# Observational Health Data Sciences and Informatics (OHDSI)

**George Hripcsak, MD, MS**
Columbia University Medical Center
NewYork-Presbyterian Hospital

Seattle Symposium on Health Care Data Analytics

# **Observational Health Data Sciences and Informatics** (OHDSI, as "Odyssey")

A multi-stakeholder, interdisciplinary, international collaborative with a coordinating center at Columbia University

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

Aiming for 1,000,000,000 patient data network

# OHDSI's global research community



- >140 collaborators from 20 different countries
- Experts in informatics, statistics, epidemiology, clinical sciences
- Active participation from academia, government, industry, providers
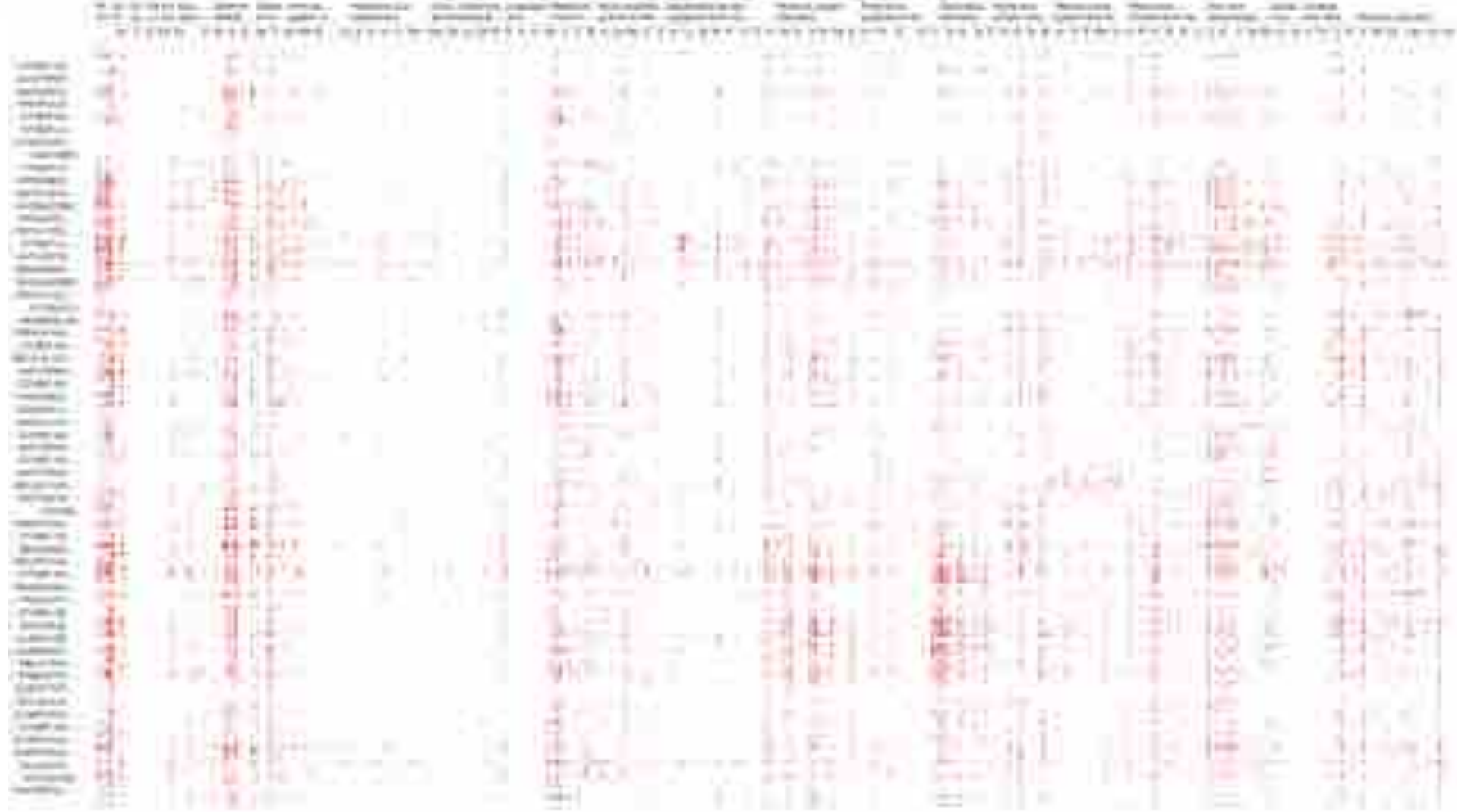- Currently 600 million patient records in 52 databases

http://ohdsi.org/who-we-are/collaborators/

# Why large-scale analysis is needed in healthcare

All health outcomes of interest



All drugs

# Patient-level predictions for personalized evidence requires big data

2 million patients seem excessive or unnecessary?

- Imagine a provider wants to compare her patient with other patients with the same gender (50%), in the same 10-year age group (10%), and with the same comorbidity of Type 2 diabetes (5%)

- Imagine the patient is concerned about the risk of ketoacidosis (0.5%) associated with two alternative treatments they are considering

- With 2 million patients, you'd only expect to observe 25 similar patients with the event, and would only be powered to observe a relative risk > 2.0

Aggregated data across a health system of 1,000 providers may contain 2,000,000 patients
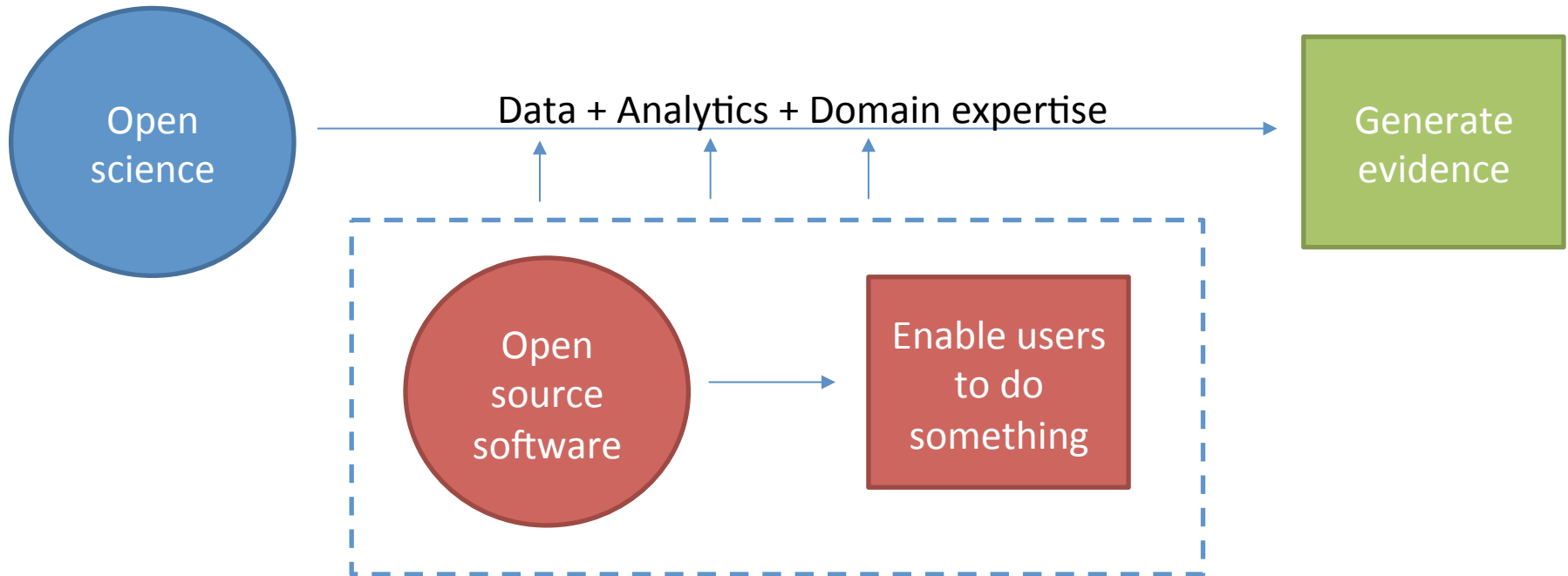
# Evidence OHDSI seeks to generate from observational data

- **Clinical characterization**
  - Natural history: Who has diabetes, and who takes metformin?
  - Quality improvement: What proportion of patients with diabetes experience complications?
- **Population-level estimation**
  - Safety surveillance: Does metformin cause lactic acidosis?
  - Comparative effectiveness: Does metformin cause lactic acidosis more than glyburide?
- **Patient-level prediction**
  - Precision medicine: Given everything you know about me, if I take metformin, what is the chance I will get lactic acidosis?
  - Disease interception: Given everything you know about me, what is the chance I will develop diabetes?
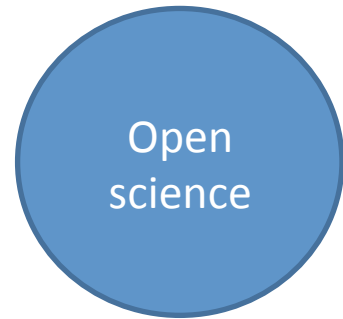
# OHDSI's approach to open science



Data + Analytics + Domain expertise

Open science → Generate evidence

Open source software → Enable users to do something

- Open science is about sharing the journey to evidence generation
- Open-source software can be part of the journey, but it's not a final destination
- Open processes can enhance the journey through improved reproducibility of research and expanded adoption of scientific best practices

# Standardizing workflows to enable transparent, reproducible research

**Open science**

Population-level estimation for comparative effectiveness research:

Is <intervention X> better than <intervention Y> in reducing the risk of <condition Z>?

**Generate evidence**

| Database summary | Cohort definition | Cohort summary | Compare cohorts | Exposure-outcome summary | Effect estimation & calibration | Compare databases |

**Defined inputs:**
- Target exposure
- Comparator group
- Outcome
- Time-at-risk
- Model specification

**Consistent outputs:**
- analysis specifications for transparency and reproducibility (protocol + source code)
- only aggregate summary statistics (no patient-level data)
- model diagnostics to evaluate accuracy
- results as evidence to be disseminated
  - static for reporting (e.g. via publication)
  - interactive for exploration (e.g. via app)

# OHDSI Distinguishing Features

- International effort (size & coverage)
  - 43 sources terminologies from around the world
- Open science (depth)
  - Infrastructure serves the science
  - Stack: Terminology, CDM, ETL, QA, Visualization, Novel analytic methods, Clinical research
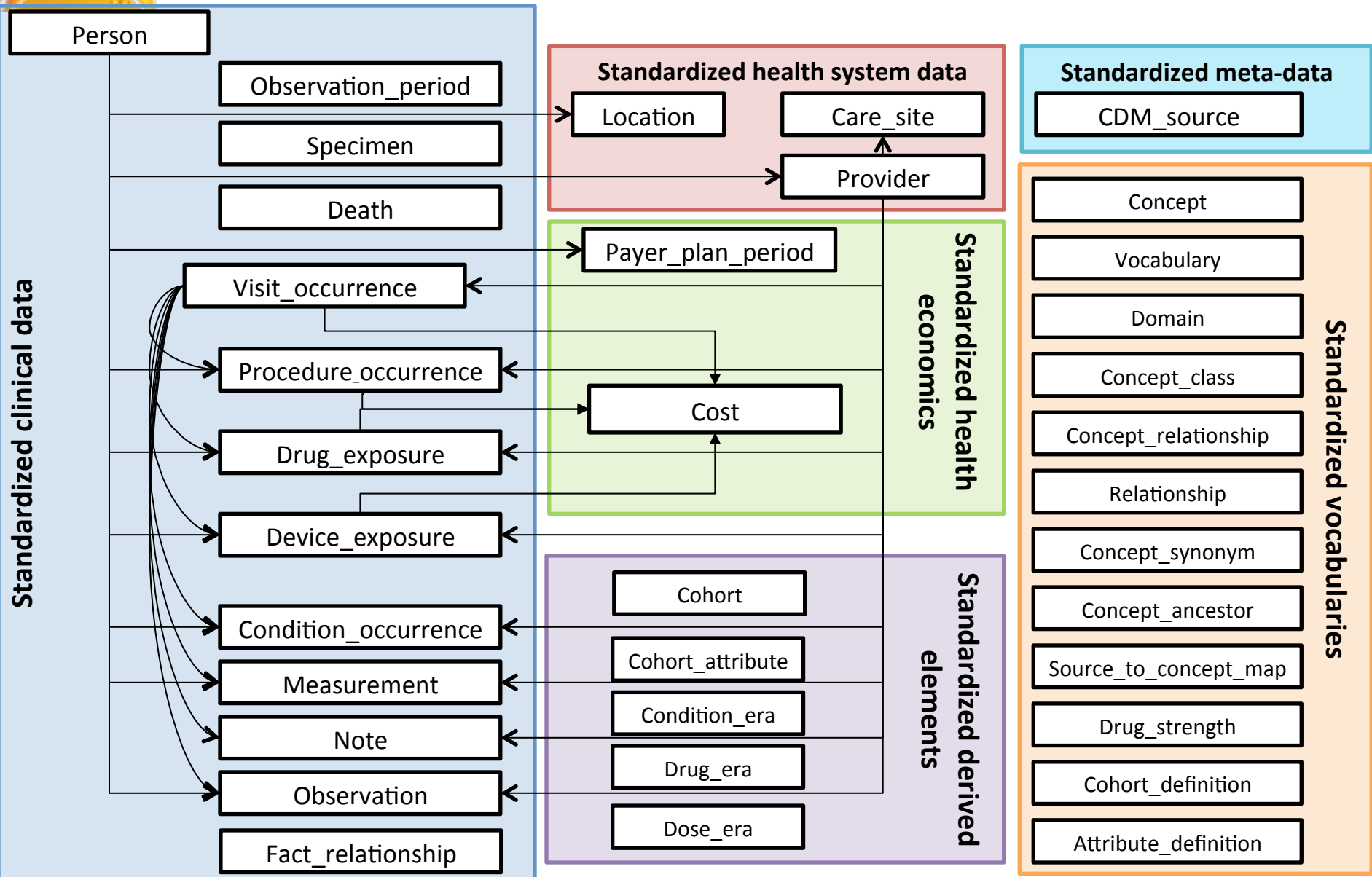- Full information model

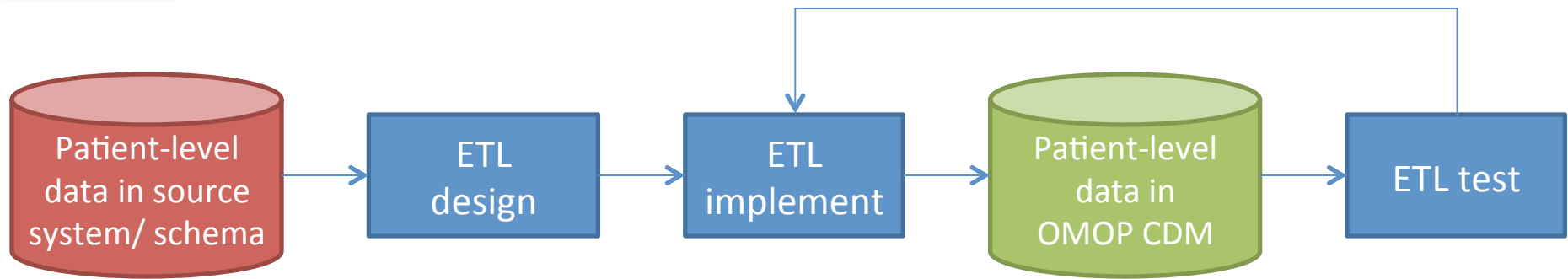# How OHDSI Works

# Deep information model
# OMOP CDM v5.0.1

# Extensive vocabularies



Breakdown of OHDSI concepts by domain, standard state, and vocabulary

# Preparing your data for analysis



**Patient-level data in source system/ schema** → **ETL design** → **ETL implement** → **Patient-level data in OMOP CDM** → **ETL test**

**OHDSI tools built to help**

**WhiteRabbit**: profile your source data

**RabbitInAHat**: map your source structure to CDM tables and fields

**ATHENA**: standardized vocabularies for all CDM domains

**Usagi**: map your source codes to CDM vocabulary

**CDM**: DDL, index, constraints for Oracle, SQL Server, PostgresSQL; Vocabulary tables with loading scripts

**ACHILLES**: profile your CDM data; review data quality assessment; explore population-level summaries

**OHDSI Forums**: Public discussions for OMOP CDM Implementers/developers

http://github.com/OHDSI

# ACHILLES Heel Data Validation

# ATLAS to build, visualize, and analyze cohorts

# Characterize the cohorts of interest

# LAERTES: Knowledge base of what we know: literature, labeling, spontaneous reporting

# OHDSI in Action

- Generate evidence
  - Randomized trial is the gold standard
  - Observational research is supporting
    - Can it become a partnership?

# Characterization

- Today we carry out RCTs without clear knowledge of actual practice
- There will be no RCTs without an observational precursor
  - It will be required to characterize a population using large-scale observational data before designing an RCT
  - Disease burden
  - Actual treatment practice
  - Time on therapy
  - Course and complication rate
  - Done now somewhat through literature and pilot studies

# Treatment Pathways

Global stakeholders

Conduits

Local stakeholders

| Public |
| --- |

Evidence

| RCT, Obs |
| --- |

| Academics |
| --- |

| Industry |
| --- |

| Regulator |
| --- |

| Social media |
| --- |

| Lay press |
| --- |

| Literature |
| --- |

| Guidelines |
| --- |

| Advertising |
| --- |

| Formulary |
| --- |

| Labels |
| --- |

| Family |
| --- |

| Patient |
| --- |

| Clinician |
| --- |

| Consultant |
| --- |

Inputs

| Indication |
| --- |

| Feasibility |
| --- |

| Cost |
| --- |

| Preference |
| --- |

# Network process

1. Join the collaborative

2. Propose a study to the open collaborative

3. Write protocol
   - http://www.ohdsi.org/web/wiki/doku.php?id=research:studies

4. Code it, run it locally, debug it (minimize others' work)

5. Publish it: https://github.com/ohdsi

6. Each node voluntarily executes on their CDM

7. Centrally share results

8. Collaboratively explore results and jointly publish findings

# OHDSI in action:
# Chronic disease treatment pathways

- Conceived at AMIA                    15Nov2014
- Protocol written, code               30Nov2014
  written and tested at 2
  sites

- Analysis submitted to                2Dec2014
  OHDSI network

- Results submitted for 7              5Dec2014
  databases

# OHDSI participating data partners

| Abbre-viation | Name | Description | Population, millions |
|---|---|---|---|
| AUSOM | Ajou University School of Medicine | South Korea; inpatient hospital EHR | 2 |
| CCAE | MarketScan Commercial Claims and Encounters | US private-payer claims | 119 |
| CPRD | UK Clinical Practice Research Datalink | UK; EHR from general practice | 11 |
| CUMC | Columbia University Medical Center | US; inpatient EHR | 4 |
| GE | GE Centricity | US; outpatient EHR | 33 |
| INPC | Regenstrief Institute, Indiana Network for Patient Care | US; integrated health exchange | 15 |
| JMDC | Japan Medical Data Center | Japan; private-payer claims | 3 |
| MDCD | MarketScan Medicaid Multi-State | US; public-payer claims | 17 |
| MDCR | MarketScan Medicare Supplemental and Coordination of Benefits | US; private and public-payer claims | 9 |
| OPTUM | Optum ClinFormatics | US; private-payer claims | 40 |
| STRIDE | Stanford Translational Research Integrated Database Environment | US; inpatient EHR | 2 |
| HKU | Hong Kong University | Hong Kong; EHR | 1 |

# Treatment pathway event flow

# Characterizing treatment pathways at scale using the OHDSI network

George Hripcsak[a,b,1], Patrick B. Ryan[c,d], Jon D. Duke[e], Nigam H. Shah[f], Rae Woong Park[g], Vojtech Huser[h], Marc A. Suchard[i,j,k], Martijn J. Schuemie[c,d], Frank J. DeFalco[c], Adler Perotte[a], Juan M. Banda[f], Christian G. Reich[l], Lisa M. Schilling[m], Michael E. Matheny[n,o], Daniella Meeker[p,q], Nicole Pratt[r], and David Madigan[s]

Observational research promises to complement experimental research by providing large, diverse populations that would be infeasible for an experiment. Observational research can test its own clinical hypotheses, and observational studies also can contribute to the design of experiments and inform the generalizability of experimental research. Understanding the diversity of populations...
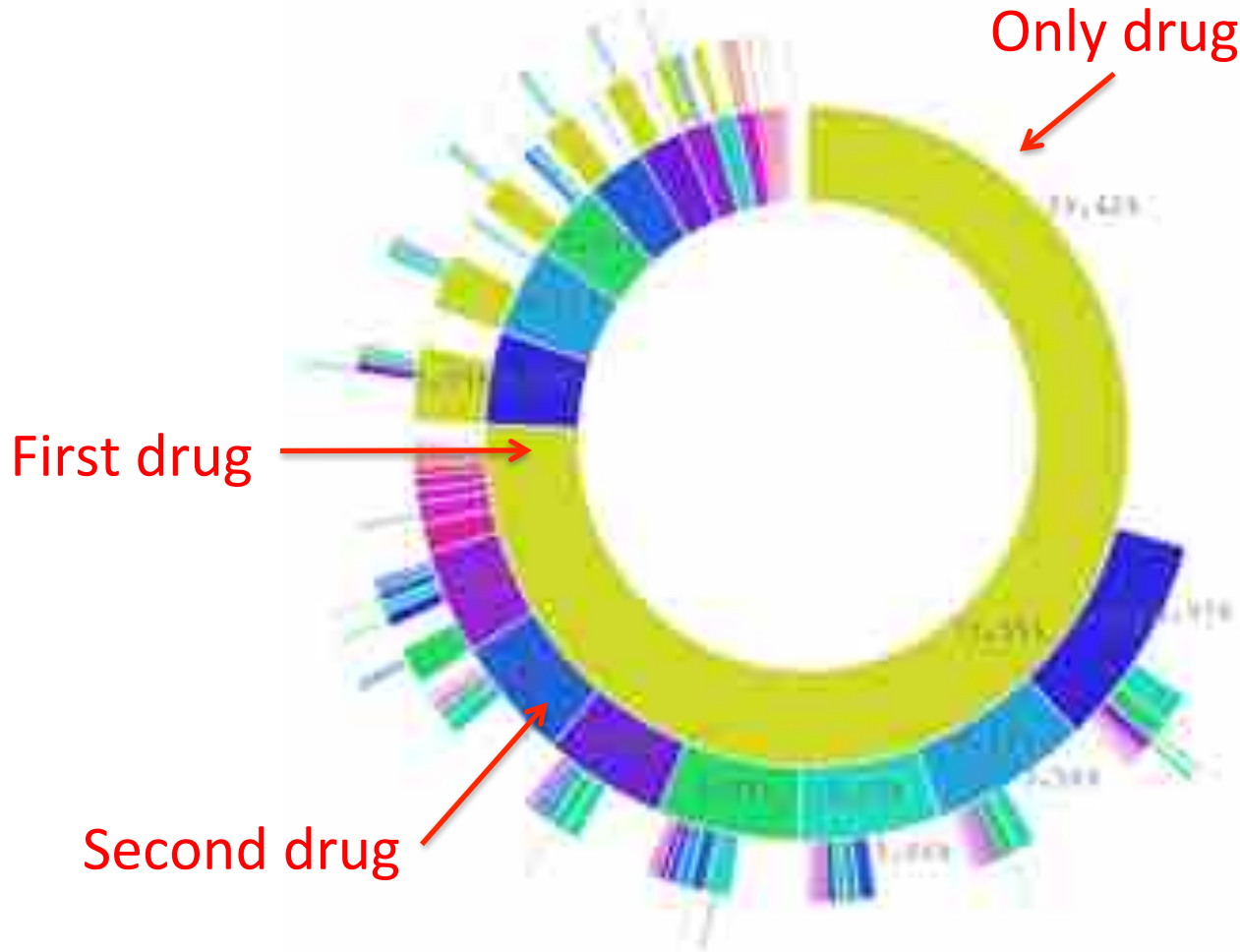
Without sufficiently broad databases available in the first stage, randomized trials are designed without explicit knowledge of actual disease states and treatment practice. Literature reviews are restricted to the population choices of previous investigations, and pilot studies usually are limited in scope. By exploiting the ClinicalTrials.gov national trial registry (9) and electronic health...
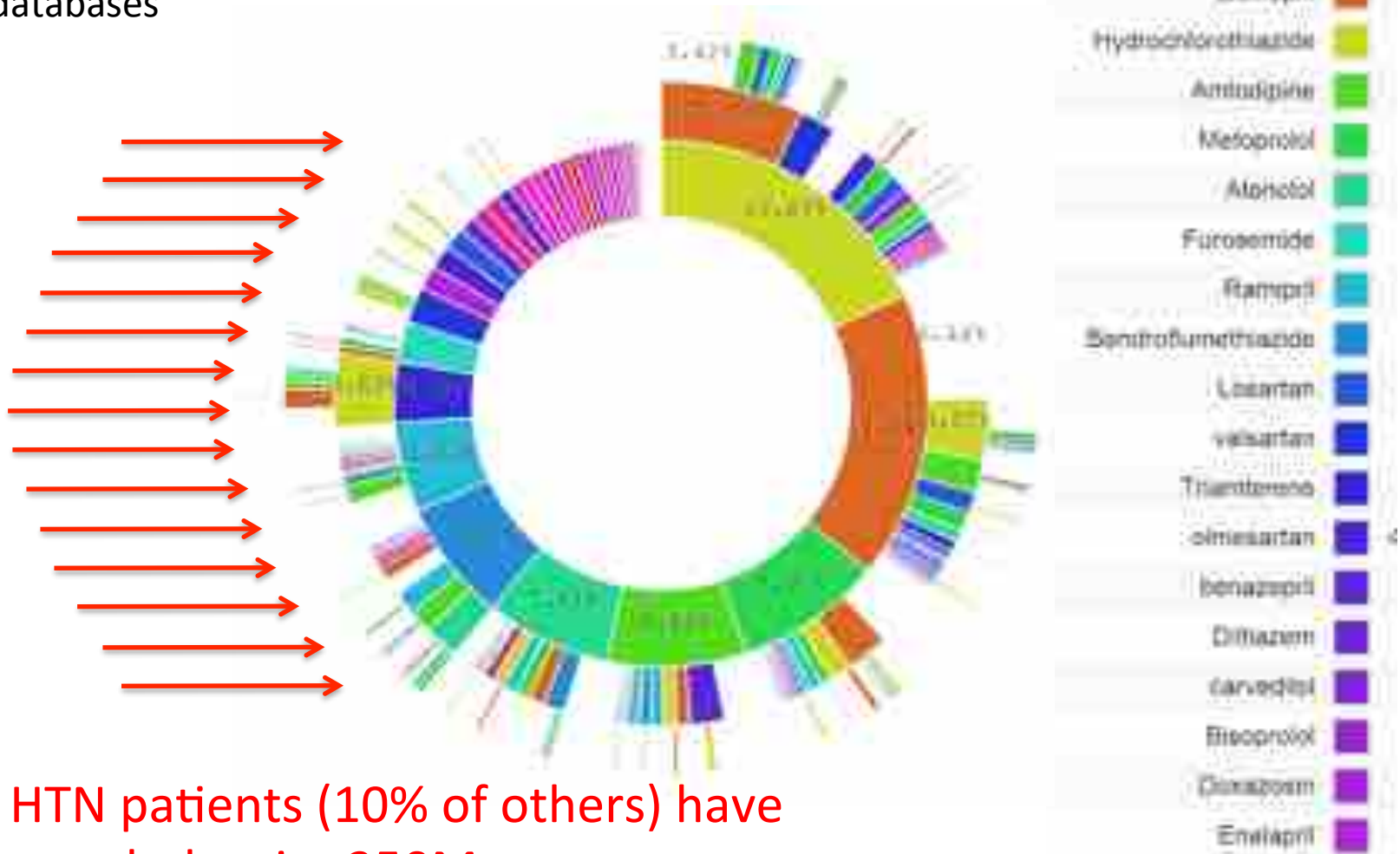
# Treatment pathways for diabetes

T2DM : All databases



Only drug

First drug

Second drug

Legend: Metformin, pioglitazone, sitagliptin, Glipzide, glimepiride, Gliclazide, Glyburide, rosiglitazone, Insulin, Glargine, Human, exenatide, Insulin, Aspart, Human, liraglutide, saxagliptin, Insulin, Lispro, Human, Glucose, insulin, isophane, Human

# Patient-level heterogeneity

HTN: All databases



25% of HTN patients (10% of others) have
a unique path despite 250M pop

# Monotherapy – diabetes

General upward trend in monotherapy

# Monotherapy – HTN

Academic medical centers differ from general practices



Legend:
- AUSOM (SKorea*)
- CCAE (US#)
- CPRD (UK*)
- CUMC (US*)
- GE (US*)
- INPC (US*#)
- JMDC (Japan#)
- MDCD (US#)
- MDCR (US#)
- OPTUM (US#)
- STRIDE (US*)
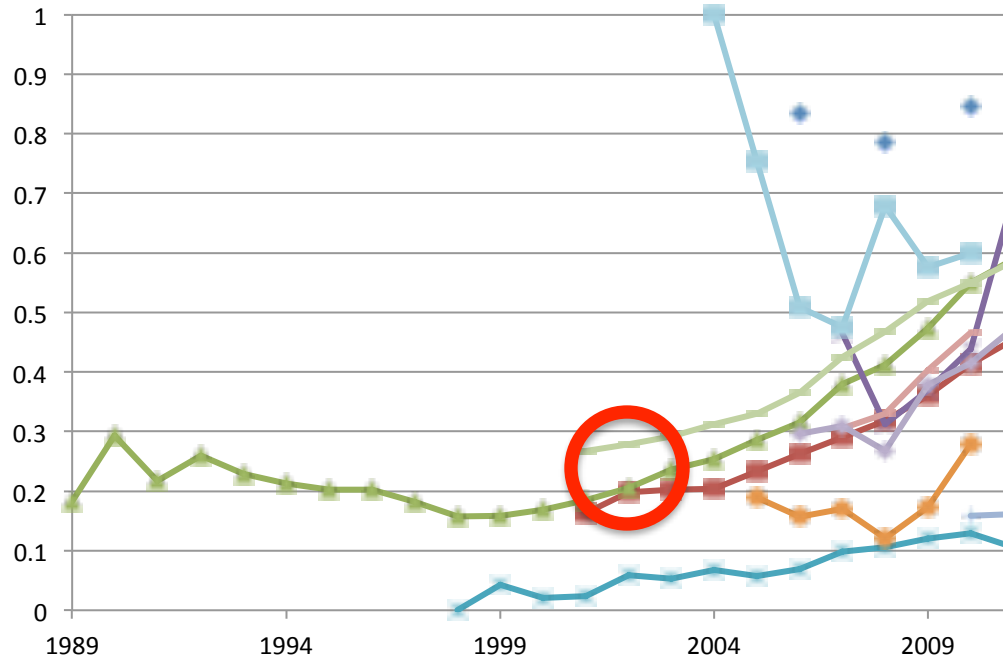
# Monotherapy – diabetes

General practices, whether EHR or claims, have similar profiles



Legend:
- AUSOM (SKorea*)
- CCAE (US#)
- CPRD (UK*)
- CUMC (US*)
- GE (US*)
- INPC (US*#)
- JMDC (Japan#)
- MDCD (US#)
- MDCR (US#)
- OPTUM (US#)
- STRIDE (US*)

# Conclusions: Network research

- It is feasible to encode the world population in a single data model
  - Over 600,000,000 records by voluntary effort (682,000,000)
- Generating evidence is feasible
- Stakeholders willing to share results
- Able to accommodate vast differences in privacy and research regulation