



Lessons from large-scale evidence generation across a network of databases (LEGEND) for hypertension

Marc A. Suchard, MD, PhD
on behalf of the LEGEND team

OHDSI Introductory Lecture
27 September 2025



News flash: new(-ish) hypertension guidelines!

CNN Health | Food | Fitness | Wellness | Parenting | Vital Signs | Live TV | US Edition | Search | More

Nearly half of Americans now have high blood pressure, based on new guidelines

By Susan Scott, CNN

Updated 2:17 AM ET, Tue November 14, 2017



Why is high blood pressure a "silent killer"? 03:29

- Lower blood pressure cutoffs
- New first-line treatment recommendations

October 2017

The New York Times

*Under New Guidelines,
Millions More Americans Will
Need to Lower Blood Pressure*





Professional guidelines

Professional society guidelines drive much of clinical practice. They synthesize

- Knowledge: RCTs, case-reports, observational studies
- Wisdom: expert opinion, collective experience

US

JOURNAL OF THE AMERICAN COLLEGE OF CARDIOLOGY
© 2016 BY THE AMERICAN COLLEGE OF CARDIOLOGY FOUNDATION AND
THE AMERICAN HEART ASSOCIATION, INC.

CLINICAL PRACTICE GUIDELINE

2017 ACC/AHA/AAPA/ABC/ACPM/
AGS/APhA/ASH/ASPC/NMA/PCNA
Guideline for the Prevention,
Detection, Evaluation, and Management
of High Blood Pressure in Adults



VOL. 71, NO. 19, 2018



ESC

European Society
of Cardiology

European Heart Journal (2018) 39, 3021–3104

doi:10.1093/eurheartj/ehy339

Europe

ESC/ESH GUIDELINES

2018 ESC/ESH Guidelines for the management
of arterial hypertension

The Task Force for the management of arterial hypertension of the
European Society of Cardiology (ESC) and the European Society of
Hypertension (ESH)

- October 2017

- June 2018



First-line agents

2017 ACC/AHA Guidelines

COR	LOE	RECOMMENDATION
I	A ^{SR}	1. For initiation of antihypertensive drug therapy, first-line agents include thiazide diuretics, CCBs, and ACE inhibitors or ARBs. (S8.1.6-1,S8.1.6-2)

"In particular, there is inadequate evidence to support the initial use of beta blockers for hypertension in the absence of specific cardiovascular comorbidities."

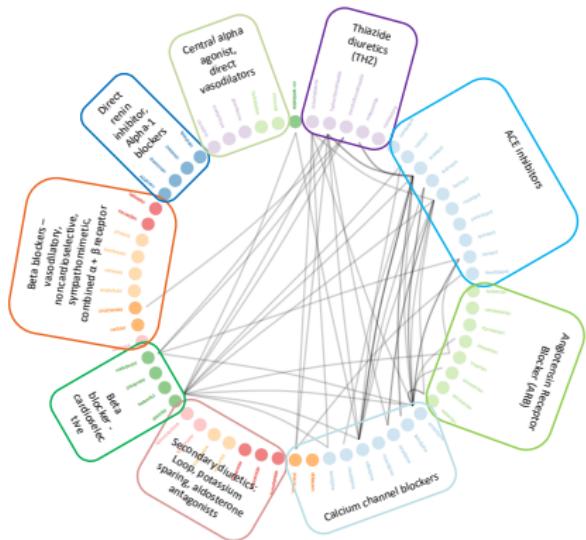
2018 ESC/ESH Guidelines: Class I, Level A

"Among all antihypertensive drugs, ACE inhibitors, ARBs, beta-blockers, CCBs, and [thiazide] diuretics ... have demonstrated effective reduction of BP and CV events in RCTs, and thus are indicated as the basis of antihypertensive treatment strategies." "...beta-blockers are usually equivalent in preventing major CV events, except for less effective prevention of stroke ..."



Current knowledge base for hypertension

Head-to-head antihypertensive drug comparisons



- Trials: 40
- $N = 102 - [1148] - 33K$

- Driven primarily by ALLHAT
 - ▶ just 3 individual drugs
- Focus: efficacy \gg safety
- New RCTs too **expensive**

Can we provide

1. reliable / reproducible – concordant extant w/ RCTs
2. rich – across “all” comparators, outcomes
3. relevant – inform practice evidence?



Odyssey (*noun*): \oh-d-si\

1. A long journey full of adventures
2. A series of experiences that give knowledge or understanding to someone





The journey to real-world evidence



Different observational data types:

- Populations, care settings, capture process, health system

Types of evidence desired:

- Cohort identification, clinical characterization, population-level effects, patient-level prediction

OHDSI's mission

- To improve health, by empowering a community to collaboratively generate the evidence that promotes better health decisions and better care

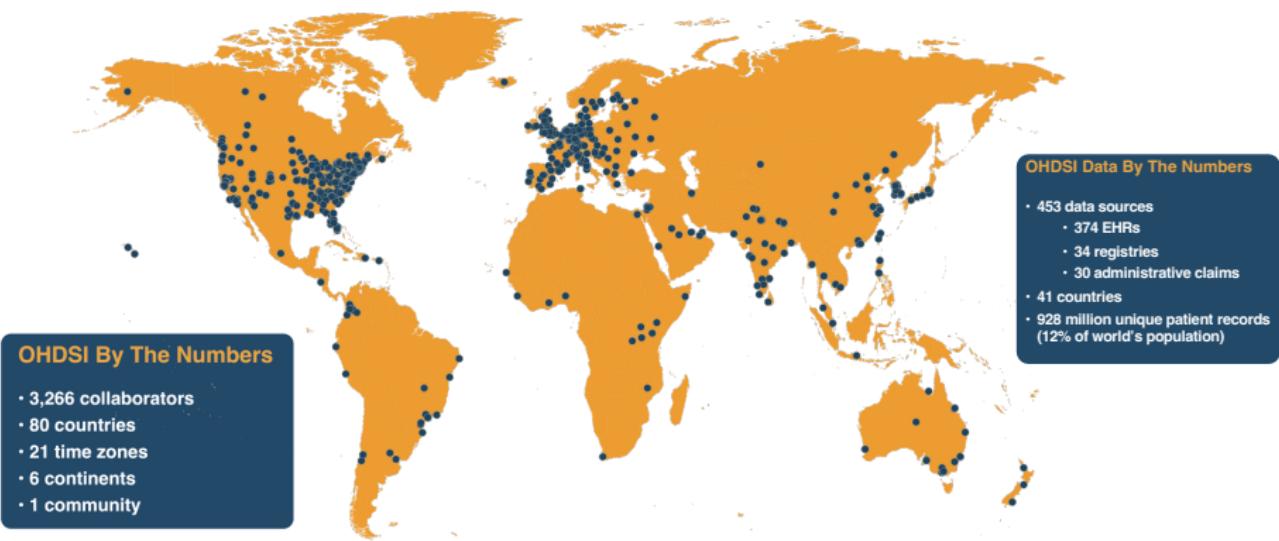


A little history about OHDSI

- Born in early 2014 out of the **Observational Medical Outcome Partnership** (OMOP) experiment
 - ▶ closed US Food & Drug Administration-guided team to evaluate the (Mini-) Sentinel initiative (short story: Sentinel was not going to work)
 - ▶ more than just a common data model (CDM)
 - ▶ often confusion between terms OMOP and OHDSI
- I was an OMOP Research Investigator, am a founding member of OHDSI and continue to serve on its executive leadership



Our journey as a community of collaborators



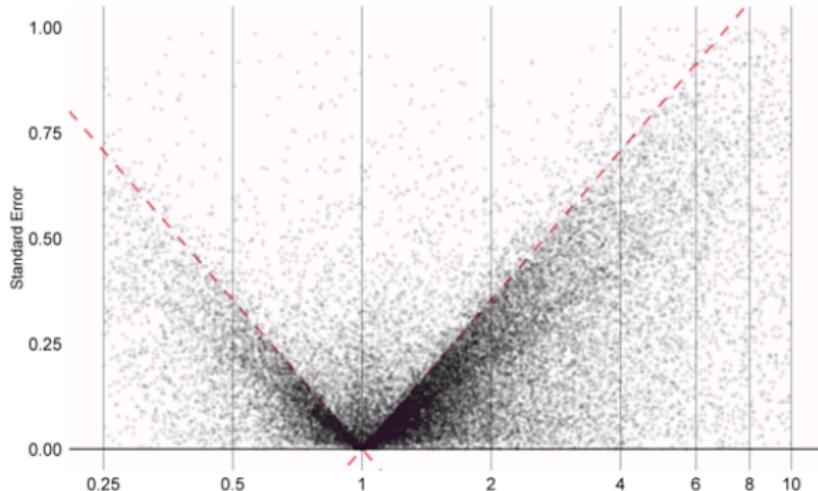
Open source, data science tool-stack
protocols, discussion, etc.



Observational research estimates in literature

29,982 drug safety estimates
from 11,758 papers

What is going **wrong**?



- 85% have reported $p < 0.05$
- Also note unusual peak along boundary

- Observational bias (confounding, selection, measurement error)
- Publication bias
- *p*-hacking (one study at a time)
- Reproducibility across populations



What is going wrong? *p*-hacking

PhD Student!

I think A may cause B,
go investigate!

Yes professor!



I ran the analysis:

$p > .05$

But did you adjust
for confounder Z?

Ehh, no

Let me get
right/back to you



After adjustment
for Z, $p < .05!$

Yay! Lets publish
a paper!

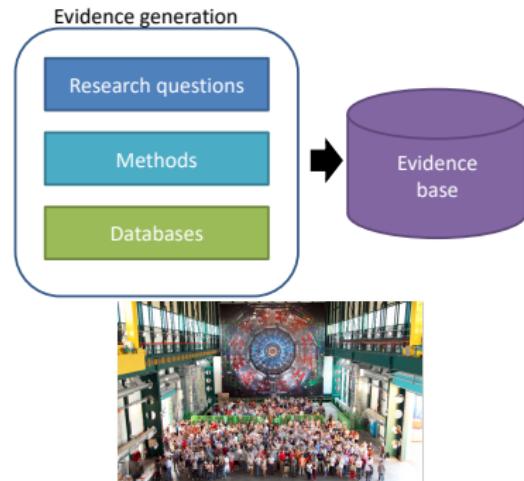




A data science solution?

Large-scale evidence generation
across a network of databases (LEGEND)

- Aims to generate reliable evidence on the effects of medical interventions using observational healthcare data
- 10 guiding principles; chief among these:
 - ▶ Generate at **large-scale** (completeness, empirical **calibration**)
 - ▶ Systematically driven by **best-practices**
 - ▶ Disseminate **everything** (open science)



No one person has all the necessary skills



Best-practices: systematic design



American Journal of Epidemiology

© The Author 2016. Published by Oxford University Press on behalf of the Johns Hopkins Bloomberg School of Public Health. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com.

Vol. 183, No. 8

DOI: 10.1093/aje/kwv254

Advance Access publication:

March 18, 2016

Practice of Epidemiology

Using Big Data to Emulate a Target Trial When a Randomized Trial Is Not Available

Miguel A. Hernán* and James M. Robins

* Correspondence to Dr. Miguel A. Hernán, Department of Epidemiology, 677 Huntington Avenue, Boston, MA 02115
(e-mail: miguel_hernan@post.harvard.edu).

Initially submitted December 9, 2014; accepted for publication September 8, 2015.

Ideally, questions about comparative effectiveness or safety would be answered using an appropriately designed and conducted randomized experiment. When we cannot conduct a randomized experiment, we analyze observational data. Causal inference from large observational databases (big data) can be viewed as an attempt to emulate a randomized experiment—the target experiment or target trial—that would answer the question of interest. When the goal is to guide decisions among several strategies, causal analyses of observational data need to be evaluated with respect to how well they emulate a particular target trial. We outline a framework for comparative effectiveness research using big data that makes the target trial explicit. This framework channels counterfactual theory for comparing the effects of sustained treatment strategies, organizes analytic approaches, provides a structured process for the criticism of observational studies, and helps avoid common methodologic pitfalls.

big data; causal inference; comparative effectiveness research; target trial

- Much older idea: William Cochran / Gertrude Cox, 1950,
Experimental Design



Target trial for comparing two initial therapies

Treatment strategies:

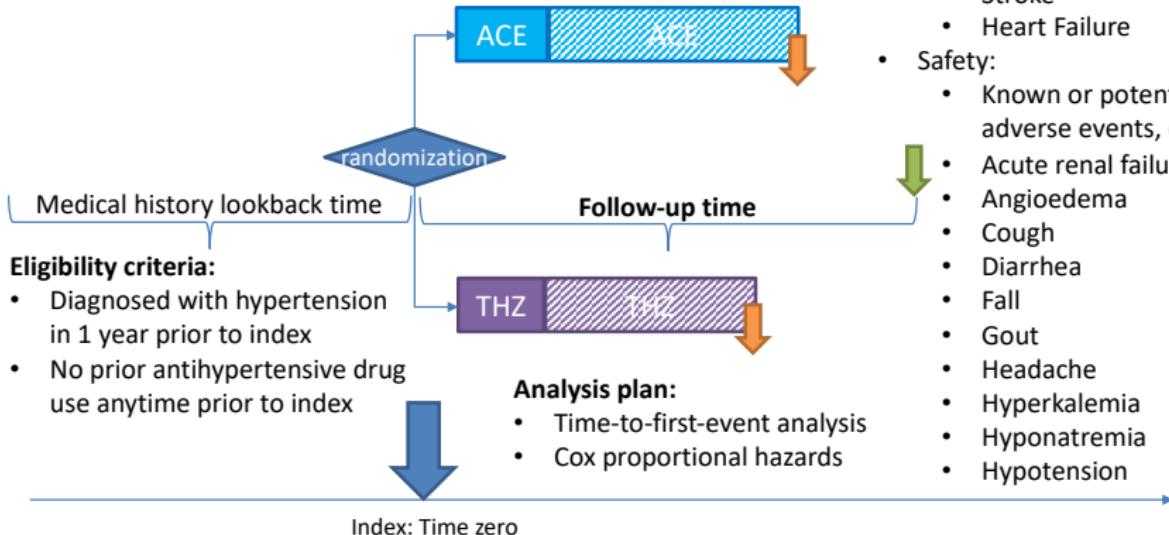
- Monotherapy with ACE
- Monotherapy with THZ

Causal contrasts of interest:

- Intent-to-treat effect
- On-treatment effect

Outcomes:

- Efficacy:
 - Myocardial infarction
 - Stroke
 - Heart Failure
- Safety:
 - Known or potential adverse events, e.g.
 - Acute renal failure
 - Angioedema
 - Cough
 - Diarrhea
 - Fall
 - Gout
 - Headache
 - Hyperkalemia
 - Hyponatremia
 - Hypotension





Observ. study for comparing two initial therapies

Treatment strategies:

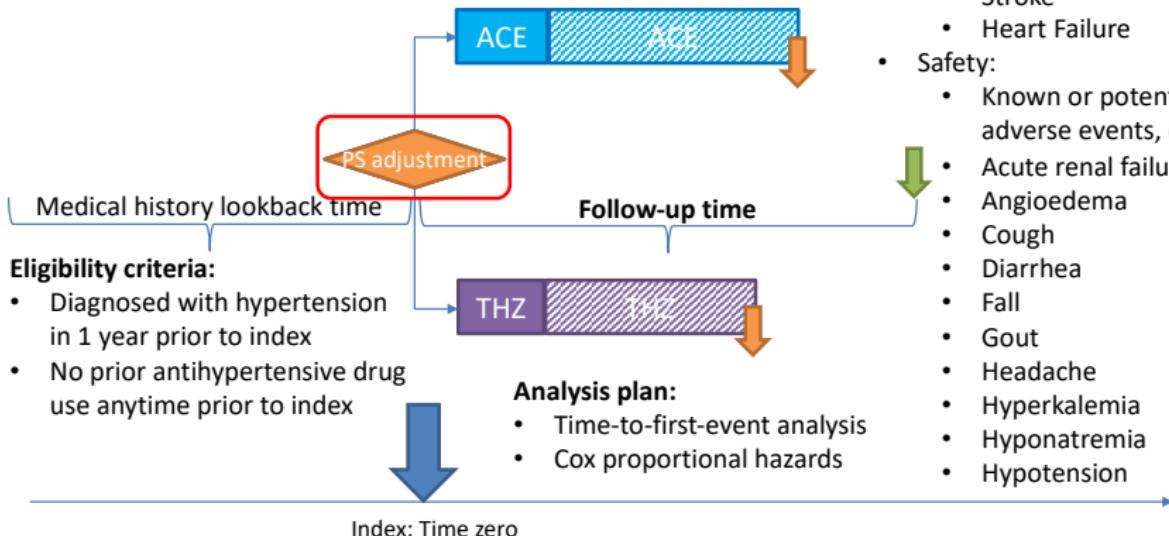
- Monotherapy with ACE
- Monotherapy with THZ

Causal contrasts of interest:

- Intent-to-treat effect
- On-treatment effect

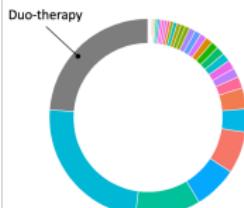
Outcomes:

- Efficacy:
 - Myocardial infarction
 - Stroke
 - Heart Failure
- Safety:
 - Known or potential adverse events, e.g.
 - Acute renal failure
 - Angioedema
 - Cough
 - Diarrhea
 - Fall
 - Gout
 - Headache
 - Hyperkalemia
 - Hyponatremia
 - Hypotension





Comparison of hypertension treatments

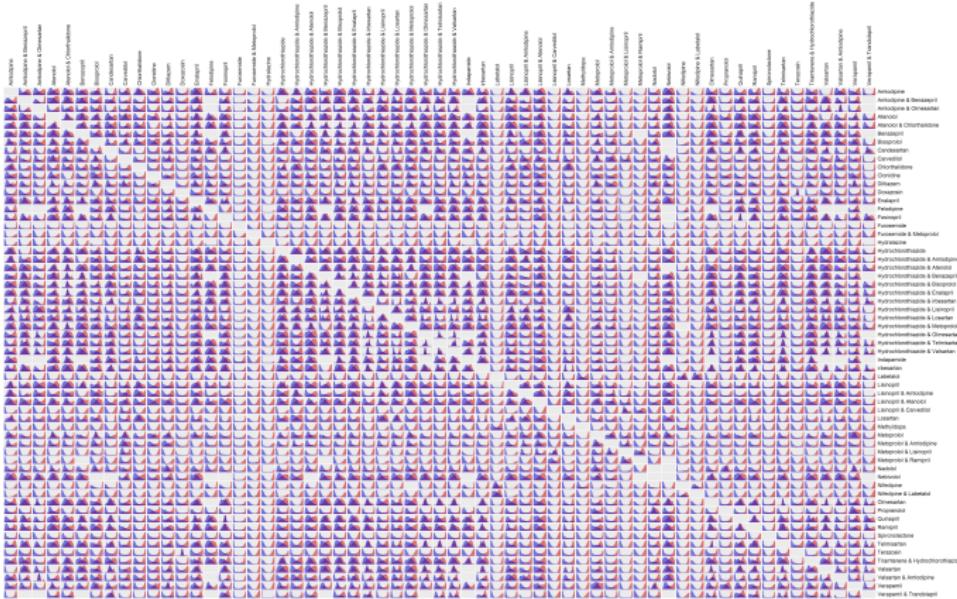


- 39 mono-drugs, 13 mono-classes
- 58 duo-drugs, 32 duo-classes
- 10,278 comparisons

	Theoretical	Observed (n > 2,500)
Single ingredients	58	39
Single ingredient comparisons	$58 * 57 = 3,306$	1,296
Single drug classes	15	13
Single class comparisons	$15 * 14 = 210$	156
Dual ingredients	$58 * 57 / 2 = 1,653$	58
Single vs duo drug comparisons	$58 * 1,653 = 95,874$	3,810
Dual classes	$15 * 14 / 2 = 105$	32
Single vs duo class comparisons	$15 * 105 = 1,575$	832
Duo vs duo drug comparisons	$1,653 * 1,652 = 2,730,756$	2,784
Duo vs duo class comparisons	$105 * 104 = 10,920$	992
...
Total comparisons	2,843,250	10,278



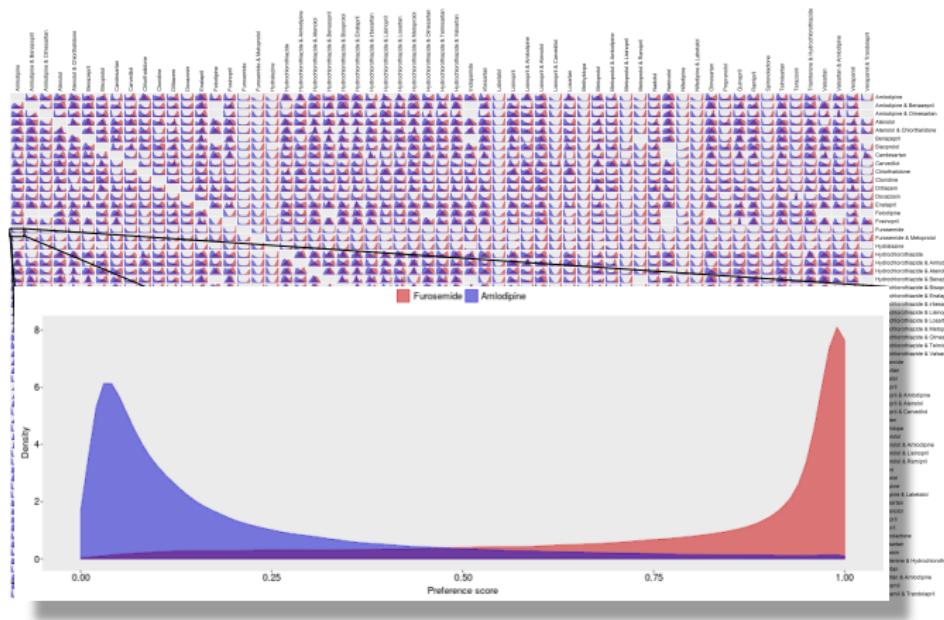
Best-practices: systematic large-scale PS



- >8,000 (regularized) baseline patient characteristics (all dx, rx, tx)
- Address observed (and some unobserved – BP control) confounding (Tian et al, 2019, IJE)



Of course, not all comparisons are valid



- Evaluation of propensity score (PS) distributions and covariate balance
- Here: poor empirical clinical equipoise



Best-practices: 58 expert-crafted outcomes

- Effectiveness (10): acute MI, heart failure, stroke
- Safety (48): known side-effects

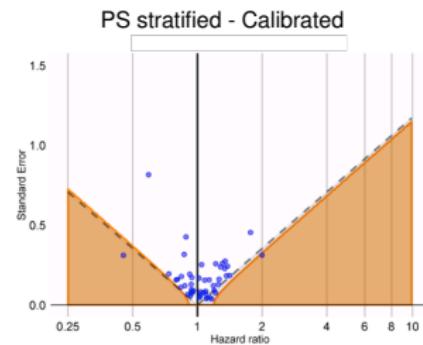
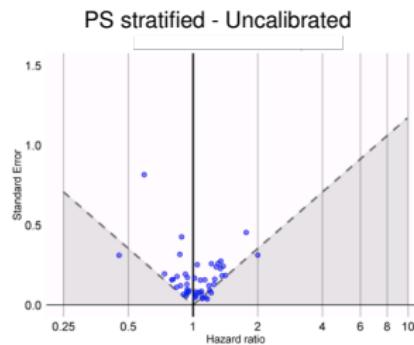
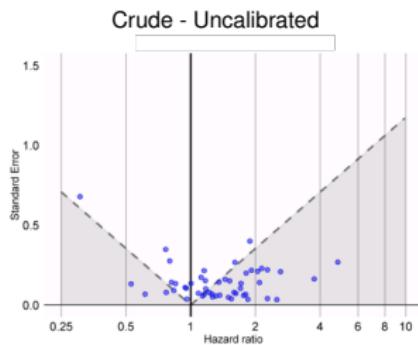
Phenotype	Logical description	Supporting references
Abdominal pain	Abdominal pain condition record of any type; successive records with > 90 day gap are considered independent episodes	4 5 6
Abnormal weight gain	Abnormal weight gain record of any type; successive records with > 90 day gap are considered independent episodes; note, weight measurements not used	7
Abnormal weight loss	Abnormal weight loss record of any type; successive records with > 90 day gap are considered independent episodes; note, weight measurements not used	8
Acute myocardial infarction	Acute myocardial infarction condition record during an inpatient or ER visit; successive records with > 180 day gap are considered independent episodes	9 10 11 12 13 14
Acute pancreatitis	Acute pancreatitis condition record during an inpatient or ER visit; successive records with >30 day gap are considered independent episodes	15 16 17 18
Acute renal failure	A diagnosis of acute renal failure in an inpatient or ER setting; must be at least 30d between inpatient/ER visits to be considered separate episodes	19 20 21 22 23 24 25 26

	Theoretical	Observed (n > 2,500)
Outcomes of interest	58	58
Target-comparator-outcomes	2,843,250 * 58 = 164,908,500	587,020



Calibrate each study (under null)

76 negative **outcome** controls (not caused by either treatment) help expose and control **residual bias**. Example: ingrown toenail



68% have $p < 0.05$

16% have $p < 0.05$

4% have $p < 0.05$

p-value empirical calibration models residual bias as exchangeable and adjusts for a (possibly) non-0 mean. (Schuemie et al, 2018, PNAS)



Calibrate each study (under alternative)

Trouble with **positive controls**:

- Few positive exemplars for a particular comparison
- Effect size never known certainty (and depends on population)
- We know they're positive and change our behavior accordingly

Proposed solution:

- Start with negative controls ($RR = 1$)
- Inject simulated outcomes during exposure until desired RR is achieved
- Simulants behave like 'real' outcomes: **preserve confounding** by injecting outcomes for patients at high risk

Yields: **Confidence intervals with (near) nominal coverage**



Network of data sources

Evidence generation

Research questions

Methods

Databases

- US insurance databases
 - ▶ IBM MarketScan CCAE
 - ▶ IBM MarketScan MDCR
 - ▶ IBM MarketScan MDCC
 - ▶ Optum Clininformatics
- Japanese insurance database: JMDC
- Korean insurance database: NHIS-NSC
- US EHR databases
 - ▶ Optum PanTher
 - ▶ Columbia University Medical Center
- German EHR database: IQVIA DA Germany



Ajou University



Columbia University

Account for population/practice heterogeneity (Madigan et al, 2013, AJE)

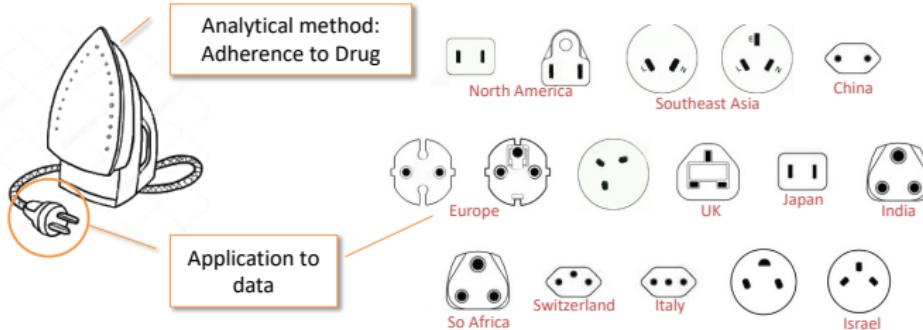
Improve generalizability



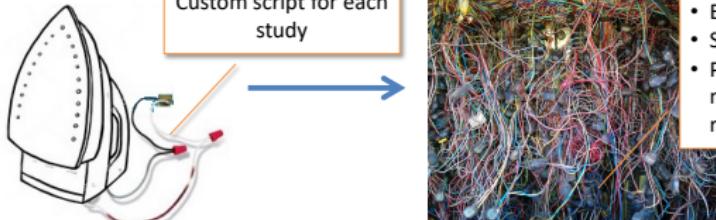
Current approach does not scale

One study question - one database - one script

"What's the adherence to my drug in the data assets I own?"



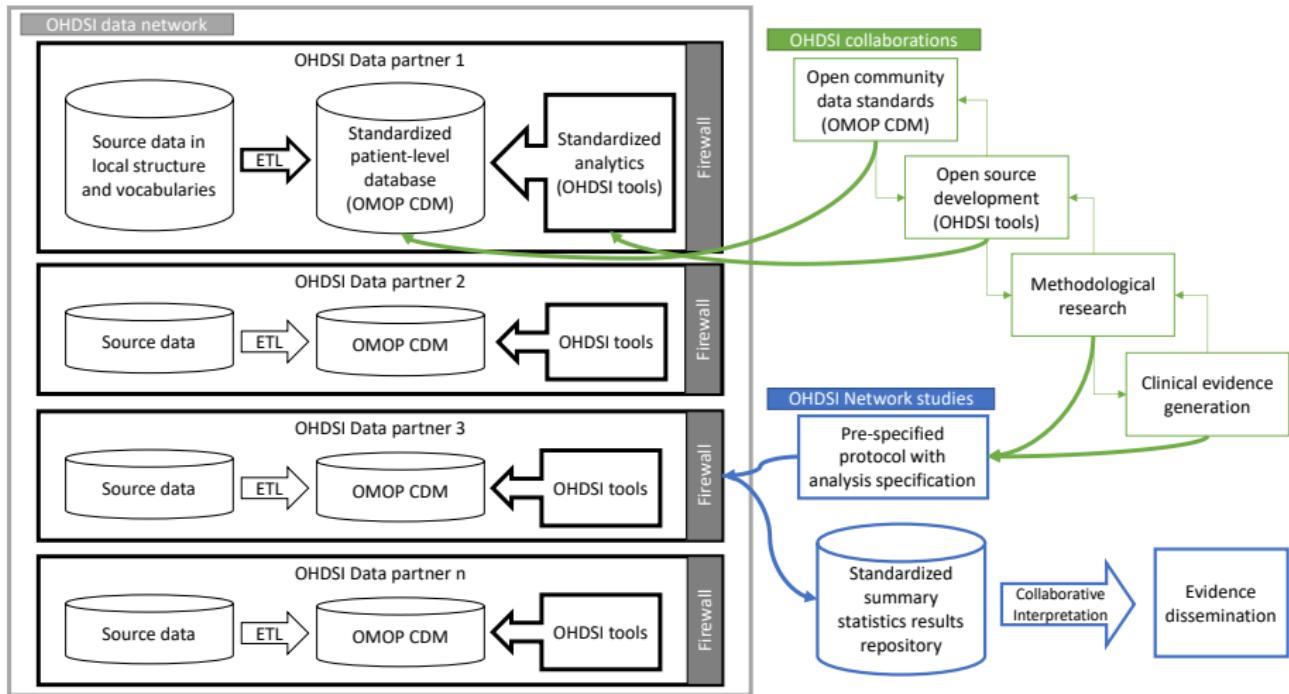
Current solution:



- Not scalable
- Expensive
- Slow
- Prohibitive to non-expert routine use

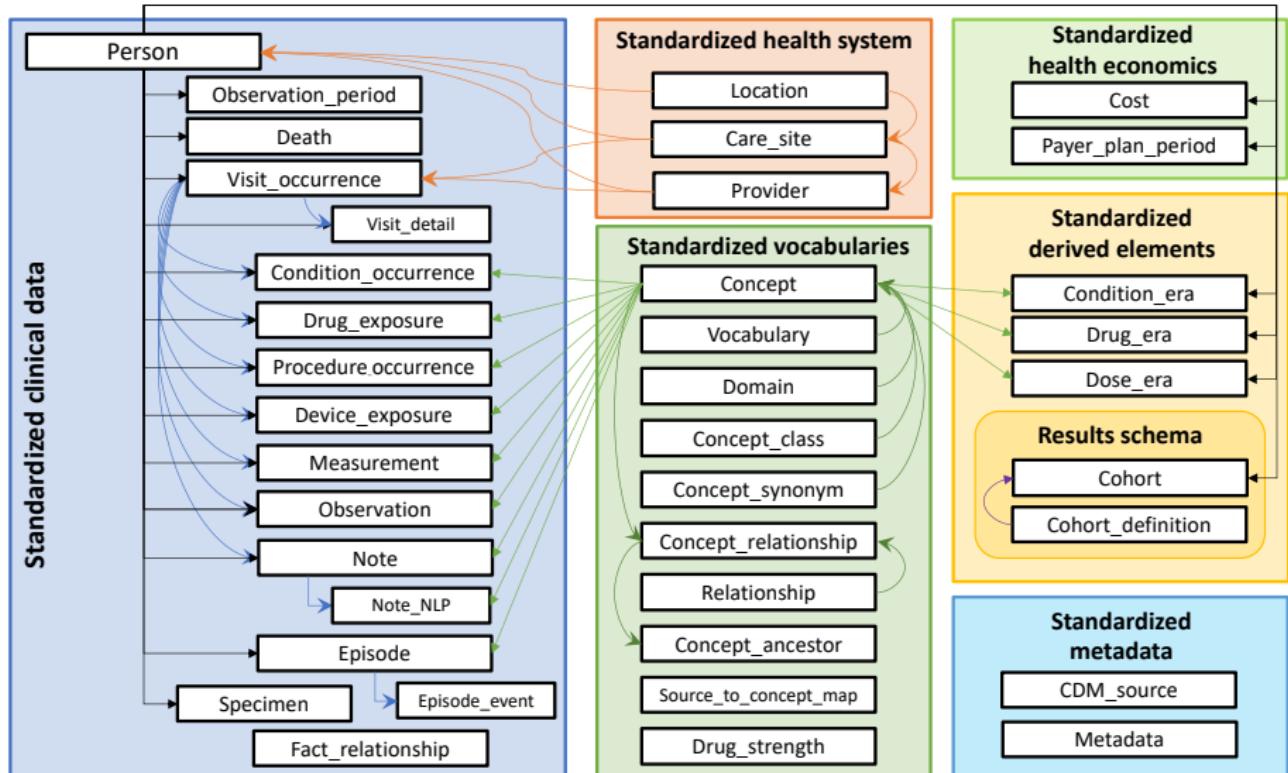


Open standards overcome the silos





OMOP common data model (CDM)





Mapping data across the global



OHDSI Vocabularies By The Numbers

as of v5.0 · Sept. 9, 2022

- 10,218,572 concepts
 - 3,549,524 standard concepts
 - 780,207 classification concepts
- 135 vocabularies
- 42 domains
- 81,243,356 concept relationships
- 85,241,004 ancestral relationships
- 3,268,183 concept synonyms

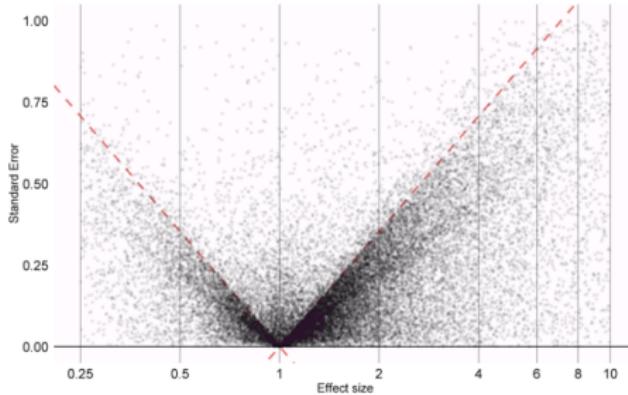
1 Shared Resource to Enable Data Standards

- Different countries (and even different hospitals in the same health network) use different “source concepts” to imply the same medical findings
 - ▶ ICD10, ICD10-CM, Snomed, RxNorm, ATC, etc.
- Need a common language (in addition to common syntax - CDM)

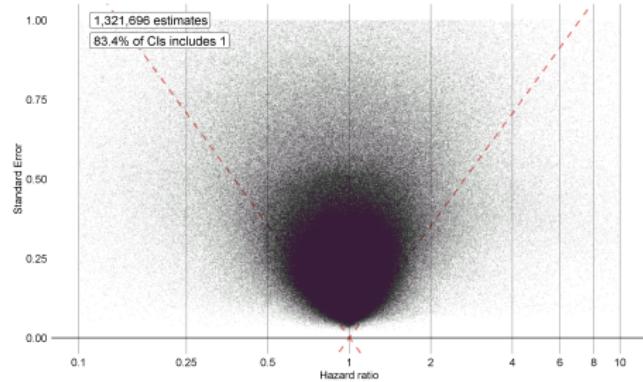


How does LEGEND perform?

Literature



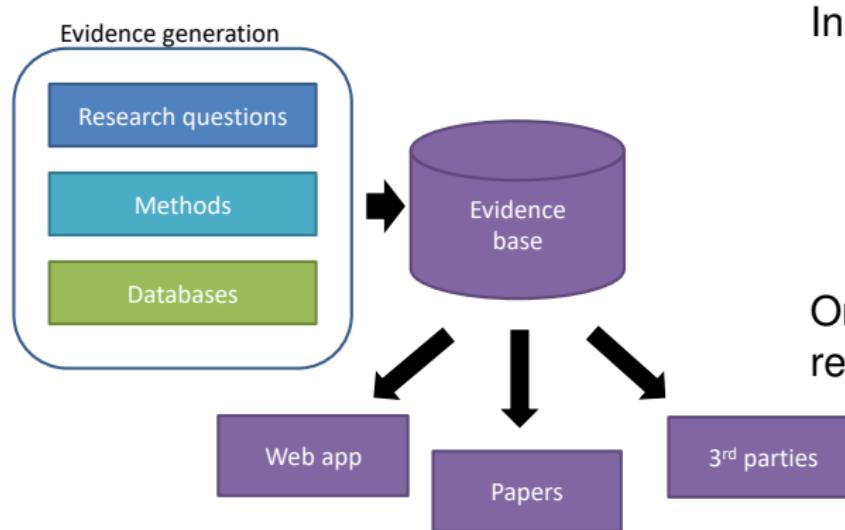
LEGEND



- Best-practices **systematic design, evaluation** and empirical **calibration** return near nominal performance
- Provide a more complete and reliable evidence basis



Unbiased LEGEND dissemination



In addition to:

- Fully specified protocol
- Open-source, end-to-end executable code

One can explore the entire result set

- <http://data.ohdsi.org/LegendBasicViewer>
(all details for each study)
- <http://data.ohdsi.org/LegendMedCentral> (**gimmick**)

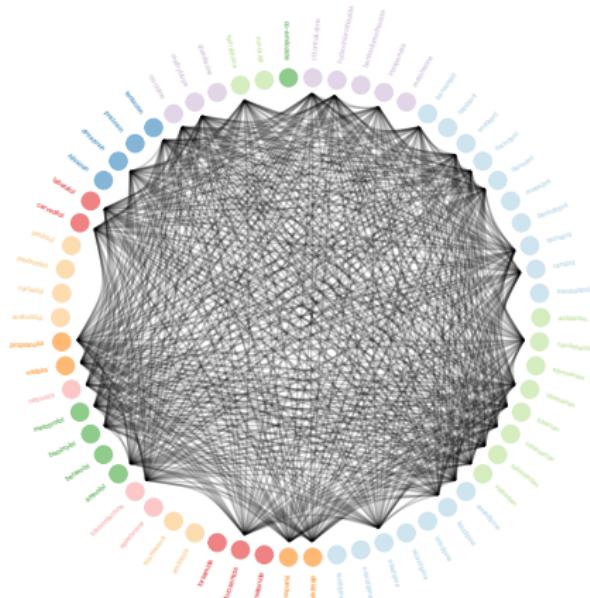


LEGEND knowledge base for hypertension

Head-to-head HTN drug comparisons



- Trials: 40
- $N = 102 - [1148] - 33K$

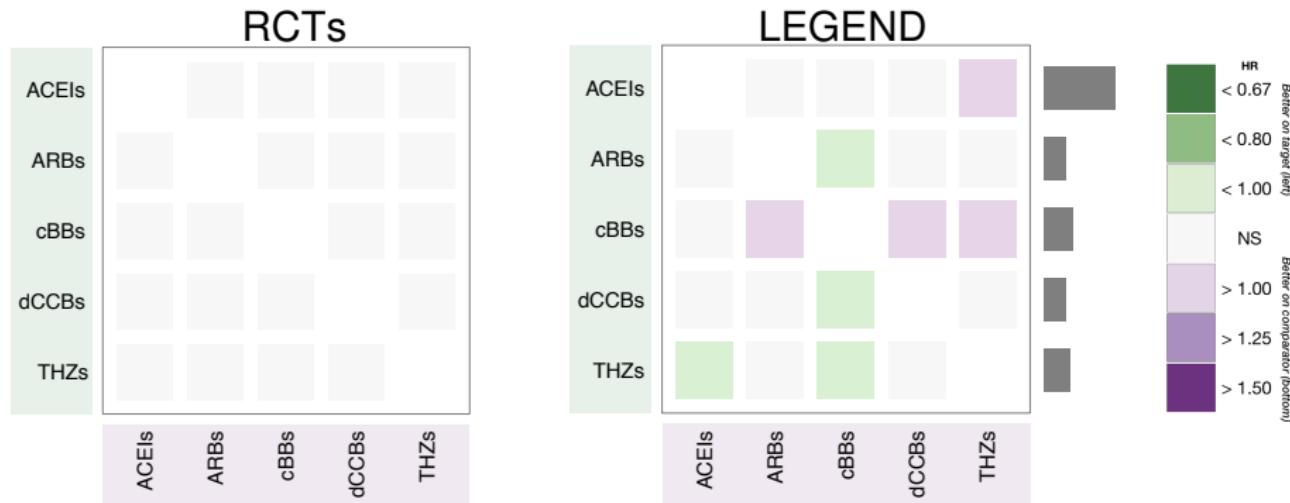


- Comparisons: 10,278
- $N = 3502 - [212K] - 1.9M$



First-line agents: comparisons from LEGEND

Efficacy outcome: **myocardial infarction**, heart failure, stroke



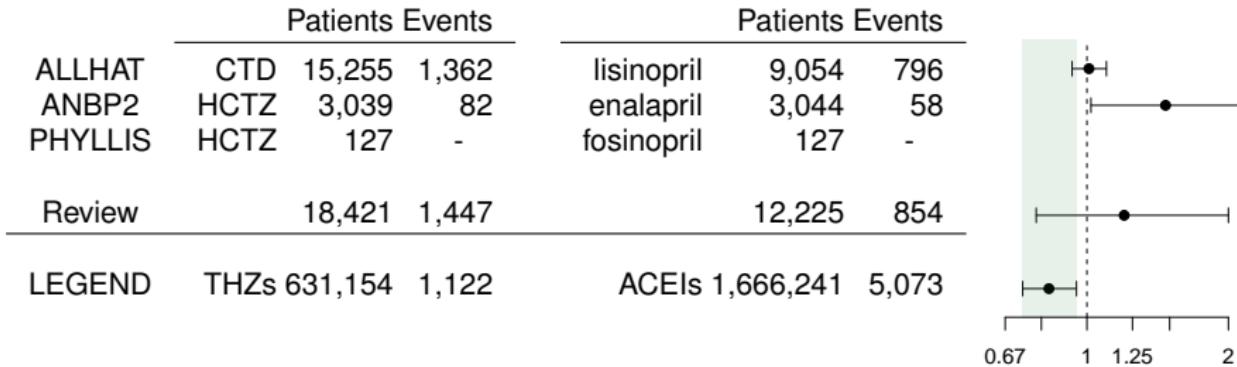
Data source: meta-analysis, ~ 1 – 2M total patients per study

- Beta blockers underperform alternatives
- Unexpected: THZs > ACEs. Reliable?



THZs vs. ACEIs: hazard ratio

Efficacy outcome: **myocardial infarction**



- THZs vs. ACEIs calibrated HR: 0.83, 95% CI 0.73 - 0.95, $p = 0.01$
- Is concordant with systematic review, but shows effect difference
- Suspect previous studies lack power



THZs vs. ACEIs: study diagnostics

Large-scale propensity score model controls for observed confounding

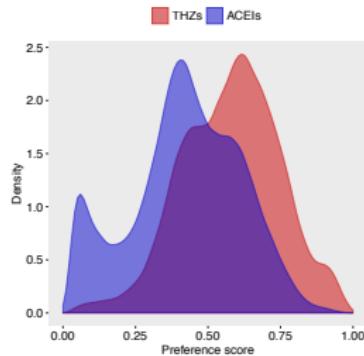


Fig. 2. Preference score distribution for TZDs and ACEIs new-users. The preference score is a transformation of the propensity score that adjusts for prevalence differences between populations. A higher overlap indicates that subjects in the two populations are more similar in terms of their predicted probability of receiving one treatment over the other.

Cohort stratification / balance:

- Achieved across all 10,868 baseline characteristics (CCAE)
- Blood pressure (pop. means in mmHg) (Panther)

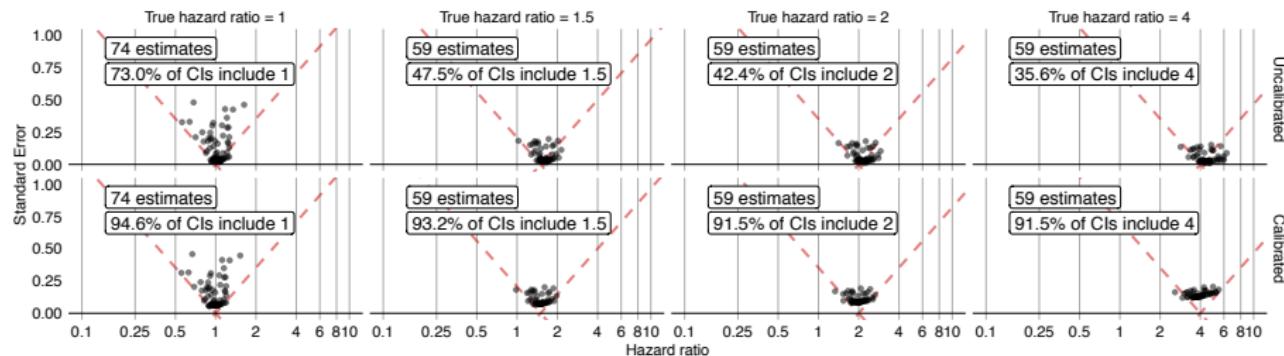
	THZs	ACEIs	$ \Delta $
before	145/89	145/87	0.13
after	145/88	145/87	0.02

No BP measurements used in PS model,
but still balanced after stratification



THZs vs. ACEIs: study outcomes

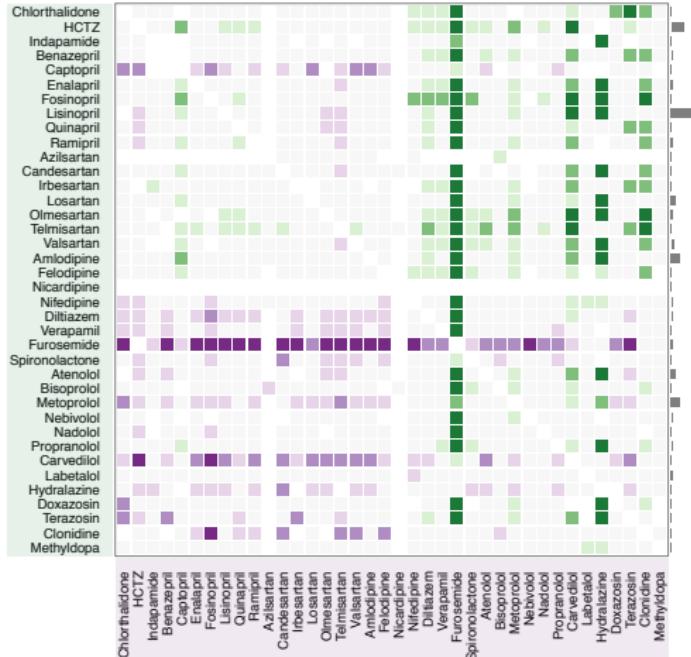
Calibration returns near nominal HR estimate coverage



- Good diagnostics → comparable cohorts (observed and unobserved); calibration → controls for residual systematic bias
- THZs are more effective than ACEIs in preventing MI
 - Could eliminate 1.3 major CV events / 1000 people



Cardiovascular efficacy by drug



Composite (MI, HF, stroke) outcome in meta-analysis

Prescriptions are not written at the class-level; must choose an individual drug for the patient

- 1st-line > 2nd-line
- Some within-class differences failed diagnostics, e.g. captoril



Chlorthalidone vs. hydrochlorothiazide

2017 ACC/AHA Guidelines

TABLE 18 Oral Antihypertensive Drugs

Class	Drug	Usual Dose, Range (mg/d)*	Daily Frequency	Comments
Primary agents				
Thiazide or thiazide-type diuretics	Chlorthalidone	12.5-25	1	■ Chlorthalidone is preferred on the basis of prolonged half-life and proven trial reduction of CVD.
	Hydrochlorothiazide	25-50	1	■ Monitor for hyponatremia and hypokalemia, uric acid and calcium levels.
	Indapamide	1.25-2.5	1	■ Use with caution in patients with history of acute gout unless patient is on uric acid-lowering therapy.
	Metolazone	2.5-5	1	

Diuretic Comparison Project (2022)

VA CSP Study No. 597: Diuretic Comparison Project

| Diuretic Comparison Project Home Page | [Information for Veterans](#) | [Information for VA Providers](#) | [Study Team and Contact Information](#) |

DCP is a national, voluntary research study funded by the [VA Cooperative Studies Program](#) (CSP), with the Department of Veterans Affairs Office of Research and Development. The goal of the DCP study is to compare the benefits of two commonly used medications, using an innovative study design. DCP uses Point of Care (POC) methodology, which embeds as much of the study procedures as possible into the routine medical care of Veterans. The result is a streamlined and efficient trial that follows clinical practice, thereby enhancing the ability to learn the answers to important questions directly within the VA healthcare system, thus supporting the goal for VA to be a Learning Healthcare System.



CTD vs. HCTZ: cardiovascular events

Figure 1. Comparability of the Populations for Commercial Claims and Encounters Database (CCAE)

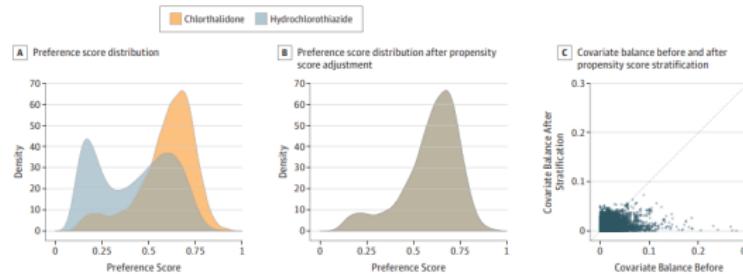
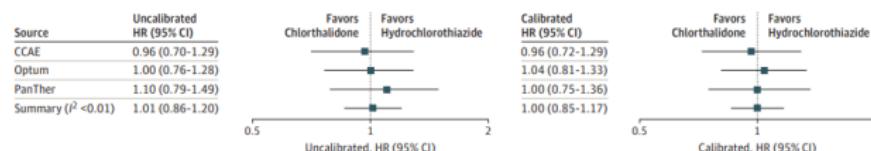
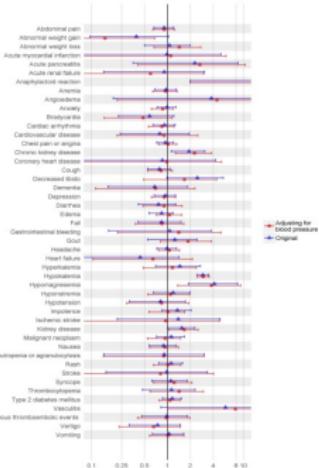


Figure 2. Homogeneity on Effectiveness



Hazard ratios (HRs) and forest plot of the 3 databases and the meta-analysis for chlorothalidone vs hydrochlorothiazide on the composite cardiovascular disease outcome. The 3 databases showed excellent agreement. CCAE indicates Commercial Claims and Encounters Database.

Figure 3. Sensitivity to balancing on baseline blood pressure in the PanTher database. We show effectiveness and safety outcomes for the PanTher database for propensity models that exclude (blue triangle) and include (red circle) baseline systolic and diastolic blood pressure in the propensity model. There are no major shifts in outcome.

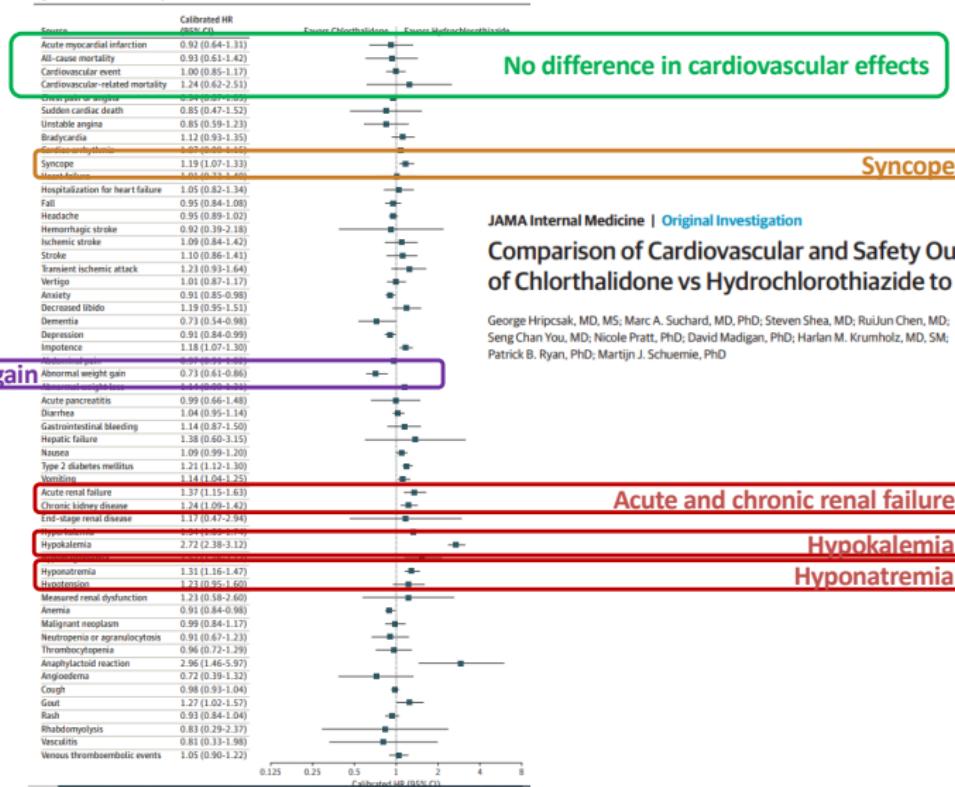


Hripcak et al. (2020) *JAMA Internal Medicine*



CTD vs. HCTZ: safety profile

Figure 3. Forest Plot of Safety and Effectiveness Outcomes



JAMA Internal Medicine | Original Investigation

Comparison of Cardiovascular and Safety Outcomes of Chlorthalidone vs Hydrochlorothiazide to Treat Hypertension

George Hripcak, MD, MS; Marc A. Suchard, MD, PhD; Steven Shea, MD; RuiJun Chen, MD; Seng Chan You, MD; Nicole Pratt, PhD; David Madigan, PhD; Harlan M. Krumholz, MD, SM; Patrick B. Ryan, PhD; Martin J. Schuemie, PhD

17 Feb 2020

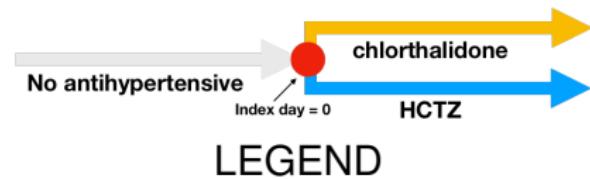
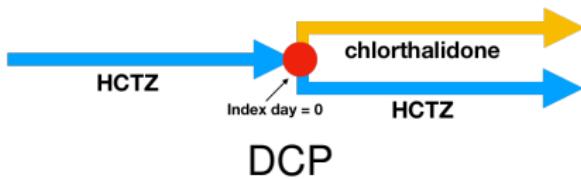


CTD vs. HCTZ: prediction \$10M = null

	DCP-RCT	LEGEND
CTD-users	6,750	25,566
HCTZ-users	6,750	528,202
avg/median follow-up	3 years	2 years
maximum follow-up	4.5 years	16 years
HR for CVE	-	1.00 (0.85 - 1.17) NS

- 25% of LEGEND patients have follow-up > 4 years

Caveat: design differences (similar vis-a-vis ALLHAT)





A dismissive rebuttal

Comment & Response

June 22, 2020

Chlorthalidone and Hydrochlorothiazide for Treatment of Patients With Hypertension

Andrew E. Moran, MD, MPH^{1,2}; Paul K. Whelton, MD, MSc³; Thomas R. Frieden, MD, MPH¹

Chlorthalidone and Hydrochlorothiazide for Treatment of Patients With Hypertension

To the Editor Hripcak et al¹ compared cardiovascular and safety outcomes of chlorthalidone and hydrochlorothiazide in the treatment of patients with hypertension. Chlorthalidone is recommended over hydrochlorothiazide because it has a longer duration of effect (24 vs 6-12 hours) and has been more extensively documented as effective in randomized clinical trials to reduce cardiovascular events and mortality.² Prior meta-analyses and observational comparisons suggest that chlorthalidone is superior in preventing cardiovascular events.^{3,4} However, to our knowledge there are no published randomized trials comparing chlorthalidone and hydrochlorothiazide; such a trial is ongoing in the US Veterans Affairs system, with results expected in 2023.⁵

Moderately strong prior evidence suggests the superiority of chlorthalidone over hydrochlorothiazide, and there is substantial likelihood that residual confounding accounts for the lack of an observed difference in cardiovascular end points in the Hripcak et al¹ study. For this reason, it is imperative to await the more definitive VA trial results in 2023⁵ before changing clinical practice recommendations on diuretic choice.

Andrew E. Moran, MD, MPH

Paul K. Whelton, MD, MSc

Thomas R. Frieden, MD, MPH



DCT reported out in November 2022

All primary outcomes and hypokalemia (only DCT-reported safety outcome) match LEGEND

	OHDSI's LEGEND in 2018/2020	Diuretic Comparison Project RCT in 2022
Cardiovascular events	1.00 (0.85-1.17)	1.04 (0.94-1.16)
Hospitalization for AMI	0.92 (0.64-1.31)	1.01 (0.80-1.28)
Hospitalization for Stroke	1.10 (0.86-1.41)	1.00 (0.74-1.36)
Hospitalization for Heart failure	1.05 (0.82-1.34)	1.04 (0.87-1.25)
Hypokalemia	2.72 (2.38-3.12)	p<0.001

- DCT published *NEJM*



How well does the collaborative model work?

Medicine:

- *JAMA Internal Medicine*
- *Hypertension* (x2)
- *The Lancet*

Clinical impact:

- Cited in **UpToDate**

Data sciences:

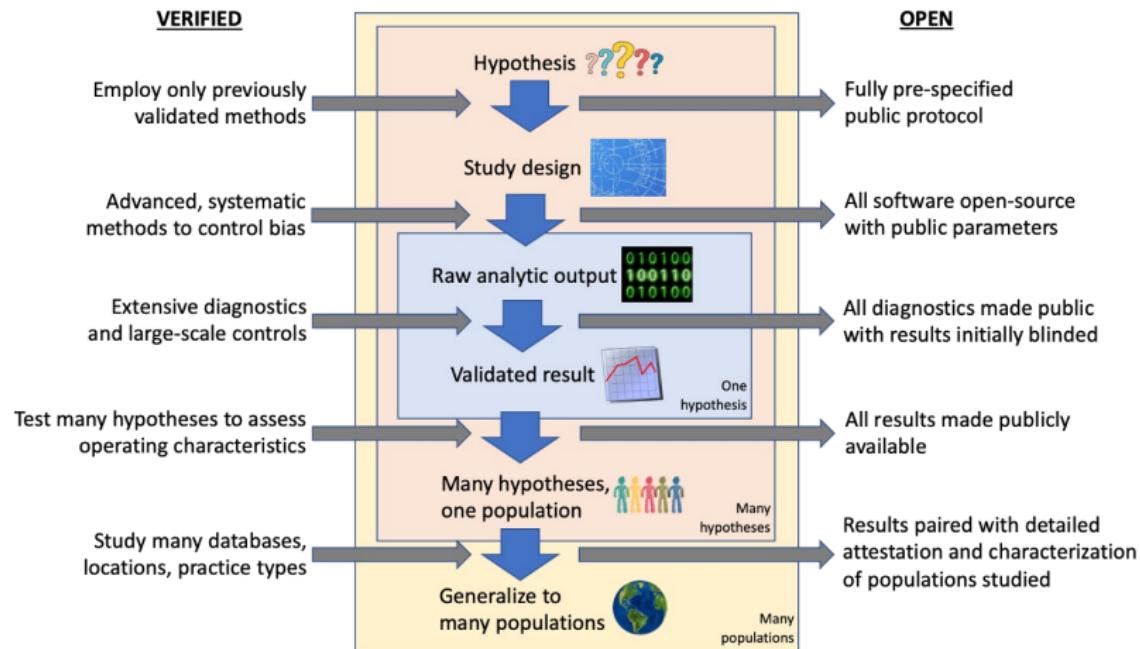
- *Bayesian Analysis*
- *JAMIA* (x2)
- *JASA*
- *Int J Epidemiol*
- *Stat Methods Med Res*

More honesty: LHC model is hard – competing goals:
research vs reliable / relevant



The BIG (legendary) open-science picture

LEGEND is ...



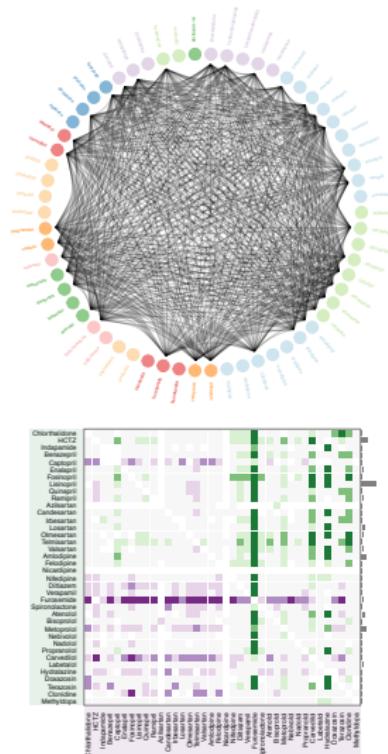


The BIG (legendary) picture

LEGEND is ...

- aimed at generating reliable, rich and relevant evidence
 - smashing silos across an extensive international network of databases
 - broadening into many other clinical conditions
 - a new way of doing (data) science

...please join the OHDSI journey!





Acknowledgments

LEGEND-HTN Scientific Group:

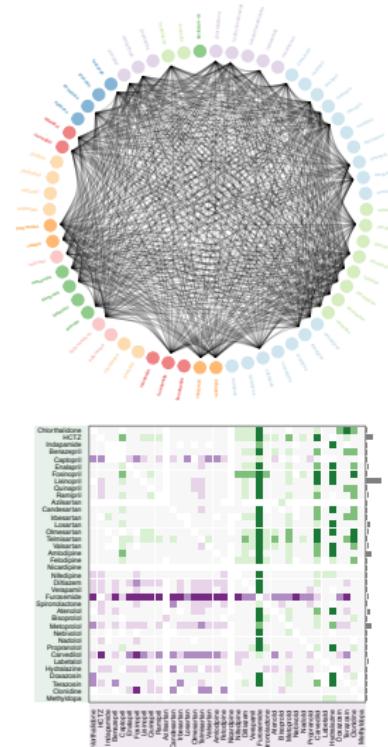
- Martijn J. Schuemie
- Patrick B. Ryan
- Seng Chan You
- Nicole Pratt
- David Madigan
- George Hripcsak
- Marc A. Suchard

Clinical Advisory Team:

- RuiJun Chen
- Jon Duke
- Harlan Krumholz
- Christian Reich

Funding:

- NSF DMS 1264153 and IIS 1251151;
NIH U19 AI135995
- Sloan Fellowship; Guggenheim Fellowship





Further reading

- Suchard, Schuemie, Krumholz et al (2019) Comprehensive comparative effectiveness and safety of first-line antihypertensive drug classes: a systemic, multinational, large-scale analysis. *The Lancet*, 394, 1816 - 1826
- Schuemie, Ryan, Pratt et al (2020) Principles of large-scale evidence generation and evaluation across a network of databases (LEGEND). *Journal of the American Medical Informatics Association*, 27, 1331 - 1337
- Rohan, Aminorroaya, Dhingra et al (2024) Comparative effectiveness of second-line antihyperglycemic agents for cardiovascular outcomes: a multinational, federated analysis of LEGEND-T2DM. *Journal of the American College of Cardiology*, 84, 904 - 917