

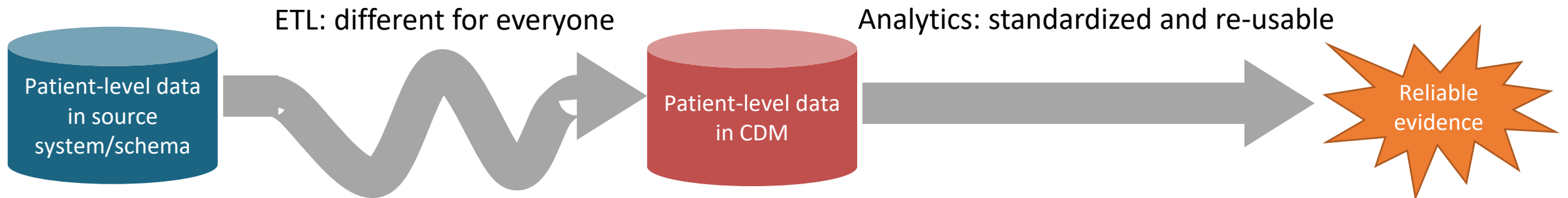


Advanced analytics with OHDSI tools



# Why convert to the Common Data Model?

- Transforming data to the OMOP CDM is a large investment
- The benefits come from being able to use the same tools and analytics across many databases

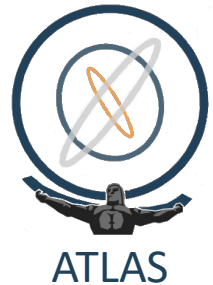




# OHDSI standardized analytics



- HADES is a set of open-source R package
- Developed and maintained by the community, for the community
- Can use cohort definitions created in ATLAS



966

*MEDINFO 2023 — The Future Is Accessible*  
*J. Bichel-Findlay et al. (Eds.)*

© 2024 International Medical Informatics Association (IMIA) and IOS Press.  
This article is published online with Open Access by IOS Press and distributed under the terms  
of the Creative Commons Attribution Non-Commercial License 4.0 (CC BY-NC 4.0).  
doi:10.3233/SHTI231108

## Health-Analytics Data to Evidence Suite (HADES): Open-Source Software for Observational Research

Martijn SCHUEMIE<sup>a,b,c,1</sup>, Jenna REPS<sup>a,b,d</sup>, Adam BLACK<sup>a,e</sup>, Frank DeFALCO<sup>a,b</sup>, Lee  
EVANS<sup>a,f</sup>, Egill FRIDGEIRSSON<sup>a,d</sup>, James P. GILBERT<sup>a,b</sup>, Chris KNOLL<sup>a,b</sup>, Martin  
LAVAILLEE<sup>a,g</sup>, Gowtham A. RAO<sup>a,b</sup>, Peter RINBECK<sup>a,d</sup>, Katarzyna SADOWSKA<sup>a,h</sup>

Analytic use case	Type	Structure
Clinical characterization	Disease Natural History	Amongst patients who are diagnosed with <b>&lt;insert your favorite disease&gt;</b> , what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <b>&lt;insert your favorite disease&gt;</b> , which treatments were patients exposed to amongst <b>&lt;list of treatments for disease&gt;</b> and in which sequence?
	Outcome incidence	Amongst patients who are new users of <b>&lt;insert your favorite drug&gt;</b> , how many patients experienced <b>&lt;insert your favorite known adverse event from the drug profile&gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> ?
Population-level effect estimation	Safety surveillance	Does exposure to <b>&lt;insert your favorite drug&gt;</b> increase the risk of experiencing <b>&lt;insert an adverse event&gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> ?
	Comparative effectiveness	Does exposure to <b>&lt;insert your favorite drug&gt;</b> have a different risk of experiencing <b>&lt;insert any outcome (safety or benefit) &gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> , relative to <b>&lt;insert your comparator treatment&gt;</b> ?
Patient level prediction	Disease onset and progression	For a given patient who is diagnosed with <b>&lt;insert your favorite disease&gt;</b> , what is the probability that they will go on to have <b>&lt;another disease or related complication&gt;</b> within <b>&lt;time horizon from diagnosis&gt;</b> ?
	Treatment response	For a given patient who is a new user of <b>&lt;insert your favorite chronically-used drug&gt;</b> , what is the probability that they will <b>&lt;insert desired effect&gt;</b> in <b>&lt;time window&gt;</b> ?
	Treatment safety	For a given patient who is a new user of <b>&lt;insert your favorite drug&gt;</b> , what is the probability that they will experience <b>&lt;insert adverse event&gt;</b> within <b>&lt;time horizon following exposure&gt;</b> ?

Analytic use case	Type	Structure
Clinical characterization	Disease Natural History	Amongst patients who are diagnosed with <insert your favorite disease>, what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <insert your favorite disease>, which treatments were patients exposed to amongst <list of treatments for disease> and in which sequence?
	Outcome incidence	Amongst patients who are new users of <insert your favorite drug>, how many patients experienced <insert your favorite known adverse event from the drug profile> within <time horizon following exposure start>?
Population effect		<insert your comparator treatment>?
Patient level prediction	Disease onset and progression	For a given patient who is diagnosed with <insert your favorite disease>, what is the probability that they will go on to have <another disease or related complication> within <time horizon from diagnosis>?
	Treatment response	For a given patient who is a new user of <insert your favorite chronically-used drug>, what is the probability that they will <insert desired effect> in <time window>?
	Treatment safety	For a given patient who is a new user of <insert your favorite drug>, what is the probability that they will experience <insert adverse event> within <time horizon following exposure>?

Standardizing the question itself helps clarify what your objective is

Analytic use case	Type	Structure
Clinical characterization	Disease Natural History	Amongst patients who are diagnosed with <b>&lt;insert your favorite disease&gt;</b> , what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <b>&lt;insert your favorite disease&gt;</b> , which treatments were patients exposed to amongst <b>&lt;list of treatments for disease&gt;</b> and in which sequence?
	Outcome incidence	Amongst patients who are new users of <b>&lt;insert your favorite drug&gt;</b> , how many patients experienced <b>&lt;insert your favorite known adverse event from the drug profile&gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> ?
Population-level effect estimation	Safety surveillance	Does exposure to <b>&lt;insert your favorite drug&gt;</b> increase the risk of experiencing <b>&lt;insert an adverse event&gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> ?
	Comparative effectiveness	Does exposure to <b>&lt;insert your favorite drug&gt;</b> have a different risk of experiencing <b>&lt;insert any outcome (safety or benefit) &gt;</b> within <b>&lt;time horizon following exposure start&gt;</b> , relative to <b>&lt;insert your comparator treatment&gt;</b> ?
Patient level prediction	Disease onset and progression	For a given patient who is diagnosed with <b>&lt;insert your favorite disease&gt;</b> , what is the probability that they will go on to have <b>&lt;another disease or related complication&gt;</b> within <b>&lt;time horizon from diagnosis&gt;</b> ?
	Treatment response	For a given patient who is a new user of <b>&lt;insert your favorite chronically-used drug&gt;</b> , what is the probability that they will <b>&lt;insert desired effect&gt;</b> in <b>&lt;time window&gt;</b> ?
	Treatment safety	For a given patient who is a new user of <b>&lt;insert your favorite drug&gt;</b> , what is the probability that they will experience <b>&lt;insert adverse event&gt;</b> within <b>&lt;time horizon following exposure&gt;</b> ?

Analytic use case	Type	Structure
Clinical characterization	Disease Natural History	Amongst patients who are diagnosed with <insert your favorite disease>, what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <insert your favorite disease>, which treatments were patients exposed to amongst <list of treatments for disease> and in which sequence?
	Outcome incidence	Amongst patients who are new users of <insert your favorite drug>, how many patients experienced <insert your favorite known adverse event from the drug profile> within <time horizon following exposure start>?
Population effect	<p>Amongst patients who are new users of <b>GLP-1s</b>, how many patients experienced <b>Acute Myocardial Infarction</b> within <b>drug exposure</b>?</p>	
Patient level prediction	Progression	For a given patient who is a new user of <insert your favorite drug>, what is the probability that they will go on to have <another disease or related complication> within <time horizon from diagnosis>?
	Treatment response	For a given patient who is a new user of <insert your favorite chronically-used drug>, what is the probability that they will <insert desired effect> in <time window>?
	Treatment safety	For a given patient who is a new user of <insert your favorite drug>, what is the probability that they will experience <insert adverse event> within <time horizon following exposure>?

Analytic use case	Type	Structure
Clinical characterization	Disease Natural History	Amongst patients who are diagnosed with <insert your favorite disease>, what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <insert your favorite disease>, which treatments were patients exposed to amongst <list of treatments for disease> and in which sequence?
	Outcome incidence	Amongst patients who are new users of <insert your favorite drug>, how many patients experienced <insert your favorite known adverse event from the drug profile> within <time horizon following exposure start>?
Population-level effect estimation	Safety surveillance	Does exposure to <insert your favorite drug> increase the risk of experiencing <insert an adverse event> within <time horizon following exposure start>?
	Comparative	Does exposure to <insert your favorite drug> have a different risk of experiencing <insert any
Patient level prediction	Treatment safety	For a given patient who is a new user of <insert your favorite drug>, what is the probability that they will experience <insert adverse event> within <time horizon following exposure>?

For a given patient who is a new user of **GLP-1s**, what is the probability that they will experience **an AMI** while **exposed to the drug**?



Analytic use case	Type	Structure
Clinical characteristics	Disease Natural History	Amongst patients who are diagnosed with <insert your favorite disease>, what are the patient's characteristics from their medical history?
	Treatment utilization	Amongst patients who have <insert your favorite disease>, which treatments were patients exposed to amongst <list of treatments for disease>, and in which sequence?
<p>Does exposure to <b>GLP-1s</b> have a different risk of experiencing <b>AMI</b> while <b>exposed to drug</b>, relative to <b>DPP-4s</b>?</p>		
Population level effect estimation	Comparative effectiveness	Does exposure to <insert your favorite drug> have a different risk of experiencing <insert any outcome (safety or benefit) > within <time horizon following exposure start>, relative to <insert your comparator treatment>?
Patient level prediction	Disease onset and progression	For a given patient who is diagnosed with <insert your favorite disease>, what is the probability that they will go on to have <another disease or related complication> within <time horizon from diagnosis>?
	Treatment response	For a given patient who is a new user of <insert your favorite chronically-used drug>, what is the probability that they will <insert desired effect> in <time window>?
	Treatment safety	For a given patient who is a new user of <insert your favorite drug>, what is the probability that they will experience <insert adverse event> within <time horizon following exposure>?



# Cohorts of our examples

Cohort: a group of people who satisfy some criteria for some period of time

- Indication cohorts:
  - Type-2 diabetes mellitus (**T2DM**)      People with T2DM, while having T2DM
- Exposures cohorts :
  - **GLP-1** agonists      People on GLP-1, while on the drug
  - **DPP-4** inhibitors      People on DPP-4, while on the drug
- Outcomes cohorts :
  - Acute myocardial infarction (**AMI**)      People with AMI, at the time of AMI

These same cohorts can be re-used to answer different questions

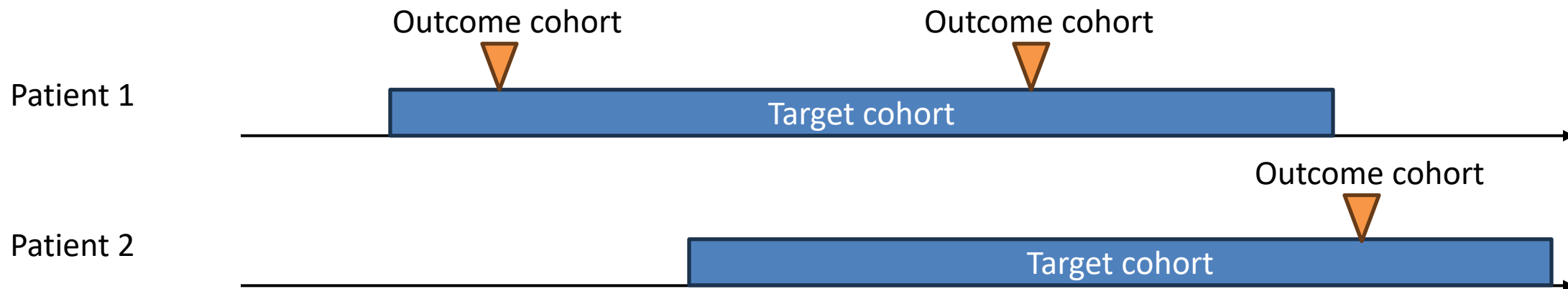
Patrick will discuss how to build these



# Characterization: CohortIncidence package

Amongst patients who are new users of **GLP-1s**, how many patients experienced **Acute Myocardial Infarction** within **drug exposure**?

- Target: **GLP-1**
- Outcome: **AMI**



Computes the incidence rate of the Outcome cohort in some Target cohort

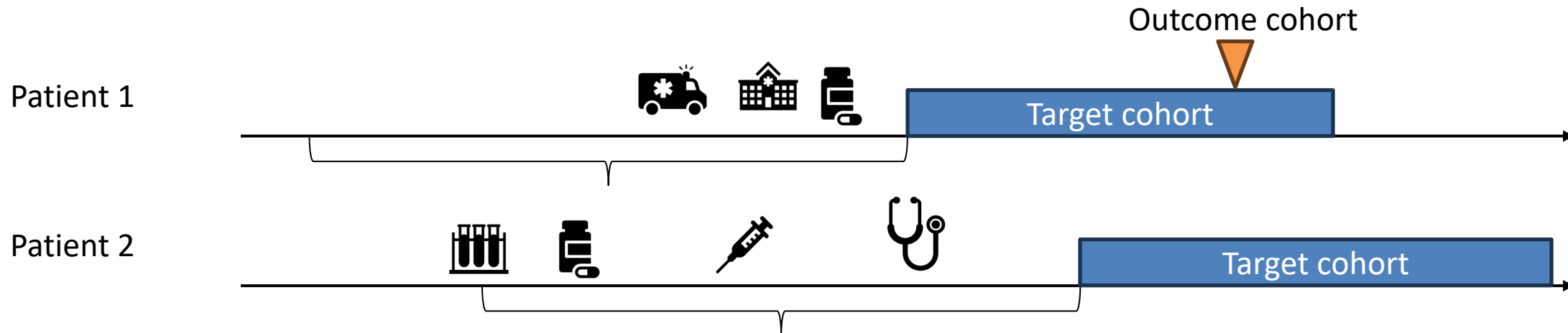
- Standardized computation of incidence rates
- Default: overall and stratified by age, sex, and calendar time



# PatientLevelPrediction package

For a given patient who is a new user of **GLP-1s**, what is the probability that they will experience **an AMI** while **exposed to the drug**?

- Target: **GLP-1**, restricted to those with **T2DM** (and first use only)
- Outcome: **AMI**



Builds a model to predict who in the Target will have the Outcome

- Uses all observed data up to Target start
- Implements many machine learning / deep learning algorithms

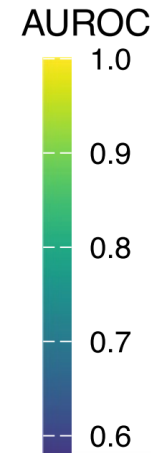
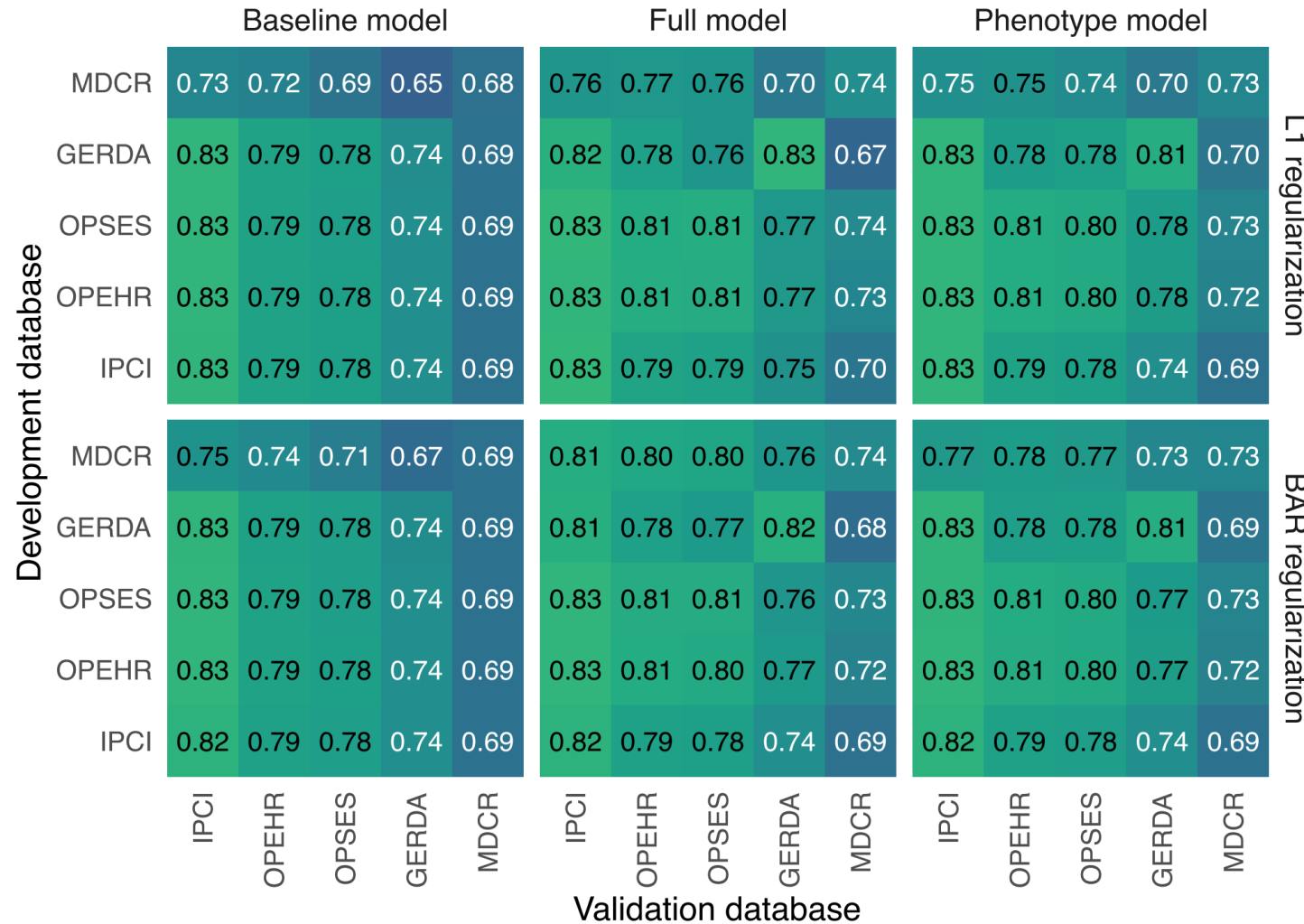


# Unique feature: external validation

- When using the PatientLevelPrediction package, models fit in one database can easily be validated in other databases
  - Portability of code
  - Standardized construction of features



# Example of external validation



John et al. *BMC Medicine* (2024) 22:308  
<https://doi.org/10.1186/s12916-024-03530-9>

BMC Medicine

RESEARCH ARTICLE

Open Access

## Development and validation of a patient-level model to predict dementia across a network of observational databases

Luis H. John<sup>1\*</sup>, Egill A. Fridgeirsson<sup>1</sup>, Jan A. Kors<sup>1</sup>, Jenna M. Reps<sup>2</sup>, Ross D. Williams<sup>1</sup>, Patrick B. Ryan<sup>2</sup> and Peter R. Rijnbeek<sup>1</sup>

Abstract

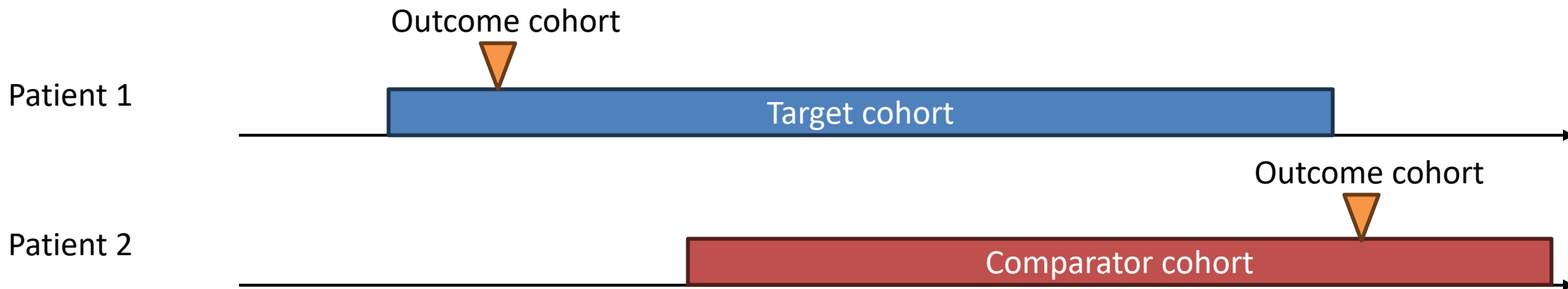
**Background** A prediction model can be a useful tool to quantify the risk of a patient developing dementia



# CohortMethod package

Does exposure to **GLP-1s** have a different risk of experiencing **AMI** while **exposed to drug**, relative to **DPP-4s**?

- Target: **GLP-1**, restricted to those with **T2DM** (and first use only)
- Comparator: **DPP-4**, restricted to those with **T2DM** (and first use only)
- Outcome: **AMI**



Computes the hazard of the Outcome cohort in the Target cohort compared to the Comparator



# Unique feature: Large-scale propensity scores

- Treatment assignment is often non-random, which can cause confounding
  - E.g. GLP-1 may be prescribed more often to obese, who already have a higher risk of AMI
- Propensity scores are an established way to address this
  - Fit a model to predict treatment assignment, and use to compute probability (propensity score)
  - Match subjects in Target to Comparator with similar propensity scores
- Traditionally, experts pick a few variables to use in the prediction model
- Large-scale propensity scores include all baseline covariates, and use machine learning



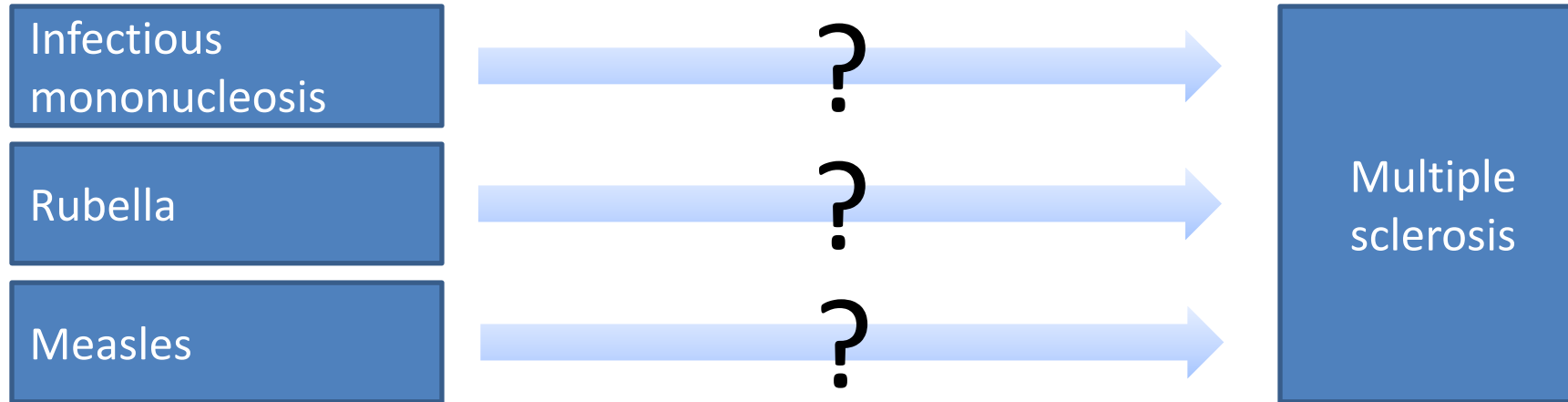


# Unique feature: objective diagnostics

- Whether study results are reliable depends on whether certain assumptions have been met
  - E.g. we assume our PS adjustment makes our treatment groups comparable
- Most of these assumptions are testable through diagnostics
  - E.g. we can test whether our PS adjustment achieved balance by computing the standardized difference of means (SDM)
- By ‘objective’ diagnostics we mean diagnostics that are evaluated while blinded to the results of the study
  - E.g. Pre-specify that we will not look at results where  $\max(|\text{SDM}|) > 0.1$
  - Unique: negative controls



## Example of a negative control



RESEARCH PAPER

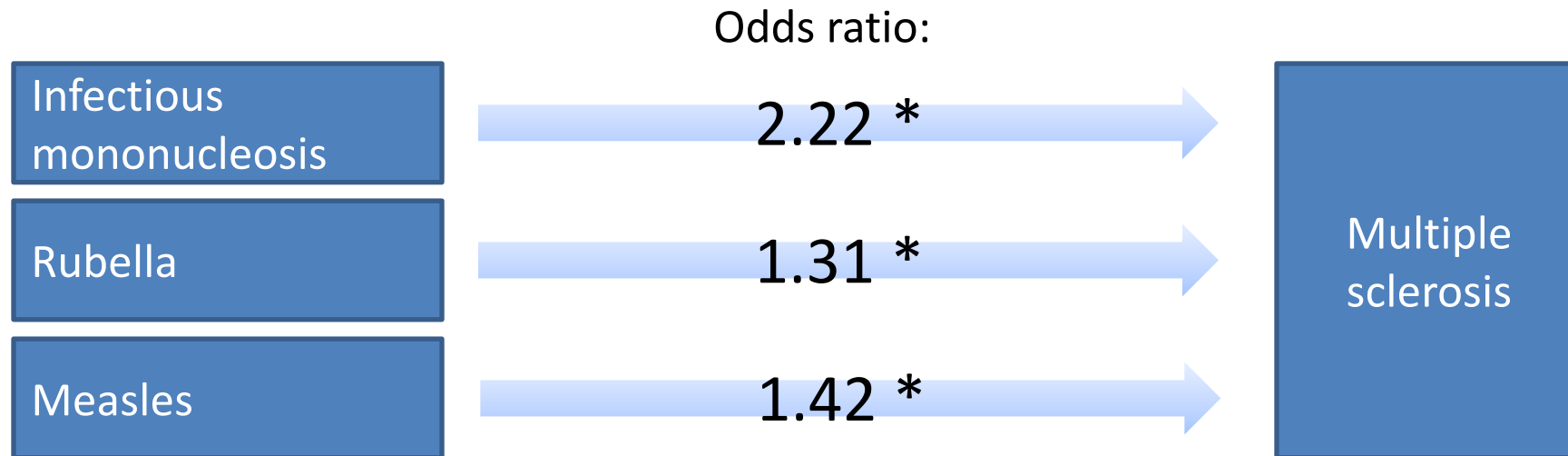
*Multiple Sclerosis* 2008; 14: 307–313

### Selective association of multiple sclerosis with infectious mononucleosis

*BM Zaadstra<sup>1,2</sup>, AMJ Chorus<sup>1</sup>, S van Buuren<sup>1,3</sup>, H Kalsbeek<sup>1</sup> and JM van Noort<sup>4</sup>*



## Example of a negative control



\* P < .05

RESEARCH PAPER

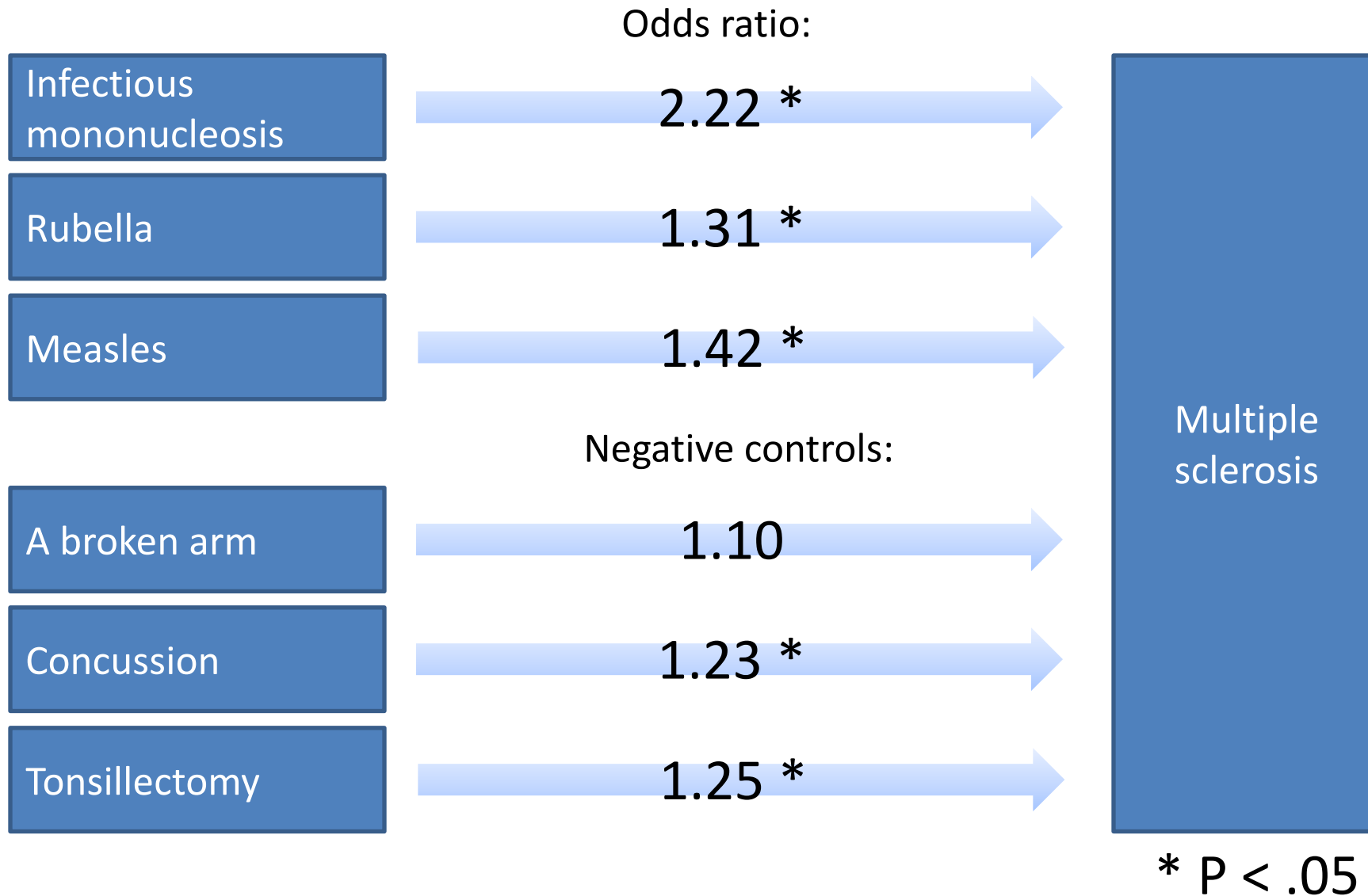
*Multiple Sclerosis* 2008; 14: 307–313

### Selective association of multiple sclerosis with infectious mononucleosis

BM Zaadstra<sup>1,2</sup>, AMJ Chorus<sup>1</sup>, S van Buuren<sup>1,3</sup>, H Kalsbeek<sup>1</sup> and JM van Noort<sup>4</sup>



## Example of a negative control



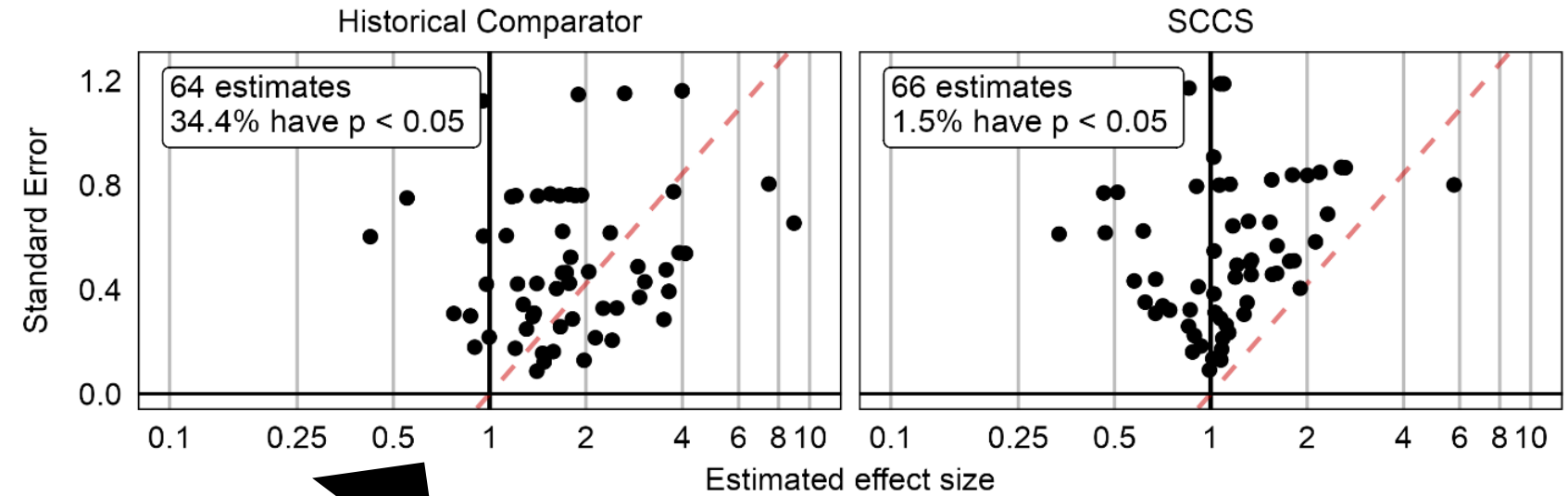


# How to interpret negative control findings?

- Unique: use a sample ( $n > 50$ ) of negative controls to understand distribution of bias
- Systematic error distribution can be used as
  - Diagnostic: if too much systematic error, we stop
  - Calibration: can adjust p-values and confidence intervals to take into account possible systematic error



# Quantifying systematic error



Need to execute estimation studies for  
66 Target-Outcome combinations.

OHDSI tools readily allow for this (simply  
swapping out the outcome cohort for the  
negative controls)

Received: 8 July 2022 | Revised: 30 September 2022 | Accepted: 8 December 2022  
DOI: 10.1002/sim.9631

## RESEARCH ARTICLE

Statistics  
in Medicine WILEY

### Adjusting for both sequential testing and systematic error in safety surveillance using observational data: Empirical calibration and MaxSPRT

Martijn J. Schuemie<sup>1,2</sup> | Fan Bu<sup>2,3</sup> | Akihiko Nishimura<sup>4</sup> | Marc A. Suchard<sup>2,3,5</sup>

<sup>1</sup>Observational Health Data Analytics, Janssen Research & Development, Titusville, New Jersey,

<sup>2</sup>Department of Biostatistics, University of California, Los Angeles, California,

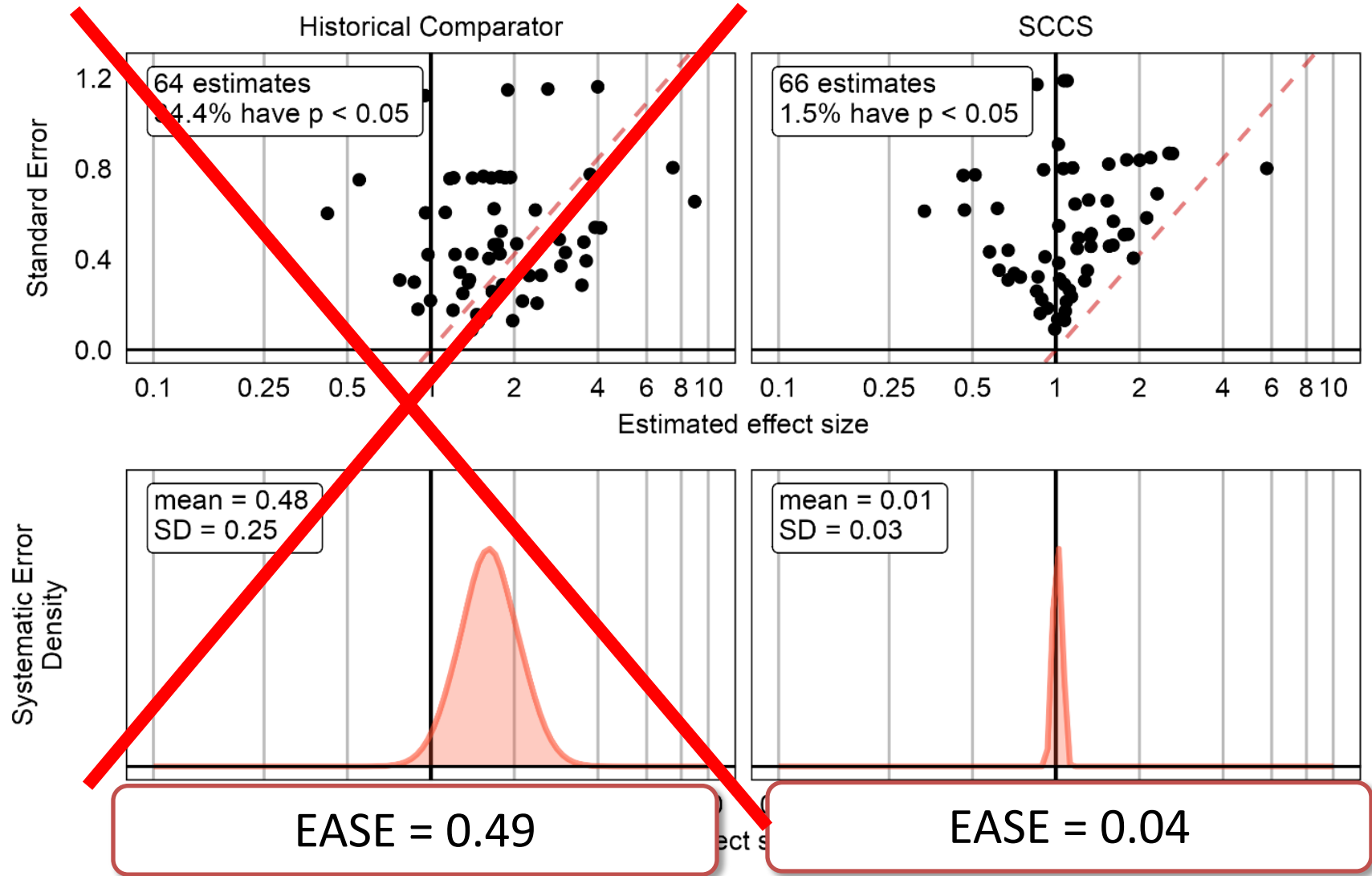
<sup>3</sup>Department of Human Genetics, University of California, Los Angeles,

Post-approval safety surveillance of medical products using observational healthcare data can help identify safety issues beyond those found in pre-approval trials. When testing sequentially as data accrue, maximum sequential probability ratio testing (MaxSPRT) is a common approach to maintaining nominal type 1 error. However, the true type 1 error may still deviate from the

# Quantifying systematic error

Expected Absolute Systematic Error (**EASE**) summarizes this distribution

We use a **prespecified** EASE threshold (EASE < 0.25) for go – no go decisions for our studies





# Distributed analyses

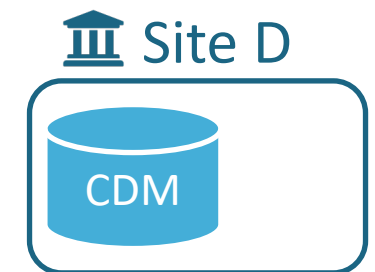
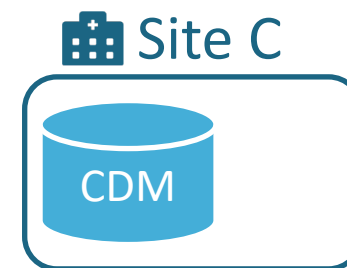
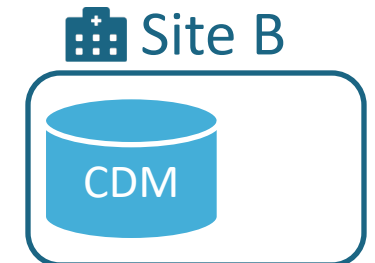
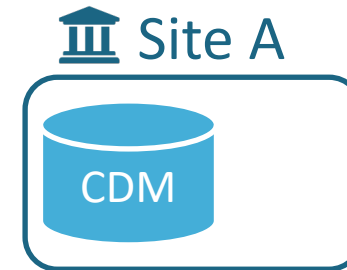
Using OHDSI tools





# Distributed Research Network

- Multiple sites with data
  - Hospital EHRs (Electronic Health Records)
  - Administrative Claims
- Patient-level data cannot be shared
- Each site uses the Common Data Model (CDM)

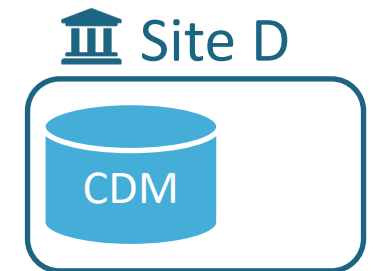
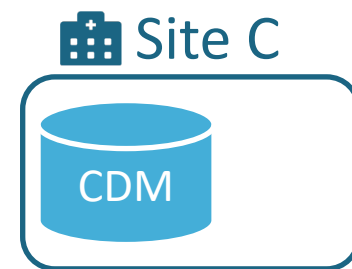
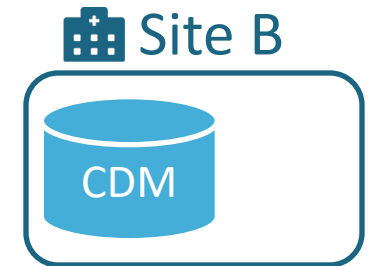
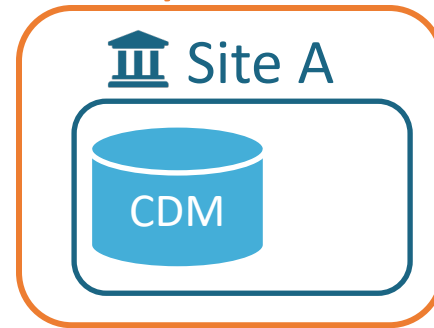




# Distributed Research Network

- A site can lead a study

Study lead

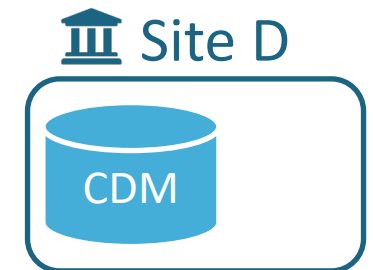
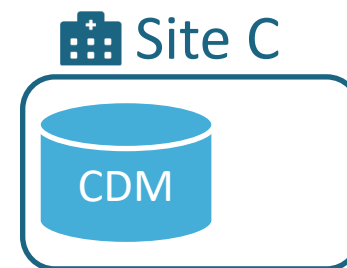
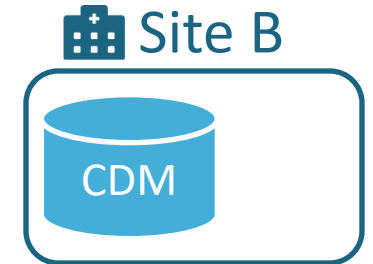
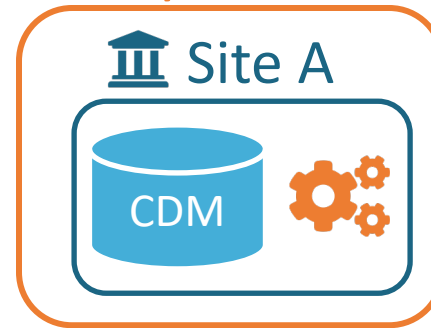




# Distributed Research Network

- A site can lead a study
- Analysis code is developed locally

Study lead

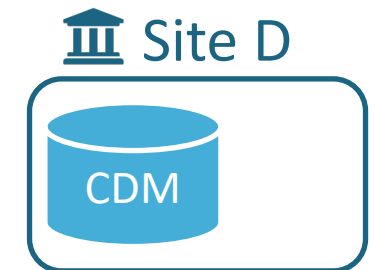
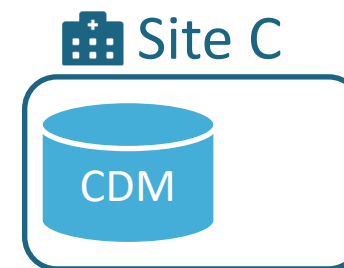
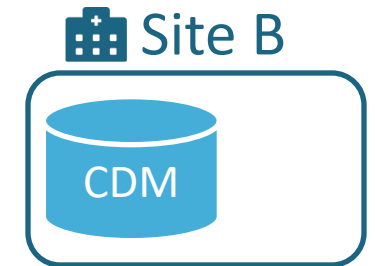
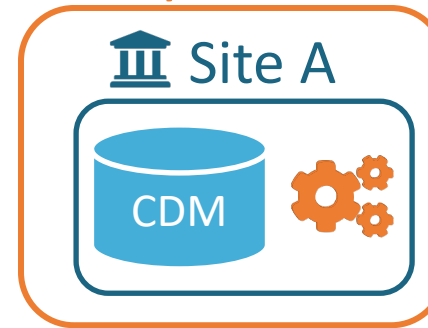




# Distributed Research Network

- A site can lead a study
- Analysis code is developed locally
- Code is distributed to study participants

Study lead

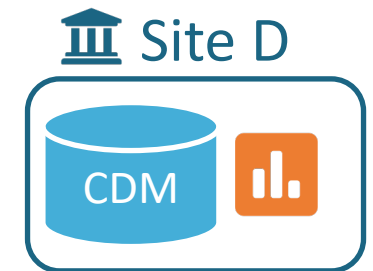
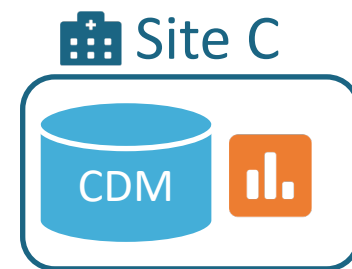
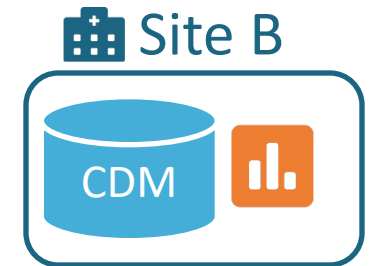
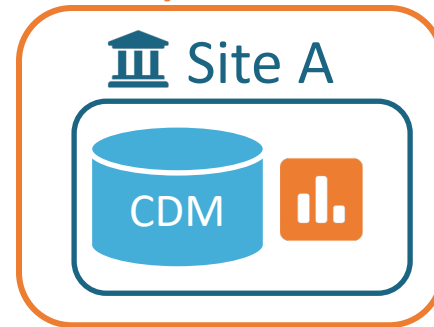




# Distributed Research Network

- A site can lead a study
- Analysis code is developed locally
- Code is distributed to study participants
- Results are generated (aggregated statistics)

Study lead

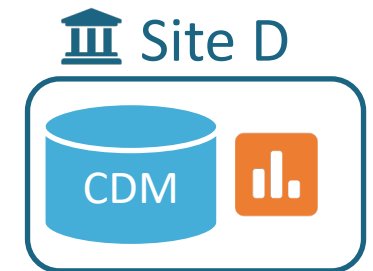
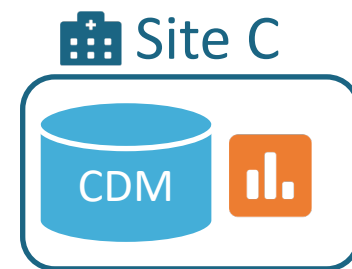
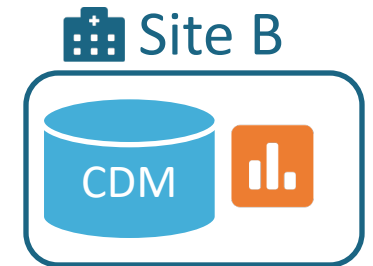
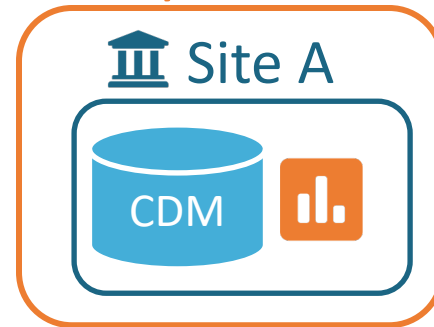




# Distributed Research Network

- A site can lead a study
- Analysis code is developed locally
- Code is distributed to study participants
- Results are generated (aggregated statistics)
- Results are sent back to lead site

Study lead

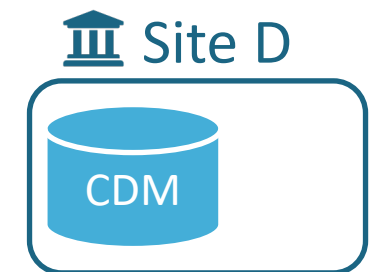
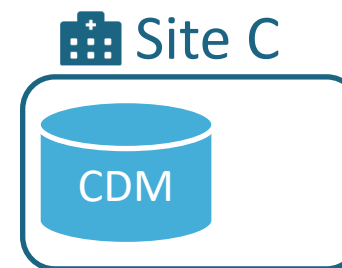
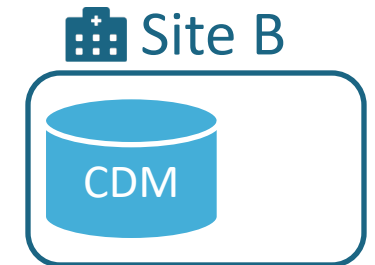
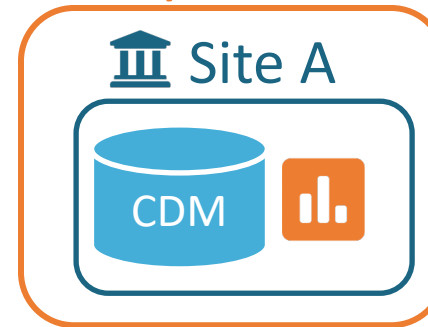




# Distributed Research Network

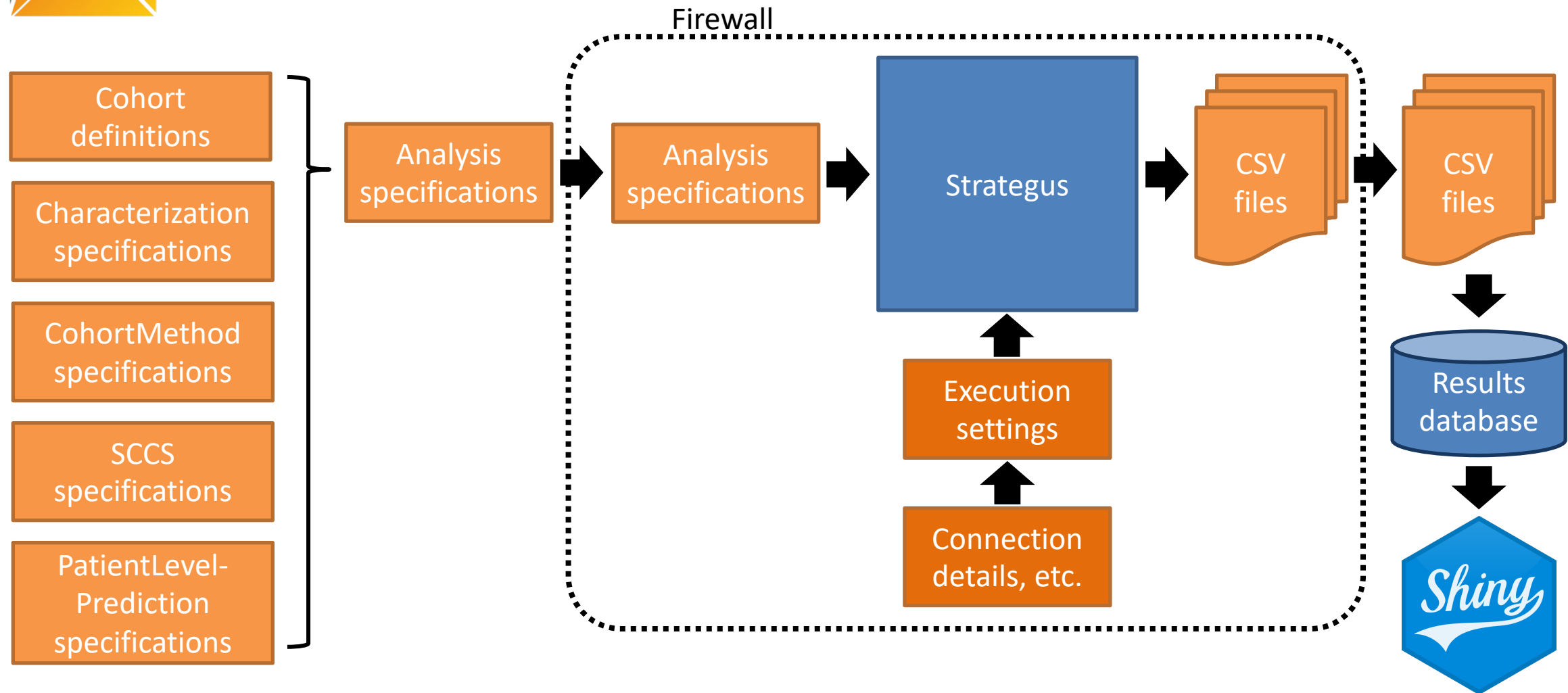
- A site can lead a study
- Analysis code is developed locally
- Code is distributed to study participants
- Results are generated (aggregated statistics)
- Results are sent back to lead site
- Evidence is synthesized

Study lead





# Strategus for study execution







## Summary



# Unique features of HADES analytics

- Re-use of cohort definitions
- Standardization of analytics in open-source software
  - Many opportunities for testing, review, fixing bugs, etc.
  - Making it hard to do the wrong thing (opinionated)
- Advanced methods to reduce bias
  - Splines for time in self-controlled case series
  - Large-scale propensity scores in cohort method
- Objective study diagnostics to improve reliability of evidence
  - Including negative controls
- Designed to run across a network of databases
  - Without sharing patient-level data