

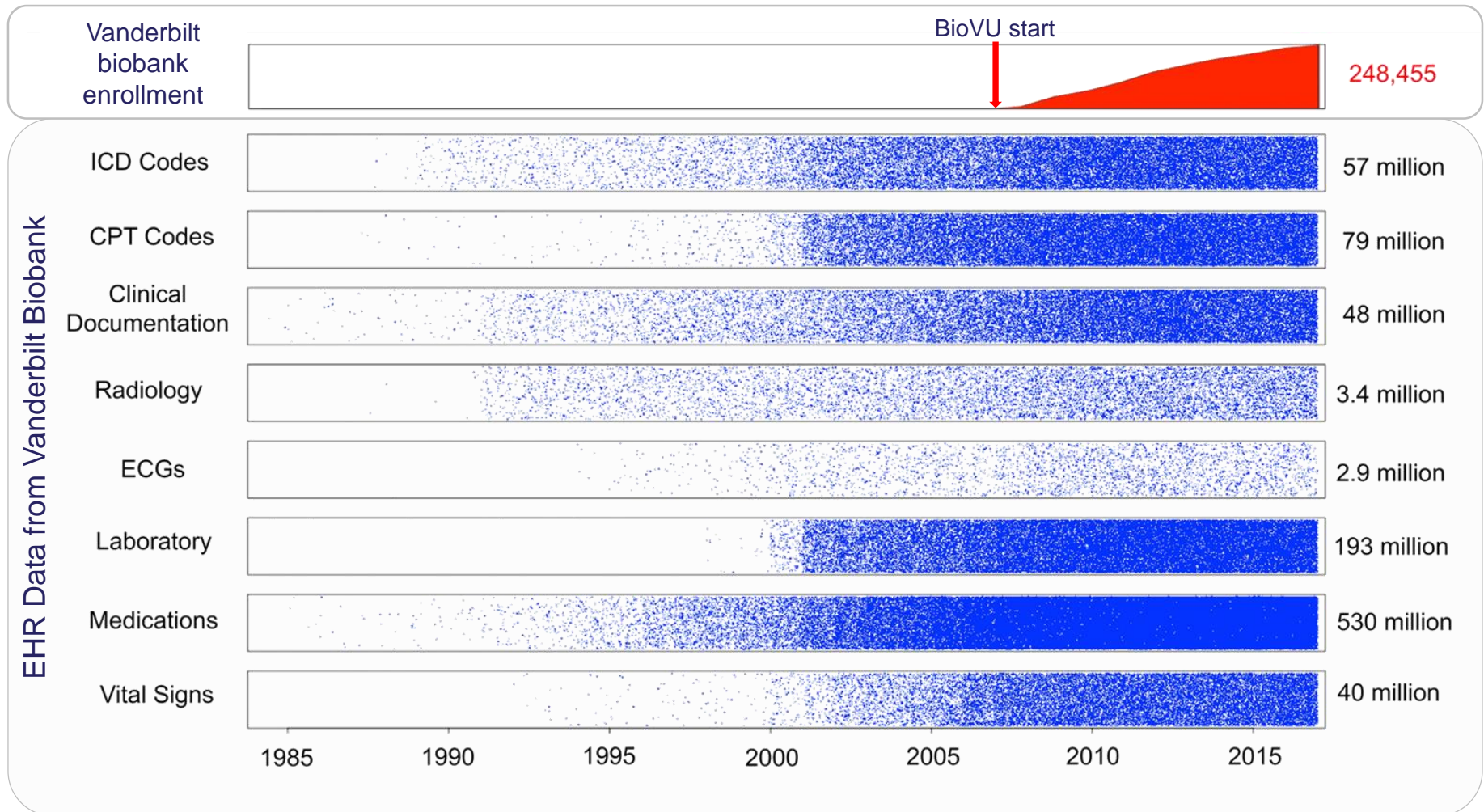
Obtaining phenotype and outcome data from EHRs

Josh Denny, MD MS

Vanderbilt University Medical Center

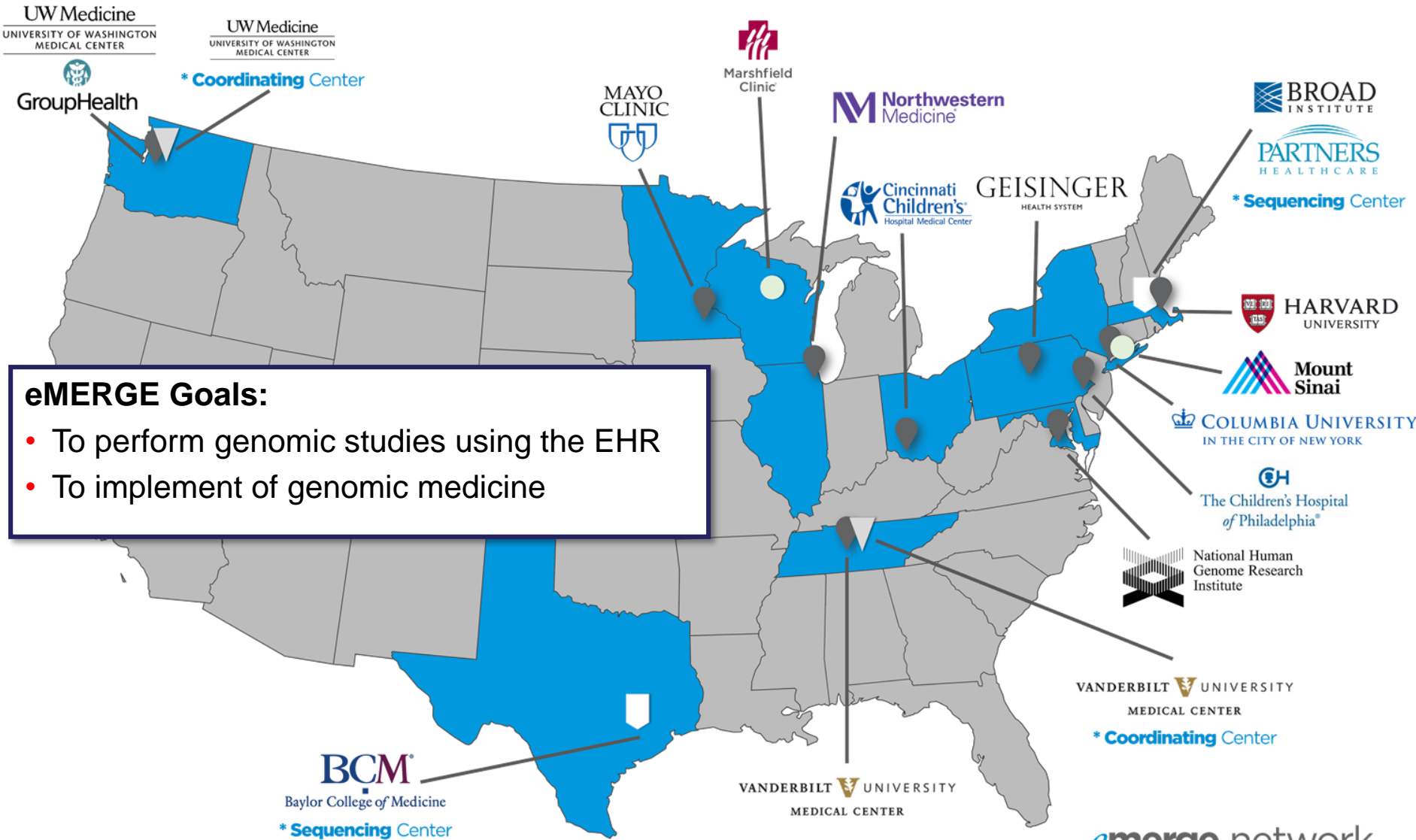
3/26/2018

EHR data are dense and efficient for discovery: Vanderbilt's experience (BioVU)



emerge network

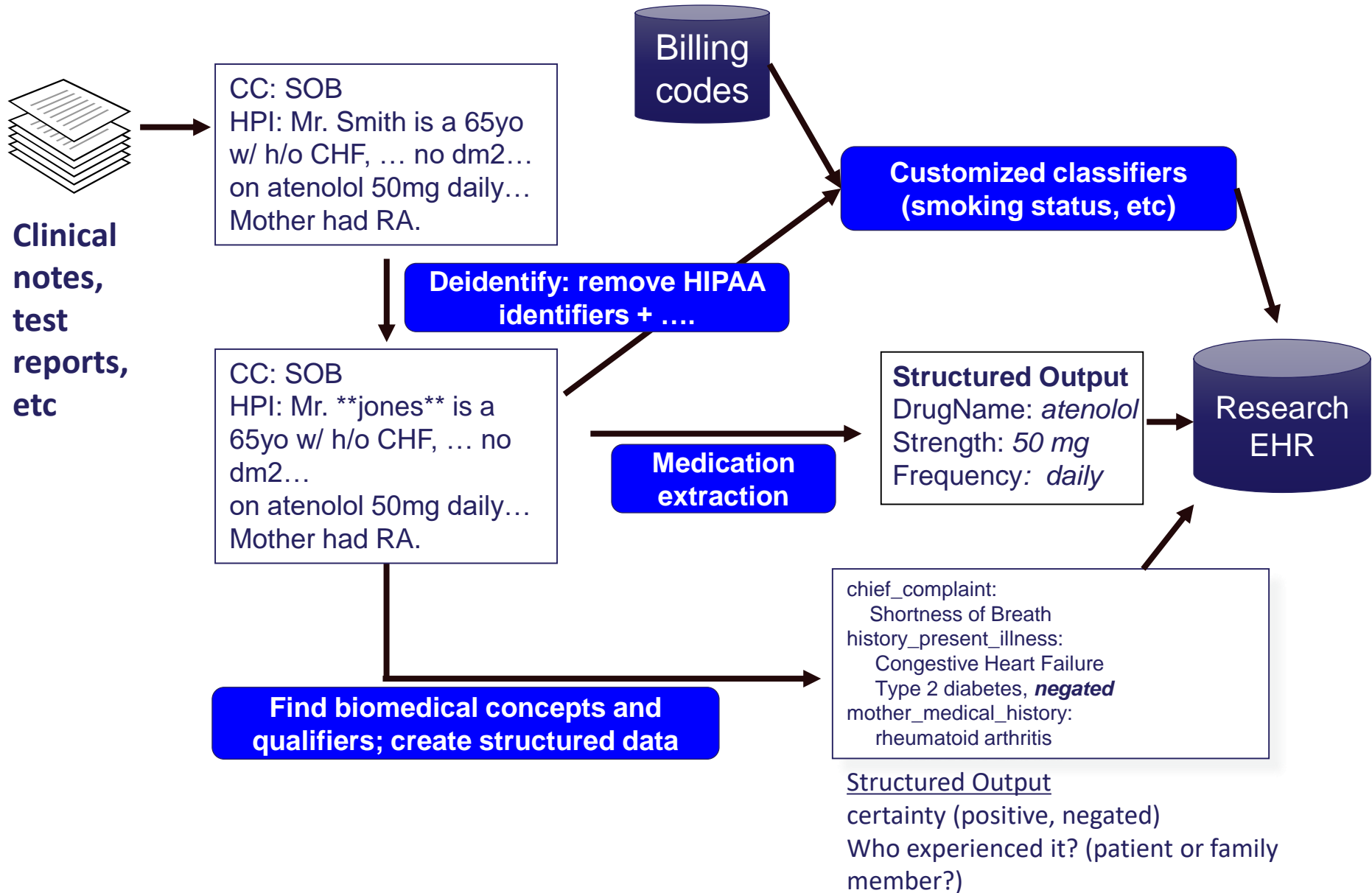
ELECTRONIC MEDICAL RECORDS & GENOMICS



eMERGE Goals:

- To perform genomic studies using the EHR
- To implement of genomic medicine

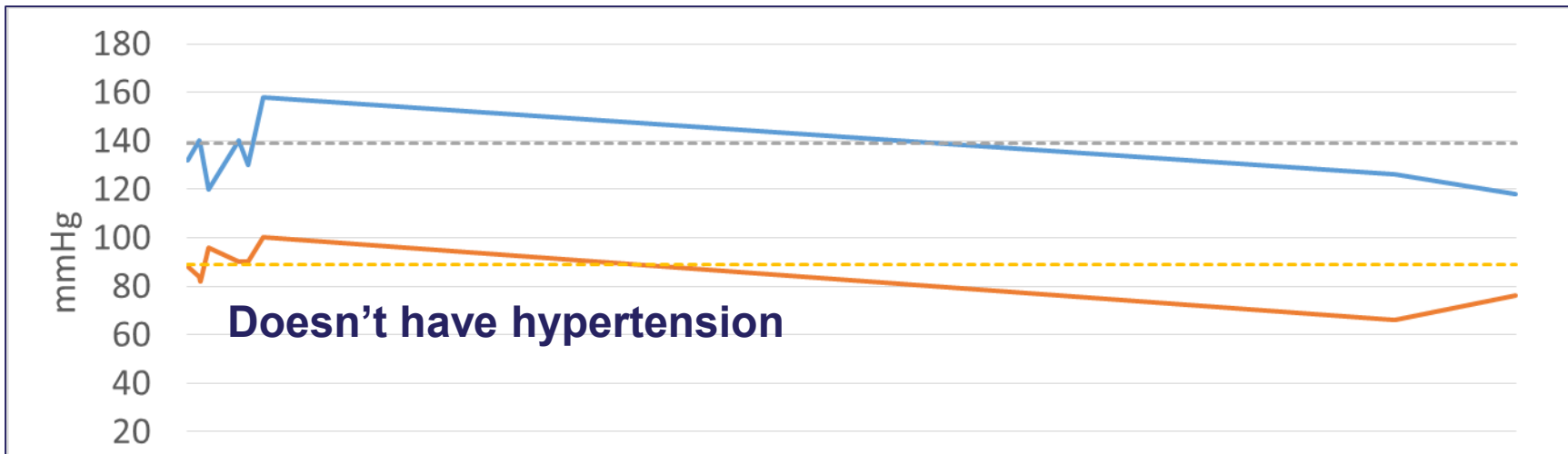
Making text documents useful for research



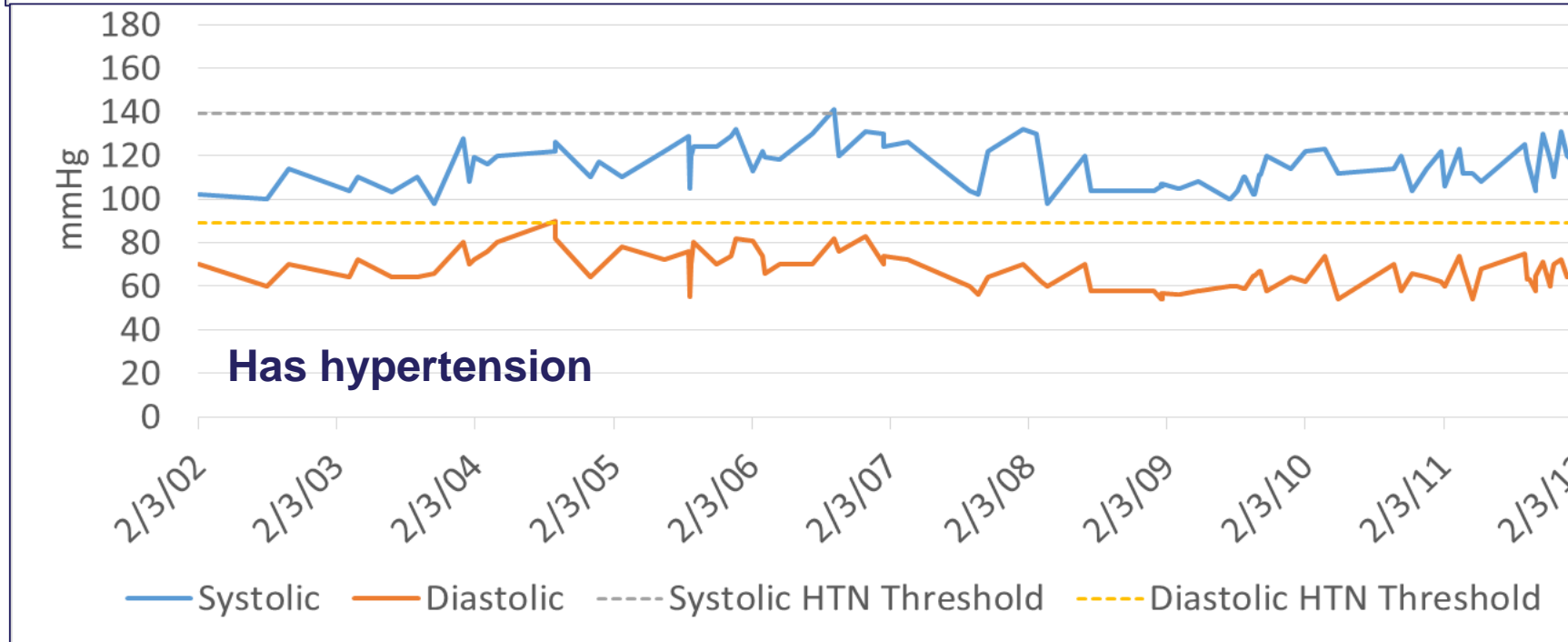
Finding a "simple" disease in the EHR: Who has hypertension?

Definition: SBP > 140 or DBP > 90

Patient 1



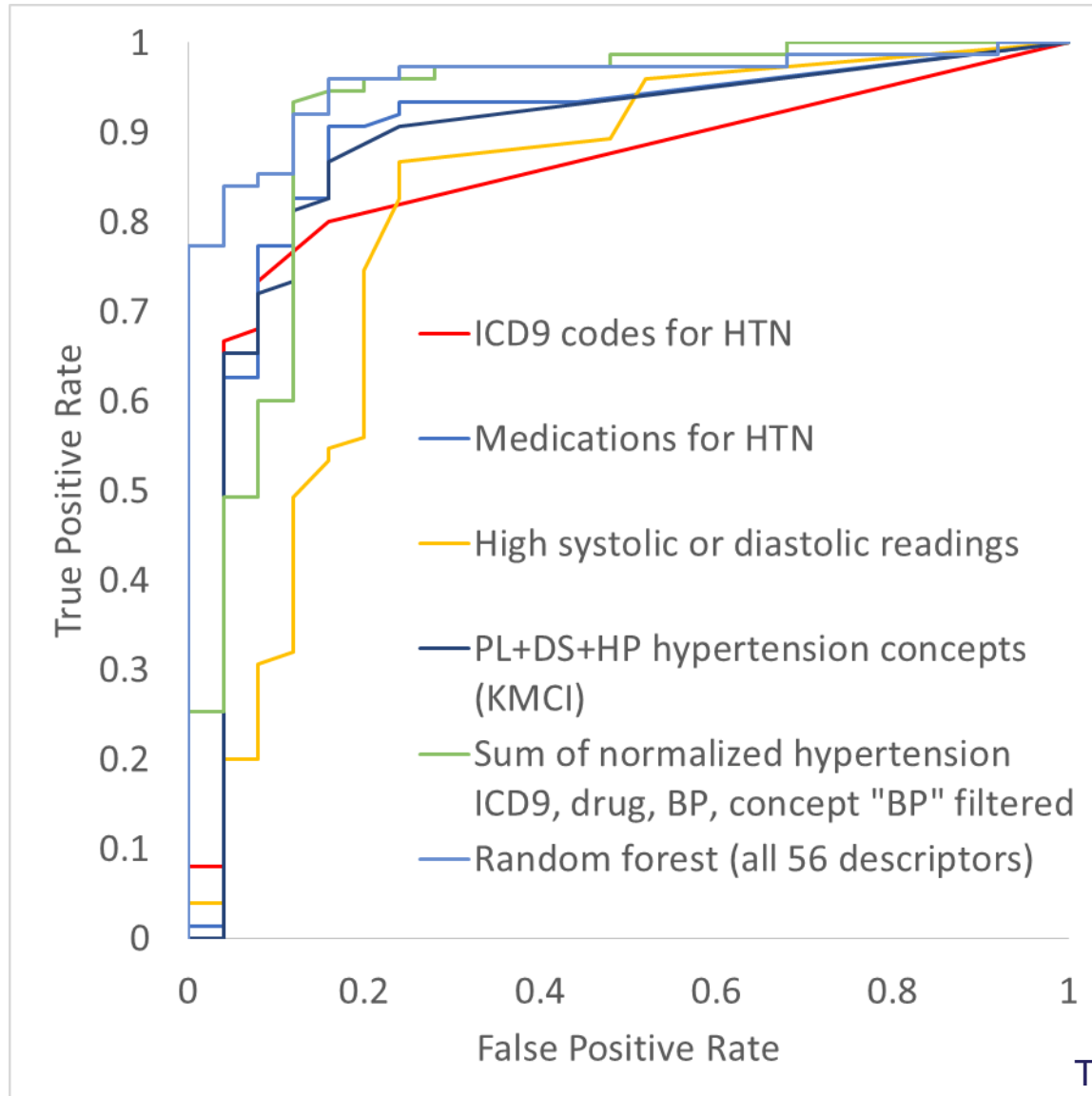
Patient 2



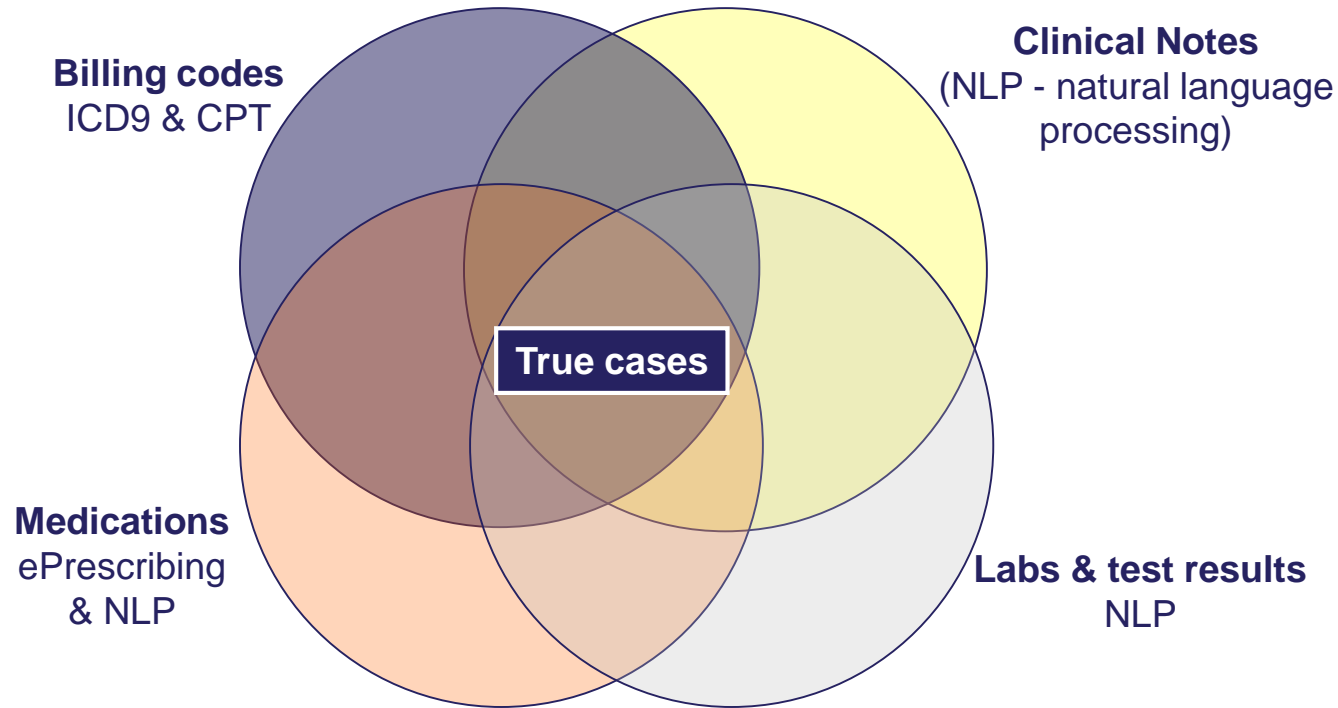
— Systolic — Diastolic - - - Systolic HTN Threshold - - - Diastolic HTN Threshold

Our “simple” example: Hypertension

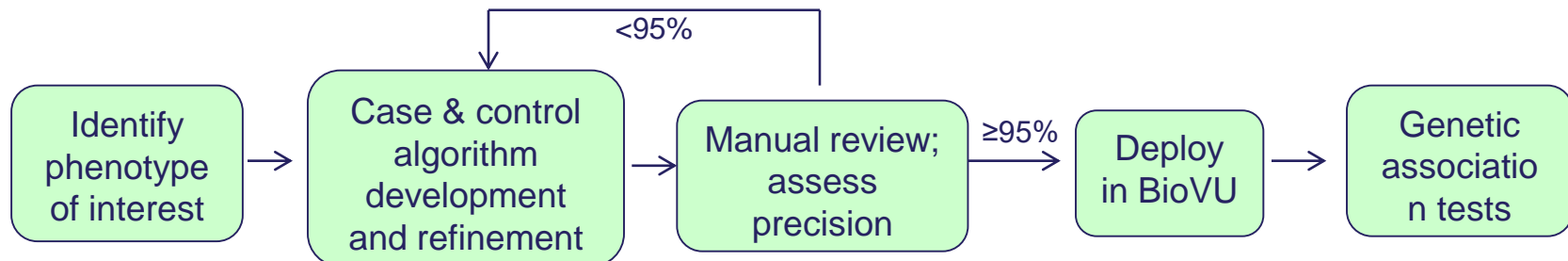
Multiple components are better
(and blood pressure is the worst)



What we learned - Finding phenotypes in the EHR



Algorithm Development and Implementation



Early discovery science in eMERGE – Hypothyroidism

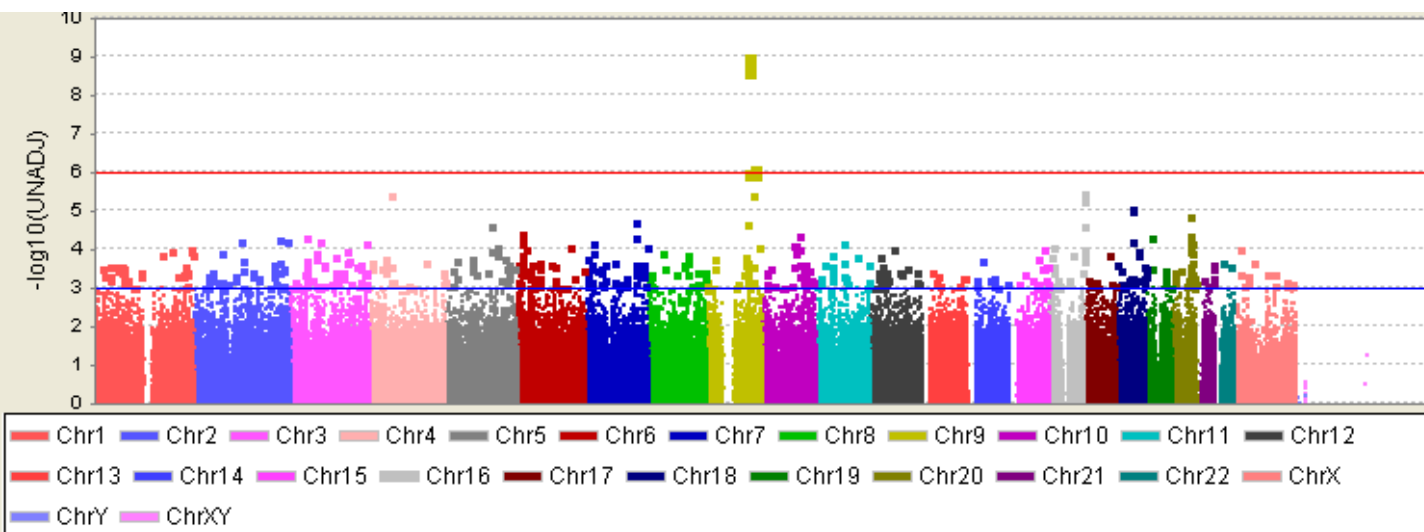
Table 1. Evaluation of Primary Hypothyroidism Algorithm at the Five eMERGE Sites

Site	Primary Phenotype	Total Genotyped Subjects	Primary Hypothyroidism			
			Cases	Controls	Case PPV (%)	Control PPV (%)
Group Health	dementia	2532	397	1,160	98	100
Marshfield	cataracts	4113	514	1,187	91	100
Mayo Clinic	peripheral arterial disease	3043	233	1,884	82	96
Northwestern	type 2 diabetes	1217	92	470	98	100
Vanderbilt	normal cardiac conduction	2712	81	352	98	100
All sites		13,617	1317	5053	92.4 ^a	98.5 ^a

Genotype counts represent all subjects who were found by the hypothyroidism algorithms at each site and who were genotyped. Counts are limited to those classified as “white” in the electronic medical record of each site. PPV = positive predictive value.

^a Average weighted for number of samples contributed to the total.

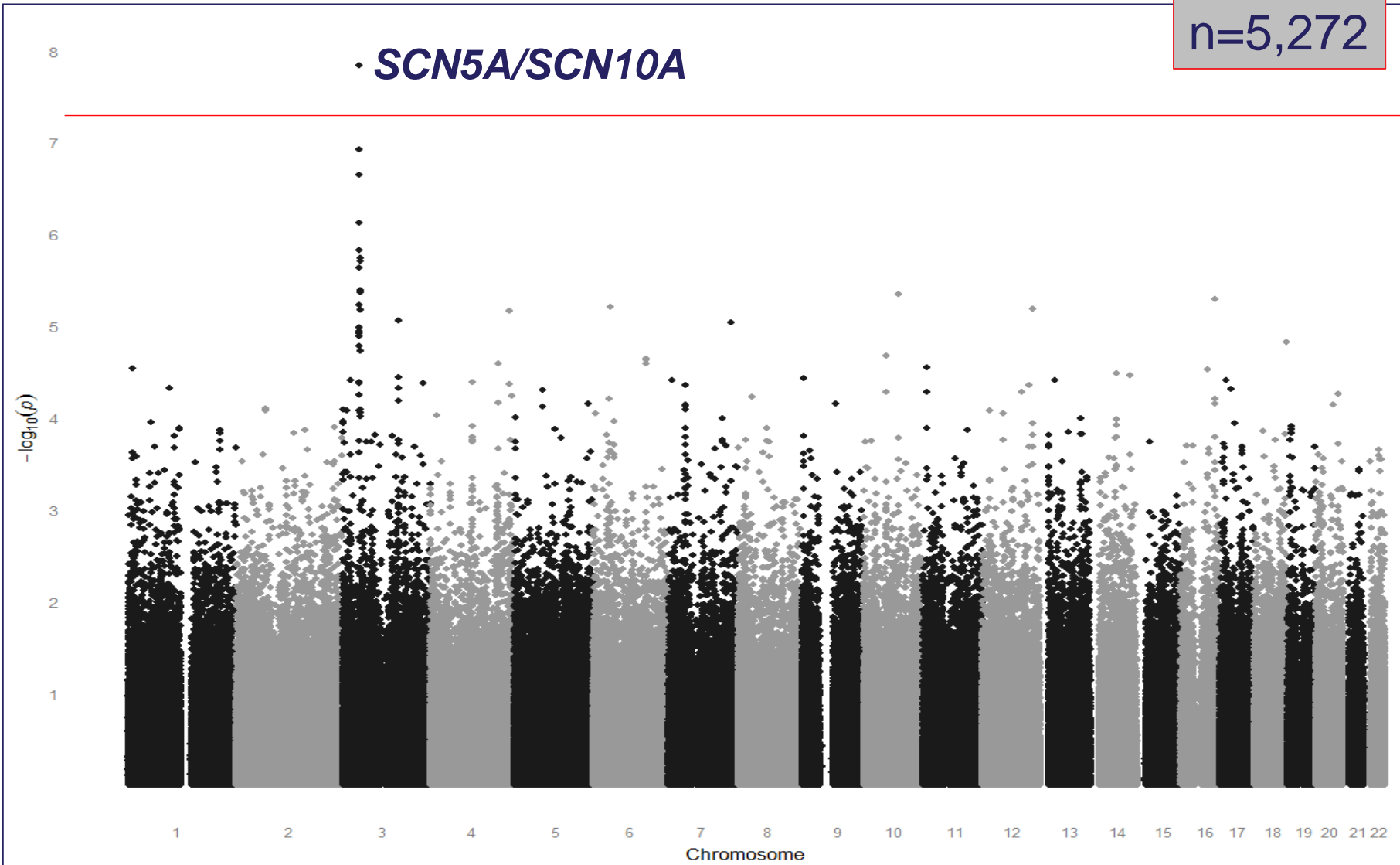
Algorithms can be deployed across multiple EHRs



Analyses can be performed using extant data

GWAS of QRS Duration in eMERGE

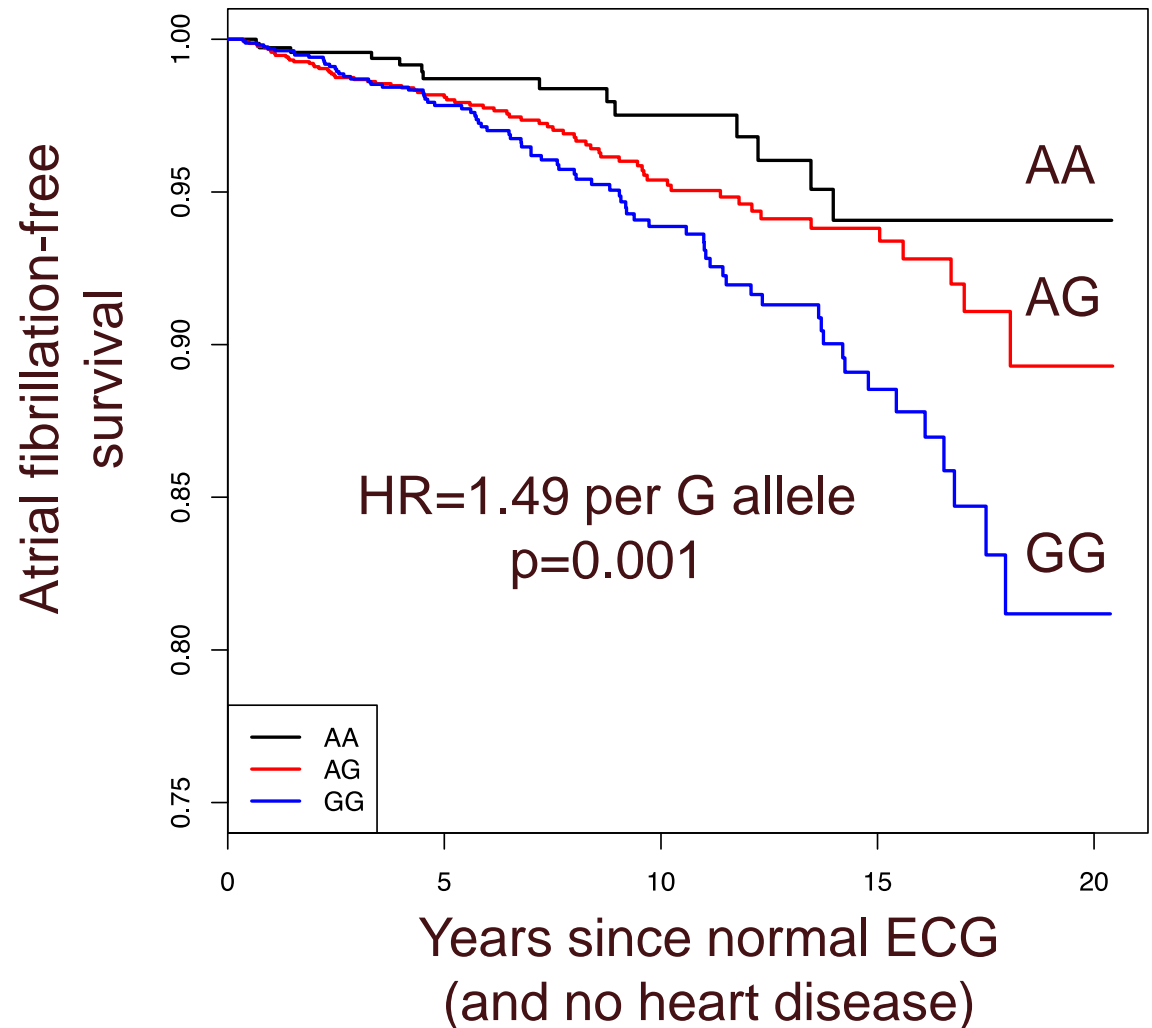
n=5,272



What happens in the “heart healthy” population?

Examined the n=5272
“heart healthy”
population

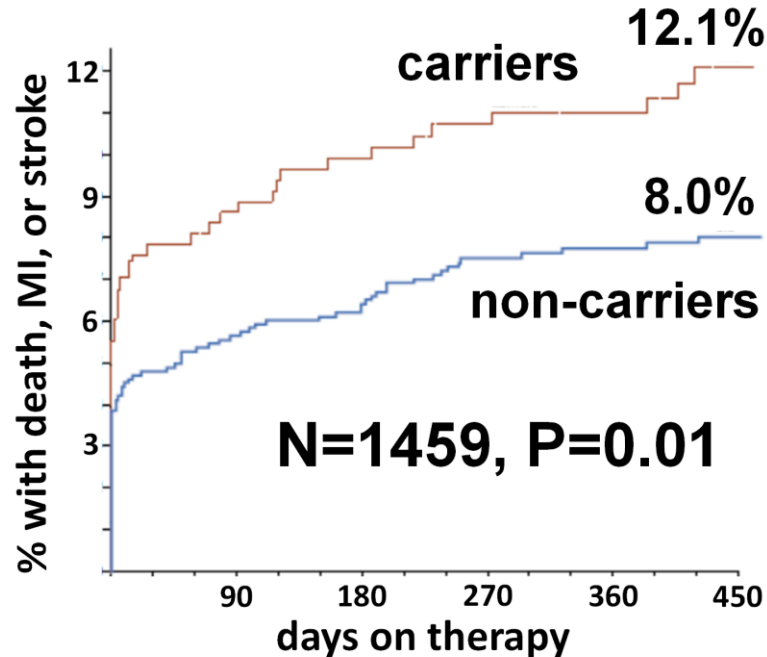
Followed for
development of **atrial
fibrillation** based on
genotype



EHRs for drug response:

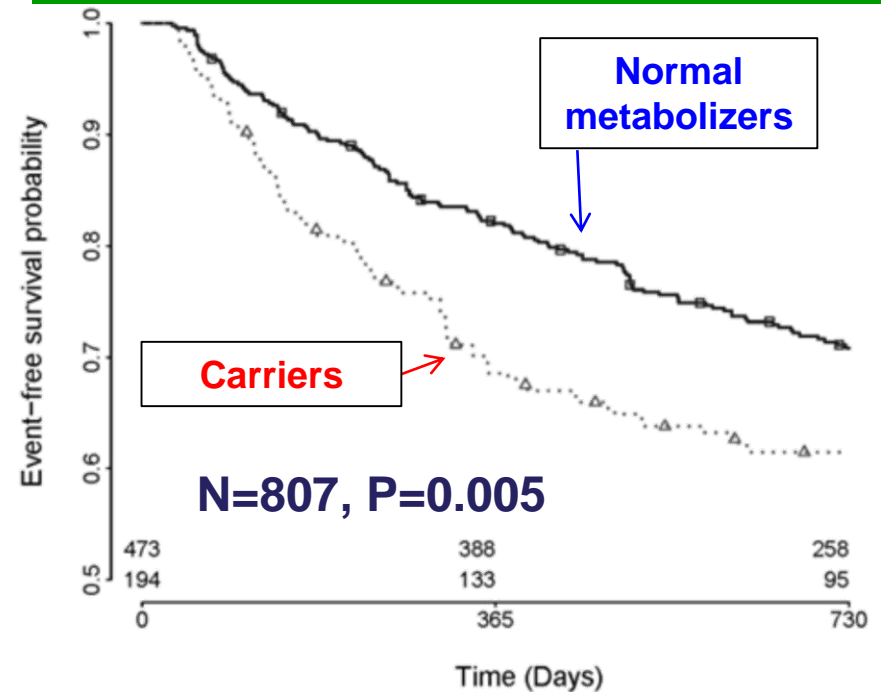
Clopidogrel adverse events associated with *CYP2C19*

From clinical trials



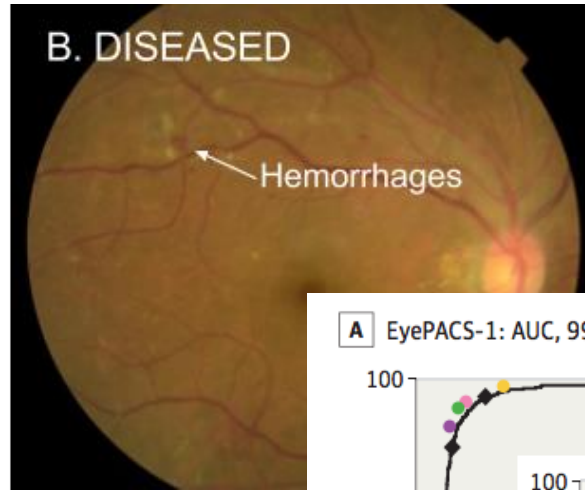
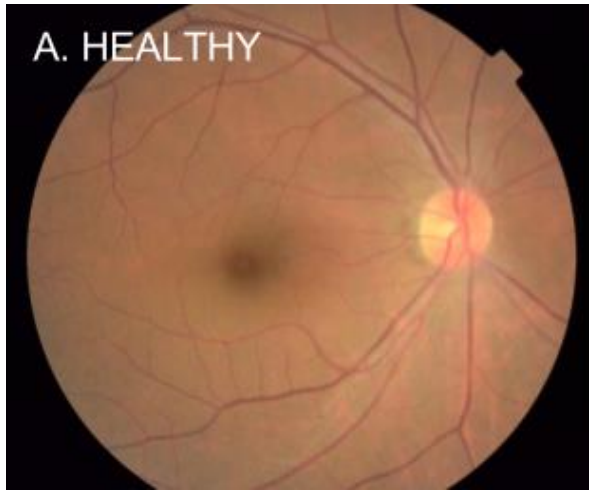
Mega et al., *NEJM* 2009

From the EHR

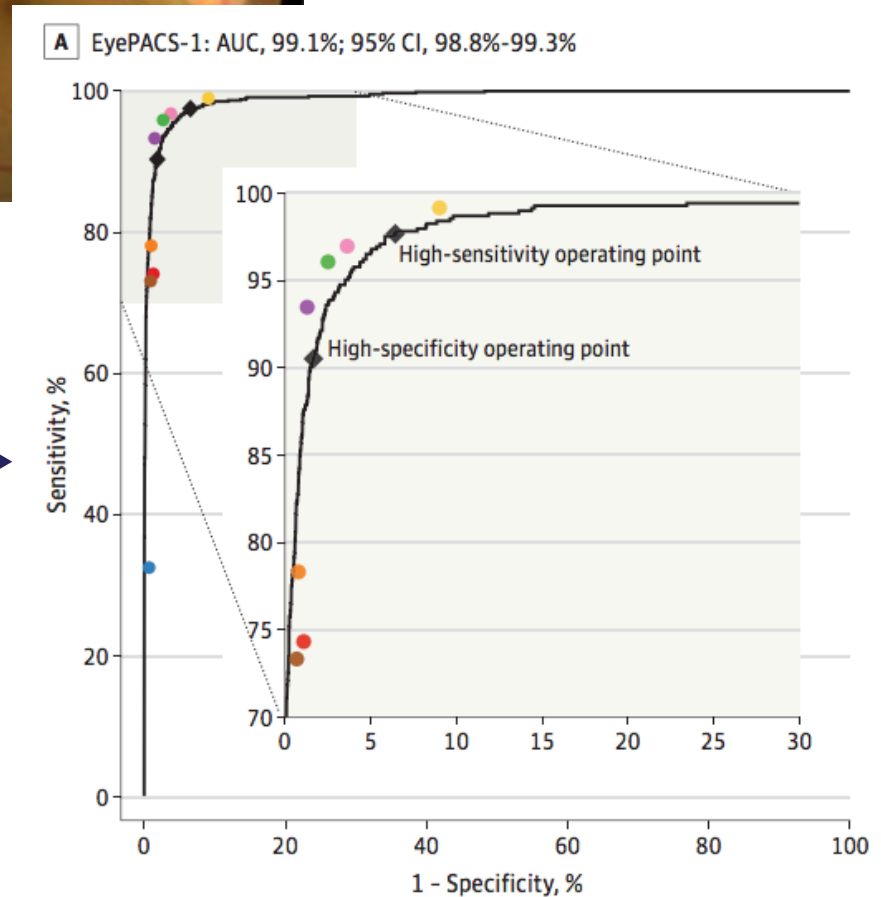


Delaney et al. *Clin Pharm Ther.* 2012

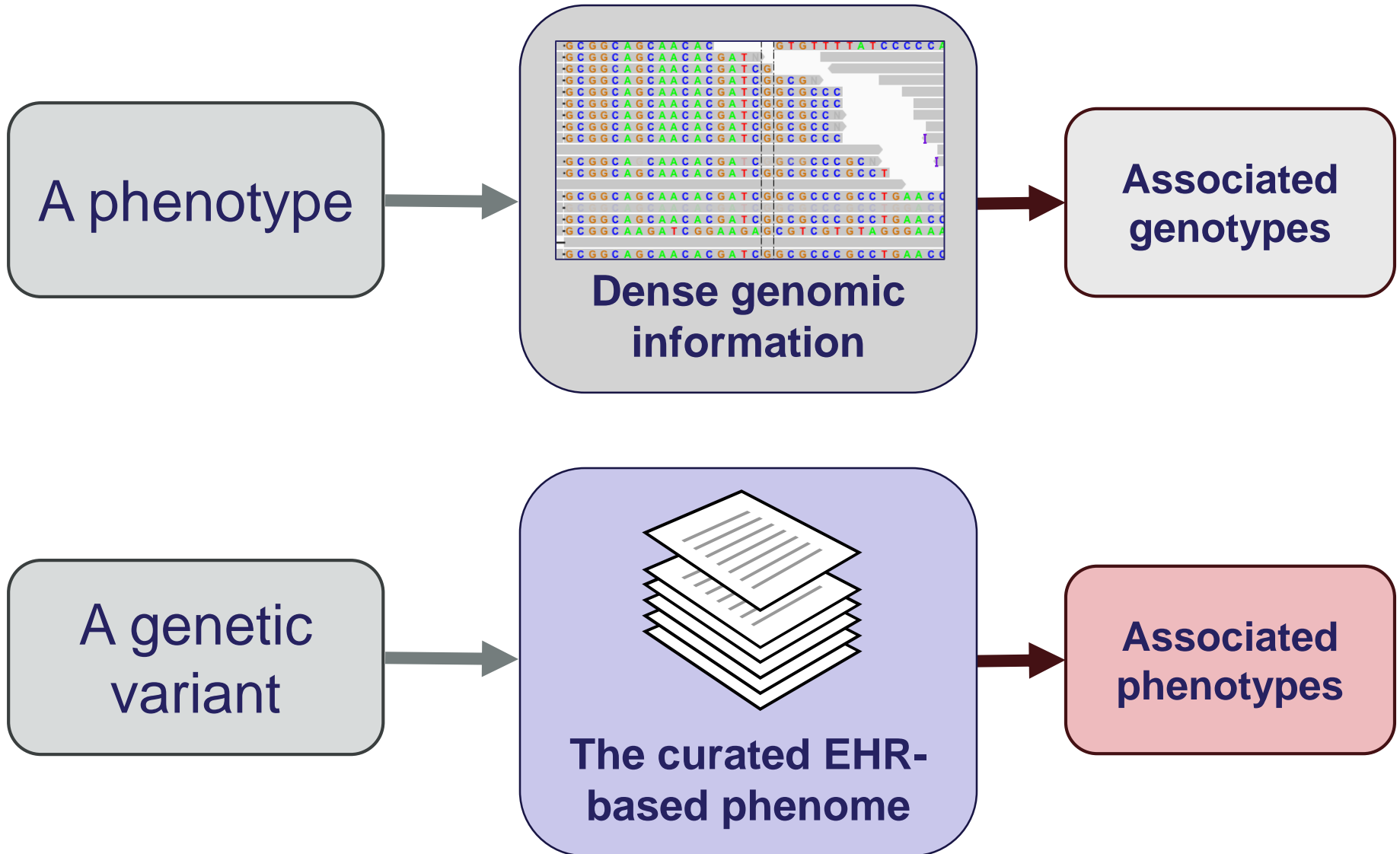
Deep learning for Diabetic Retinopathy



Train a machine learning algorithm over >128k images



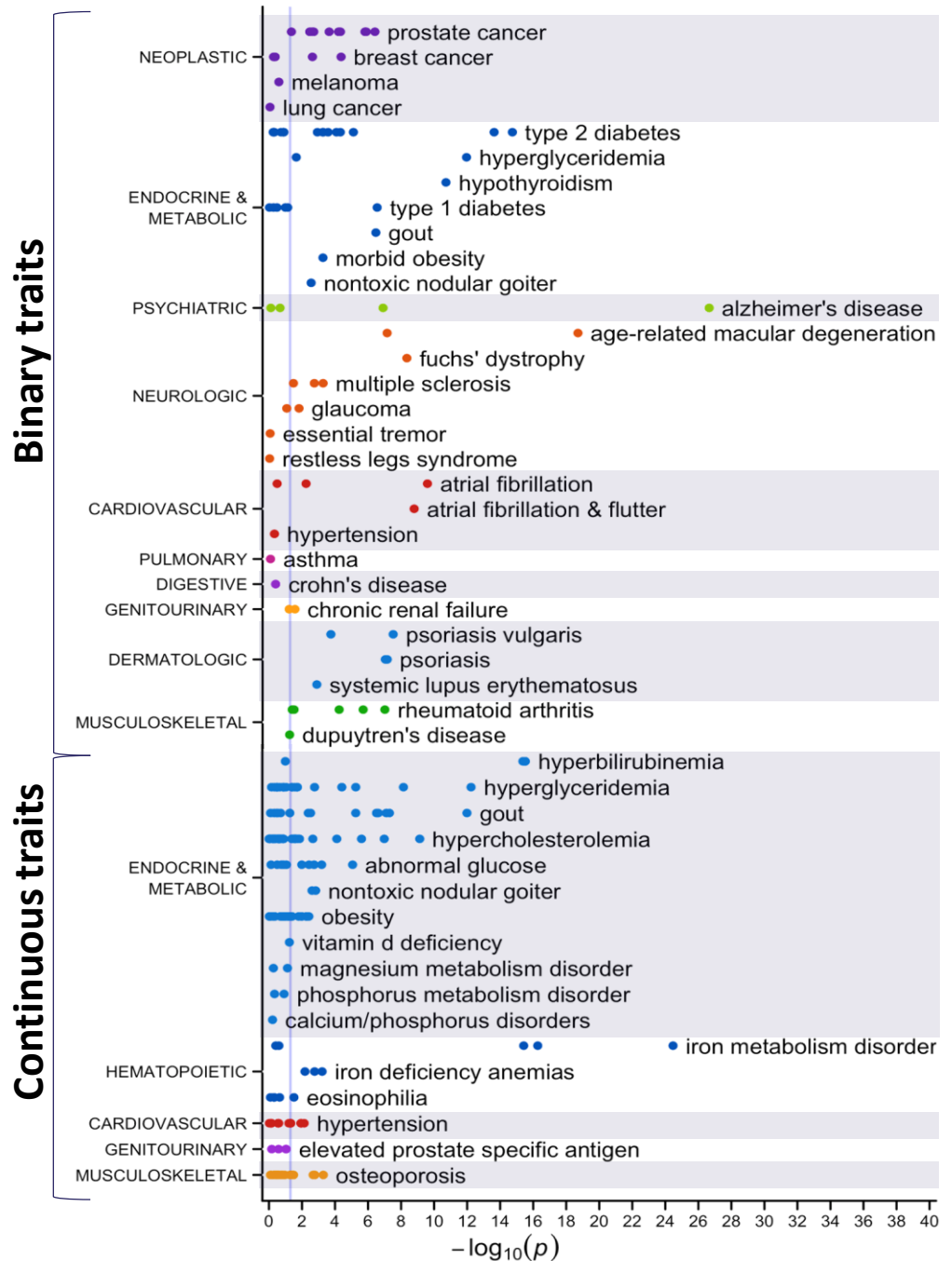
Phenome scanning (*PheWAS*) in the EHR



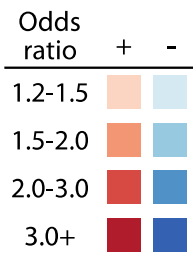
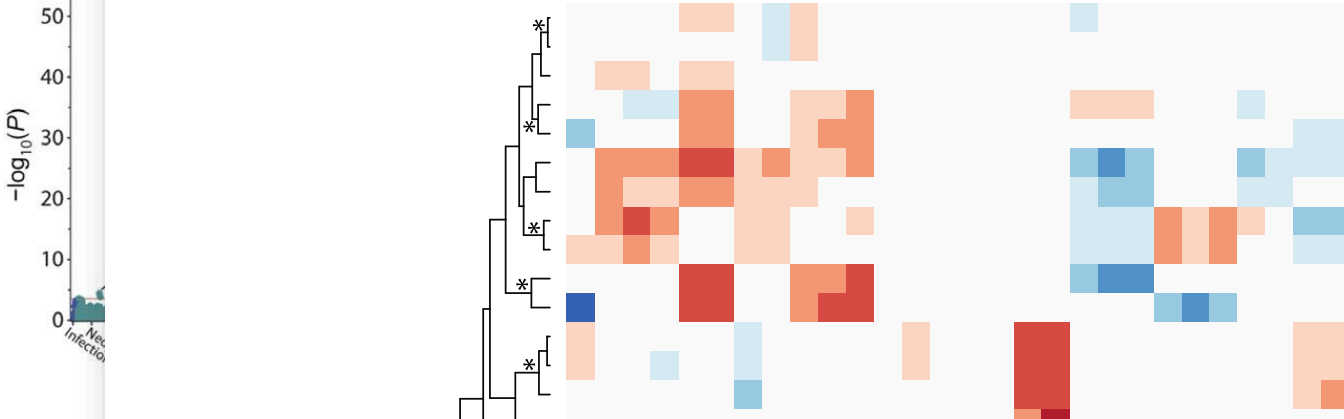
Replications of GWAS associations via PheWAS

P-value for replication:

- All - 210/751: 2×10^{-98}
- Powered - 51/77: 3×10^{-47}

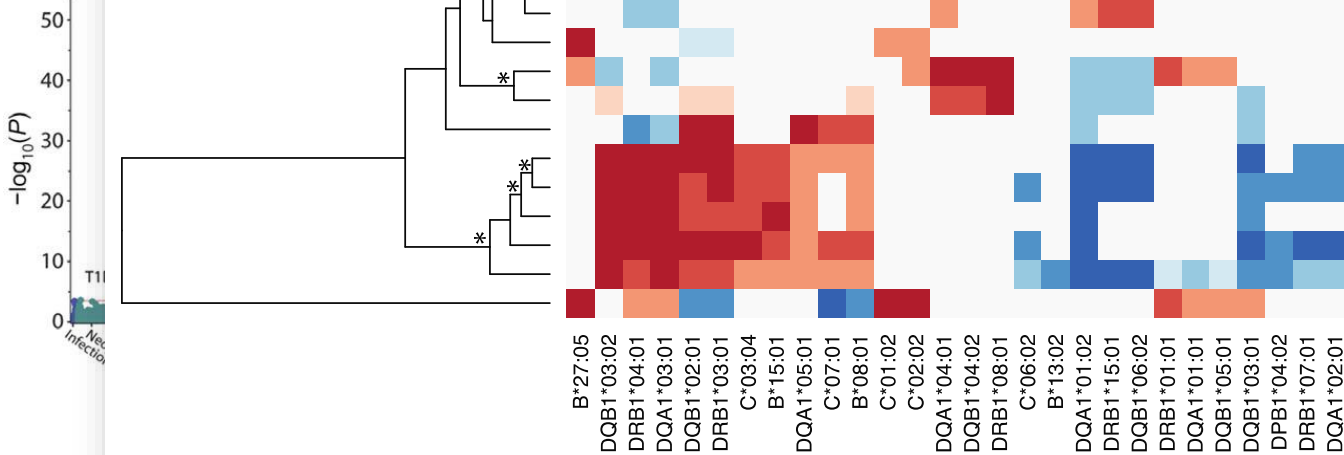


A



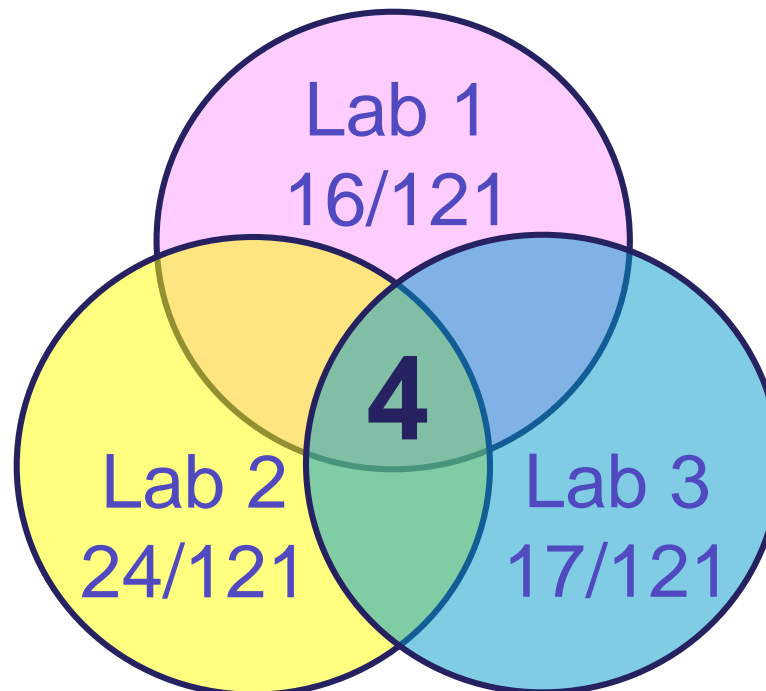
- Hypothyroidism NOS
- Hypothyroidism
- Diabetes mellitus
- Lupus (localized and systemic)
- Diffuse diseases of connective tissue
- Diabetic retinopathy
- Insulin pump user
- Rheumatoid arthritis
- Rheumatoid arthritis and other infl. polyarth.
- Dermatomyositis
- Sicca syndrome
- Psoriasis
- Psoriasis and related disorders
- Psoriasis vulgaris
- Loss of teeth or edentulism
- Multiple sclerosis
- Other demyelinating diseases of CNS
- Other inflammatory spondylopathies
- Juvenile rheumatoid arthritis
- Primary biliary cirrhosis
- Celiac disease
- Type 1 diabetes with ophthalmic manif.
- Type 1 diabetes with renal manifestations
- Type 1 diabetes with neurological manif.
- Type 1 diabetes with ketoacidosis
- Type 1 diabetes
- Ankylosing spondylitis

E

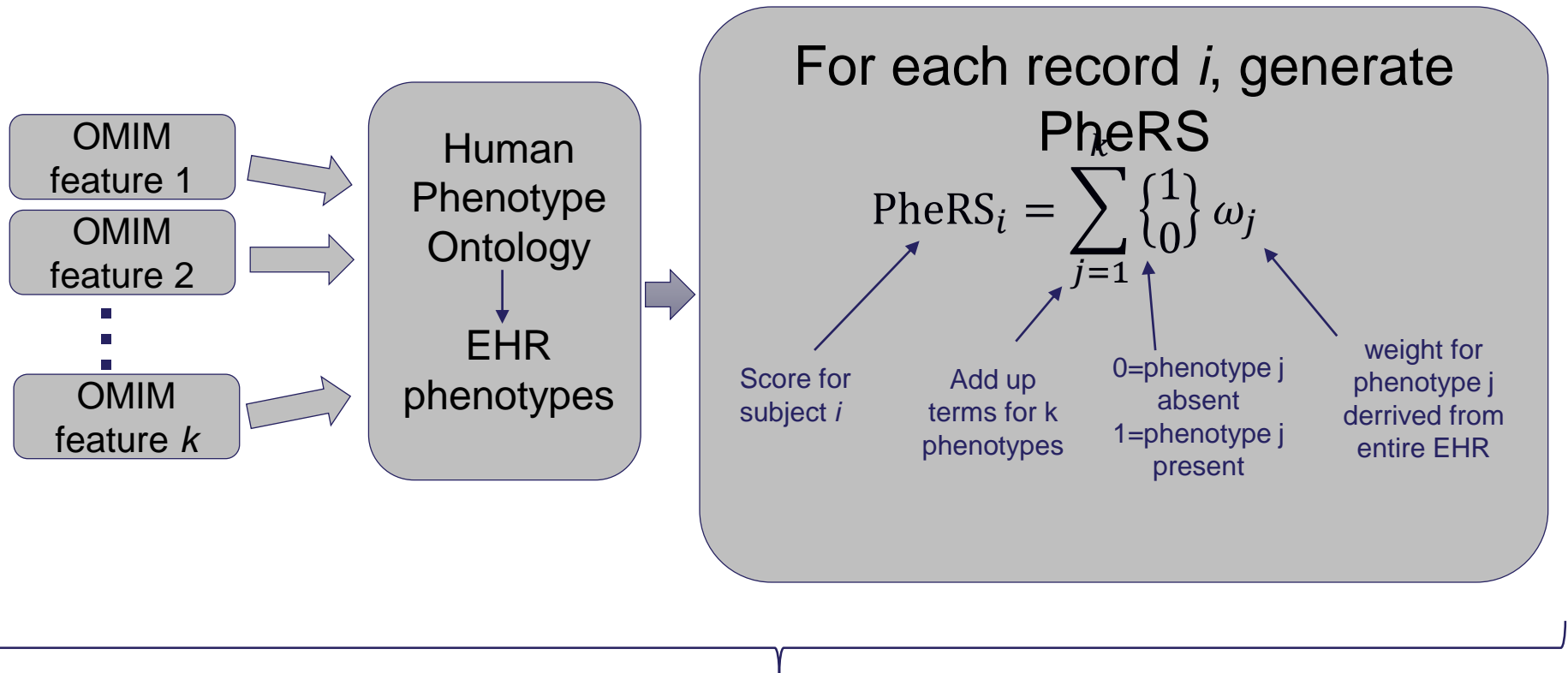


The potential for “call back” deeper phenotyping: Long QT genes (*SCN5A* and *KCNH2*) in 2,200 sequenced patients in eMERGE

- 83 rare (MAF < 1%) in *SCN5A*, 45 in *KCNH2*
- 121/128 MAF < 0.5%, 92 singletons
- Three labs assessed known/likely pathogenicity



Calculating a Phenotype Risk Score (PheRS)



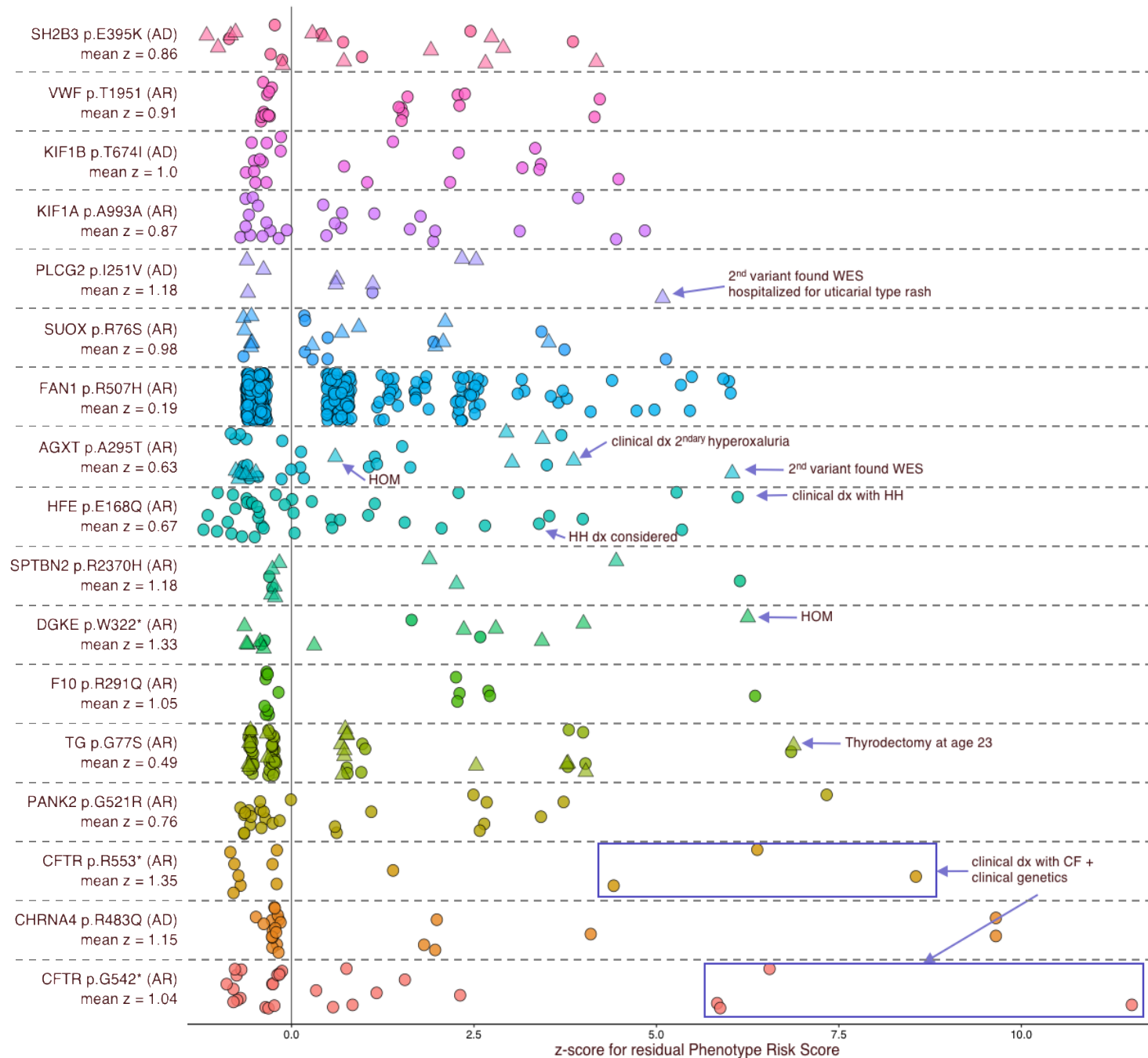
Repeat this for all Mendelian diseases

Example: a phenotype risk score in Cystic Fibrosis

	CF cases				CF controls			
Age/Sex	18F	26M	29F	29M	18F	26M	29F	29M
Chronic airway obstruction	Red	Red	Red	Red	White	White	White	White
Pneumonia	Red	Red	Red	Red	Red	White	White	White
Diseases of pancreas	Red	Red	Red	Red	White	White	White	White
Hypovolemia	Red	Red	Red	White	White	White	White	White
Acute upper respiratory infections	Red	White	Red	Red	Red	White	White	White
Asthma	Red	White	White	Red	Red	White	White	White
Bronchiectasis	White	White	Red	Red	White	White	White	White
Intestinal malabsorption	Red	White	White	Red	White	White	White	White
Hepatomegaly	Red	White	White	White	White	White	White	White
Acute pulmonary heart disease	White	White	White	White	White	White	White	White
Phenotype Risk Score	9.8	4.4	6.3	7.8	2.5	0.7	0.0	0.7

PheRS identified potentially pathogenic SNVs

N=21k on exome chip
6k SNVs



The *All of Us* Research Program – Breaking Down Data Silos



Overview of the *All of Us* approach and protocol

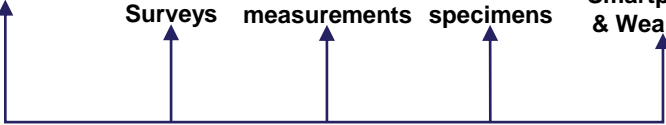


**Direct
Volunteers**

**Health Care Provider
Organizations**

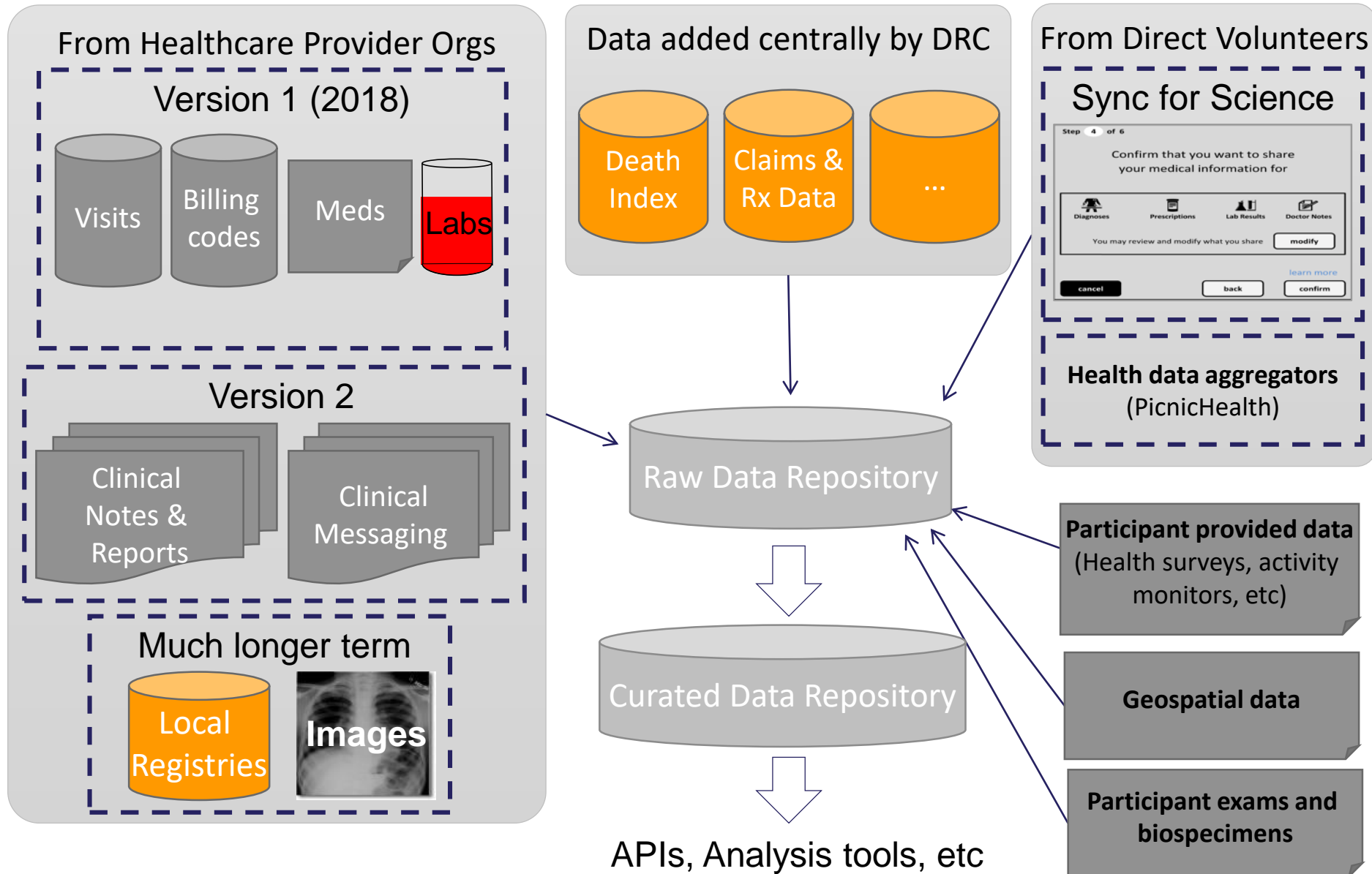


EHR data Health Surveys Baseline measurements Bio-specimens Smartphones & Wearables



Multiple data types linked together by semantic standards

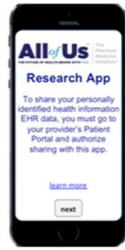
All of Us will aggregate data from many sources



Sync 4 Science (S4S) – a technology to share health data

Research App

Launch Portal



Return to App



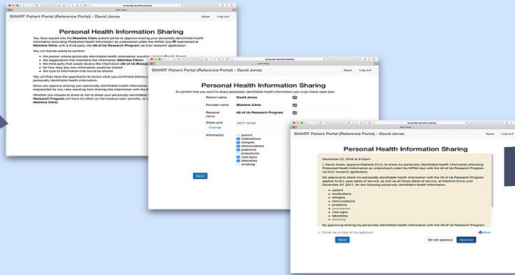
Workflow

Login

Patient's Confidential Credentials

Approve

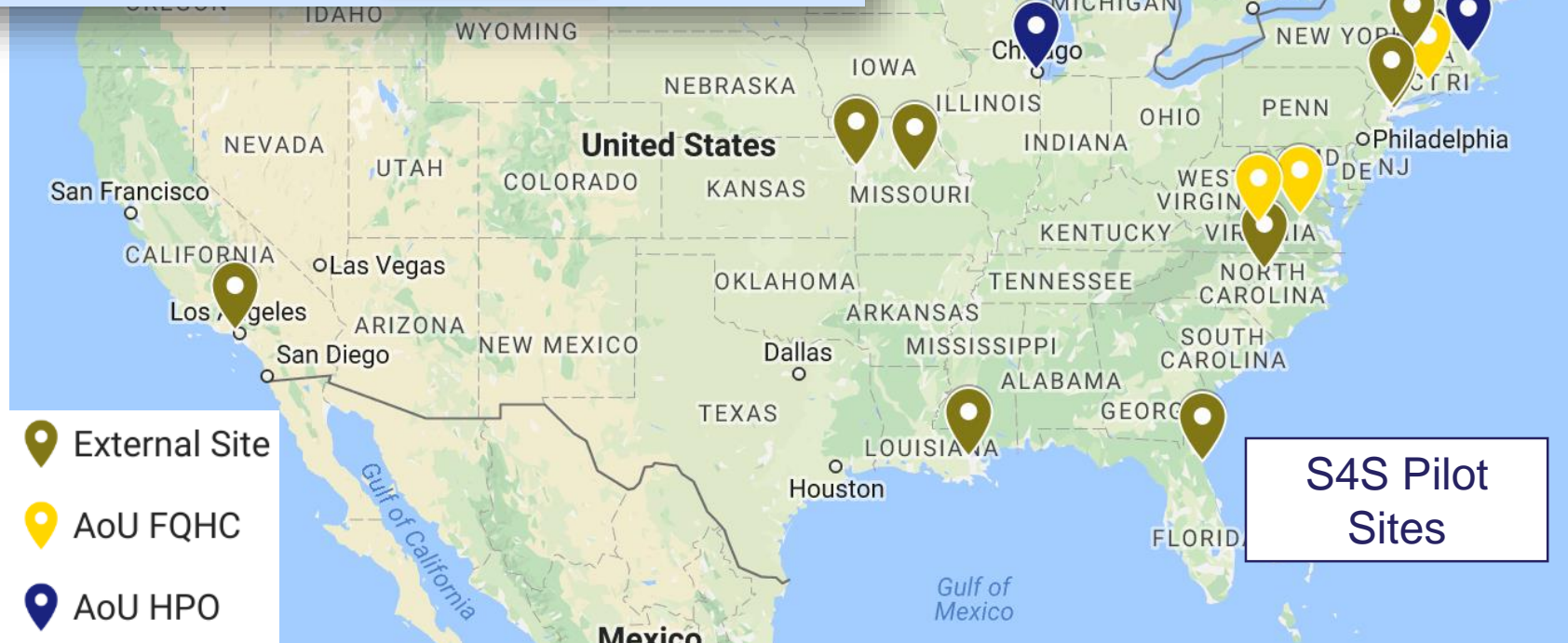
or NOT!



Patient Portal

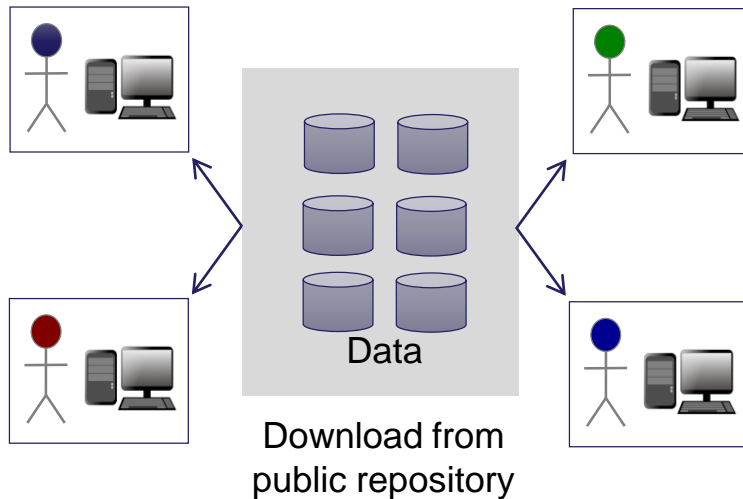
S4S:

- FHIR-based
- Starting with MU Common Clinical Data set



Data Access is centralized in *All of Us*

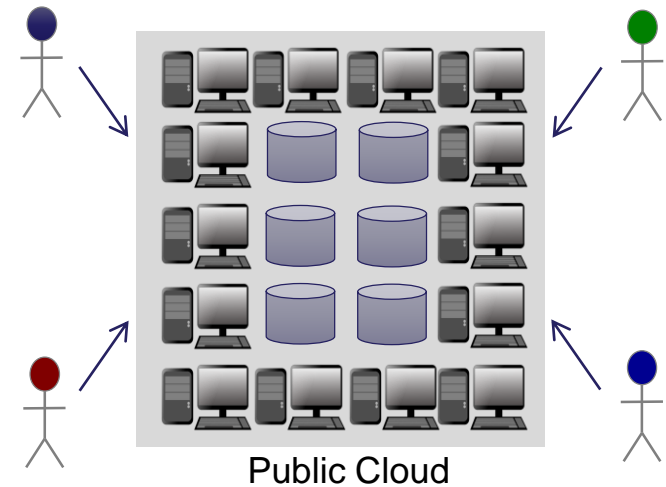
Traditional Approach: Bring data to researchers



Problems

- **Data sharing = data copying**
- **Security (data handoffs)**
- **Huge infrastructure needed**
- **Siloed compute**

AoU Approach: Bring researchers to the data



Advantages

- **Cost**
- **Threat detection and auditing**
- **Increased Accessibility**
- **Shared compute**

The power of a data biosphere of common semantics and APIs

