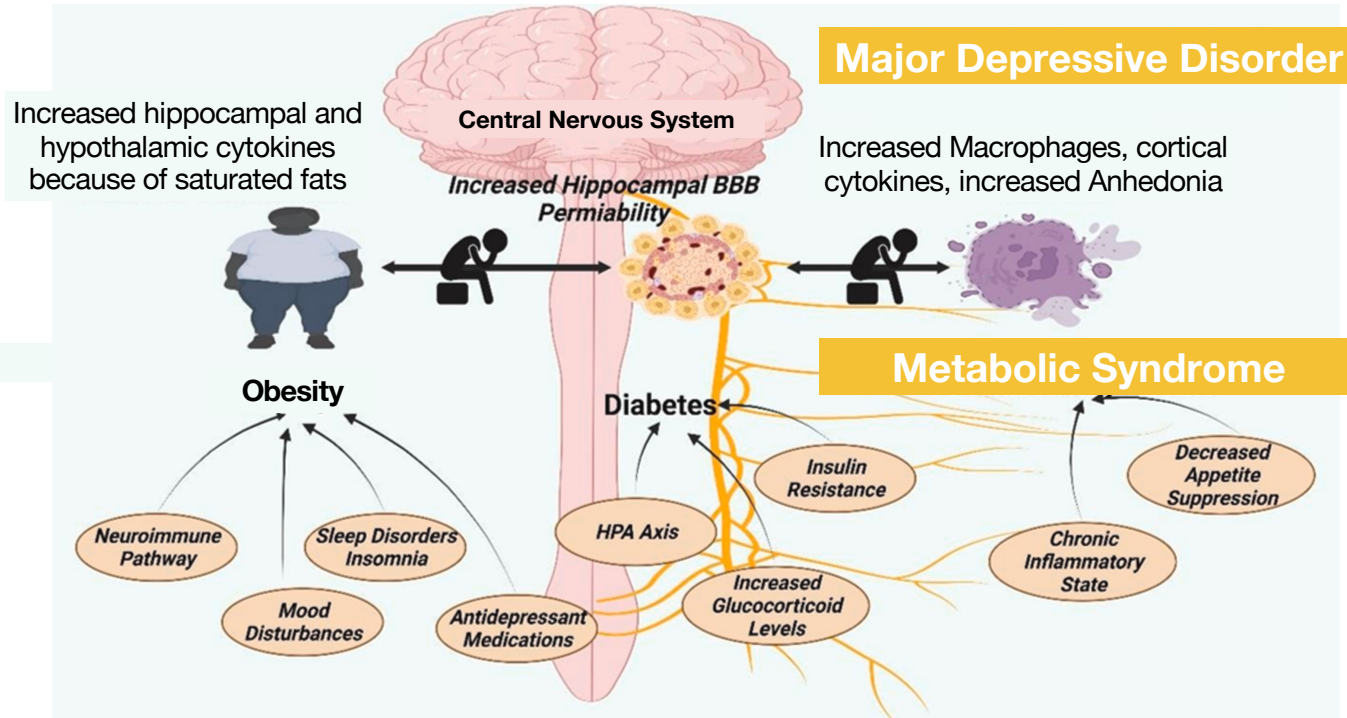


Lee, D.Y., Kim N.R., et al, *Psychiatry Research*, 2024



Lee, D.Y., Byeon, G, et al. *Transl Psychiatry*, 2024

# Bi-directional relationship between Major Depressive Disorder and Metabolic Syndrome

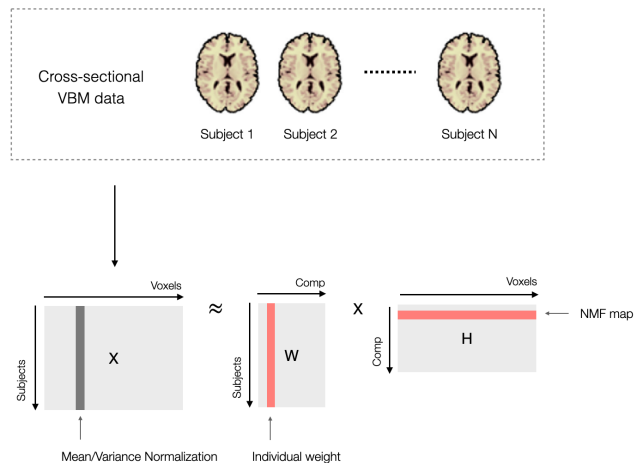


# Research Hypothesis

- **The bidirectional interaction** between depression and metabolic syndrome is **mediated** by specific **brain volume components** and **peripheral laboratory markers**.
- These components can serve as biomarkers **to classify the presence of metabolic syndrome** in patients with major depressive disorder (MDD).

# Overview

## Dimension Reduction in VBM



## Canonical Correlation Analysis

Brain features

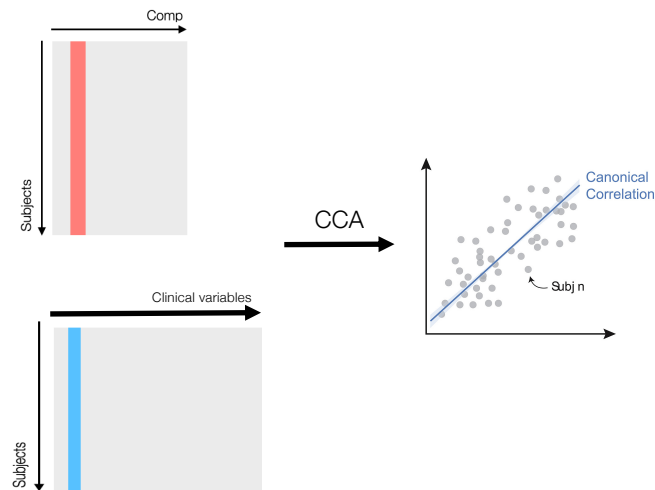


Clinical variables

Metabolic Syndrome variables

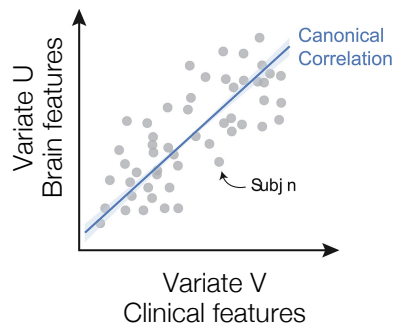
SBP	DBP	BMI
HDL	BST	Triglyceride

- 38 Blood laboratory measurements (CRP, Albumin, LDL, AST, etc...)



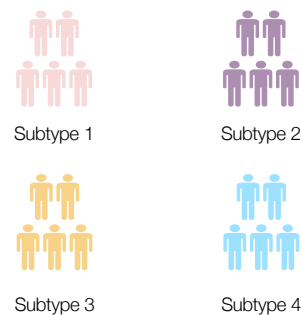
## Classification

Significant Mode



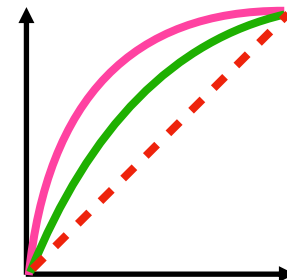
K means

MDD subtypes



XGB

Classification



# Study population



OMOP CDM

- Electronic Health Records database from Ajou University School of Medicine (AUSOM)
- January 1994 to July 2023 (OMOP-CDM 534)



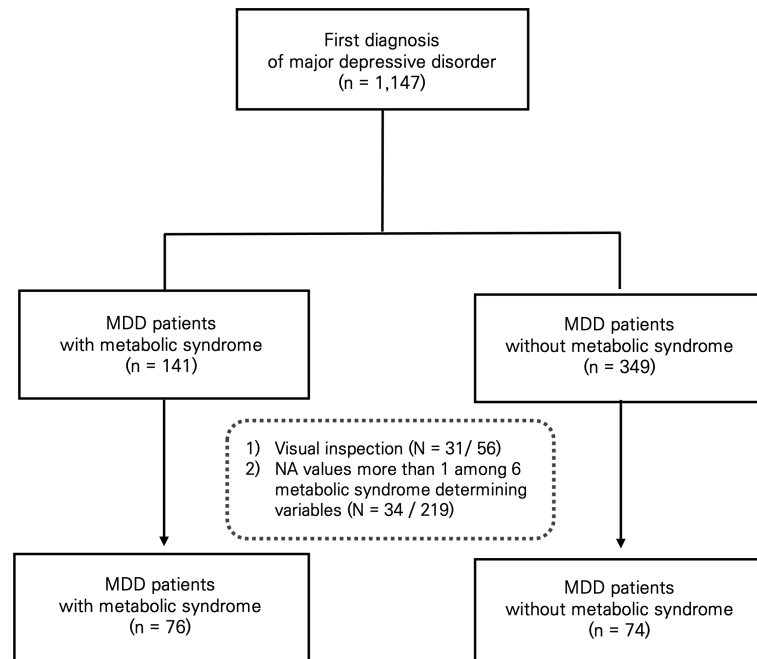
Cohort

- Major Depressive Disorder diagnosis, for the first time
- Brain MRI procedures (- 365 d ~ + 30 d from the index date)
- No history of psychiatric comorbidities (bipolar disorder, schizophrenia, or dementia), substance disorders, brain injury, hydrocephalus

## Definition for Metabolic syndrome

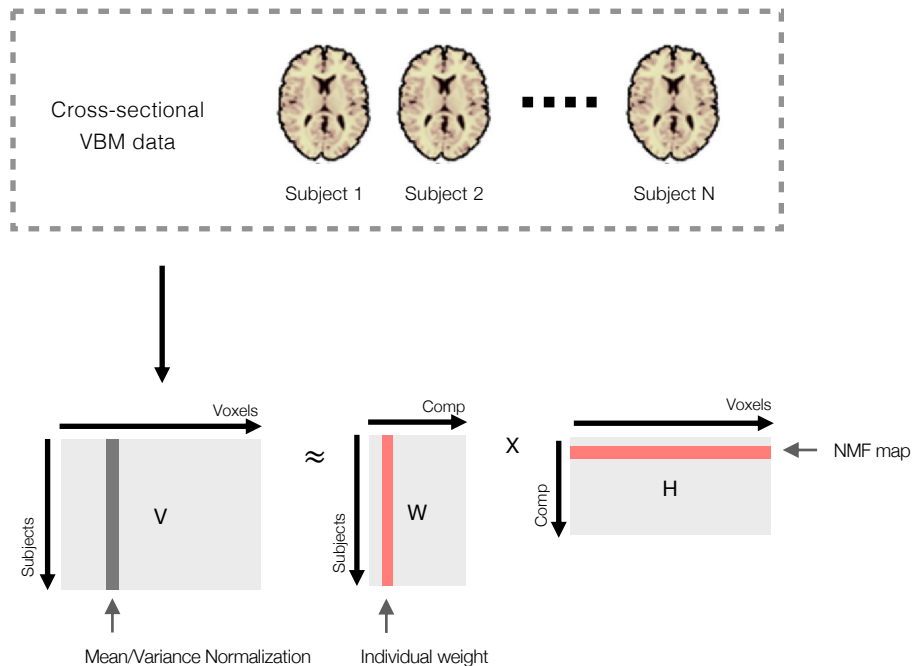
- More than 2 out of the following MetS criteria
  - 1) Antihypertensive drugs uses or **systolic blood pressure (SBP)** > 130
  - 2) Antidiabetic drugs uses or **blood sugar level (BST)** > 100
  - 3) Low **HDL-cholesterol** levels (<40 mg/dL for men and <50 mg/dL for women)
  - 4) Hypertriglyceridemia (**Triglyceride** ≥ 150 mg/dL)
  - 5) **BMI** ≥ 30

## Study population flow chart



# Dimension Reduction

## Non-negative Matrix Factorization



$$V \approx W \cdot H$$

The matrix  $V$  is represented by the two smaller matrices  $W$  and  $H$ , which, when multiplied, approximately reconstruct  $V$

$$H_{[i,j]}^{n+1} \leftarrow H_{[i,j]}^n \frac{((W^n)^T V)_{[i,j]}}{((W^n)^T W^n H^n)_{[i,j]}}$$

and

$$W_{[i,j]}^{n+1} \leftarrow W_{[i,j]}^n \frac{(V(H^{n+1})^T)_{[i,j]}}{(W^n H^{n+1} (H^{n+1})^T)_{[i,j]}}$$

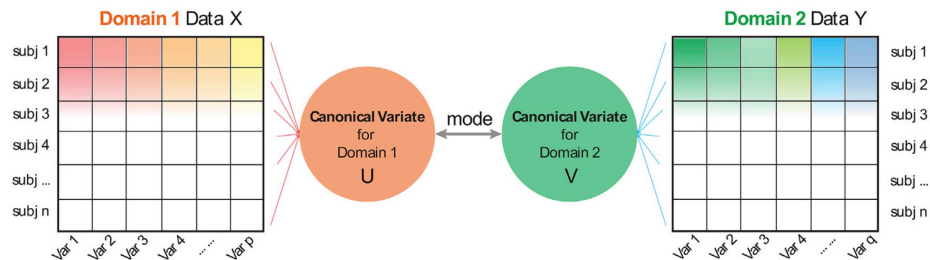
Initialize:  $W$  and  $H$  non negative.

Until  $W$  and  $H$  are stable.



# Dimension Reduction

## Canonical Correlation Analysis



Original Variables  
Canonical Vector

X or Y a or b

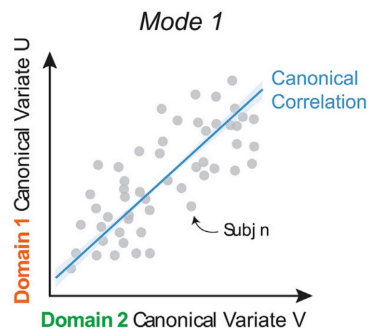
Canonical Variate

U or V

Var 1	X	0.4
Var 2	X	0.2
Var 3	X	0
Var 4	X	-0.1
⋮	⋮	⋮
Var p	X	0.3

=

Var 1
+
Var 2
+
Var 4
+
Var p



$$\Sigma_{XX} = \text{Cov}(X, X) = X^T X \text{ and } \Sigma_{YY} = \text{Cov}(Y, Y) = Y^T Y$$

$$\Sigma_{XY} = \text{Cov}(X, Y) = X^T Y$$



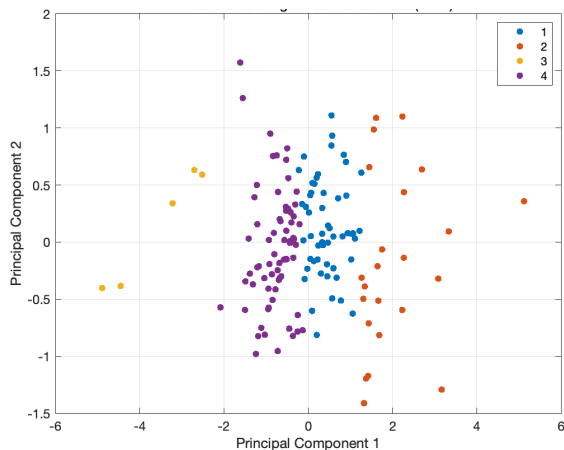
$$U = c^T \Sigma_{xx}^{-1/2} X = a^T X$$

$$V = d^T \Sigma_{yy}^{-1/2} Y = b^T Y.$$

# Clustering and Classification

## K means clustering analysis

- Using **only the first component in rCCA**
- Determination of the optimal number of k (1:10)
  - based on Cophenetic Correlation, Silhouette Coefficient, Residual Sum of Squares (RSS)



Optimal Number of K = 4  
Mean silhouette: 0.523



**4 MDD  
Subtypes**

## XGBoost

- Dataset split
  - Train set : Test set (75 : 25)
  - 5 folds cross validation
- Parameters
  - Numbers of estimators: 50, 100, 200
  - Learning rate: 0.01, 0.1, 0.2
  - Max depth: 3, 4, 5
  - Colsample by tree: 0.8, 0.9, 1.0
  - Subsample: 0.8, 0.9, 1.0

## Model description

---

**Model 1** Demographics + MetS

---

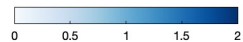
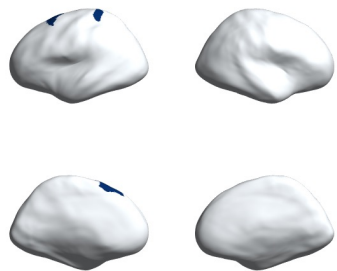
**Model 2** Demographics + MetS + VBM NMF

---

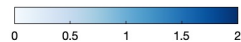
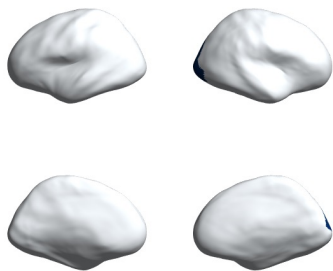
**Model 3** Demographics + MetS + VBM NMF + Subtypes

# NMF-derived structural networks

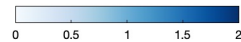
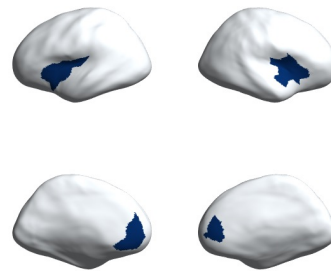
NMF Component 1



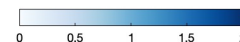
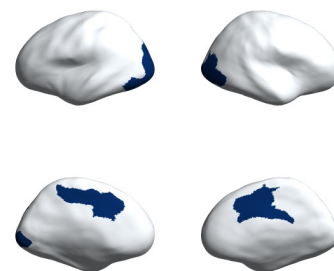
NMF Component 2



NMF Component 3

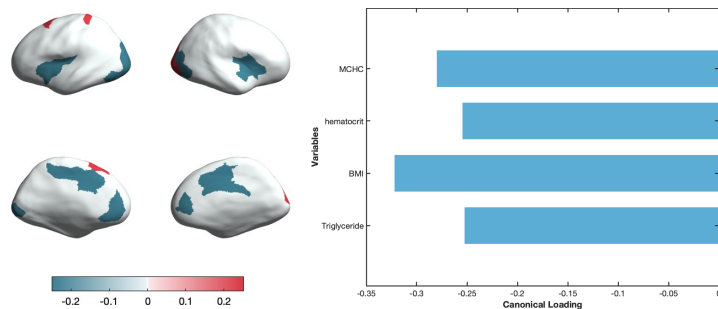


NMF Component 4

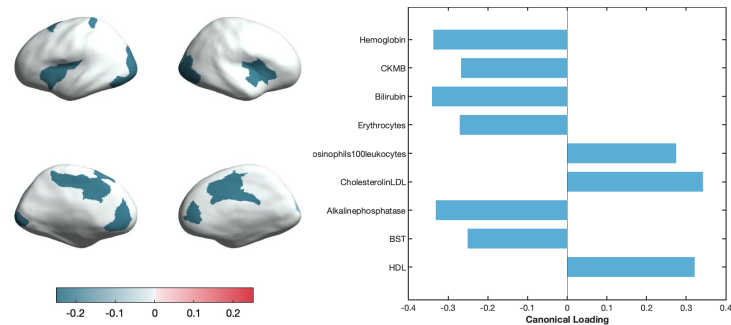


# CCA Multivariate patterns of brain imaging and clinical variables

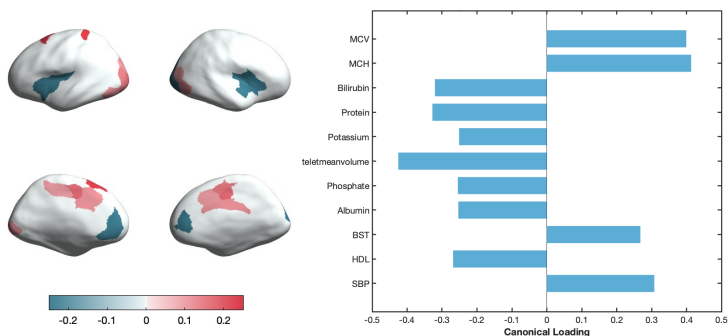
## CCA Component 1



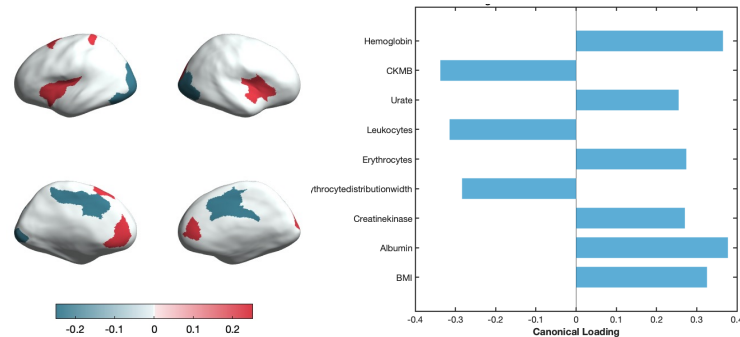
## CCA Component 2



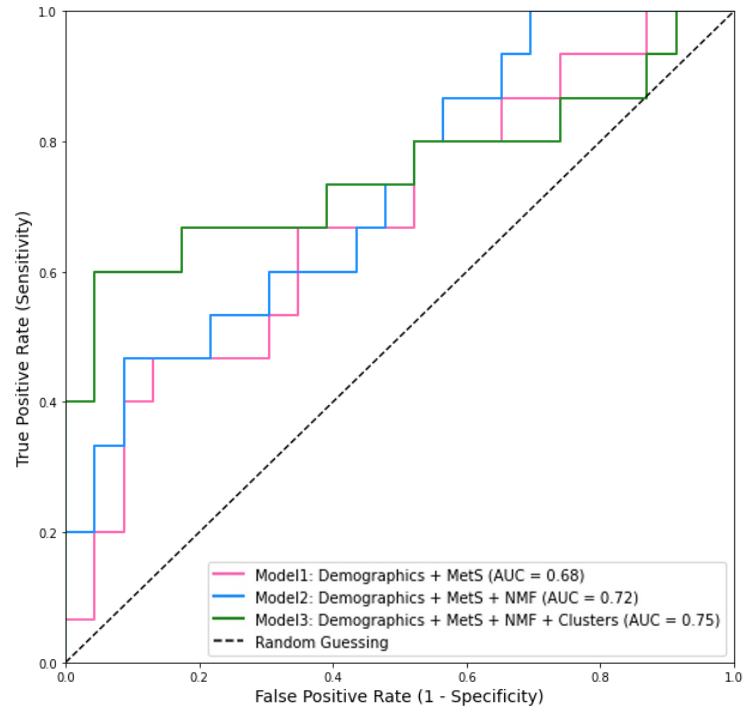
## CCA Component 3



## CCA Component 4



# Classification Model Performance



	Model 1	Model 2	Model 3
AUROC	0.680	0.720	0.750
Accuracy	0.605	0.580	0.632
Precision	0.500	0.610	0.640
Recall	0.667	0.610	0.733
F1-score	0.571	0.580	0.611

## Research Summary

- Through **the initial exploration of brain-body interactions**, we applied OMOP CDM to integrate brain imaging and clinical data for **classifying metabolic syndrome in MDD patients**
- Combined brain imaging data (reduced via NMF) with clinical variables using CCA
- Demonstrates **the potential of CDM** in pioneering **multimodal studies** and future scalability in **multi-organ interactions research**

# Thank you



## Authors

**Sujin Gan, R.N**

**Narae Kim, M.S**

## Advisors

**Rae Woong Park, M.D., PhD**

**Bumhee Park, PhD**



# **CohortConstructor – an R package to support cohort building pipelines**



Ed Burn





# Cohorts

- Cohorts are a key building block in research studies – people fulfilling some criteria for some amount of time
- Established OHDSI tools (such as ATLAS/ Capr and CIRCE) allow us to define cohorts that can be instantiated in a database and stored in a library for future re-use
- **However,**
  - computational challenges remain when making many cohorts, and
  - bespoke cohort logic may not be supported by current tools



# Many cohorts, all at once



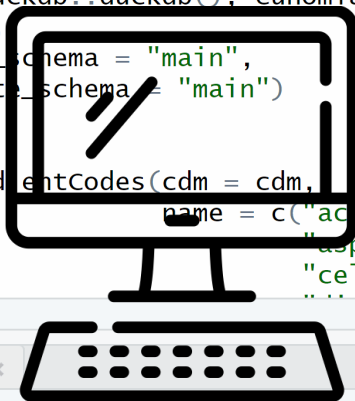
# Building cohorts by domain

- In OHDSI studies cohorts are typically defined independently and instantiated sequentially
- *CohortConstructor* builds cohorts by domain instead





```
1
2 # libraries|
3 library(CDMConnector)
4 library(CodelistGenerator)
5 library(CohortConstructor)
6
7 # connect to eunomia data
8 con <- DBI::dbConnect(duckdb::duckdb(), eunomia_dir())
9 cdm <- cdm_from_con(con,
10                     cdm_schema = "main",
11                     write_schema = "main")
12
13 # get drug codes
14 meds_cs <- getDrugIngredientCodes(cdm = cdm,
15                                   name = c("acetaminophen",
16                                           "aspirin",
17                                           "celecoxib",
```



2:12 (Top Level) ▾

Console Terminal x Background Jobs x

R 4.4.0 · ~/

	<int>	<int>	<int>
1	1	137	137
2	2	4379	1927
3	3	1800	1800
4	4	13908	2679
5	5	830	830
6	6	35	35

```
> cdm_disconnect(cdm)
```

```
> |
```

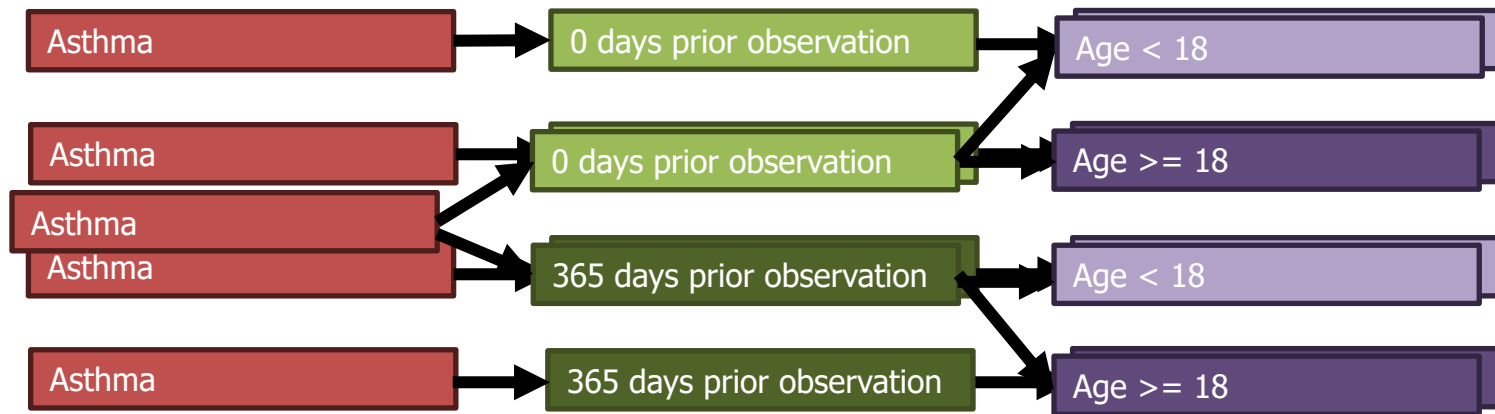


# From one cohort to many



# Deriving cohorts from other cohorts

- Often, we build many cohorts that only vary slightly
- *CohortConstructor* encourages creating base cohorts from which study cohorts can be derived





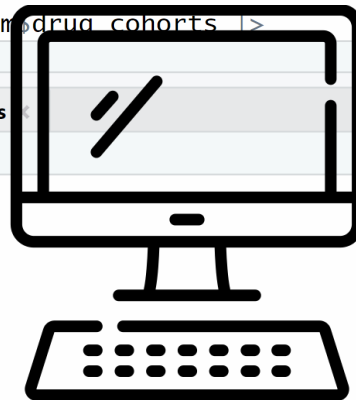
```
33 |  
34 cdm$drug_cohorts <- cdm$drug_cohorts |>  
35   requirePriorObservation(c(0, 365))  
36  
37 settings(cdm$drug_cohorts)  
38  
39 cohortCount(cdm$drug_cohorts)  
40  
41 cdm$drug_cohorts <- cdm$drug_cohorts |>
```

33:1 (Top Level) ▾

Console Terminal x Background Jobs x

R 4.4.0 · ~/ ↗

> |





# Flexible cohort pipelines



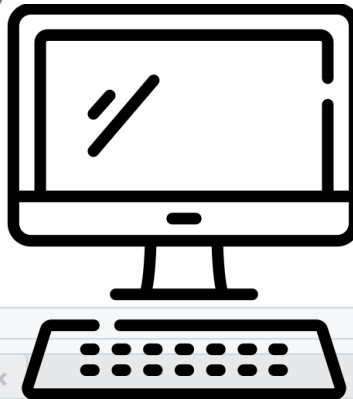


# Cohort utilities

- Often, custom study-specific requirements need to be applied to a cohort
- CohortConstructor provides various utilities to support building bespoke cohorts
  - Add or subtract days from cohort entry and exit
  - Create cohorts based on age (e.g. entry on 18<sup>th</sup> birthday)
  - Reset entry/ exit on first/ last of set of date variables
  - Require cohort subjects are present in (or absence from) another cohort or table in some time window
  - Take a random sample of each cohort
  - Persist cohort entry across multiple observation periods (*next release*)



```
1 |
2 cdm$hip_fx <- conceptCohort(cdm,
3                             conceptSet = list("hip_fracture" = 4230399L),
4                             name = "hip_fx",
5                             exit = "event_start_date") |>
6   padCohortEnd(180) |>
7   requireIsFirstEntry() |>
8   sampleCohorts(n = 100)
9
10 settings(cdm$hip_fx)
11
12 cohortCount(cdm$hip_fx)
13
14
```



1:1 (Top Level) ▾

Console

Terminal x

Background Jobs x

R 4.4.0 · ~/ ↗

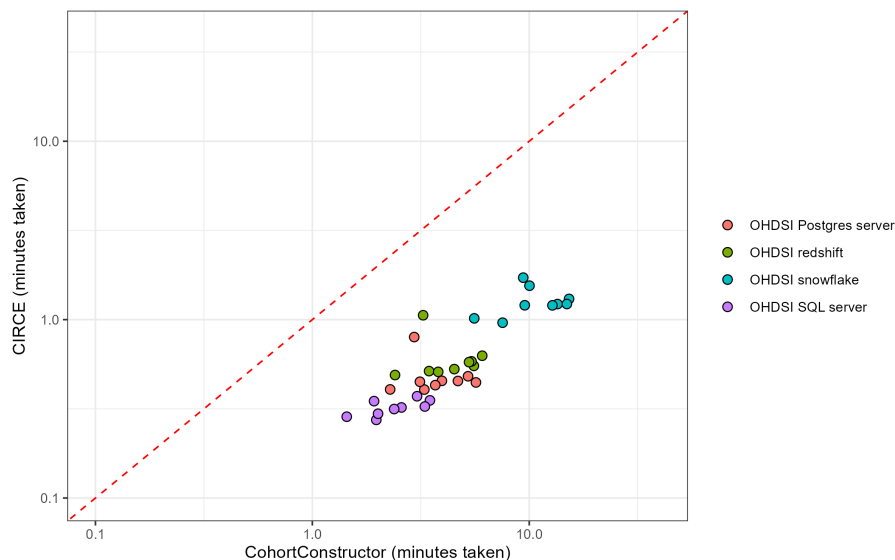
> |



# Benchmarking

# Benchmark results

- Selected 9 cohorts from the OHDSI phenotype library

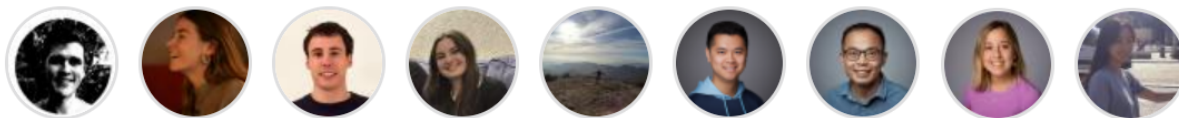




# More information

Package website	<a href="https://ohdsi.github.io/CohortConstructor">ohdsi.github.io/CohortConstructor</a>
GitHub	<a href="https://github.com/OHDSI/CohortConstructor">github.com/OHDSI/CohortConstructor</a>
Benchmarking code	<a href="https://github.com/oxford-pharmacoepi/BenchmarkCohortConstructor">github.com/oxford-pharmacoepi/BenchmarkCohortConstructor</a>

## Contributors 9



# Unlocking Efficiency in Real-world Collaborative Studies

A Multi-site International Study with **COLA-GLMM**  
(Collaborative One-shot Lossless Algorithm for  
Generalized Linear Mixed Model)

Jiayi (Jessie) Tong, Assistant Professor

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health

Oct 23, 2024, OHDSI 2024 Global Symposium

## Map of Collaborators

The OHDSI community brings together volunteers from around the world to establish open community data standards, develop open-source software, conduct methodological research, and apply scientific best practices to both answer public health questions and generate reliable clinical evidence.

Our community is ALWAYS seeking new collaborators. Do you want to focus on data standards or methodological research? Are you passionate about open-source development or clinical applications? Do you have data that you want to be part of global network studies? Do you want to be part of a global community that truly values the benefits of open science? Add a dot to the map below and JOIN THE JOURNEY!

### OHDSI By The Numbers

- 3,266 collaborators
- 80 countries
- 21 time zones
- 6 continents
- 1 community



## A **Primary Challenge** in Multi-site International Study

### **Individual Patient-level Data (IPD) cannot be shared across sites**

- Regulatory Approval Processes
- Country-Specific Laws (e.g., HIPAA in USA, PIPEDA in Canada)
- Institutional Policies Data Sharing Restrictions



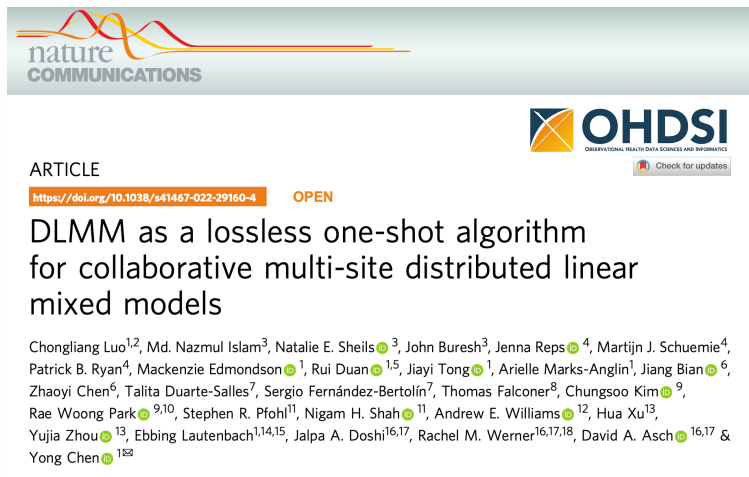
### **Privacy-Preserving Federated Learning Algorithms**

- Enables fitting statistical models in a federated manner
- Requires summary statistics, instead of IPD
- Ensures data privacy and security



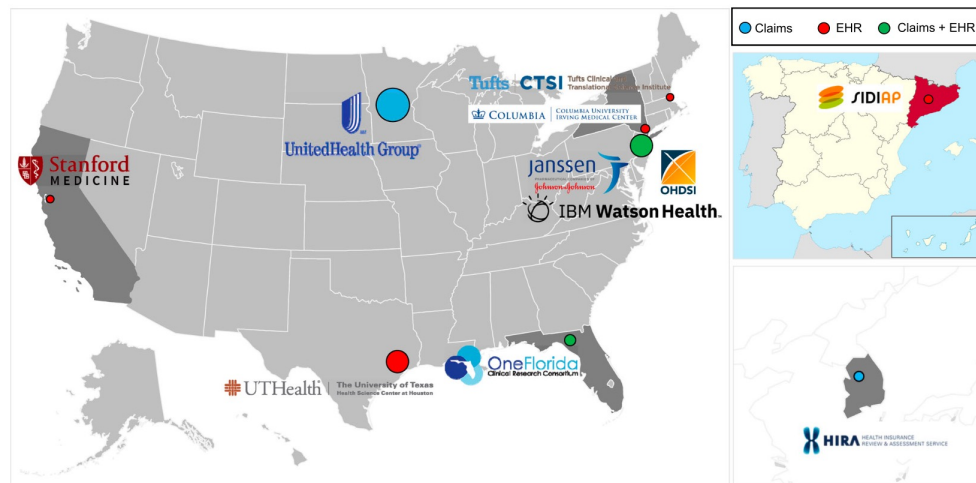
# Multi-site Collaborative Study with Observational Data

-- An OHDSI Study Using Distributed **Linear** Mixed Model (DLMM)



Luo et al, 2022, Nature Communications

**Outcome of interest:** Length of Stay (continuous outcome)



11 databases from 3 countries

**Investigation on the associations between demographic and clinical characteristics and length of hospital stay in COVID-19 patients**

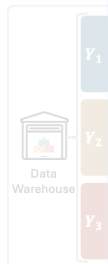
Luo C, Islam MN, Sheils NE, Buresh J, Reps J, Schuermie MJ, Ryan PB, Edmondson M, Duan R, Tong J, Marks-Anglin A. DLMM as a lossless one-shot algorithm for collaborative multi-site distributed linear mixed models.

Nature communications. 2022 Mar 30;13(1):1678.

# Two Ideal Properties of Federated Learning Algorithms

Lossless

One-shot



Pool

To date, only a few algorithms have successfully achieved both lossless and one-shot properties simultaneously:

- Linear Regression (i.e., Chen et al., 2006, IEEE)
- Linear mixed models (i.e., Luo et al., 2022, Nature Communications)



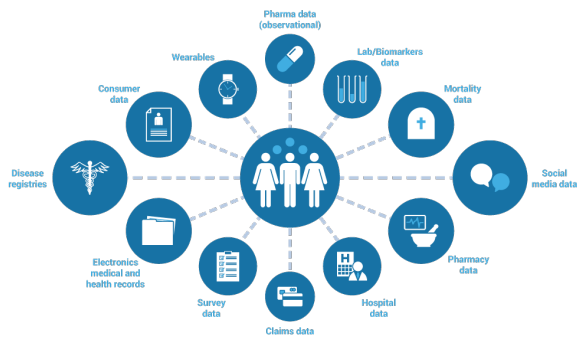
Results

Identical Results

No accuracy loss due to data sharing constraints

Only a single round of communication is needed

# Challenges in Real-world Data



## Non-continuous outcomes

- Binary outcome
- Categorical outcome
- Count outcome
- ...



## Between-site heterogeneity

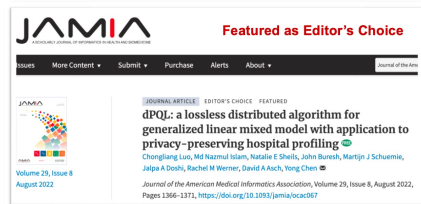
We need **Federated Learning Algorithms for Generalized Linear Mixed Model (GLMM)**

**GLMM Model:**  $y_{ij} = b_i + \beta x_{ij} + \epsilon_{ij}$ , where  $\epsilon_{ij} \sim N(0, \sigma^2)$ ,  $b_i \sim N(0, \tau^2)$  is site-specific random effect for k-th site,  $\beta$  is fixed effect.

# Existing Works on Federated Learning Algorithms for GLMM



Zhu et al, 2020, *Bioinformatics*



Luo et al, 2022, *JAMIA*



Yan et al, 2022, *arxiv*



Lossless



One-shot



Communication  
Round

Iterative  
(500~1000 rounds)

< 5 rounds

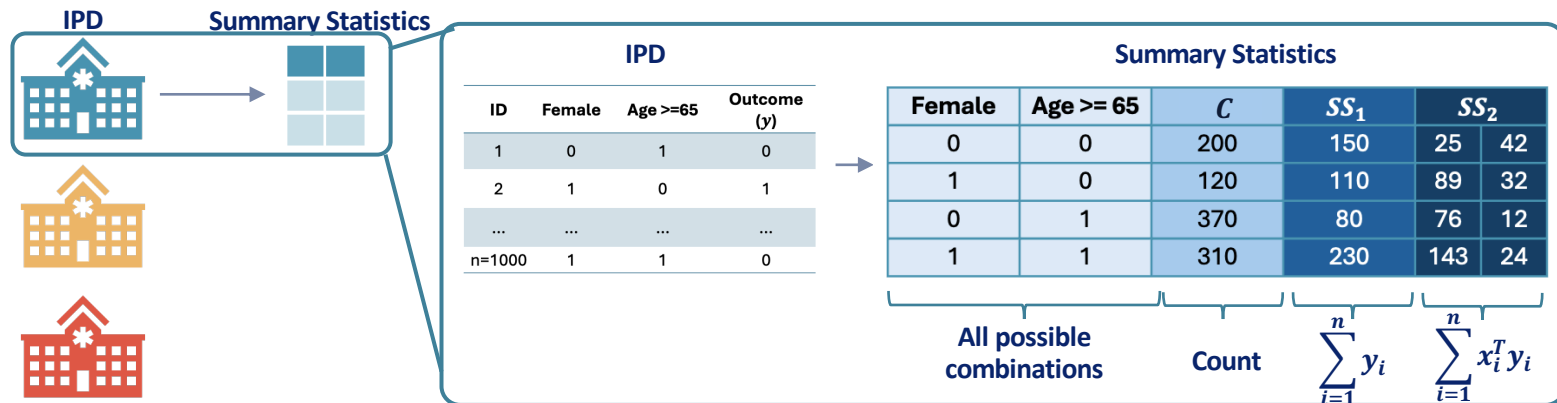
1 or 2 rounds\*  
(Depends on initialization)

1 round

## Proposed Method – COLA-GLMM

### Collaborative One-shot Lossless Algorithm for Generalized Linear Mixed Model

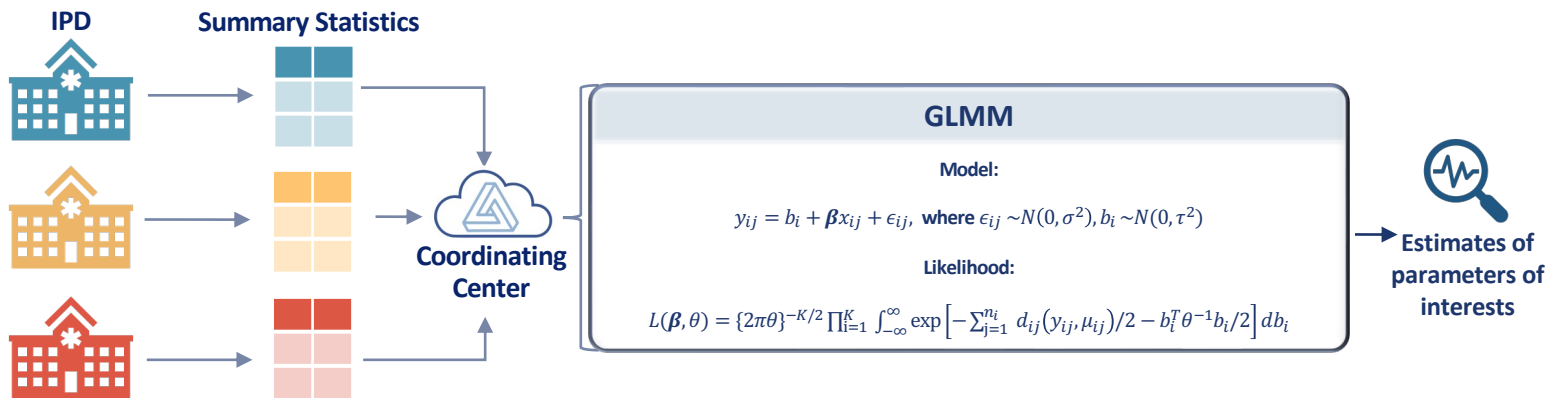
- Suppose that a **common** set of covariates are available at all collaborating sites.
- The covariates have been **standardized into categorical variables**.
- Pipeline:



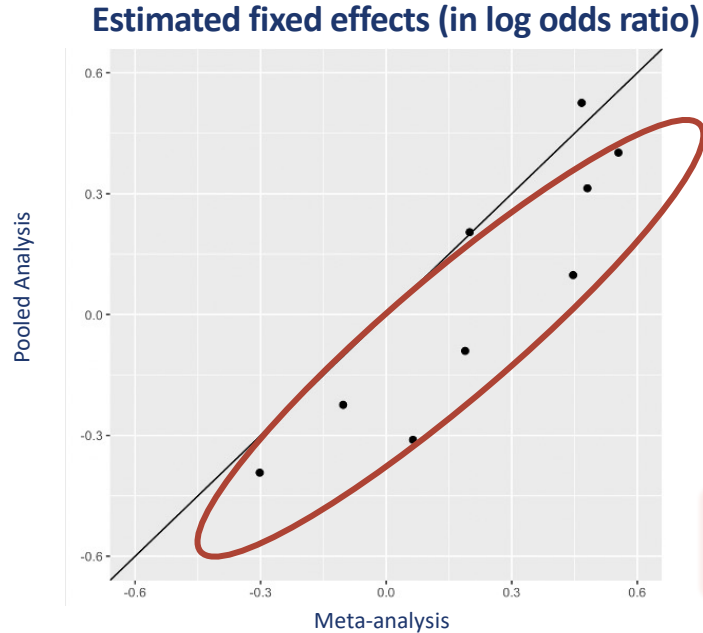
## Proposed Method – COLA-GLMM

### Collaborative One-shot Lossless Algorithm for Generalized Linear Mixed Model

- Suppose that a **common** set of covariates are available at all collaborating sites.
- The covariates have been **standardized into categorical variables**.
- Pipeline:



# Simulation Study – Meta-analysis vs Pooled analysis



## Simulation setting:

- 8 sites in total
- 9 risk factors
- Binary outcome
- Heterogeneous site-level random effects

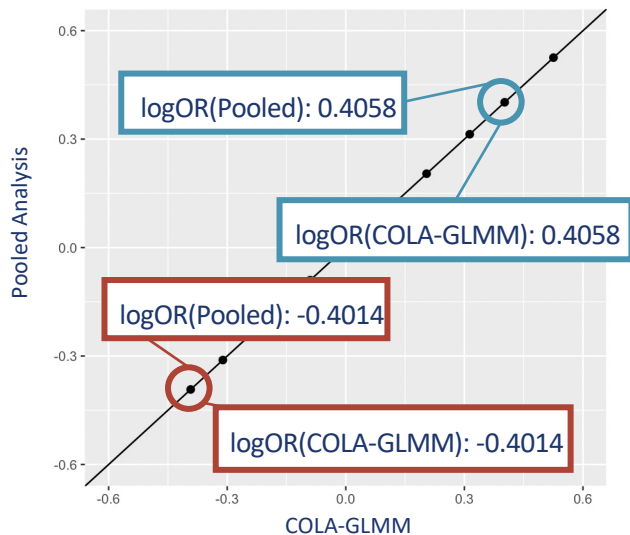
## Methods to compare:

- Pooled Analysis
- Meta-analysis

**Meta-analysis has low accuracy**

# Simulation Study – Compare Pooled Analysis and COLA-GLMM

## No cell suppression



U.S. Dept. of Health & Human Services  
**Guidance Portal**

[Return to Search](#)

## CMS Cell Suppression Policy

Guidance for CMS Cell Suppression Policy Web Page

Final

Issued by: Centers for Medicare & Medicaid Services (CMS)

Raw value	Report
0	0
1-10	<11
>=11	as it

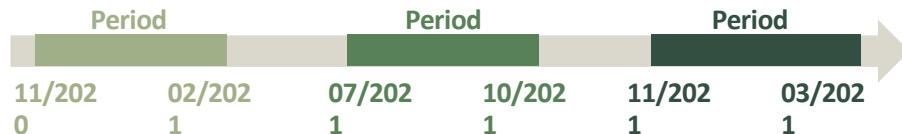


## Real-world Case Study

- **Scientific Question:**

Identify COVID-19 mortality **risk factors**  
over **three time periods** among hospitalized patients

- **Study Period:**



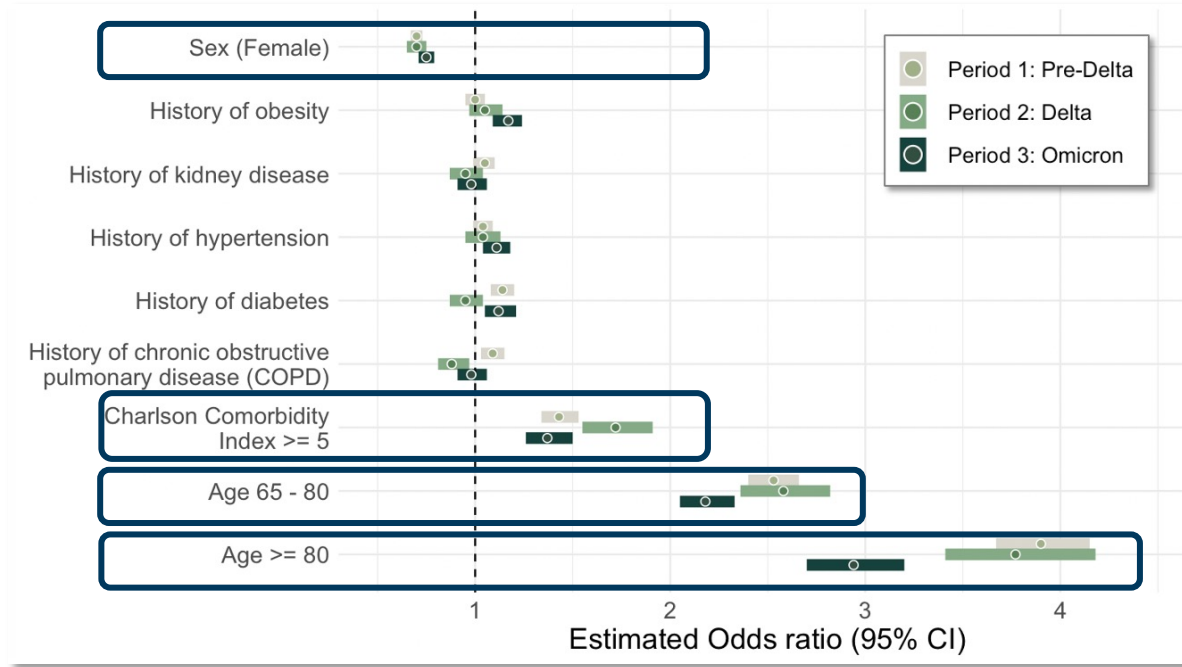
- **Databases (3 countries):**

- Optum® de-identified Electronic Health Record Dataset (Optum EHR);
- Optum's Clinformatics® Data Mart (CDM or Clinformatics®);
- IQVIA Hospital CDM;
- University of Florida Health;
- Department of Veterans Affairs;
- Integrated Primary Care Information (IPCI), The Netherlands;
- Columbia University Irving Medical Center (CUIMC);
- Parc Salut Mar Barcelona (PSMAR), Spain.

**Inclusion criteria:**

- Patients aged 18 years and older
- Had an inpatient visit with either a diagnosis of COVID-19 or a positive test for COVID-19 between 21 days prior to the inpatient visit and the end of the inpatient visit

## Real-world Case Study Results



- Sex (female):**

- Reference group: Male
- Female patients consistently exhibit a lower risk of mortality compared to males across all periods

- Charlson Comorbidity Index (CCI):**

- Reference group: CCI < 5
- Higher CCI scores are statistically associated with an increased risk of mortality.

- Age:**

- Reference group: Age < 65
- Higher age indicates significantly increased risk of mortality

## Summary – COLA-GLMM

Collaborative One-shot Lossless Algorithm for Generalized Linear Mixed Model

- **Lossless One-Shot**
- **Summary Statistics Only**
- **Heterogeneity-Aware**
- **Scalable, Applicable, and Implementation-Ready in OHDSI Network**



PDA R Package: 13300+ downloads since 2020



PDA Github Page: <https://github.com/Penncil/pda>



PDA website: <https://pdamethods.org/>



Penn security office certified

**PDA-OTA**

PDA-OTA: <https://pda-ota.pdamethods.org/>

# Acknowledgements



## Poster #117

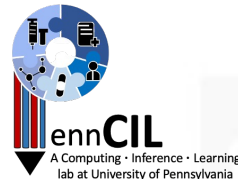
- Yong Chen, University of Pennsylvania
- David A. Asch, University of Pennsylvania
- Jenna Repts, Janssen Research and Development
- Chongliang Luo, Washington University in St. Louis
- Yiwen Lu, University of Pennsylvania
- Milou T. Brand, Real World Solutions, IQVIA
- Scott L. DuVall, VA Informatics and Computing Infrastructure
- Thomas Falconer, Columbia University
- Juan Manuel Ramirez-Anguita, Hospital del Mar Research Institute (HMRIB)
- Miguel A. Mayer, Hospital del Mar Research Institute (HMRIB)
- Michael E Matheny, VA Informatics and Computing Infrastructure
- Alex Mayer Fuentes, Parc Taulí Hospital Universitari
- Xing He, University of Florida
- Bhavnisha K Patel, VA Informatics and Computing Infrastructure
- Katherine R Simon, VA Informatics and Computing Infrastructure
- Marc A. Suchard, University of California, Los Angeles
- Guojun Tang, University of Calgary
- Benjamin Viernes, VA Informatics and Computing Infrastructure
- Fei Wang, Weill Cornell Medicine
- Ross D. Williams, Erasmus University Medical Center
- Mui van Zandt, Real World Solutions, IQVIA
- Jiang Bian, University of Florida
- Jiayu Zhou, Michigan State University

### Correspondence to:

- **Jessie Tong**, [jtong20@jhu.edu](mailto:jtong20@jhu.edu)
- **Yong Chen**, [ychen123@upenn.edu](mailto:ychen123@upenn.edu)



Penn Medicine



OHDSI  
OBSERVATIONAL HEALTH DATA SCIENCES AND INFORMATICS

Department of Biostatistics, Epidemiology and Informatics

# NCO-Calibrated DID Analysis: Addressing Unmeasured Confounding in Difference-in-Differences Analyses Using Negative Control Outcomes Experiments

Dazheng Zhang, Ph.D. candidate in Biostatistics at the University of Pennsylvania

2024 OHDSI Symposium

Advisor: Yong Chen, Ph.D., Professor of Biostatistics

Director of Center for Health AI & Synthesis of Evidence (CHASE), University of Pennsylvania

Joint work with Bingyu Zhang, Dr. Huiyuan Wang, Dr. Charles J. Wolock, Yiwen (Iris) Lu, Dr. Yong Chen





# Real-World Case Study

- **Racial/Ethnic disparities** long lasting in healthcare.



# Real-World Case Study

- **Racial/Ethnic disparities** long lasting in healthcare.
- Does the pandemic worsen racial/ethnic disparities?

**This Issue** Views **164,791** Citations **1,505** Altmetric **940**

**Viewpoint** FREE

May 11, 2020

**COVID-19 and Racial/Ethnic Disparities**

Monica Webb Hooper, PhD<sup>1</sup>; Anna María Nápoles, PhD, MPH<sup>1</sup>; Eliseo J. Pérez-Stable, MD<sup>1</sup>

[» Author Affiliations](#) | [Article Information](#)

JAMA. 2020;323(24):2466-2467. doi:10.1001/jama.2020.8598

 Editorial Comment

 Related Articles



# Real-World Case Study

- **Racial/Ethnic disparities** long lasting in healthcare.
- Does the pandemic worsen racial/ethnic disparities?
  - **Difference-in-difference (DiD) approach** finds racial/ethnic disparities attributable to the pandemic while controlling for pre-existing disparities.

**This Issue** Views **164,791** Citations **1,505** Altmetric **940**

**Viewpoint** FREE

May 11, 2020

## COVID-19 and Racial/Ethnic Disparities

Monica Webb Hooper, PhD<sup>1</sup>; Anna María Nápoles, PhD, MPH<sup>1</sup>; Eliseo J. Pérez-Stable, MD<sup>1</sup>

[» Author Affiliations](#) | [Article Information](#)

JAMA. 2020;323(24):2466-2467. doi:10.1001/jama.2020.8598

Editorial Comment

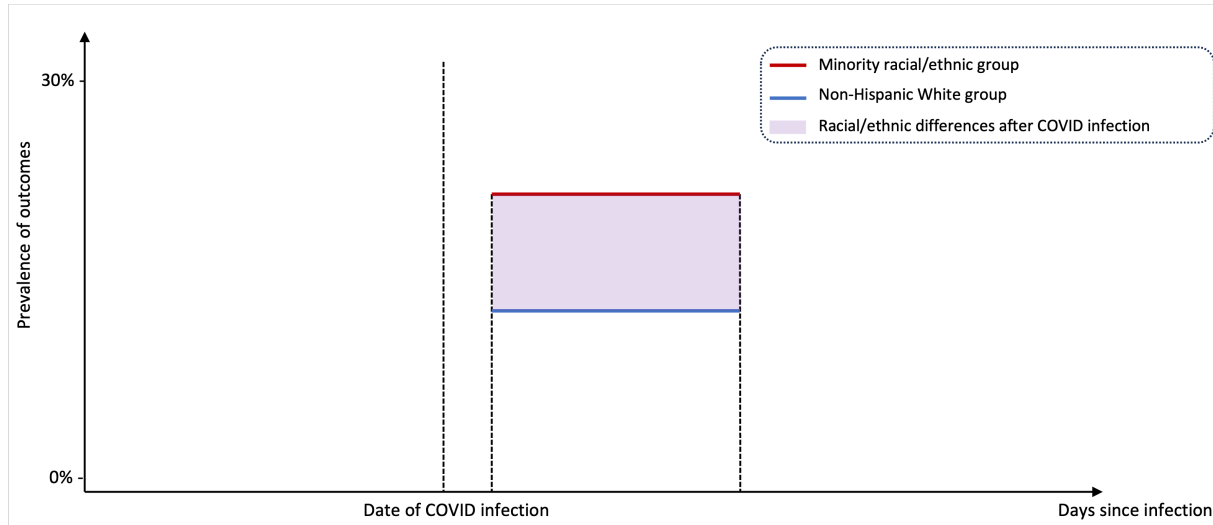
Related Articles





# Real-World Case Study

- **Racial/Ethnic disparities** long lasting in healthcare.
- Does the pandemic worsen racial/ethnic disparities?
  - **Difference-in-difference (DiD) approach** finds racial/ethnic disparities attributable to the pandemic while controlling for pre-existing disparities.



This Issue Views 164,791 Citations 1,505 Altmetric 940

Viewpoint

May 11, 2020

## COVID-19 and Racial/Ethnic Disparities

Monica Webb Hooper, PhD<sup>1</sup>; Anna María Nápoles, PhD, MPH<sup>1</sup>; Eliseo J. Pérez-Stable, MD<sup>1</sup>

[» Author Affiliations](#) | [Article Information](#)

JAMA. 2020;323(24):2466-2467. doi:10.1001/jama.2020.8598

Editorial  
Comment

Related  
Articles

FREE



# Real-World Case Study

This Issue Views 164,791 Citations 1,505 Altmetric 940

Viewpoint

May 11, 2020

## COVID-19 and Racial/Ethnic Disparities

Monica Webb Hooper, PhD<sup>1</sup>; Anna María Nápoles, PhD, MPH<sup>1</sup>; Eliseo J. Pérez-Stable, MD<sup>1</sup>

[» Author Affiliations](#) | [Article Information](#)

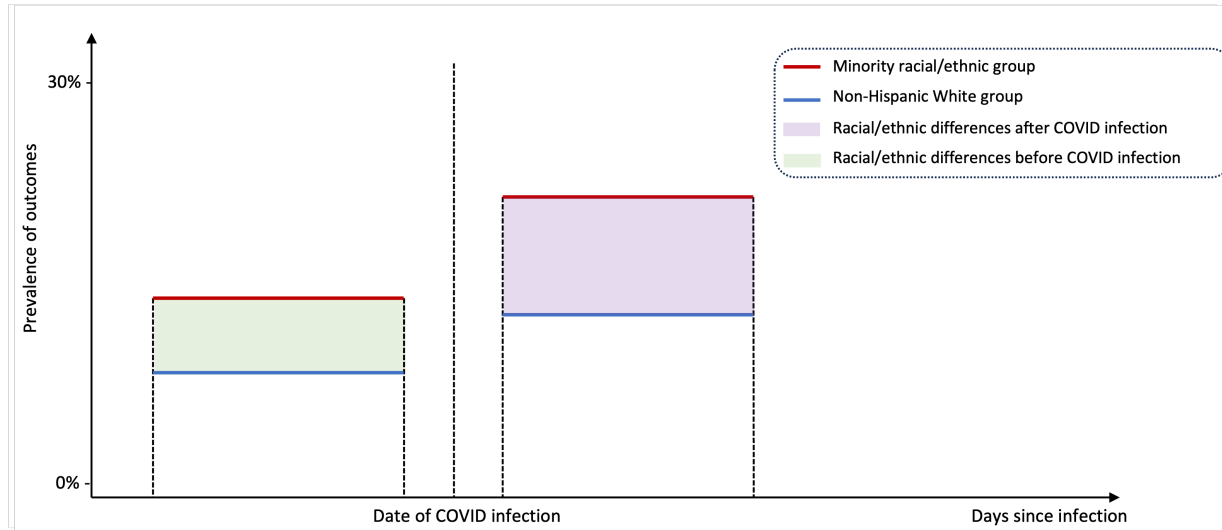
JAMA. 2020;323(24):2466-2467. doi:10.1001/jama.2020.8598

Editorial  
Comment

Related  
Articles

FREE

- **Racial/Ethnic disparities** long lasting in healthcare.
- Does the pandemic worsen racial/ethnic disparities?
  - **Difference-in-difference (DiD) approach** finds racial/ethnic disparities attributable to the pandemic while controlling for pre-existing disparities.





# Real-World Case Study

This Issue Views 164,791 Citations 1,505 Altmetric 940

Viewpoint

May 11, 2020

## COVID-19 and Racial/Ethnic Disparities

Monica Webb Hooper, PhD<sup>1</sup>; Anna María Nápoles, PhD, MPH<sup>1</sup>; Eliseo J. Pérez-Stable, MD<sup>1</sup>

[» Author Affiliations](#) | [Article Information](#)

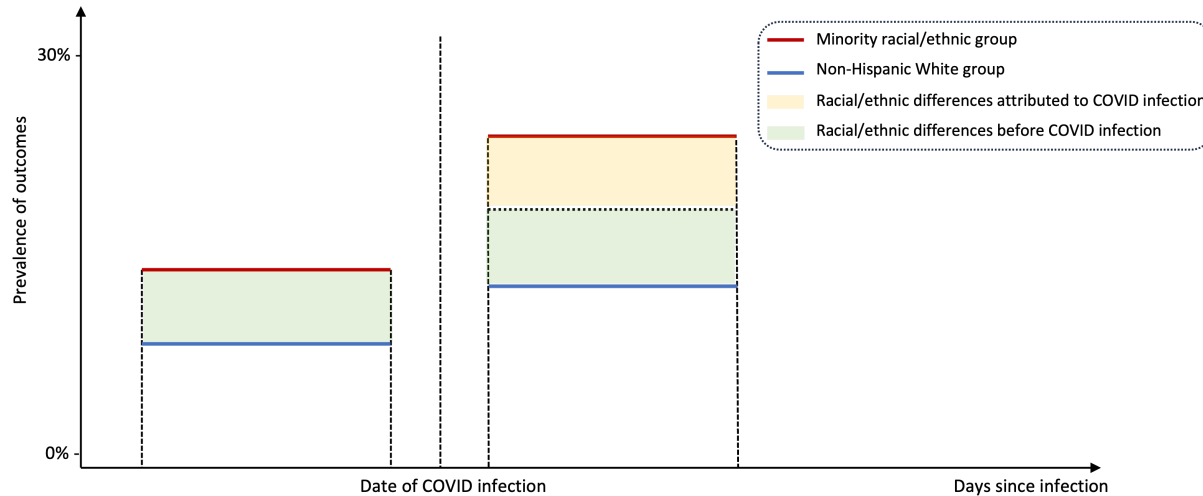
JAMA. 2020;323(24):2466-2467. doi:10.1001/jama.2020.8598

Editorial  
Comment

Related  
Articles

FREE

- **Racial/Ethnic disparities** long lasting in healthcare.
- Does the pandemic worsen racial/ethnic disparities?
  - **Difference-in-difference (DiD) approach** finds racial/ethnic disparities attributable to the pandemic while controlling for pre-existing disparities.





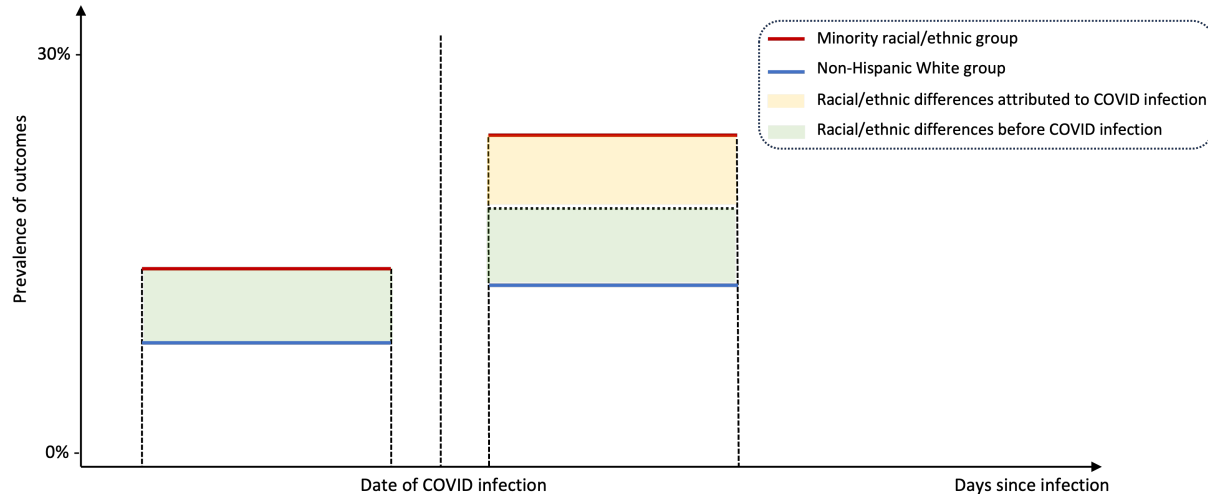
# Parallel Trends Assumption in DiD Model

- **Parallel trends assumption:** in the absence of intervention, the before-intervention difference and the after-intervention should be the same.



# Parallel Trends Assumption in DiD Model

- **Parallel trends assumption:** in the absence of intervention, the before-intervention difference and the after-intervention should be the same.





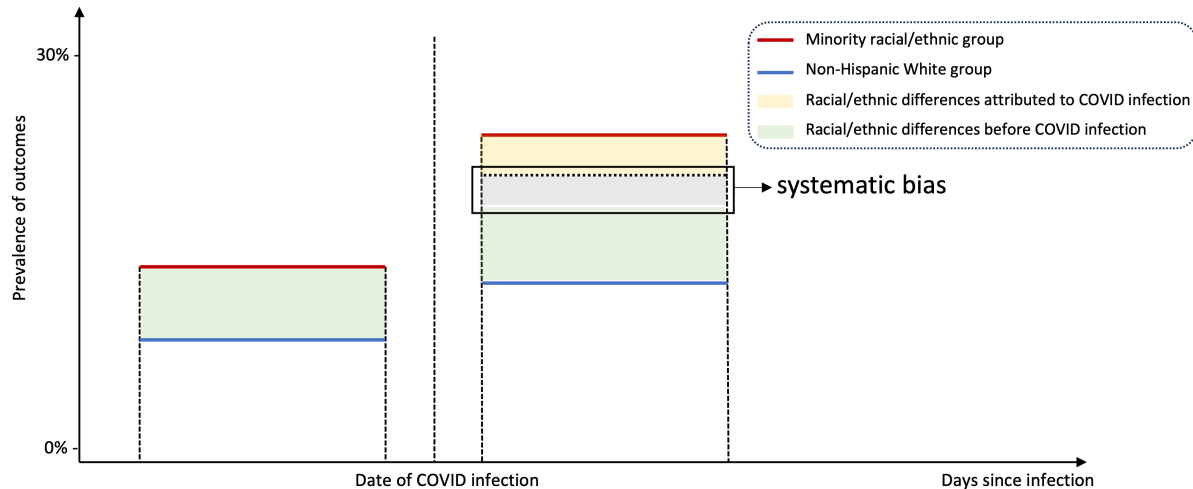
# Parallel Trends Assumption in DiD Model

- **Parallel trends assumption:** in the absence of intervention, the before-intervention difference and the after-intervention should be the same.
  - **Violation of parallel trends assumption:** systematic bias from unmeasured confounding variables makes the non-parallel trends for two groups.
-



# Parallel Trends Assumption in DiD Model

- **Parallel trends assumption:** in the absence of intervention, the before-intervention difference and the after-intervention should be the same.
- **Violation of parallel trends assumption:** systematic bias from unmeasured confounding variables makes the non-parallel trends for two groups.





# Negative Control Experiments Empirical Calibration

Negative control outcome (NCO), known **in priori** to be unrelated to exposure.

## Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data

Martijn J. Schuemie<sup>a,b,1</sup>, George Hripcsak<sup>a,c,d</sup>, Patrick B. Ryan<sup>a,b,e</sup>, David Madigan<sup>a,e</sup>, and Marc A. Suchard<sup>a,f,g,h</sup>

<sup>a</sup>Observational Health Data Sciences and Informatics, New York, NY 10032; <sup>b</sup>Epidemiology Analytics, Janssen Research & Development, Titusville, NJ 08560; <sup>c</sup>Department of Biomedical Informatics, Columbia University, New York, NY 10032; <sup>d</sup>Medical Informatics Services, New York-Presbyterian Hospital, New York, NY 10032; <sup>e</sup>Department of Statistics, Columbia University, New York, NY 10027; <sup>f</sup>Department of Biomathematics, University of California, Los Angeles, CA 90095; <sup>g</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; and <sup>h</sup>Department of Microbiology, University of California, Los Angeles, CA 90095.





# Negative Control Experiments Empirical Calibration

Negative control outcome (NCO), known **in priori** to be unrelated to exposure.

## Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data

Martijn J. Schuemie<sup>a,b,1</sup>, George Hripcsak<sup>a,c,d</sup>, Patrick B. Ryan<sup>a,b,c</sup>, David Madigan<sup>a,e</sup>, and Marc A. Suchard<sup>a,f,g,h</sup>

<sup>a</sup>Observational Health Data Sciences and Informatics, New York, NY 10032; <sup>b</sup>Epidemiology Analytics, Janssen Research & Development, Titusville, NJ 08560; <sup>c</sup>Department of Biomedical Informatics, Columbia University, New York, NY 10032; <sup>d</sup>Medical Informatics Services, New York-Presbyterian Hospital, New York, NY 10032; <sup>e</sup>Department of Statistics, Columbia University, New York, NY 10027; <sup>f</sup>Department of Biomathematics, University of California, Los Angeles, CA 90095; <sup>g</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; and <sup>h</sup>Department of Microbiology, University of California, Los Angeles, CA 90095

Legend study (Suchard et al. 2021 Lancet) used “ingrown nail” as an adverse event that is known to be unrelated to the antihypertension

### JOURNAL ARTICLE

## Large-scale evidence generation and evaluation across a network of databases (LEGEND): assessing validity using hypertension as a case study

Martijn J Schuemie , Patrick B Ryan, Nicole Pratt, RuiJun Chen, Seng Chan You, Harlan M Krumholz, David Madigan, George Hripcsak, Marc A Suchard

*Journal of the American Medical Informatics Association* Volume 27 Issue 8 August 2020



# Negative Control Experiments Empirical Calibration

Negative control outcome (NCO), known **in priori** to be unrelated to exposure.

## Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data

Martijn J. Schuemie<sup>a,b,1</sup>, George Hripcsak<sup>a,c,d</sup>, Patrick B. Ryan<sup>a,b,c</sup>, David Madigan<sup>a,e</sup>, and Marc A. Suchard<sup>a,f,g,h</sup>

<sup>a</sup>Observational Health Data Sciences and Informatics, New York, NY 10032; <sup>b</sup>Epidemiology Analytics, Janssen Research & Development, Titusville, NJ 08560; <sup>c</sup>Department of Biomedical Informatics, Columbia University, New York, NY 10032; <sup>d</sup>Medical Informatics Services, New York-Presbyterian Hospital, New York, NY 10032; <sup>e</sup>Department of Statistics, Columbia University, New York, NY 10027; <sup>f</sup>Department of Biomathematics, University of California, Los Angeles, CA 90095; <sup>g</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; and <sup>h</sup>Department of Medicine, University of California, Los Angeles, CA 90095

Legend study (Suchard et al. 2021 Lancet) used “ingrown nail” as an adverse event that is known to be unrelated to the antihypertension

### JOURNAL ARTICLE

## Large-scale evidence generation and evaluation across a network of databases (LEGEND): assessing validity using hypertension as a case study

Martijn J Schuemie ✉, Patrick B Ryan, Nicole Pratt, RuiJun Chen, Seng Chan You, Harlan M Krumholz, David Madigan, George Hripcsak, Marc A Suchard

Journal of the American Medical Informatics Association Volume 27 Issue 8 August 2020

An example list of 40 NCOs for a vaccine study.

### Related Features



Original Research | 9 January 2024

## Real-World Effectiveness of BNT162b2 Against Infection and Severe Diseases in Children and Adolescents

Authors: Qiong Wu, PhD, Jiayi Tong, MS, Bingyu Zhang, MS, Dazheng Zhang, MS, Jiajie Chen, PhD, Yuqing Lei, MS, Yiwen Lu, BS, ... [SHOW ALL](#) ... and Yong Chen, PhD | [AUTHOR, ARTICLE, & DISCLOSURE INFORMATION](#)

Publication: Annals of Internal Medicine • Volume 177, Number 2 • <https://doi.org/10.7326/M23-1754>

Categories	Examples
Infectious and parasitic diseases	Impetigo, Tinea capitis, Tinea corporis, Insect bite
Diseases of the skin tissue	Contact dermatitis, Diaper rash, Acne
Diseases of the musculoskeletal system and connective tissue	Dislocations, Closed fracture of distal end of radius, Sprain of ankle, Scoliosis, Foot pain, Injury of free lower limb, Injury of upper extremity, Injury of right leg, Injury of left leg, Injury of right foot
Diseases of the nervous system	Seizure, Epilepsy, Concussion, Closed injury of head



# Negative Control Experiments Empirical Calibration

Negative control outcome (NCO), known **priori** to be unrelated to exposure.

## Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data

Martijn J. Schuemie<sup>1,2</sup>, George Hripcsak<sup>3,4,5</sup>, Patrick B. Ryan<sup>1,2,3</sup>, David Madigan<sup>6,7</sup>, and Marc A. Suchard<sup>8,9,10</sup>

<sup>1</sup>Observational Health Data Sciences and Informatics, New York, NY 10032; <sup>2</sup>Epidemiology Analytics, Janssen Research & Development, Titusville, NJ 08560; <sup>3</sup>Department of Biomedical Informatics, Columbia University, New York, NY 10032; <sup>4</sup>Medical Informatics Services, New York-Presbyterian Hospital, New York, NY 10032; <sup>5</sup>Department of Statistics, Columbia University, New York, NY 10027; <sup>6</sup>Department of Biomathematics, University of California, Los Angeles, CA 90095; <sup>7</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; <sup>8</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; <sup>9</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095; <sup>10</sup>Department of Biostatistics, University of California, Los Angeles, CA 90095

An example list of 40 NCOs for a vaccine study.

## Related Features

[VISUAL ABSTRACT](#)

Original Research | 9 January 2024

## Real-World Effectiveness of BNT162b2 Against Infection and Severe Diseases in Children and Adolescents

Authors: Qiong Wu, PhD , Jialy Tong, MS , Bingyu Zhang, MS , Dazheng Zhang, MS , Jialie Chen, PhD , Yuying Lei, MS , Yiyen Lu, BS , ... [SHOW ALL](#) ... and Yong Chen, PhD  [AUTHOR ARTICLE & DISCLOSURE INFORMATION](#)

Publication: Annals of Internal Medicine • Volume 177, Number 2 • <https://doi.org/10.7326/M23-1754>

# Can we make empirical calibration for DiD model?

“ingrown nail” as an adverse event that is known to be unrelated to the antihypertension

## JOURNAL ARTICLE

## Large-scale evidence generation and evaluation across a network of databases (LEGEND): assessing validity using hypertension as a case study

Martijn J Schuemie , Patrick B Ryan, Nicole Pratt, RuiJun Chen, Seng Chan You, Harlan M Krumholz, David Madigan, George Hripcsak, Marc A Suchard

Journal of the American Medical Informatics Association Volume 27 Issue 8 August 2020

Categories	Examples
Parasitic diseases	Impetigo, Tinea capitis, Tinea corporis, Insect bite
Diseases of the skin tissue	Contact dermatitis, Diaper rash, Acne
Diseases of the musculoskeletal system and connective tissue	Dislocations, Closed fracture of distal end of radius, Sprain of ankle, Scoliosis, Foot pain, Injury of free lower limb, Injury of upper extremity, Injury of right leg, Injury of left leg, Injury of right foot
Diseases of the nervous system	Seizure, Epilepsy, Concussion, Closed injury of head



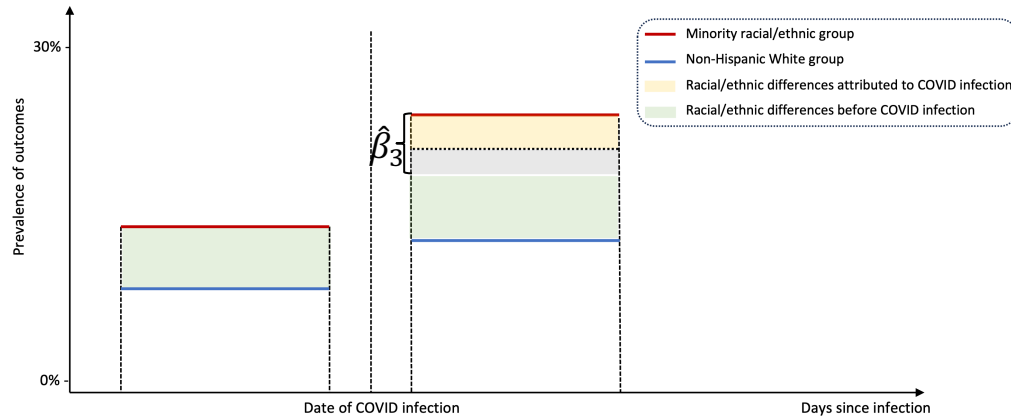
## Proposed Method: NCO-DiD





# Proposed Method: NCO-DiD

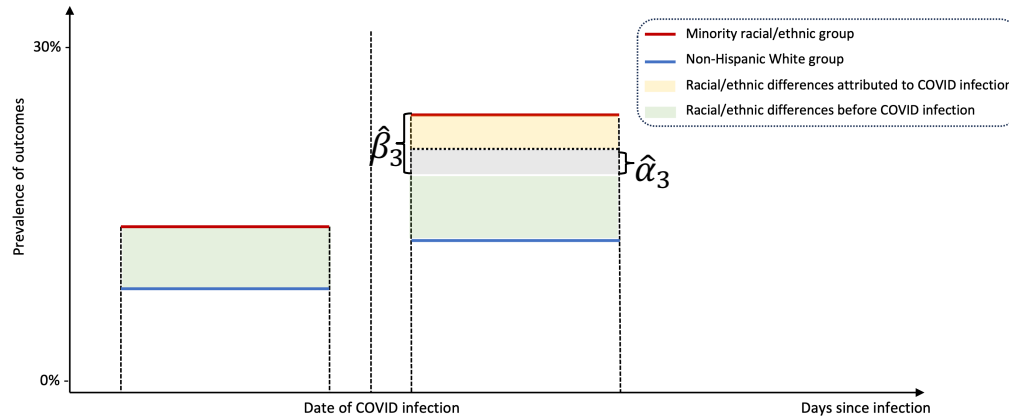
- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$





# Proposed Method: NCO-DiD

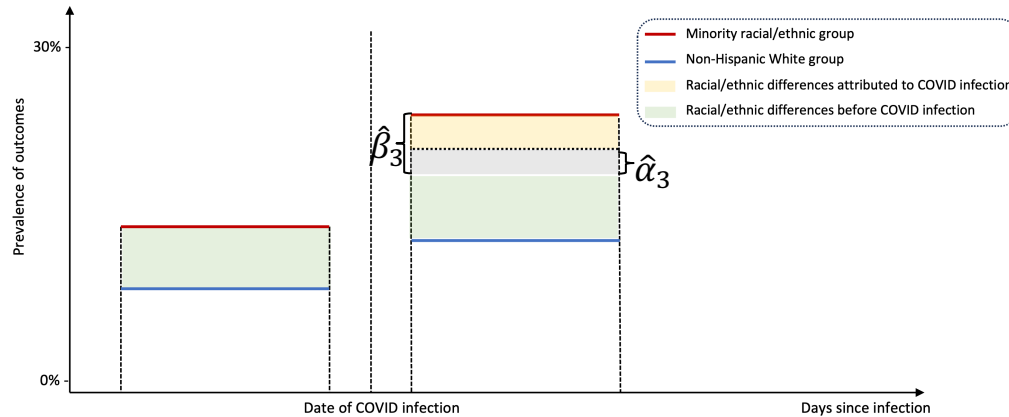
- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
NCO  $\sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$





# Proposed Method: NCO-DiD

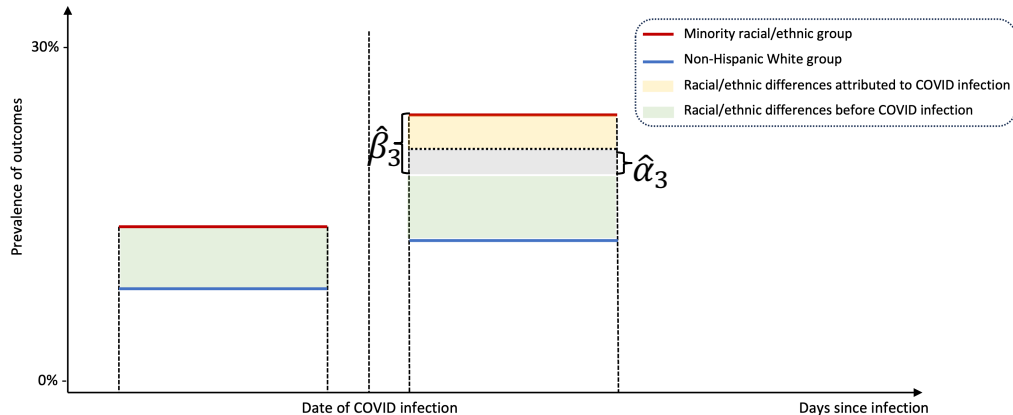
- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
NCO  $\sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$





# Proposed Method: NCO-DiD

- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
NCO  $\sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$



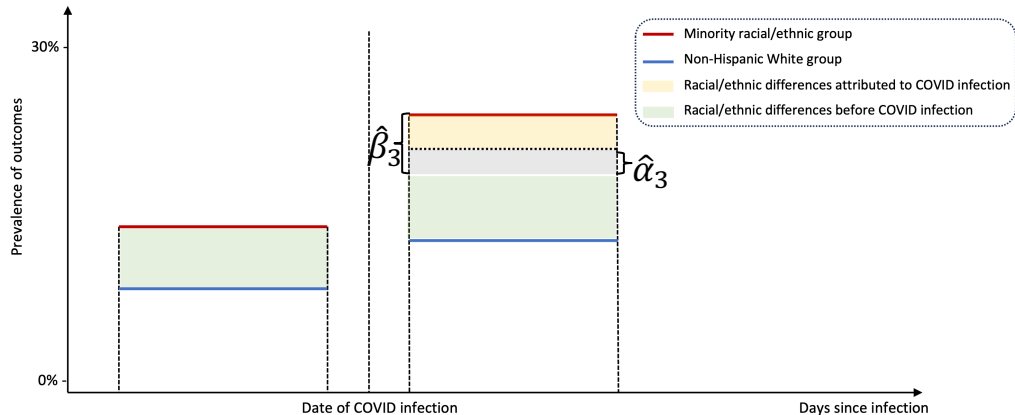
Estimation for outcome of interest  
 $\hat{\beta}_3 = \text{systematic bias} + \beta_3$





# Proposed Method: NCO-DiD

- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
NCO  $\sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$



Estimation for outcome of interest

$$\hat{\beta}_3 = \text{systematic bias} + \beta_3$$

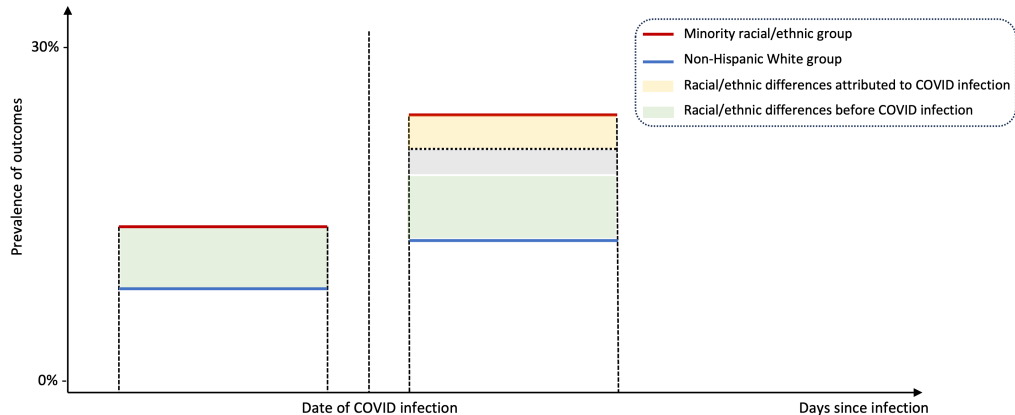
Estimation for NCO:

$$\hat{\alpha}_3 = \text{systematic bias} + \alpha_3$$



# Proposed Method: NCO-DiD

- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
Outcome  $\sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
NCO  $\sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$



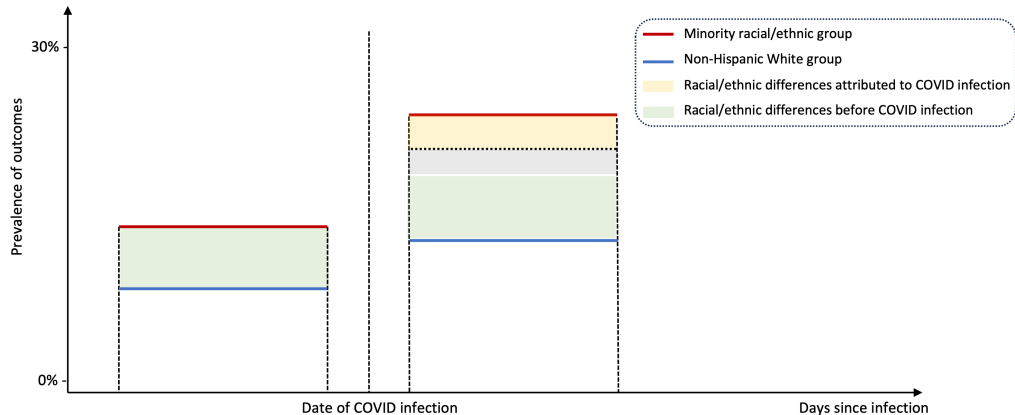
Estimation for outcome of interest  
 $\hat{\beta}_3 = \text{systematic bias} + \beta_3$

Estimation for NCO:  
 $\hat{\alpha}_3 = \text{systematic bias} + 0$



# Proposed Method: NCO-DiD

- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
$$\text{Outcome} \sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$$
- Step 2: Estimate systematic bias from an NCO ( $\hat{\alpha}_3$ ; true  $\alpha_3 = 0$ ).  
$$\text{NCO} \sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$



Estimation for outcome of interest

$$\hat{\beta}_3 = \text{systematic bias} + \beta_3$$

$$\hat{\alpha}_3 \approx \text{systematic bias}$$

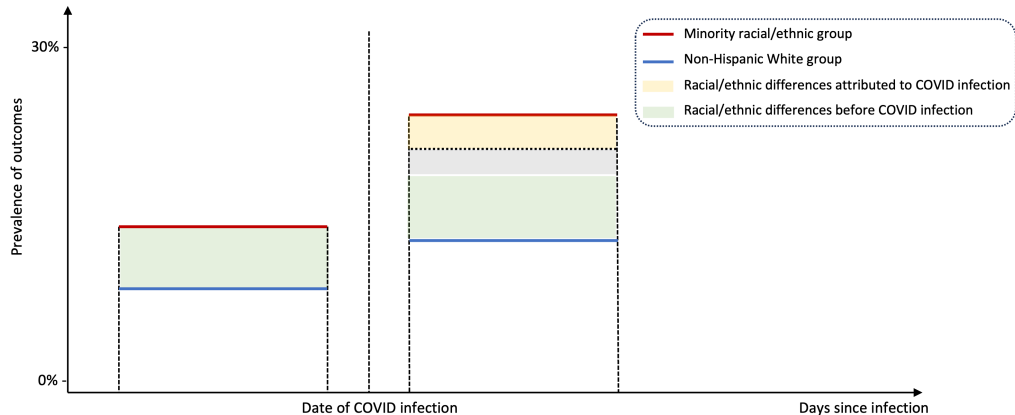
Estimation for NCO:

$$\hat{\alpha}_3 = \text{systematic bias} + 0$$



# Proposed Method: NCO-DiD

- Step 1: Fit a DiD model on the cohort (matched on measured confounders) to find estimate ( $\hat{\beta}_3$ )  
$$\text{Outcome} \sim \beta_0 + \beta_1 \text{Time} + \beta_2 \text{Intervention} + \beta_3 \text{Time} \times \text{Intervention}$$
- Step 2: Estimate systematic bias from **M NCOs**.  
$$\text{NCO} \sim \alpha_0 + \alpha_1 \text{Time} + \alpha_2 \text{Intervention} + \alpha_3 \text{Time} \times \text{Intervention}$$
- Step 3: Empirically calibrate  $\hat{\beta}_3 - \hat{\alpha}_3$



Estimation for outcome of interest

$$\hat{\beta}_3 = \text{systematic bias} + \beta_3$$

$$\hat{\alpha}_3 \approx \text{systematic bias}$$

Estimation for NCO:

$$\hat{\alpha}_3 = \text{systematic bias} + \mathbf{0}$$



# Baseline Method vs Proposed Method

## Baseline method (DiD)

### Matching

- ☐ Learn propensity score.
- ☐ Construct matched cohort.



### Inference for intervention effect

- ☐ Fit DiD model.

## Proposed method (NCO-DiD)

### Matching

- ☐ Learn propensity score.
- ☐ Construct matched cohort.



### Inference for intervention effect

- ☐ Fit DiD model

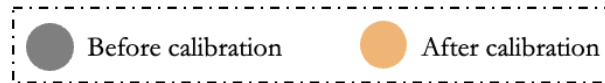
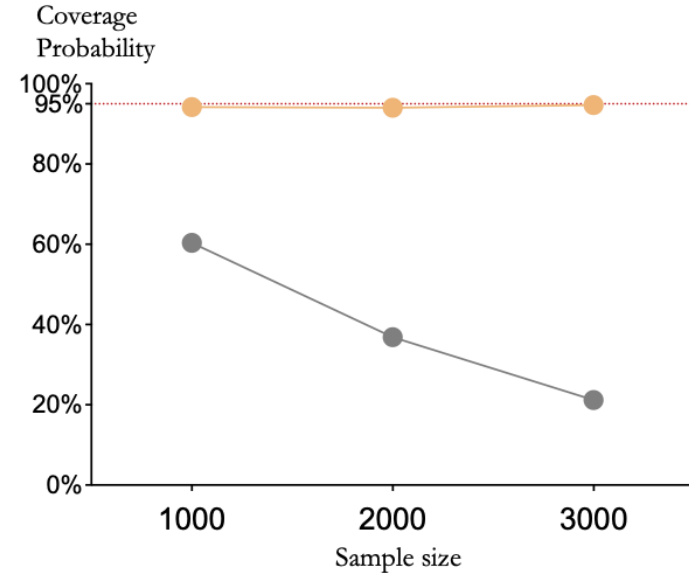
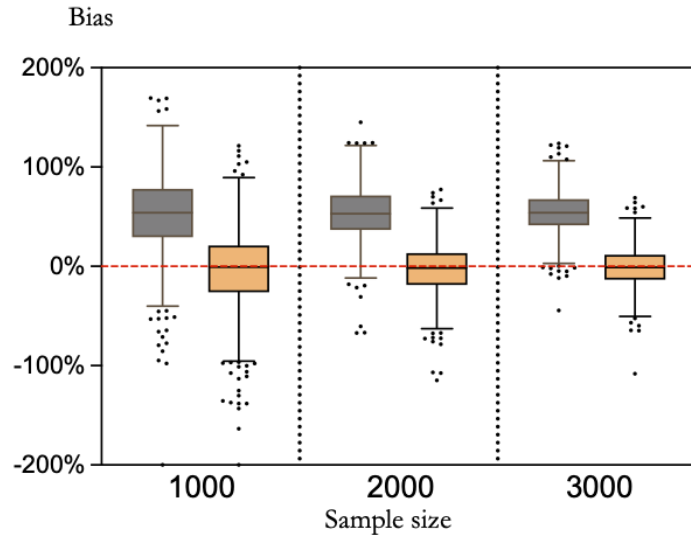


### Empirical Calibration

- ☐ Use NCOs to estimate systematic bias.
- ☐ Calibrate systematic bias.



# Results from NCO-DiD



- Key message: across all scenarios, the proposed NCO-DiD model successfully calibrates the systematic bias.



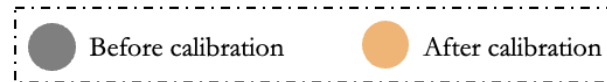
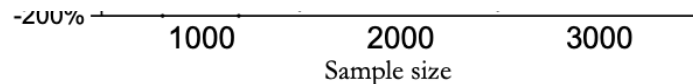
# Results from NCO-DiD

Coverage  
Bias

## Racial/Ethnic Differences in Long-COVID-Associated Symptoms among Pediatrics Population: Findings from Difference-in-differences Analyses in RECOVER Program

Author list:

Dazheng Zhang, MS<sup>1^</sup>, Bingyu Zhang, MS<sup>1,2^</sup>, Qiong Wu, PhD<sup>1</sup>, Ting Zhou, PhD<sup>1</sup>, Jiayi Tong, MS<sup>1</sup>, Yiwen Lu<sup>1,2</sup>, Jiajie Chen, PhD<sup>1</sup>, Huiyuan Wang, PhD<sup>1</sup>, Deena J Chisolm, PhD<sup>3</sup>, Ravi Jhaveri, MD<sup>4</sup>, Rachel C Kenney, PhD<sup>5,6</sup>, Russell L Rothman, MD, MPP<sup>7</sup>, Suchitra Rao, MD<sup>8</sup>, David A Williams, MD<sup>9</sup>, Mady Hornig, MA, MD<sup>10</sup>, Linbo Wang, PhD<sup>11</sup>, Jeffrey S Morris, PhD<sup>12</sup>, Christopher B Forrest, MD, PhD<sup>13</sup>, and Yong Chen, PhD<sup>1,2</sup> on behalf of the RECOVER consortium.



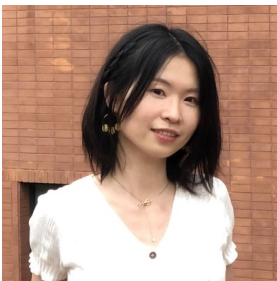
- Key message: across all scenarios, the proposed NCO-DiD model successfully calibrates the systematic bias.



# Takeaway Messages

- Extend OHDSI empirical calibration framework.
- Control systematic bias from unmeasured confounding variables for DiD model.
- Help understand racial/ethnic differences in long COVID conditions after COVID-19 infection among children and adolescents.





**Bingyu Zhang**



**Huiyuan Wang, Ph.D.**



**Charles J. Wolock, Ph.D.**



**Yiwen Lu**



**Yong Chen, Ph.D.**

**Acknowledgments:**

Jiayi (Jessie) Tong, Ph.D.

Qiong Wu, Ph.D.,

Jiajie Chen, Ph.D.,

Lu Li,

Yuqing Lei,

and other lab members for your  
help with the project.



**Poster: # 118**



**Oct 23, 3:00 pm - 5:00 pm**



Code



Google scholar

**Correspondence to:**

**Dazheng Zhang**

Email: [dazheng.zhang@penndicine.upenn.edu](mailto:dazheng.zhang@penndicine.upenn.edu)

**Yong Chen, Ph.D.**

Email: [ychen123@upenn.edu](mailto:ychen123@upenn.edu)

# Health Trends Across Communities in Minnesota (HTAC-MN):

a Statewide Dashboard Leveraging the OMOP CDM to Monitor the Prevalence of Health Conditions

2024 OHDSI Symposium  
October 23, 2024

Sam Patnoe, HealthPartners Institute  
on behalf of the Minnesota EHR Consortium

HTAC-MN is funded through a Minnesota Public Health Infrastructure Grant from the Minnesota Department of Health.





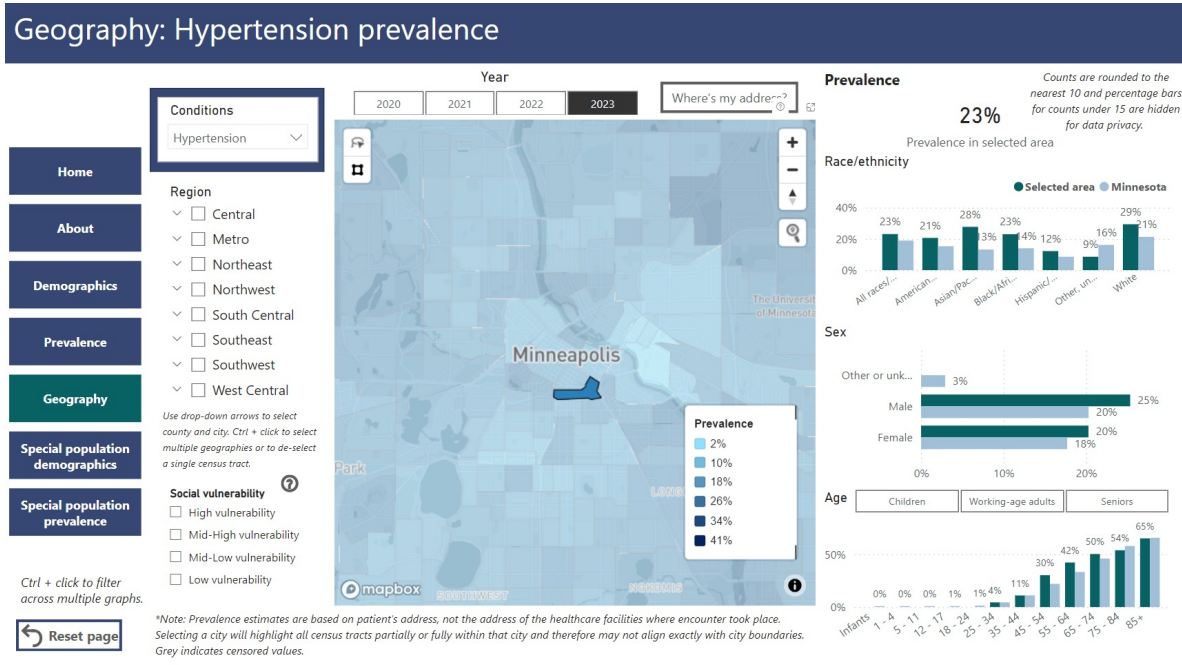
# MN EHR Consortium

- Formed in March 2020 to address gaps in COVID-19 data sharing and communication
- 11 largest health systems in Minnesota
- >90% of residents in MN, all regions of state
- Federated data model - OMOP CDM (v5.3) adoption across all health systems from 2022-2023
- Health Trends Across Communities in Minnesota (HTAC-MN) began in 2023
  - Goal: Build a comprehensive statewide dashboard to support public health surveillance, inform community health assessments, and promote health equity



# The HTAC-MN Dashboard

- Prevalence estimates for 30+ health conditions
  - 5.4M+ people (>90% of MN population)
- Race/ethnicity, age, sex
- State → census tract
- Data for 2020-2023
- Timely – refreshed annually



[mnehrconsortium.org/htac](https://mnehrconsortium.org/htac)

# Selecting and defining health conditions

## 1. Identified and prioritized health conditions for dashboard:

- Public health significance
- Potential for action
- Lack of/limitations of existing data
- Emerging conditions
- Alignment with current public health priorities
- Detailed EHR data could support assessment work



## 2. Selected health conditions:

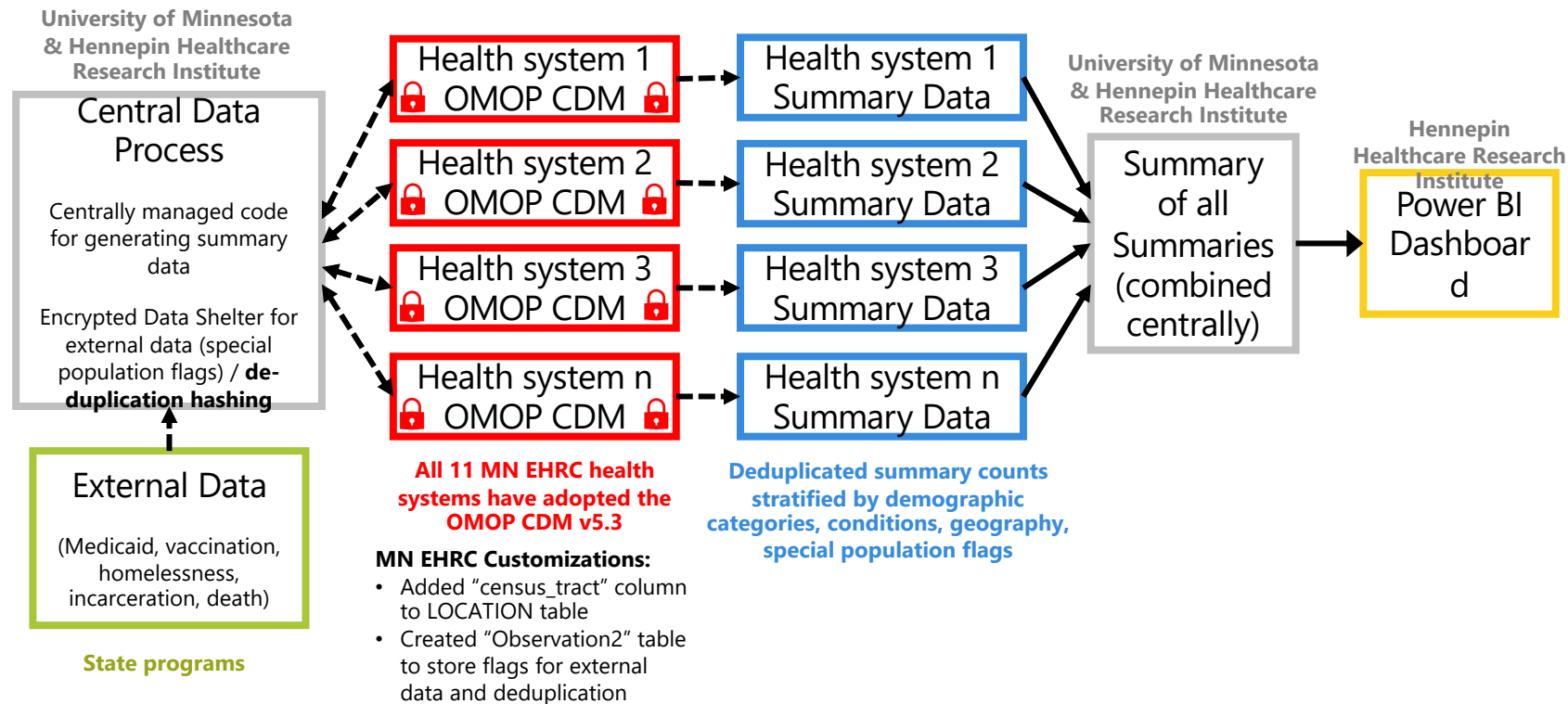
Chronic Conditions	Mental Health
<ul style="list-style-type: none"><li>• Asthma</li><li>• COPD</li><li>• Chronic kidney disease</li><li>• Diabetes, Type 2</li><li>• Heart failure</li><li>• Hyperlipidemia</li><li>• Hypertension</li><li>• Ischemic heart disease</li><li>• Obesity</li><li>• Peripheral vascular disease</li></ul>	<ul style="list-style-type: none"><li>• Anxiety</li><li>• Bipolar disorder</li><li>• Depression</li><li>• PTSD</li><li>• Psychotic disorders</li><li>• Suicidal ideation or recent attempt</li></ul>
Substance Use	Maternal & Child Health
<ul style="list-style-type: none"><li>• Alcohol</li><li>• Cannabis</li><li>• Cocaine</li><li>• Hallucinogens</li><li>• Inhalants</li><li>• Opioids</li><li>• Psychostimulants</li><li>• Sedatives</li></ul>	<ul style="list-style-type: none"><li>• Obstetrical deliveries</li><li>• Severe maternal morbidity</li><li>• Maternal opioid use</li></ul>
	Other
	<ul style="list-style-type: none"><li>• Acute myocardial infarction</li><li>• Firearm injury</li><li>• Lung cancer</li><li>• Stroke</li></ul>



## 3. Developed standardized OMOP concept sets for each condition:

- Mapped existing ICD-10-CM diagnostic code sets to OMOP concepts (incl. SNOMED & ICD concepts)
- Accounts for metadata across MN EHR Consortium health systems
- Reviewed by clinicians
- Centrally managed
- Each condition is defined uniformly across all health systems

# HTAC-MN Data Infrastructure



# Data flow – health system perspective

University of Minnesota  
& Hennepin Healthcare  
Research Institute

## Central Data Process

Centrally managed code for generating summary data

Encrypted Data Shelter for external data (special population flags) / **de-duplication hashing**

## External Data

(Medicaid, vaccination, homelessness, incarceration, death)

State programs

4. **R code** for generating summary data (incl. concept set table for conditions, Observation2 exclusions for de-duplication) is sent to all systems

Health system n  
OMOP CDM

Deduplicated summary counts stratified by demographic categories, conditions, geography, special population flags

Health system n  
Summary Data

5. System submits summary data file to be combined centrally

1. Each system sends de-identified patient-level file with de-duplication information for their patient population to Encrypted Data Shelter using one-way encryption

2. In Encrypted Data Shelter, patient-level files are deduplicated and linked to external data using privacy-preserving hashing functions

3. Matched, encrypted, de-duplicated, person-level file is returned to each health system and used for populating an "Observation2" table w/ de-duplication & external data flags

# Conclusion

- The HTAC-MN Dashboard demonstrates how the OMOP CDM can facilitate sharing summary EHR data across an entire state to monitor community health
- For more information, visit me at poster 119 and check out the website:

## HEALTH TRENDS ACROSS COMMUNITIES IN MINNESOTA DASHBOARD

Knowledge is power when working to improve the health of our communities. Better information can lead to a better understanding of community health needs and more effective

[HTAC Project Summary](#) | [Health Condition Descriptions](#) | [Glossary of Technical Terms](#) | [FAQs](#)

Please contact [MNEHTAC@mninstitute.org](mailto:MNEHTAC@mninstitute.org) with questions.

[HTAC Project Summary](#) | [Health Condition Descriptions](#) | [Glossary of Technical Terms](#) | [FAQs](#)

[View Additional Resources >](#)



[mnehrconsortium.org/htac](https://mnehrconsortium.org/htac)





# How Often: Characterizing Heterogeneity in Drug-Outcome Incidence Rate Estimates Attributed to Drug Indication

Results from the 2023 OHDSI Global Symposium

Hsin Yi Chen, BS, Christopher Knoll, BS, Elise Ruan, MD, MPH, Adam Black, BA,  
Sarah Seager, BA, Patrick Ryan, PhD, George Hripcsak, MD, MS



# Incidence Rates

- Incidence rate calculations is one of the most common analyses in pharmacoepidemiology
  - Comparative background estimation
  - Drug Adverse Events

$$\text{Incidence Rate} = \frac{\begin{array}{c} \text{\# persons in the target cohort who have} \\ \text{new outcome occurrence during the time-at-risk} \end{array}}{\begin{array}{c} \text{person-time-at-risk for persons in the target cohort} \\ \text{with time at risk} \end{array}}$$



# Why Incidence Rates?

- "Simple" (fewer assumptions)
- Not causal, but still useful
  - If incidence is low and side effect is not serious, then we're good
  - If incidence is high, then need to look out for it even if not caused by drug



# ...but incidence rate calculations can still be complex

- Heterogeneity in incidence rates estimates can be influenced by factors such as age, sex, calendar time, and differences in databases
  - Magnitude of potential impact... up to 1000 fold

## Factors Influencing Background Incidence Rate Calculation: Systematic Empirical Evaluation Across an International Network of Observational Databases

Anna Ostropolets<sup>1†</sup>, Xintong Li<sup>2†</sup>, Rupa Makadia<sup>3</sup>, Gowtham Rao<sup>3</sup>, Peter R. Rijnbeek<sup>4</sup>, Talita Duarte-Salles<sup>5</sup>, Anthony G. Sena<sup>3,4</sup>, Azza Shaabi<sup>5</sup>, Marc A. Suchard<sup>6,7</sup>, Patrick B. Ryan<sup>1,3</sup>, Daniel Prieto-Alhambra<sup>2</sup> and George Hripcsak<sup>1,8\*</sup>

### OPEN ACCESS

#### Edited by:

Elisabetta Poluzzi,  
University of Bologna, Italy

#### Reviewed by:

Michèle Fusaroli,  
University of Bologna, Italy  
Angela Accosta,  
ICESI University, Colombia  
Raquel Herrera Comoglio,  
National University of Córdoba,  
Argentina

#### \*Correspondence:

George Hripcsak  
gh13@cumc.columbia.edu

<sup>†</sup>These authors have contributed  
equally to this work and share first  
authorship

#### Specialty section:

This article was submitted to  
Pharmacoepidemiology,  
a section of the journal  
Frontiers in Pharmacology

Received: 12 November 2021

Accepted: 17 March 2022

Published: 26 April 2022

#### Citation:

Ostropolets A, Li X, Makadia R, Rao G,  
Rijnbeek PR, Duarte-Salles T,  
Sena AG, Shaabi A, Suchard MA,  
Ryan PB, Prieto-Alhambra D and  
Hripcsak G (2022) Factors Influencing Background Incidence Rate Calculation: Systematic Empirical Evaluation Across an International Network of Observational Databases. *Front. Pharmacol.* 13:811111. doi: 10.3389/fphar.2022.811111

<sup>1</sup>Columbia University Medical Center, New York, NY, United States, <sup>2</sup>Centre for Statistics in Medicine, NDOIMS, University of Oxford, Oxford, United Kingdom, <sup>3</sup>Janssen Research and Development, Titusville, NJ, United States, <sup>4</sup>Department of Medical Informatics, Erasmus University Medical Center, Rotterdam, Netherlands, <sup>5</sup>Fundació Institut Universitari per a la Recerca a L'Atenció Primària de Salut Jordi Gol i Gurina (IDIAPJGol), Barcelona, Spain, <sup>6</sup>Department of Biostatistics, Fielding School of Public Health, University of California, Los Angeles, Los Angeles, CA, United States, <sup>7</sup>Department of Human Genetics, David Geffen School of Medicine at UCLA, University of California, Los Angeles, Los Angeles, CA, United States, <sup>8</sup>New York-Presbyterian Hospital, New York, NY, United States

**Objective:** Background incidence rates are routinely used in safety studies to evaluate an association of an exposure and outcome. Systematic research on sensitivity of rates to the choice of the study parameters is lacking.

**Materials and Methods:** We used 12 data sources to systematically examine the influence of age, race, sex, database, time-at-risk, season and year, prior observation and clean window on incidence rates using 15 adverse events of special interest for COVID-19 vaccines as an example. For binary comparisons we calculated incidence rate ratios and performed random-effect meta-analysis.

**Results:** We observed a wide variation of background rates that goes well beyond age and database effects previously observed. While rates vary up to a factor of 1,000 across age groups, even after adjusting for age and sex, the study showed residual bias due to the other parameters. Rates were highly influenced by the choice of anchoring (e.g., health visit, vaccination, or arbitrary date) for the time-at-risk start. Anchoring on a healthcare encounter yielded higher incidence comparing to a random date, especially for short time-at-risk. Incidence rates were highly influenced by the choice of the database (varying by up to a factor of 100), clean window choice and time-at-risk duration, and less so by secular or seasonal trends.



# How does indication affect incidence rates?

- Drug indication is an obvious but usually ignored source of error
  - Beta blockers can be indicated for acute myocardial infarction (AMI) and hypertension, and of course, those taking a beta blocker for AMI are at higher risk of subsequent AMI
- We know incidence of different health outcomes differ by indication.
  - **What is the extent of the variation?**
  - **What is its relative contribution to heterogeneity compared to age, biological sex, and database?**



# OHDSI Symposium October 2023

- How Often: Large scale characterization of incidence of outcomes following drug exposure
- Pre-Symposium
  - Draft protocol
  - Develop and evaluate phenotypes
  - Gathered research questions from OHDSI community
  - Release analysis package that has all the targets and outcomes of interest
- During Symposium (October 2023)
  - Execute How Often Analysis Package across OHDSI Network
  - Deploy viewer to allow for exploration of results
  - Collaborate on appropriate use of evidence
    - How to ensure reliability of results?
    - How to improve user interface to disseminate results?
    - What have we learned that can fill evidence gaps and improve decision making?



# Method

- Analysis was conducted in October 2023 on 13 databases
- Study Design:
  - Target cohorts: First occurrence of drug exposure (12 different classes, stratified by indication)
  - Outcome cohorts: 73 different outcomes (defined in the OHDSI phenotype library)
  - Time at risk: 1 day to 365 day after cohort start (Intent to treat)
  - Stratifications: Age and gender

$$\text{Incidence Rate} = \frac{\begin{array}{c} \text{\# persons in the target cohort who have} \\ \text{new outcome occurrence during the time-at-risk} \end{array}}{\begin{array}{c} \text{person-time-at-risk for persons in the target cohort} \\ \text{with time at risk} \end{array}}$$



# Method

Target cohorts:  
12 Drug classes,  
nested by  
indication

	Indications
Beta Blockers	1) hypertension, 2) heart failure, 3) acute myocardial infarction
Cephalosporins	1) Urinary tract infection, 2) pneumonia
Calcium Channel Blockers	1) Hypertension
DPP-4 Inhibitors	1) Type 2 diabetes mellitus
Fluoroquinolones	1) Urinary tract infection, 2) pneumonia
GLP-1 antagonists	1) Type 2 diabetes mellitus, 2) obesity
IL-23 Inhibitors	1) Psoriasis
JAK inhibitors	1) Rheumatoid arthritis, 2) Ulcerative colitis
SGLT2 Inhibitors	1) Type 2 diabetes mellitus, 2) heart failure
Thiazide Diuretics	1) Hypertension
Trimethoprim	1) Urinary tract infection, 2) pneumonia
TNF-alpha inhibitors	1) Rheumatoid arthritis, 2) Psoriatic Arthritis, 3) Crohns disease, 4) Ulcerative colitis, 5) Psoriasis





# Method

## Outcomes Cohort examples (73 total)

### Cardiovascular

- 3 and 4-point major adverse cardiovascular event (MACE) outcomes
- Cardiac death
- Torsades de Pointes
- Hospitalization with heart failure events

### Neurologic

- Stroke
- Headache
- Guillen-Barre Syndrome (GBS)

### Gastrointestinal

- Abdominal Pain
- Acute Liver Injury
- Diarrhea
- GI Bleed



# Analysis

- Random effect meta-analysis of incidence rates across the 13 databases
- For drug classes with >1 indication: Variance components analysis to quantify relative heterogeneity between age, biological sex, database, and indication
- R `metafor` package (`rma`)



# Results

- 77,631 total incidence rates calculated
- 8 different drug classes had at least 2 indications

Drug class	Indications	Median VC
Beta Blockers	1) Essential Hypertension, 2) Left Heart Failure, 3) Acute Myocardial Infarction	0.1013
SGLT2 Inhibitors	1) Type 2 Diabetes Mellitus, 2) Left Heart Failure	0.2642
<b>GLP-1 antagonists</b>	<b>1) Type 2 Diabetes Mellitus, 2) Obesity</b>	<b>&lt;0.001</b>
Cephalosporins	1) Urinary Tract Infection, 2) Acute Typical Pneumonia	0.0397
Fluoroquinolones	1) Urinary Tract Infection, 2) Acute Typical Pneumonia	0.0983
<b>Trimethoprim</b>	<b>1) Urinary tract infection, 2) Pneumonia</b>	<b>0.4887</b>
JAK inhibitors	1) Rheumatoid Arthritis, 2) Ulcerative Colitis	0.0383
TNF-alpha inhibitors	1) Plaque Psoriasis, 2) Rheumatoid Arthritis, 3) Ulcerative Colitis, 4) Psoriatic Arthritis, 5) Crohn's Disease	0.0332



# Relative Variance Components

Median Variance Components attributed to...

Drug class	Indications	Database	Age	Biological Sex
Beta Blockers	0.1013	0.1537	<b>0.3102</b>	0.0204
SGLT2 Inhibitors	0.2642	<b>0.3170</b>	0.2779	0.0155
GLP-1 antagonists	<0.001	<b>0.6117</b>	0.3678	0.0289
Cephalosporins	0.0397	0.7230	<b>1.4631</b>	0.0515
Fluoroquinolones	0.0983	<b>0.994</b>	0.8573	0.0696
Trimethoprim	0.4887	0.2772	<b>1.5228</b>	0.1219
JAK inhibitors	0.0383	0.1792	<b>0.2055</b>	0.0937
TNF-alpha inhibitors	0.0332	0.1675	<b>0.1815</b>	0.0221



## Key Takeaways & Next Steps

- Among the drug classes we looked at, Trimethoprim is the drug class that is *most* sensitive to indications; GLP-1 the least
- Relative heterogeneity:
  - Database/Age > Indications > Biological Sex
- Next Steps: How Often All x All



Thank you! 😊

