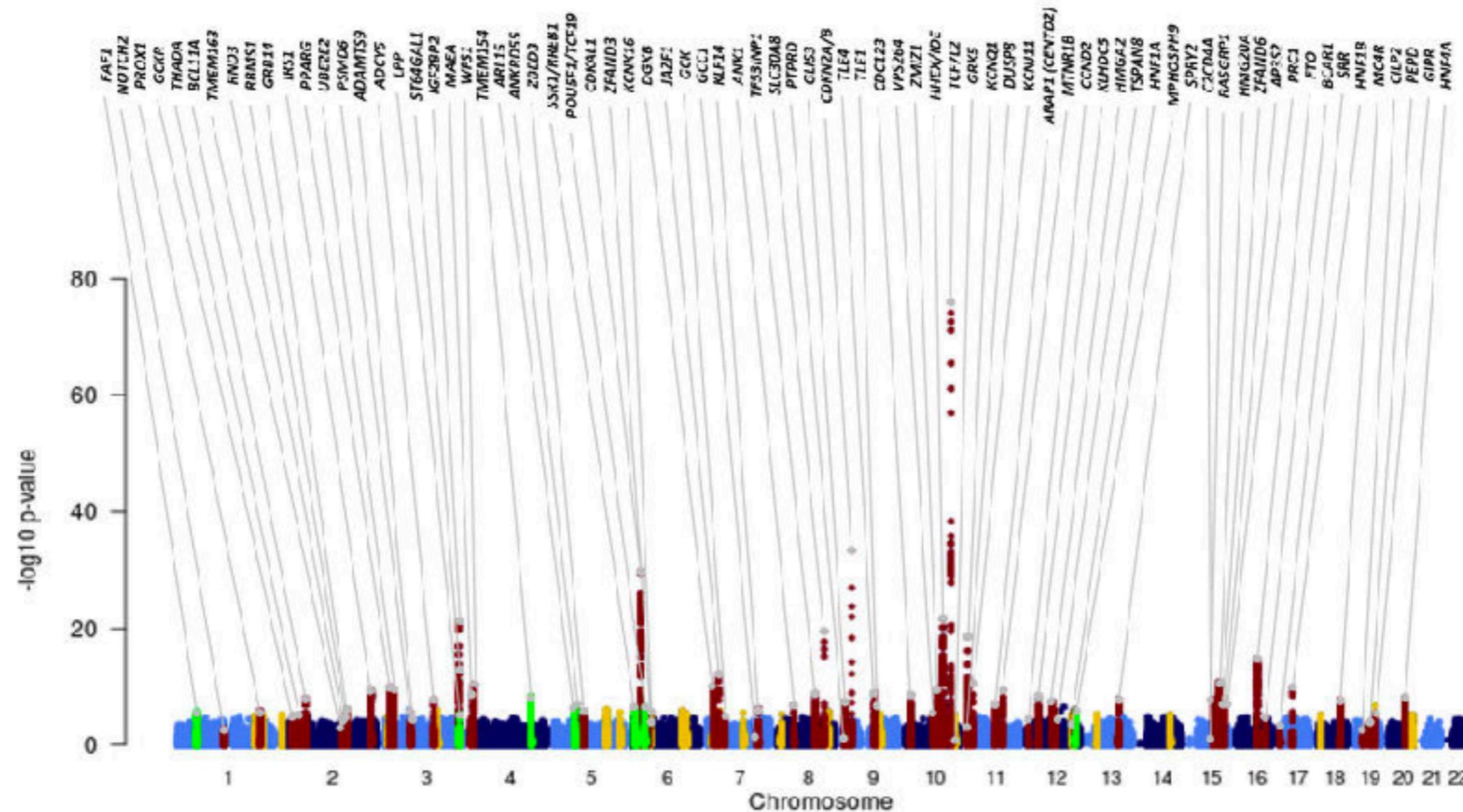


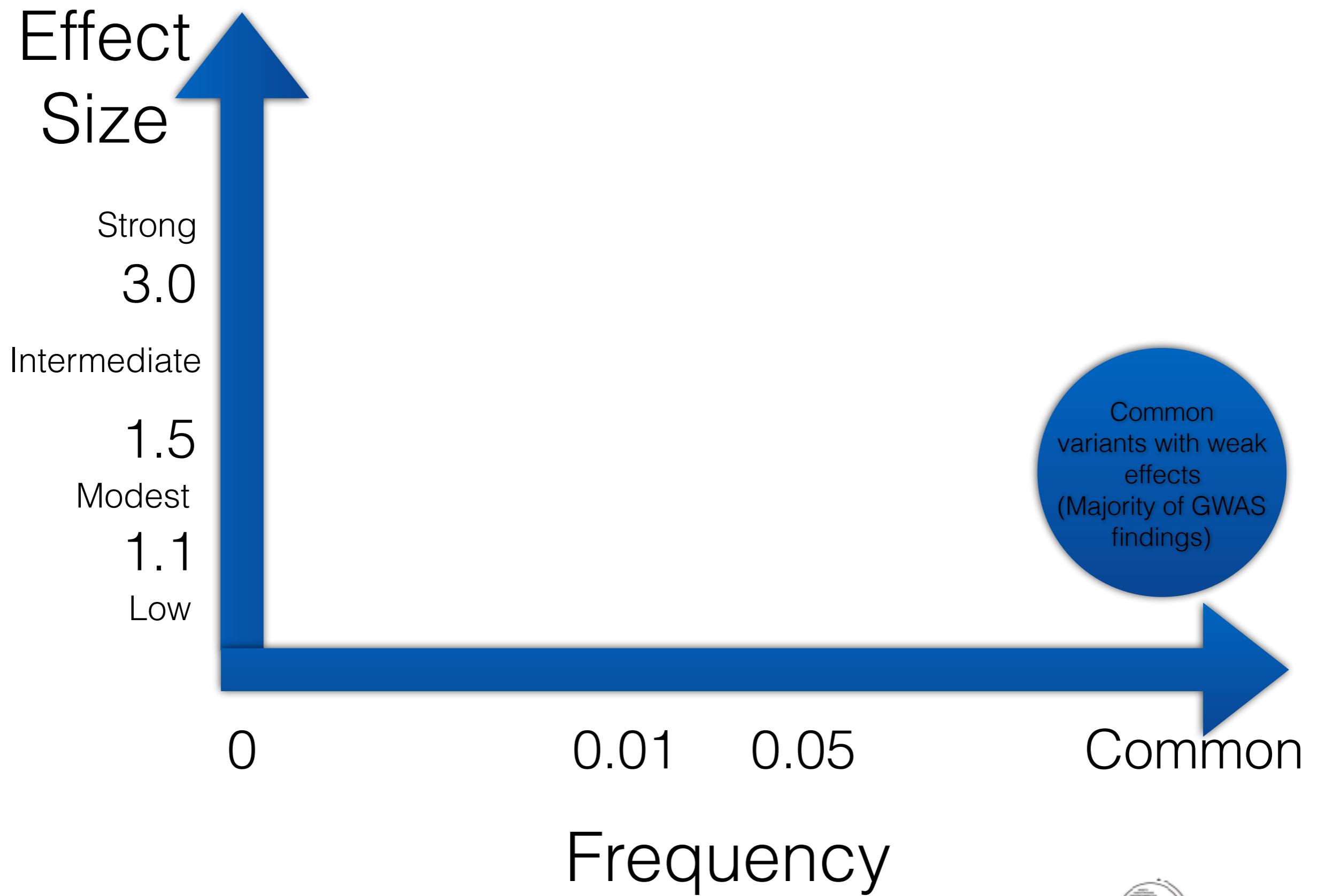
High-dimensional inference from summary statistics

Manuel A. Rivas
Department of Biomedical Data Science
Stanford University
rivaslab.stanford.edu

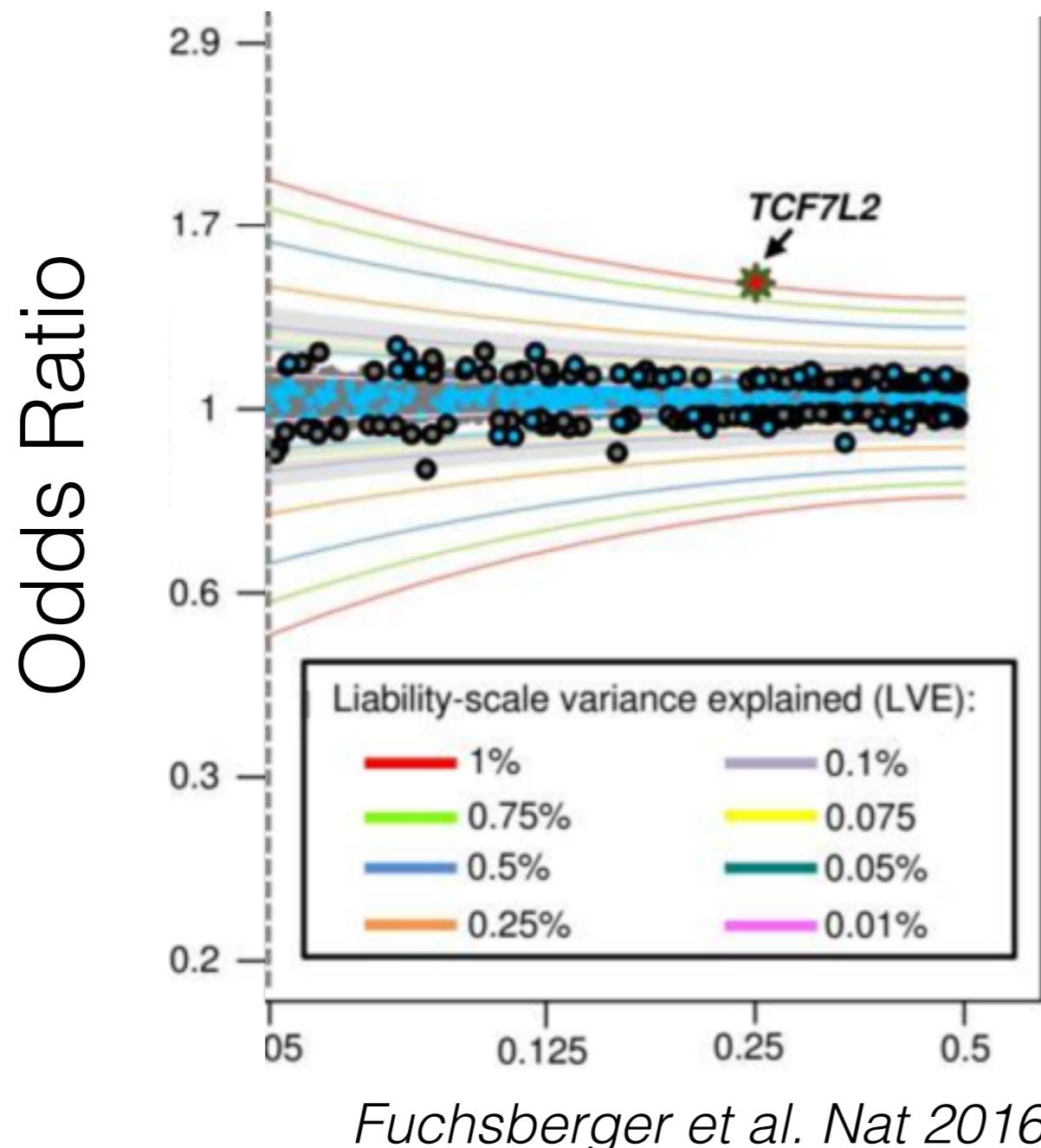


Genome-wide association studies have been very successful*

Why the “*”?

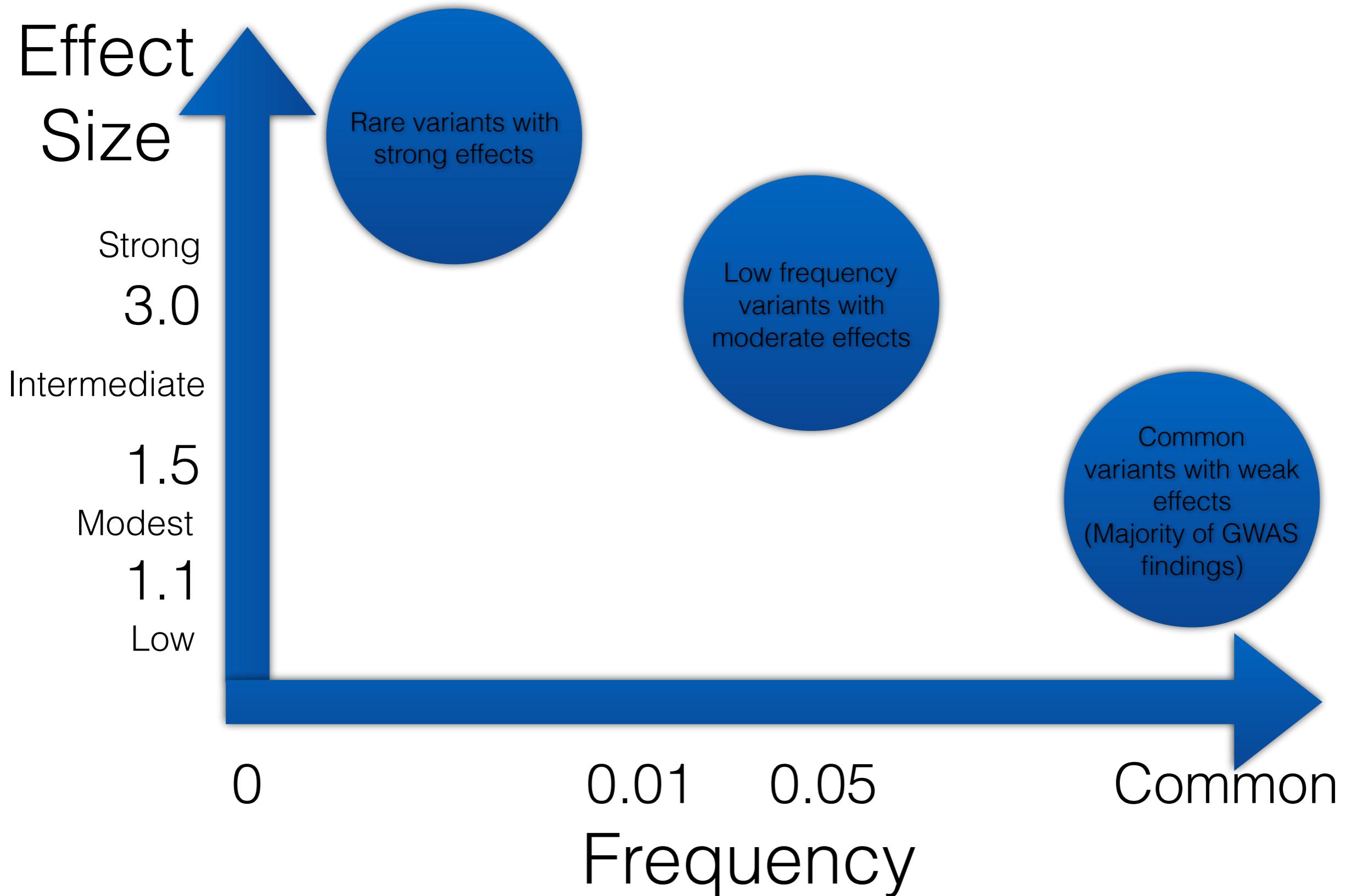


In the context of **type 2 diabetes** all common variant associations were tiny



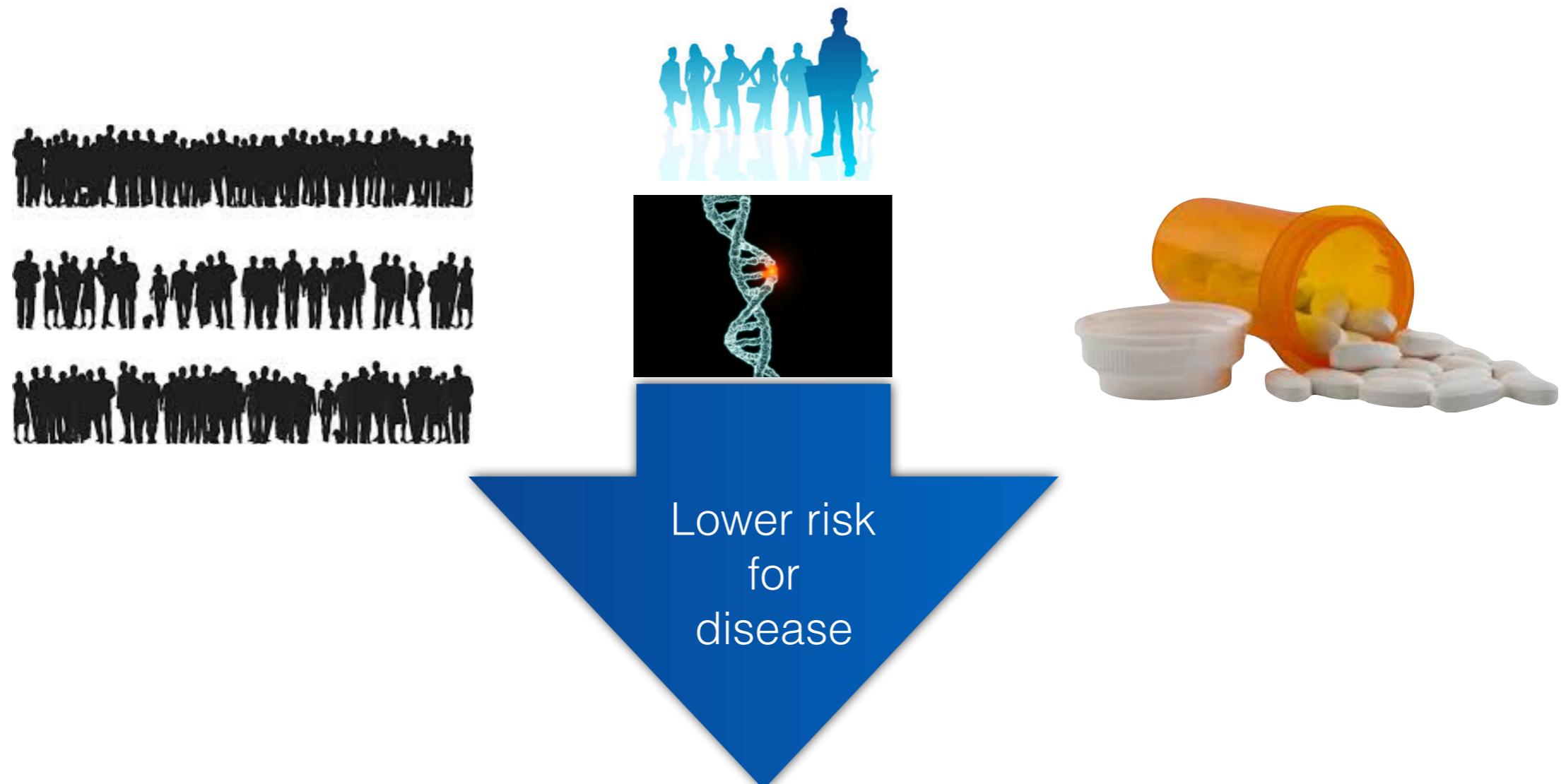
Tiny effect sizes for all associated variants

Minor allele frequency

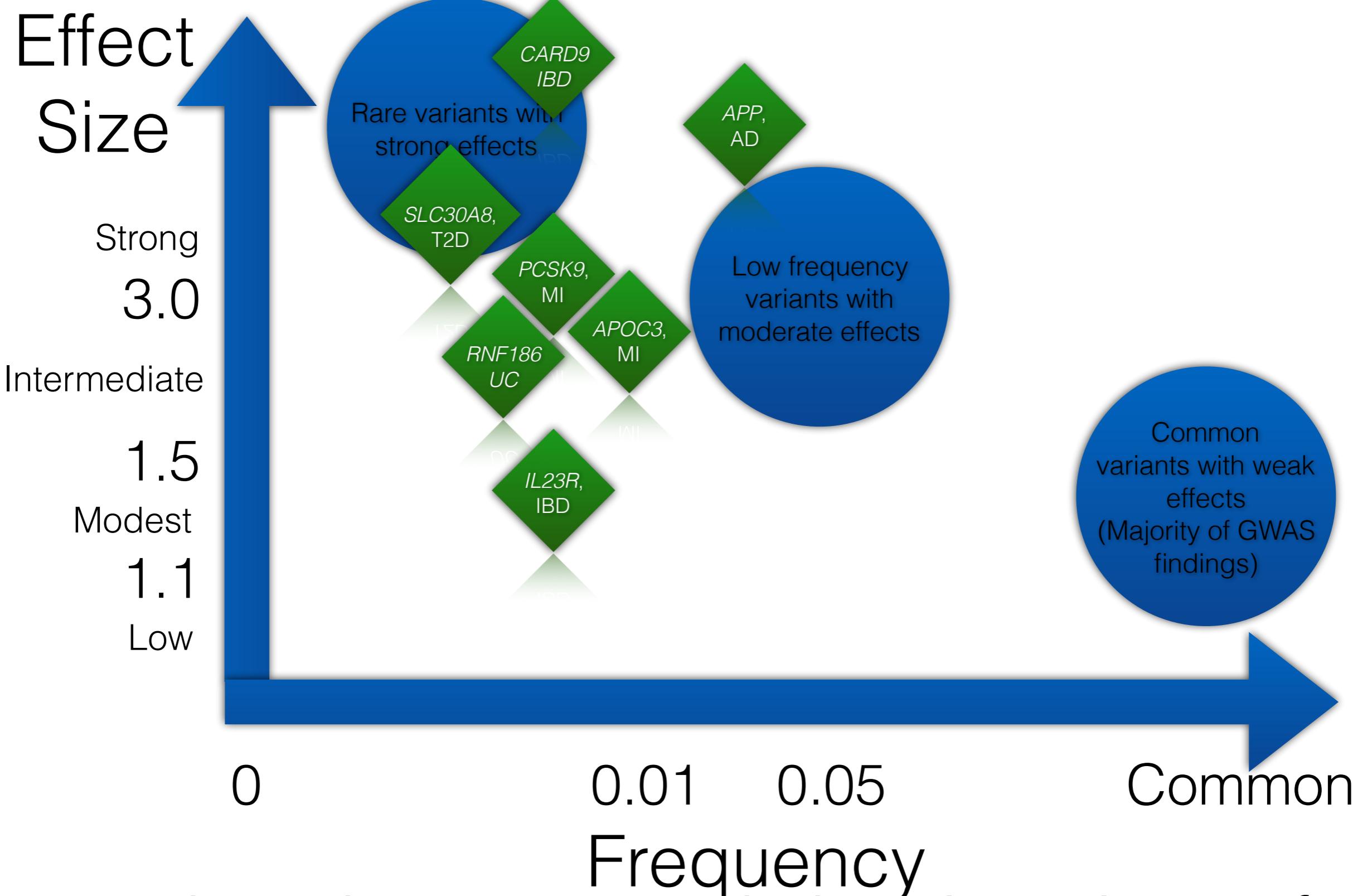


Additional signals started emerging from sequencing studies of rare variants

Experiments of Nature



Can Guide Selection of Targets



Human knockouts started showing signs of protection against disease

Population Biobank Studies

UK Biobank



UK Biobank



Health records

$n=500,000$



Web-based questionnaire
 $n=200,000$



Physical activity monitor $n=100,000$



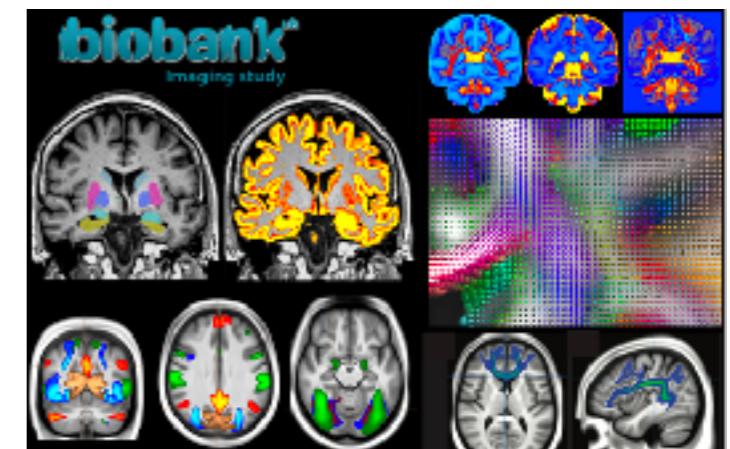
Genotyping $n=500,000$



Baseline Biochemistry

n=500,000

Q1 2018



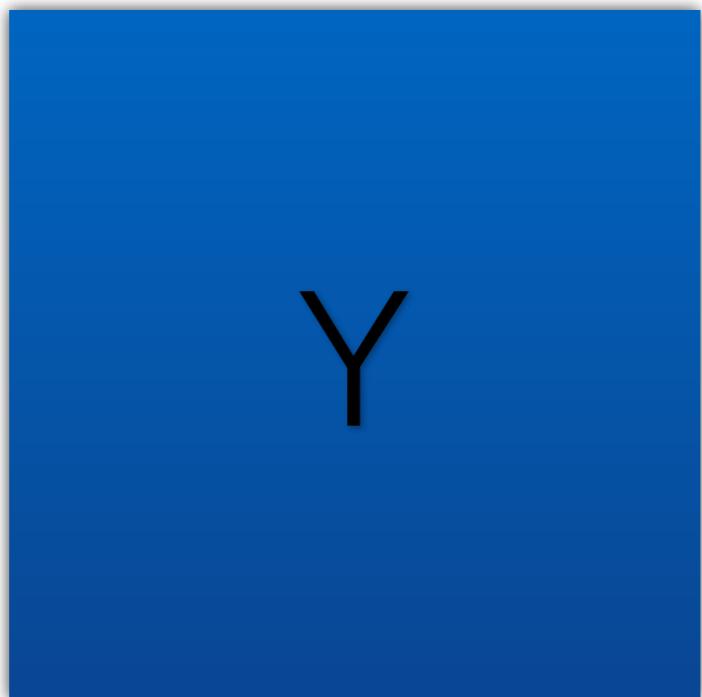
Imaging $n=100,000$ 2015-2023



Massive Datasets

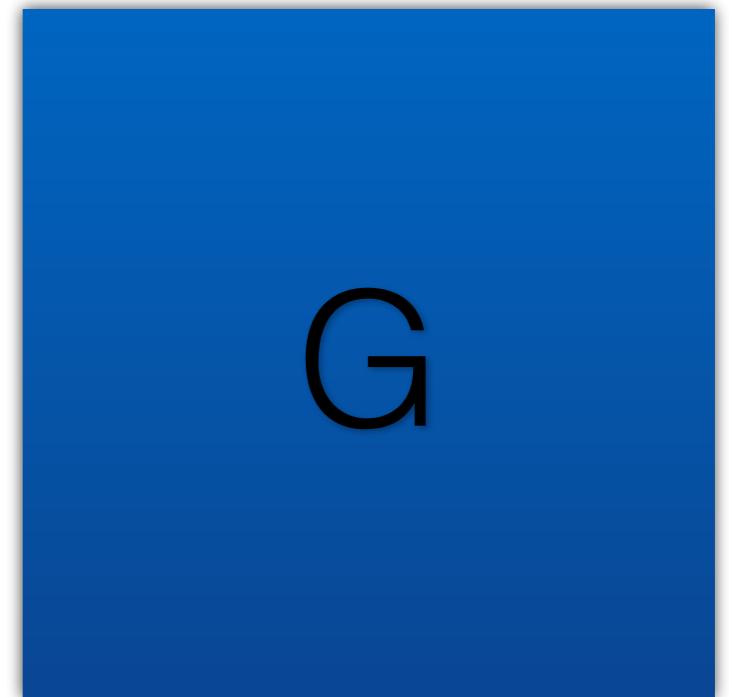
$q \sim 10^4$

$N \sim 10^6$



$p \sim 10^7$

$N \sim 10^6$

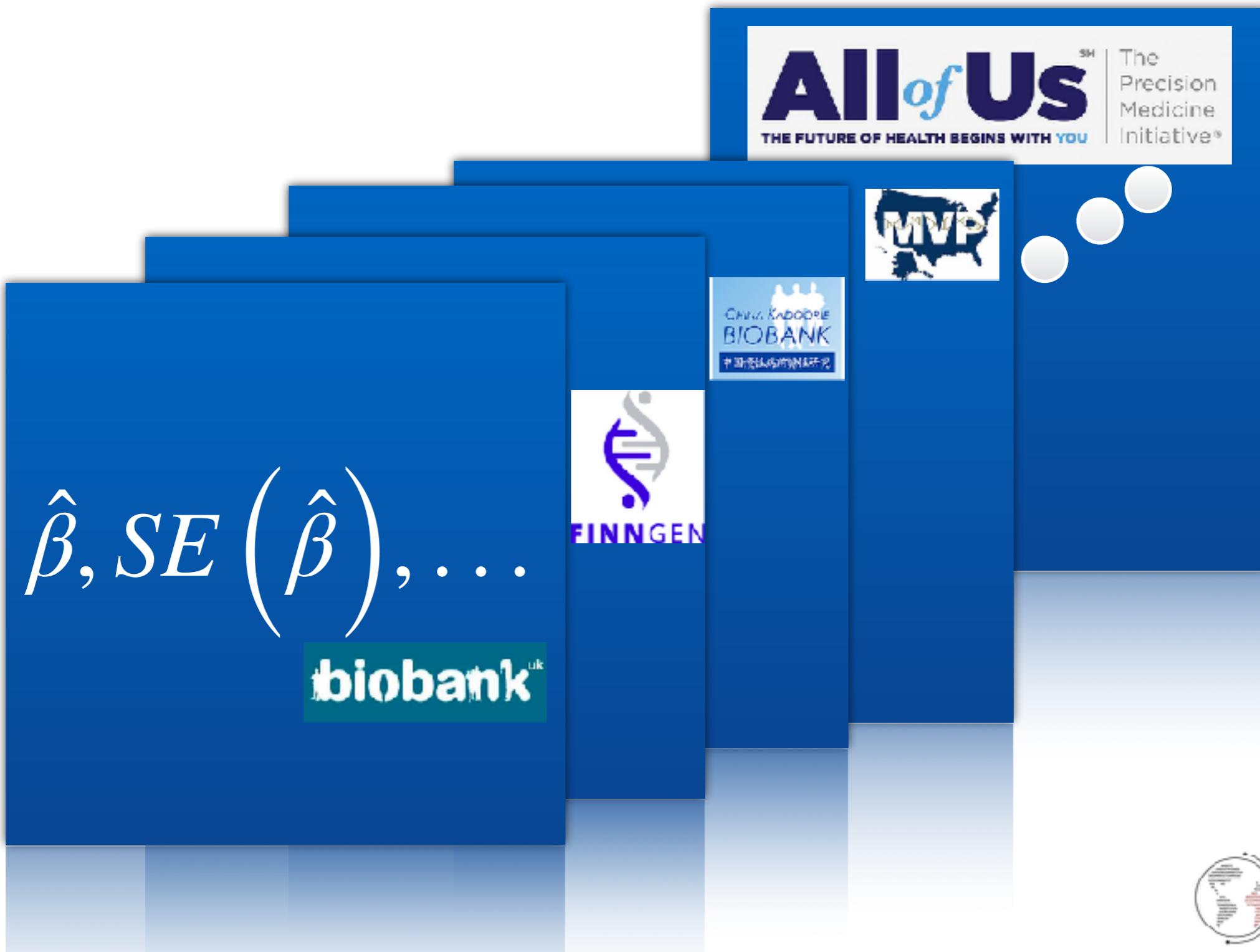


Going from ‘BIG DATA’ to ‘SMALL DATA’

$K \sim 10^4$

$M \sim 10^7$ $\hat{\beta}, SE(\hat{\beta}), \dots$

Improving inference



Global Biobank Engine

Global Biobank Engine About Downloads Terms Contact HLA Alleles Power Genetic correlation Decomposition FAQ
G Select Language ▾

Global Biobank Engine (pre-alpha)

Search for a gene or variant or region or phenotype coding (coming soon)

Examples - Gene: [F5](#), Transcript: [ENST00000367797](#), Variant: [1:169519049](#), RS ID: [rs6025](#), Region: [10:114696614-114706614](#)

Genetic Association Results

Note: We present summary statistic results from the UK Biobank hospital in-patient health-related outcomes summary information data ([Data-Field 41202](#)); computational grouping of phenotypes with cancer ([Category 100092](#)) registry, death registry data ([Category 100093](#)), algorithmically-defined outcomes ([Category 42](#)), and verbal questionnaire data ([Category 100071](#)); and manually curated grouping of phenotypes.

Browseable phenotypes

Cancer

Recent News

Up next

October 10-18, 2017

- Imputation results upload with ([Neale Lab](#)).

November 21, 2017

- HLA Disease Map added ([Here](#)).

October 25, 2017

- Update of gene-based results: Top 100

Bayesian model comparison for rare variant association studies of multiple phenotypes

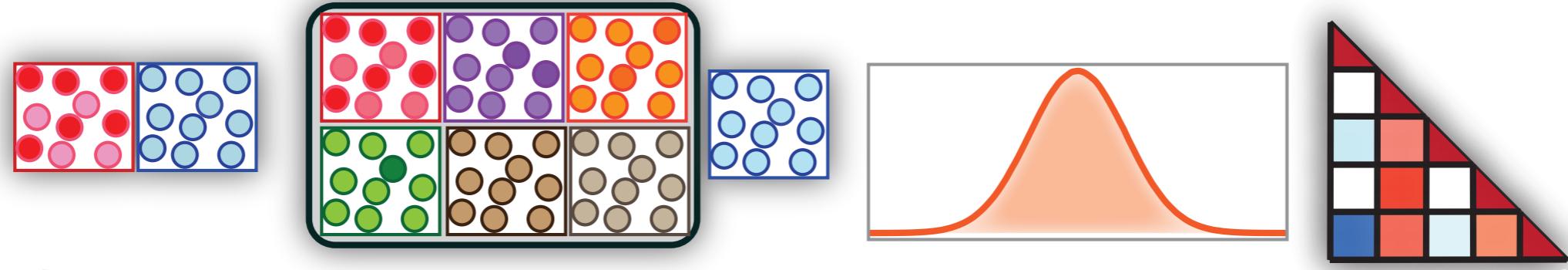
Introduction

- How can we detect associations for rare variants ($MAF < 1\%$)?
- Common approach: aggregate rare variants by gene or other genomic “unit”
- Can we incorporate multiple genetically-related phenotypes to improve power?
- **M**ultiple **R**are variants and **P**henotypes (**MRP**) Bayesian model comparison:
 - Compare genetic effects for rare variants in a “unit” (e.g. gene) using null model of no effects or alternative model that captures correlation, scale, and location of genetic effects

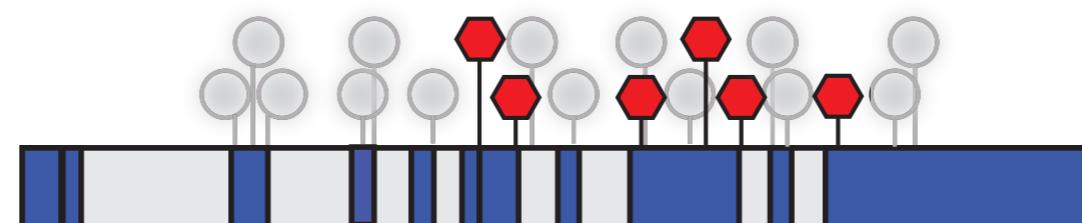
Study designs for multiple phenotypes

- MRP can be applied to different study designs:
 - Case-control
 - Multiple diseases and shared controls
 - Quantitative phenotypes
 - Mixture of case-control and quantitative phenotypes

Study design



*Rare variant
genetic analysis*



Bayesian model comparison

Defined as a Bayes factor between the alternative model and the null model.

Bayesian model comparison

Defined as a Bayes factor between the alternative model and the null model.

Where the null model is simply

$$\beta = 0$$

Bayesian model comparison

The Bayes Factor (BF) is obtained as a ratio of the marginal likelihoods of the data for the two models:

$$\text{BF} = \frac{\int_{\beta_1} p(\text{Data}|\beta_1) p(\beta_1) d\beta_1}{\int_{\beta_0} p(\text{Data}|\beta_0) p(\beta_0) d\beta_0}$$

Bayesian model comparison

The Bayes Factor (BF) is obtained as a ratio of the marginal likelihoods of the data for the two models:

$$\text{BF} = \frac{\int_{\beta_1} p(\text{Data}|\beta_1) p(\beta_1) d\beta_1}{\int_{\beta_0} p(\text{Data}|\beta_0) p(\beta_0) d\beta_0}$$

Data can correspond to the effect size estimates and the estimated variance-covariance matrix:

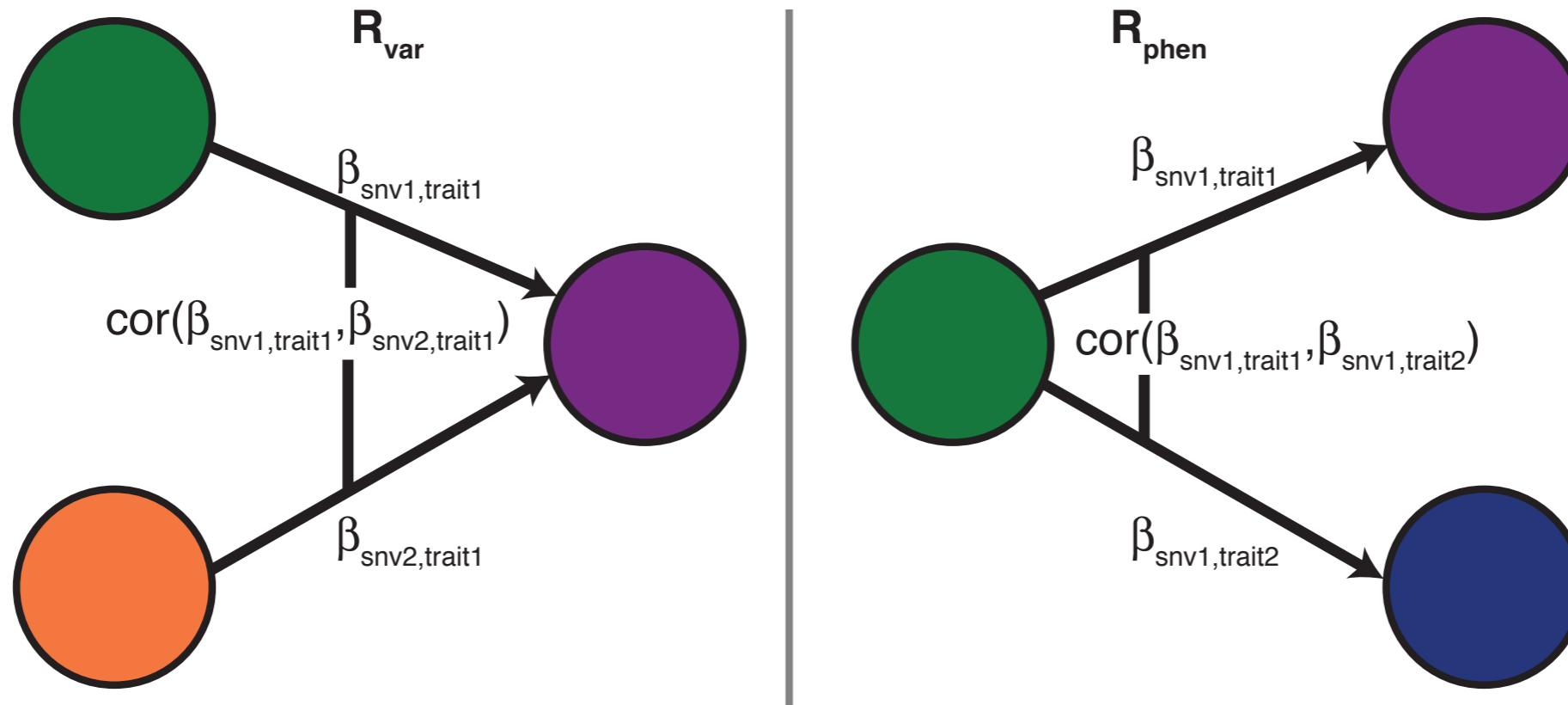
$$\widehat{\beta}, \widehat{\mathbf{V}_B}$$

Bayesian model comparison

In MRP we extend it so that the effect size estimates is across all variants analyzed, across all studies, and all phenotypes. We specify the prior distribution for the alternative model in three steps:

1. Single study
2. Multiple studies
3. Prior mean of genetic effects

Bayesian model comparison: Single study



The prior density for β incorporates the expected correlation of genetic effects among a group of variants (R_{var}) and among a group of phenotypes (R_{phen}).

Bayesian model comparison: Multiple studies

We introduce the matrix

$$\mathbf{R}_{\text{study}}$$

to specify prior on the similarity in effect size across studies

Bayesian model comparison: Multiple studies

We introduce the matrix

$$\mathbf{R}_{\text{study}}$$

to specify prior on the similarity in effect size across studies

$$\mathbf{U} = \mathbf{R}_{\text{study}} \otimes (\mathbf{S}_{\text{var}} \otimes \mathbf{R}_{\text{phen}})$$

and the prior is

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \mathbf{U})$$

Bayesian model comparison: Multiple studies

Attractive property of Bayesian model comparison is ability to obtain a measure of relative support across a number of models

Additional use case example

- Quality control

Heterogeneity of effects across studies may be representative of QC/technical artifacts!

Bayesian model comparison: prior mean of genetic effects

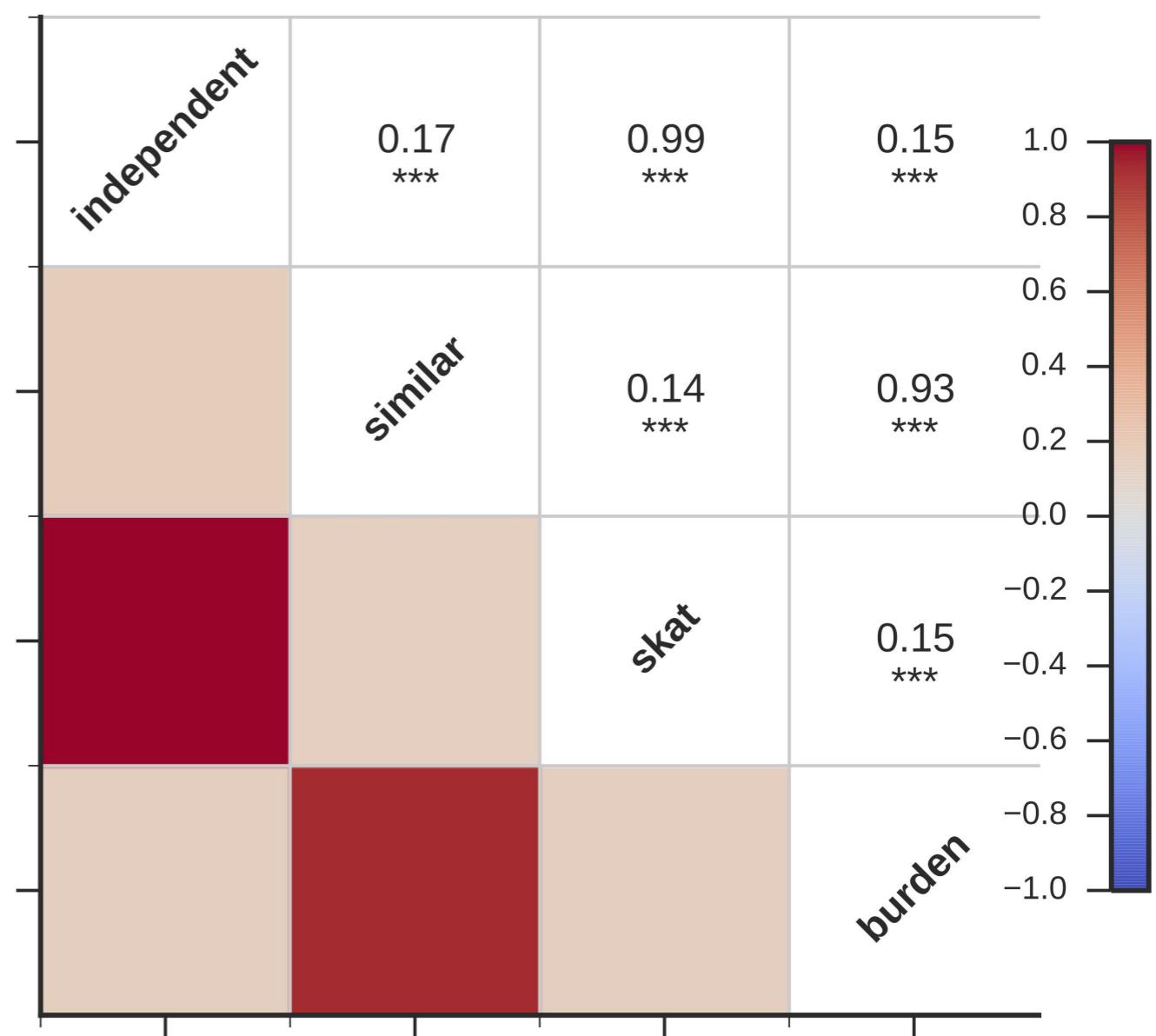
It is straightforward to incorporate prior mean of genetic effects

$$\beta \sim \mathcal{N}(\mu, \mathbf{U})$$

Useful for prioritizing protective modifiers of disease risk

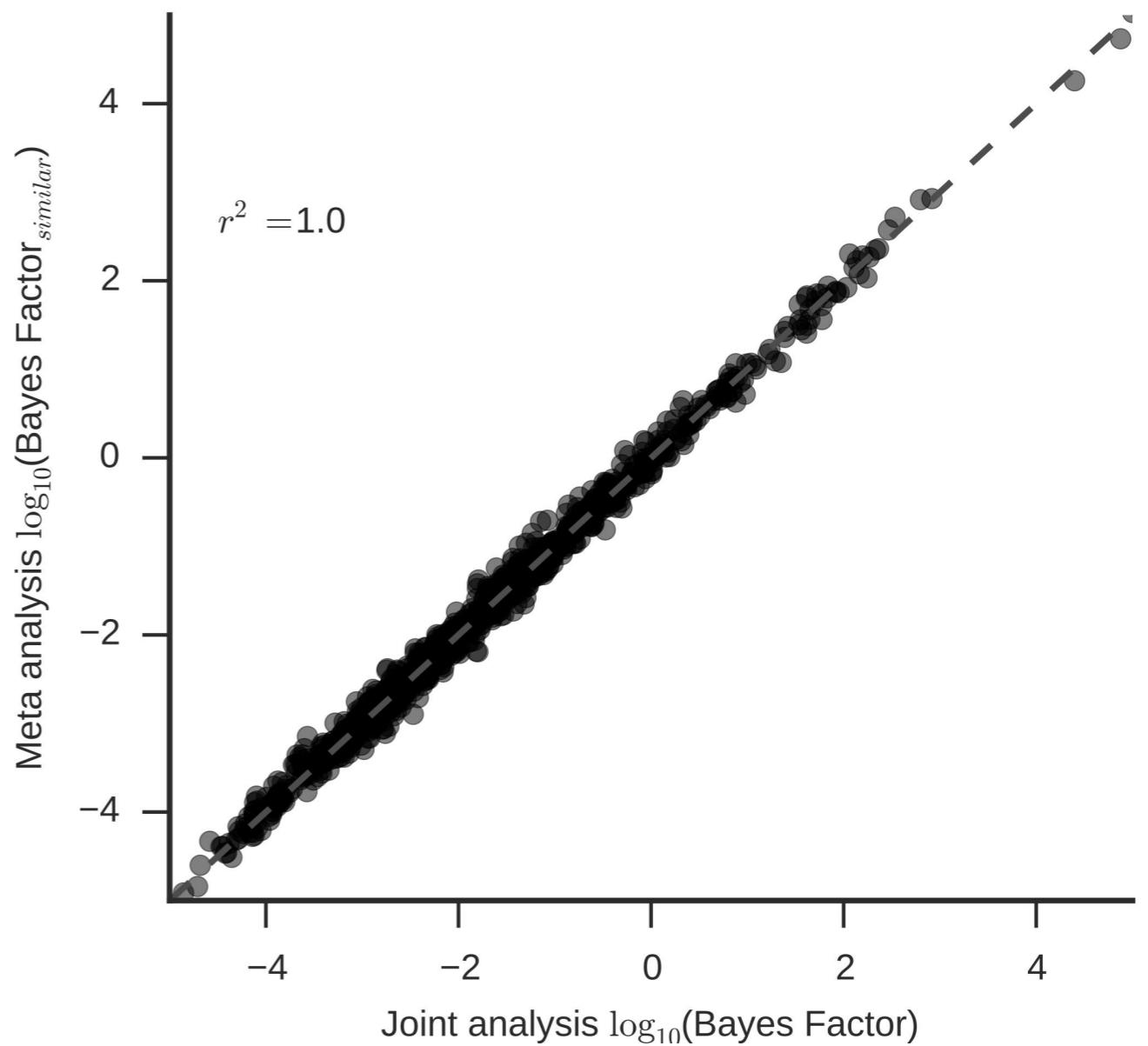
Comparison to other methods

- MRP independent effects:
different effects for
variants in group
- MRP similar effects:
similar effects for variants
in group
- MRP independent or
similar effects model
recover SKAT or burden
results

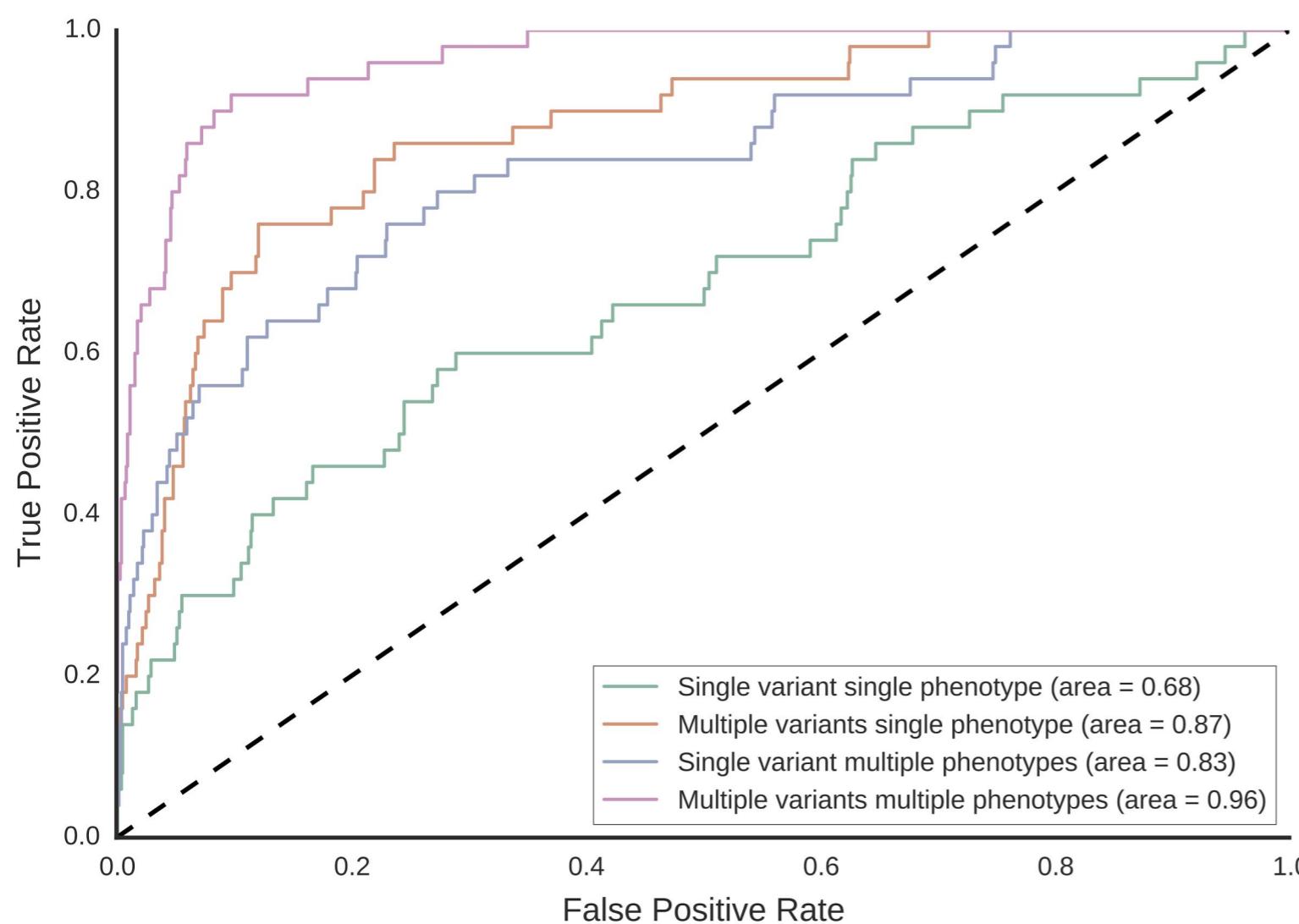


MRP can use GWAS summary statistics

- Results from MRP using raw genotype-phenotype data (x-axis) or GWAS summary statistics (y-axis) are highly correlated

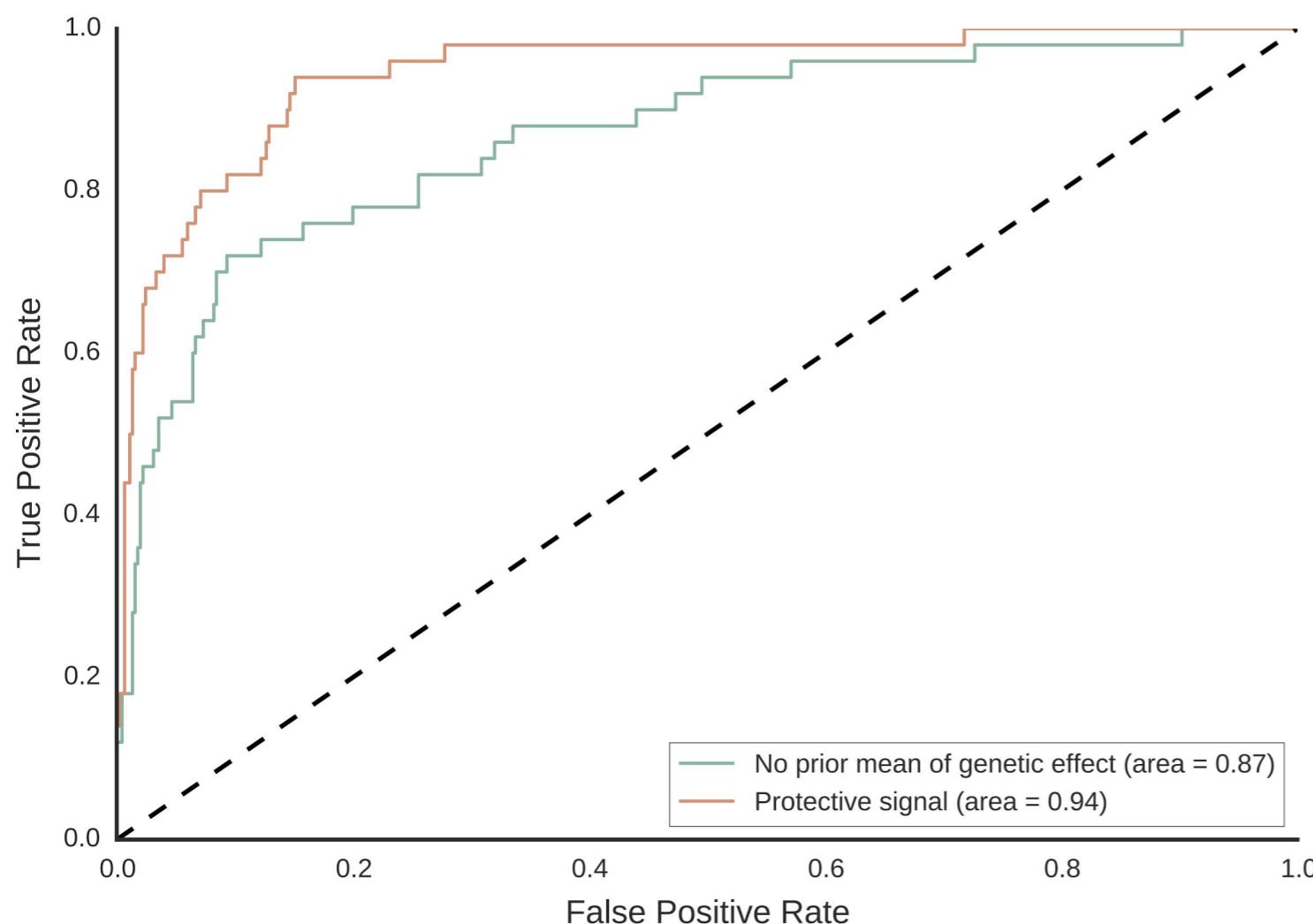


Univariate vs. multivariate analyses



- Using multiple variants and multiple phenotypes (purple line) improves power to detect associations

Protective associations



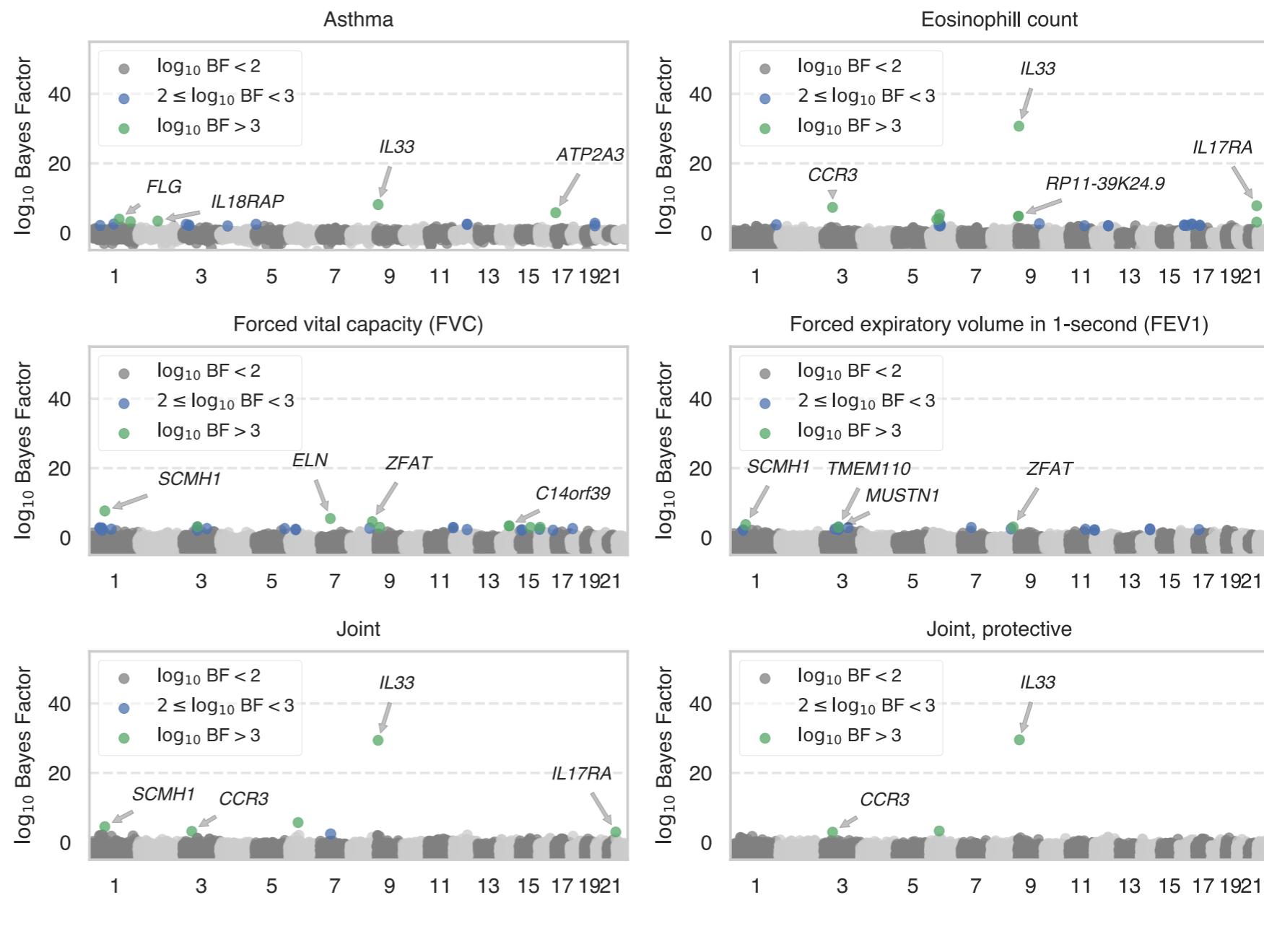
- Adding a prior mean for genetic effects can boost power to identify protective associations

Asthma application

- Applied MRP to asthma and related phenotypes from UK Biobank

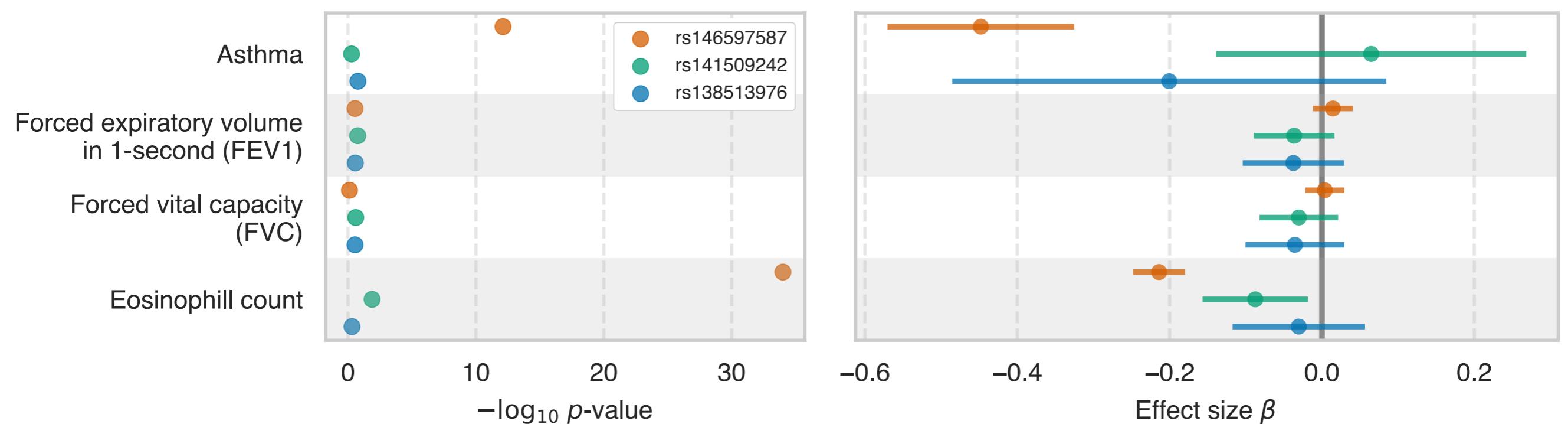
- Identified association for rare variants in *IL33* that protect against asthma

- Also observed moderate evidence for protective associations for drug target *CCR3*



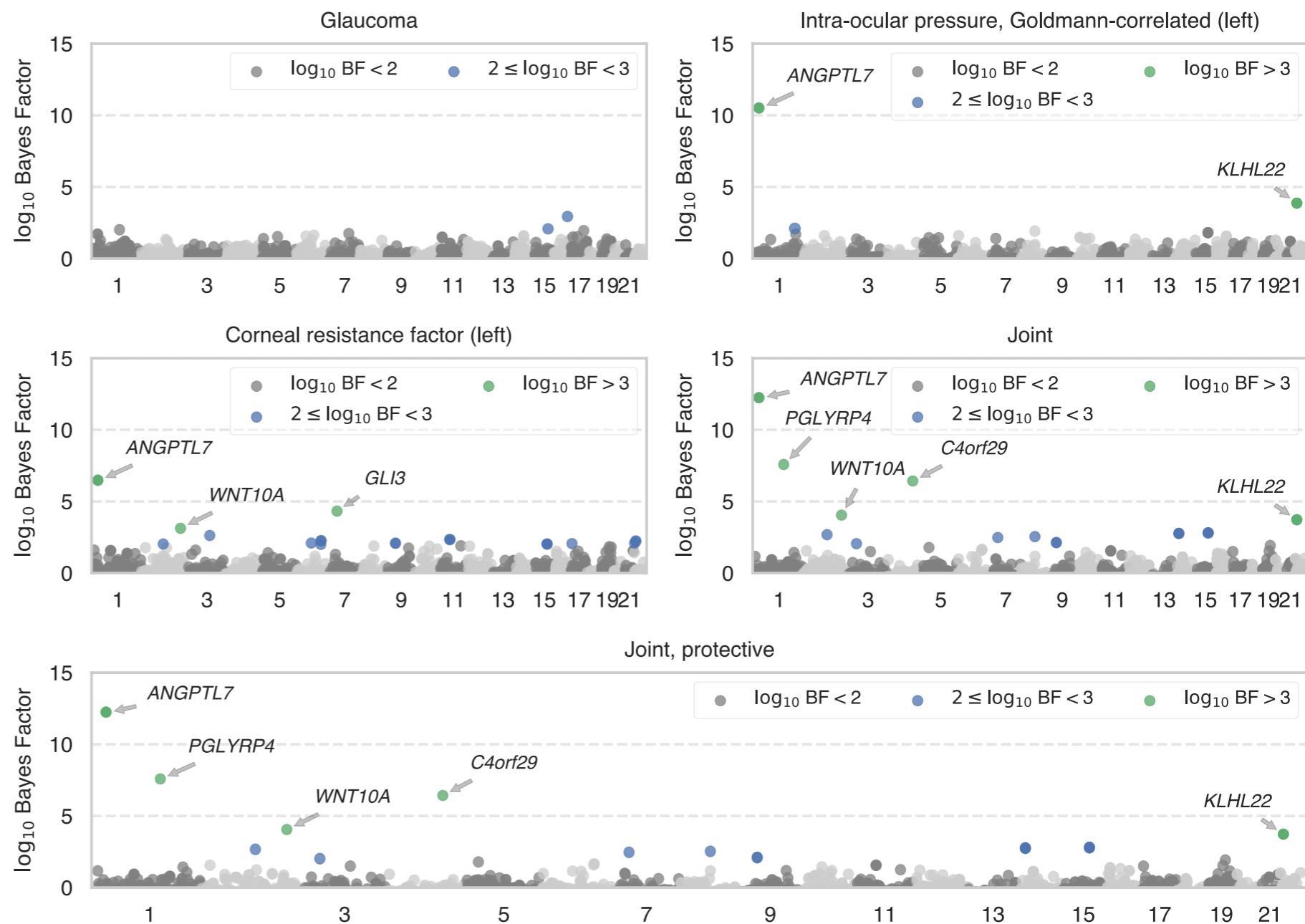
Asthma application

- Protective association in *IL33* driven by PTV rs146597587 (Smith *et al.* 2017, DeBoever *et al.* 2017)
- PTV also has strong effect on eosinophil counts



Glaucoma application

- Applied MRP to glaucoma and two related traits from the UK Biobank
- Identified protective association for *ANGPTL7*
- *ANGPTL7* is unregulated in glaucoma and the protein expression is unique to the vitreous humor of the eye



Conclusions

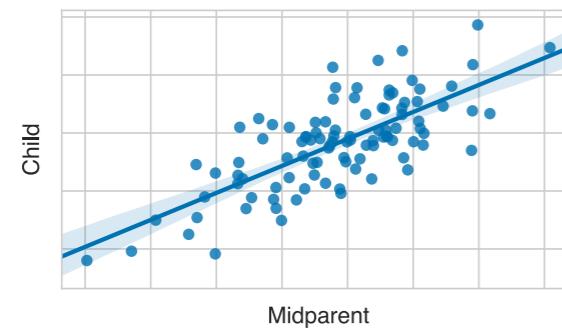
- MRP
 - Can use raw genotype/phenotype data or GWAS summary statistics
 - Uses information from multiple variants and phenotypes to boost power to detect rare variant associations
 - Can prioritize associations that protect against disease

Multivariate Polygenic Mixture Models

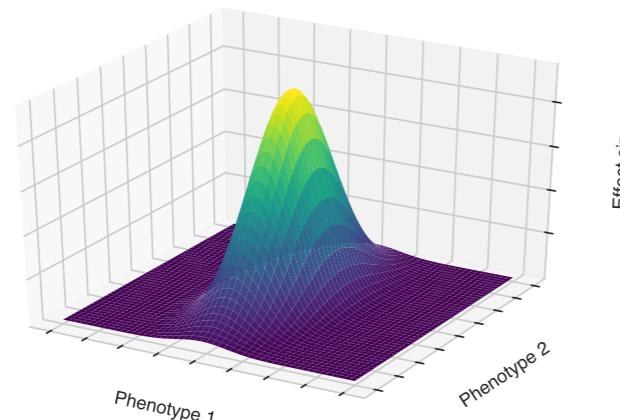
Motivation

- Most common diseases are polygenic
- Large number of phenotypes and loci tested by GWAS
- Opportunity to estimate genetic parameters

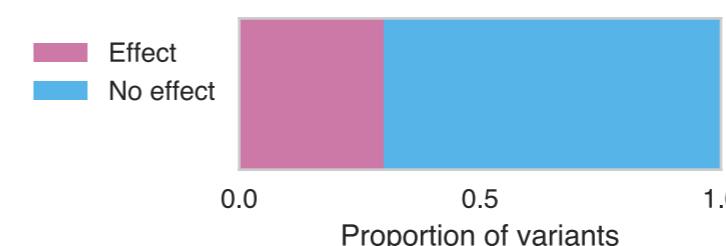
- h^2 : heritability



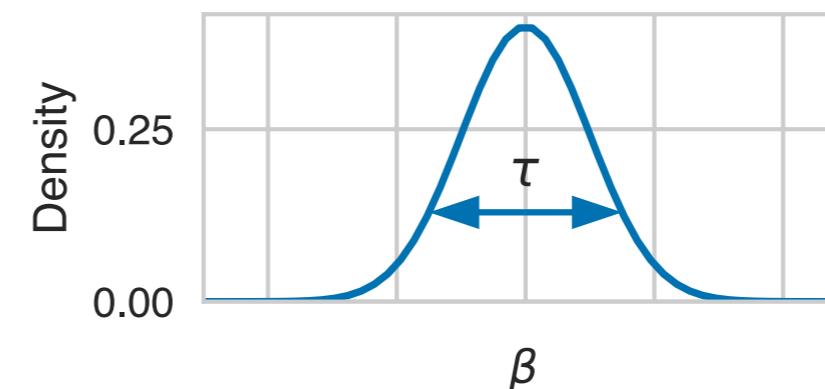
- Ω : genetic correlation



- π : membership



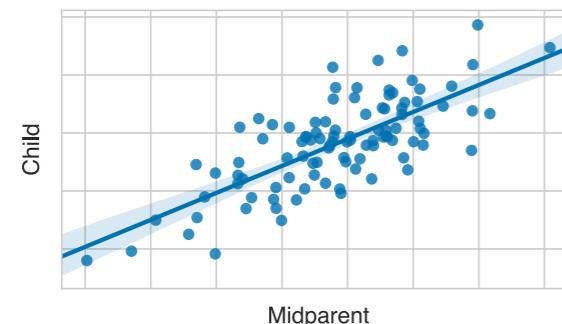
- τ : spread/scale of effects



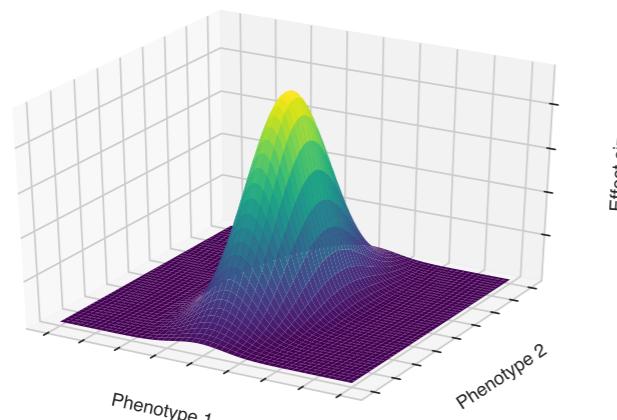
RIVASLAB

Motivation

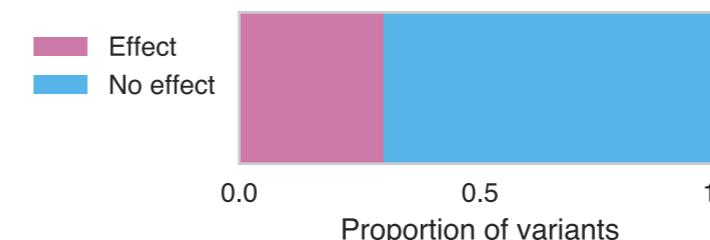
- Most common diseases are polygenic
- Large number of phenotypes and loci tested by GWAS
- Opportunity to estimate genetic parameters
 - h^2 : heritability



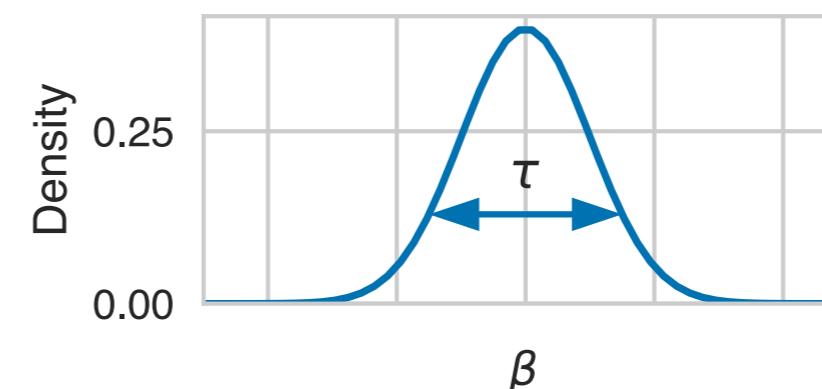
- Ω : genetic correlation



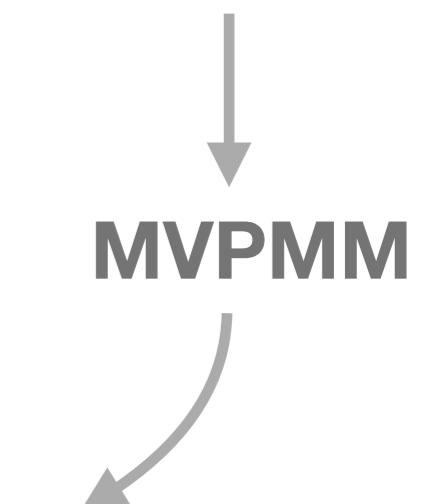
- π : membership



- τ : spread/scale of effects



Phenotype 1		Phenotype 2	
BETA	SE	BETA	SE
1.2	0.03	1.15	0.04
1.5	0.9	1.3	0.1
-1.1	0.006	-1.2	0.05
...



Estimating genetic parameters

- Model GWAS summary statistics as generated from one of two components
 - Null component with correlated errors
 - Non-null component with correlated genetics effects and errors

$\hat{\beta}_i$: regression effect size for locus i

$\hat{\sigma}_i$: regression SE for locus i

Null component

$$\hat{\beta}_i \sim \text{MVN}(0, \Sigma_{\Theta i})$$

Non-null component

$$\hat{\beta}_i \sim \text{MVN}(0, \Sigma_{\Theta i} + \Sigma_{\Omega})$$

Correlated errors

$$\Sigma_{\Theta i} = \text{diag}(\hat{\sigma}_i) \cdot \Theta \cdot \text{diag}(\hat{\sigma}_i)$$

Correlated genetic effects

$$\Sigma_{\Omega} = \text{diag}(\tau) \cdot \Omega \cdot \text{diag}(\tau)$$



RIVASLAB

Can we use genetic parameters to evaluate different phenotyping methods for population biobanks?

What are the estimates of genetic parameters across diseases?

UK Biobank

- Diverse data collected on ~500,000 participants
- Array genotyping for ~800,000 variants
- Hospital in-patient records, cancer registry, death registry

⌚📁 J45 Asthma	-
J45.0 Predominantly allergic asthma	79
J45.1 Nonallergic asthma	1
J45.8 Mixed asthma	2
J45.9 Asthma, unspecified	2659
J46 Status asthmaticus	300
J47 Bronchiectasis	631

<http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=41202>

- Verbal questionnaire data

UK Biobank. General Everywhere : Interview. Cancer * TRAINING/DEMONSTRATION

Select/Search for Cancer

Cancer

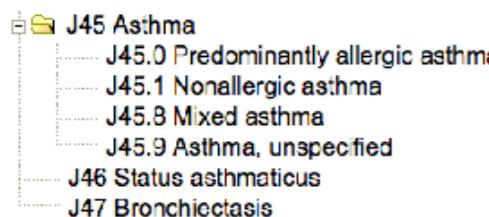
Cancer Type	Count
breast cancer	+ 3
ear/nose/throat cancer	+ 3
gastrointestinal cancer	+ 8
genital tract cancer	+ 2
haematological malignancy	+ 5
neurological system cancer	+ 5
other cancer	+ 9
respiratory/ intrathoracic cancer	+ 6
skin cancer	+ 2
urinary tract cancer	+ 3

Search Enter text to search

< Prev Help ----- Lock Next >

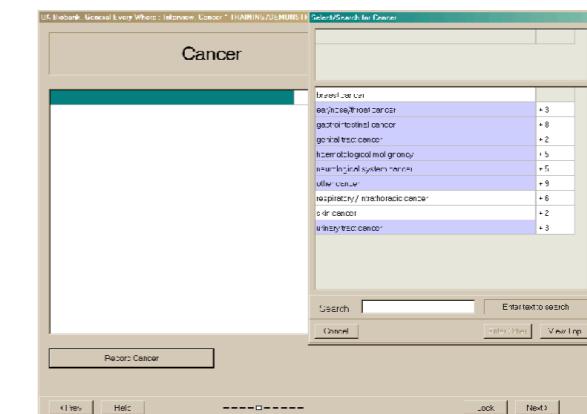
<http://biobank.ctsu.ox.ac.uk/crystal/label.cgi?id=100074>

How can we use genetic parameters to evaluate different phenotyping methods for population biobanks?



-
79
1
2
2659
300
631

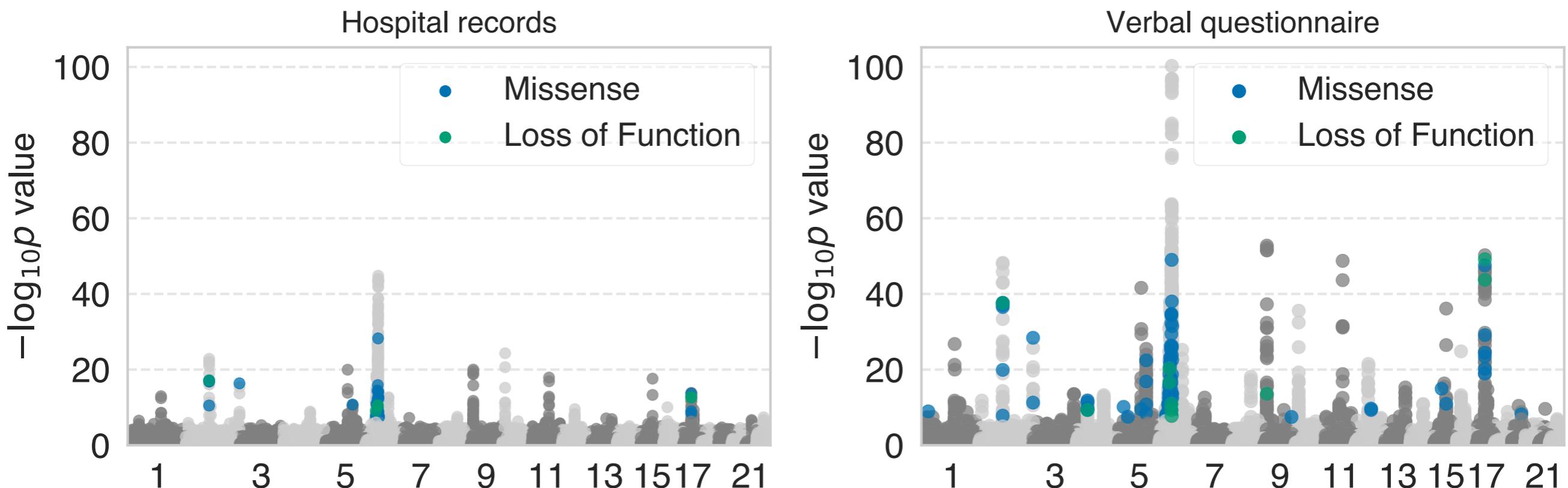
VS.



What are the estimates of genetic parameters across diseases?

Verbal questionnaire vs. hospital records

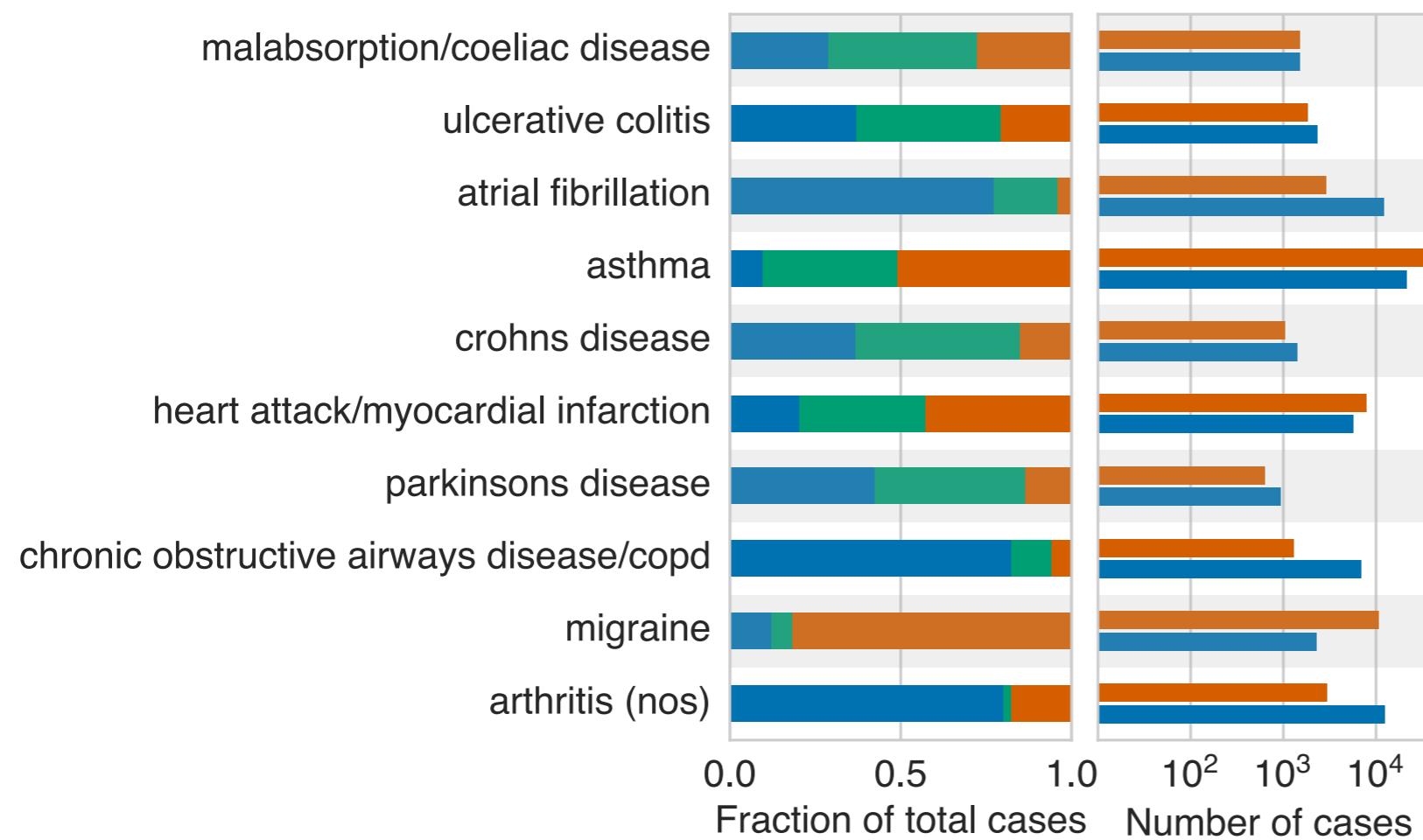
- Asthma cases defined according to either hospital records **or** questionnaire data



Verbal questionnaire vs. hospital records

- Genetic correlation for phenotypes defined using either hospital and registry records **or** questionnaire data

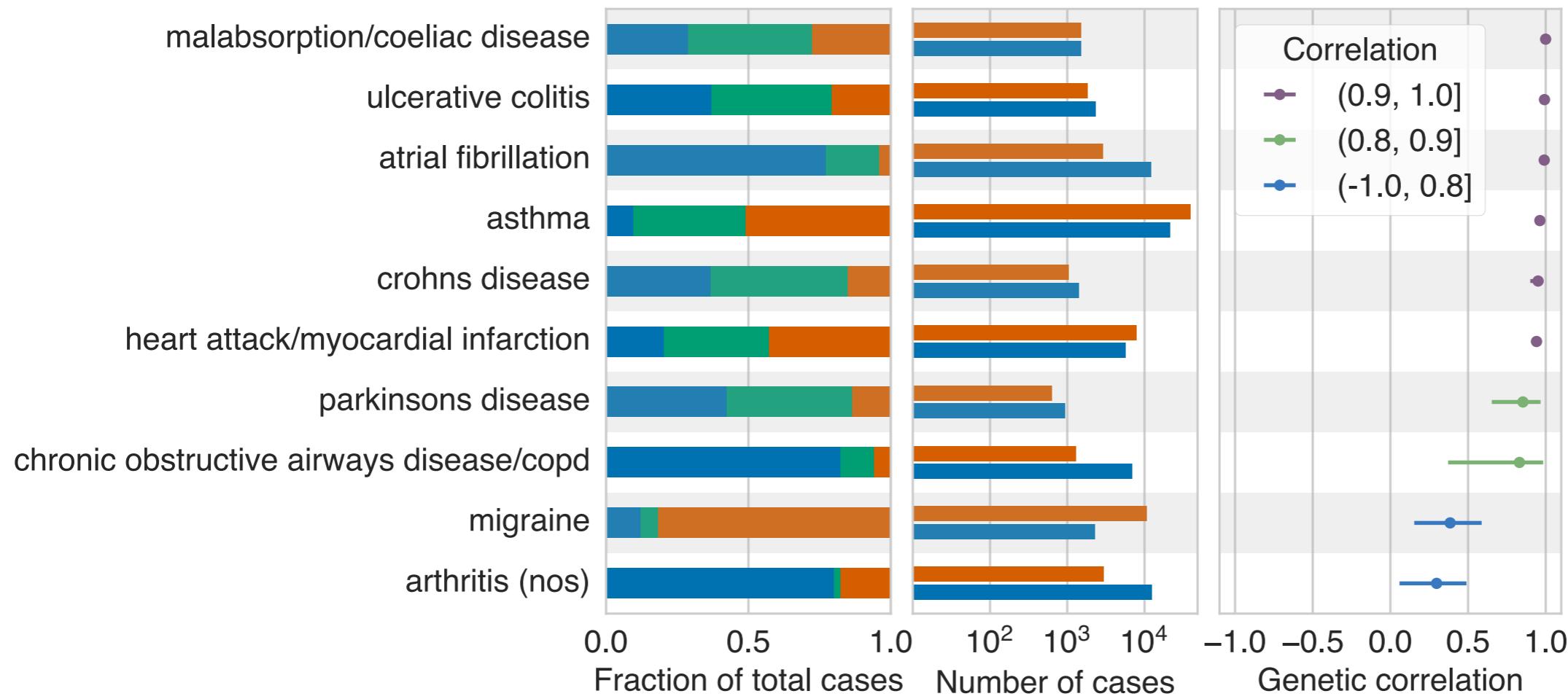
Hospital Shared Questionnaire



Verbal questionnaire vs. hospital records

- Genetic correlation for phenotypes defined using either hospital and registry records **or** questionnaire data

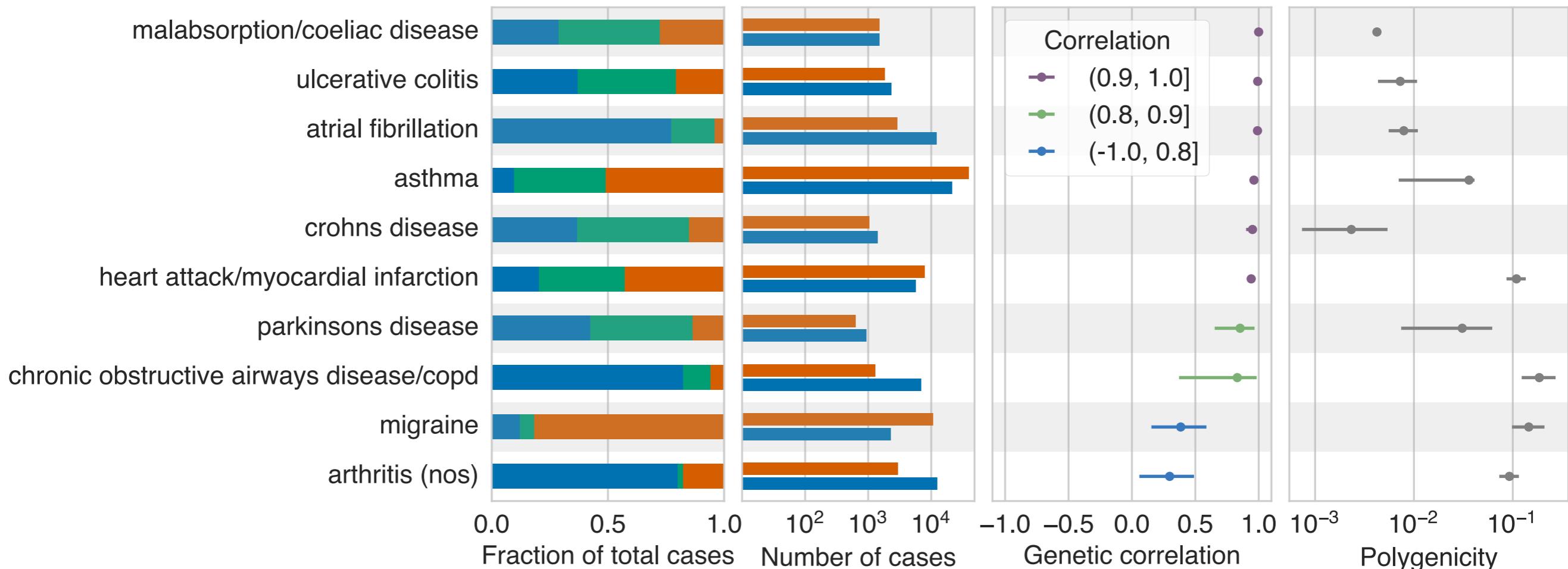
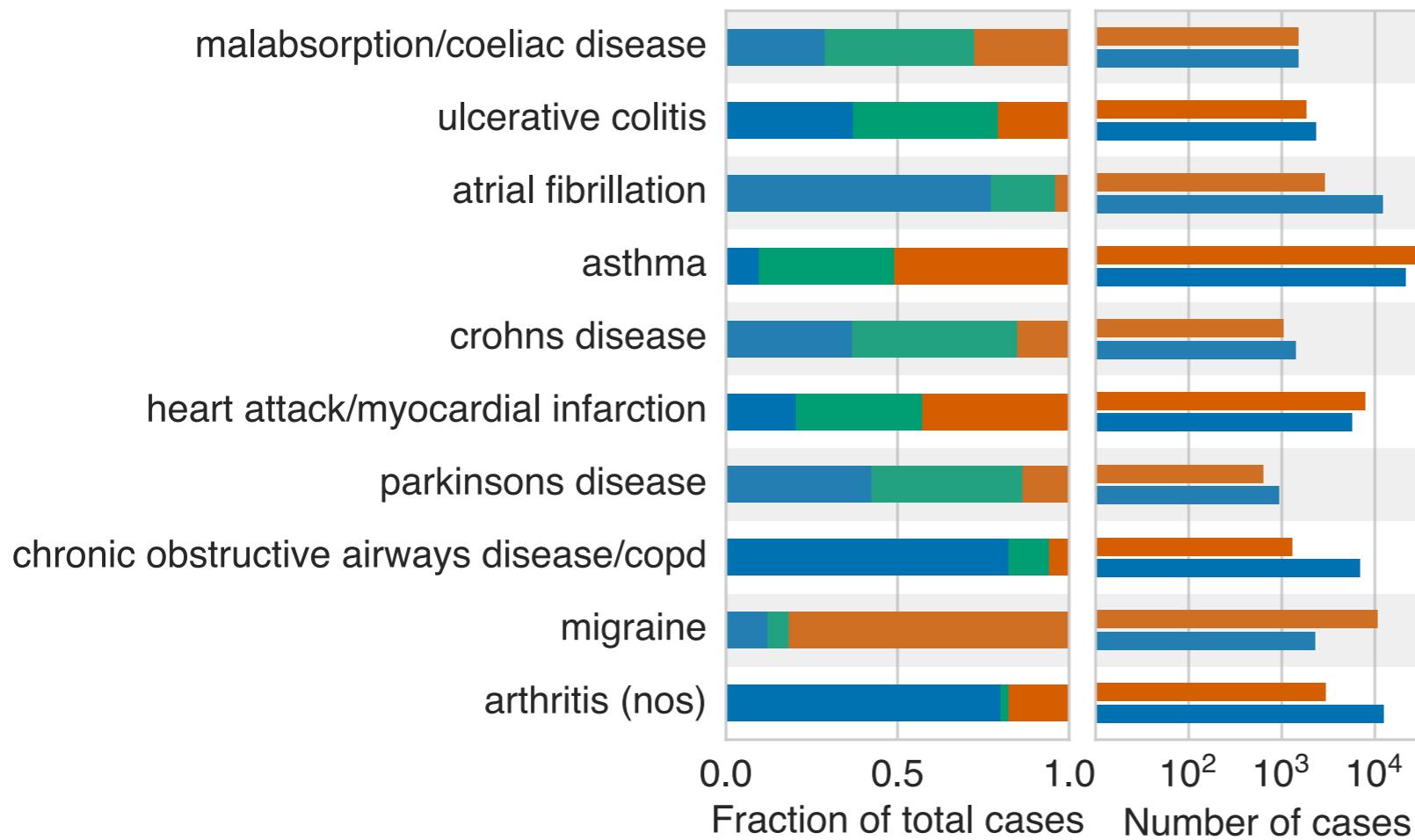
Hospital Shared Questionnaire



Verbal questionnaire vs. hospital records

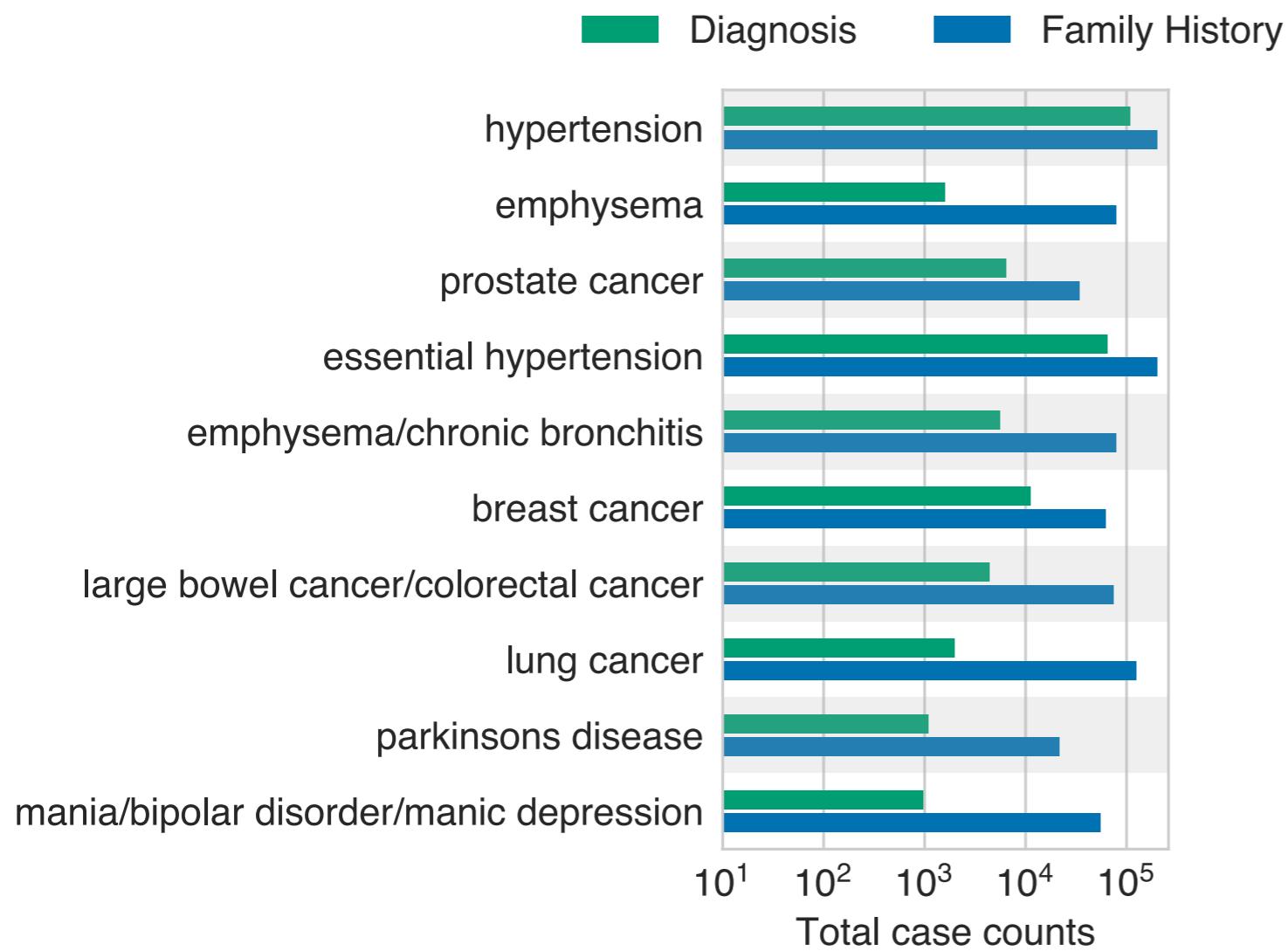
- Genetic correlation for phenotypes defined using either hospital and registry records **or** questionnaire data

Hospital Shared Questionnaire



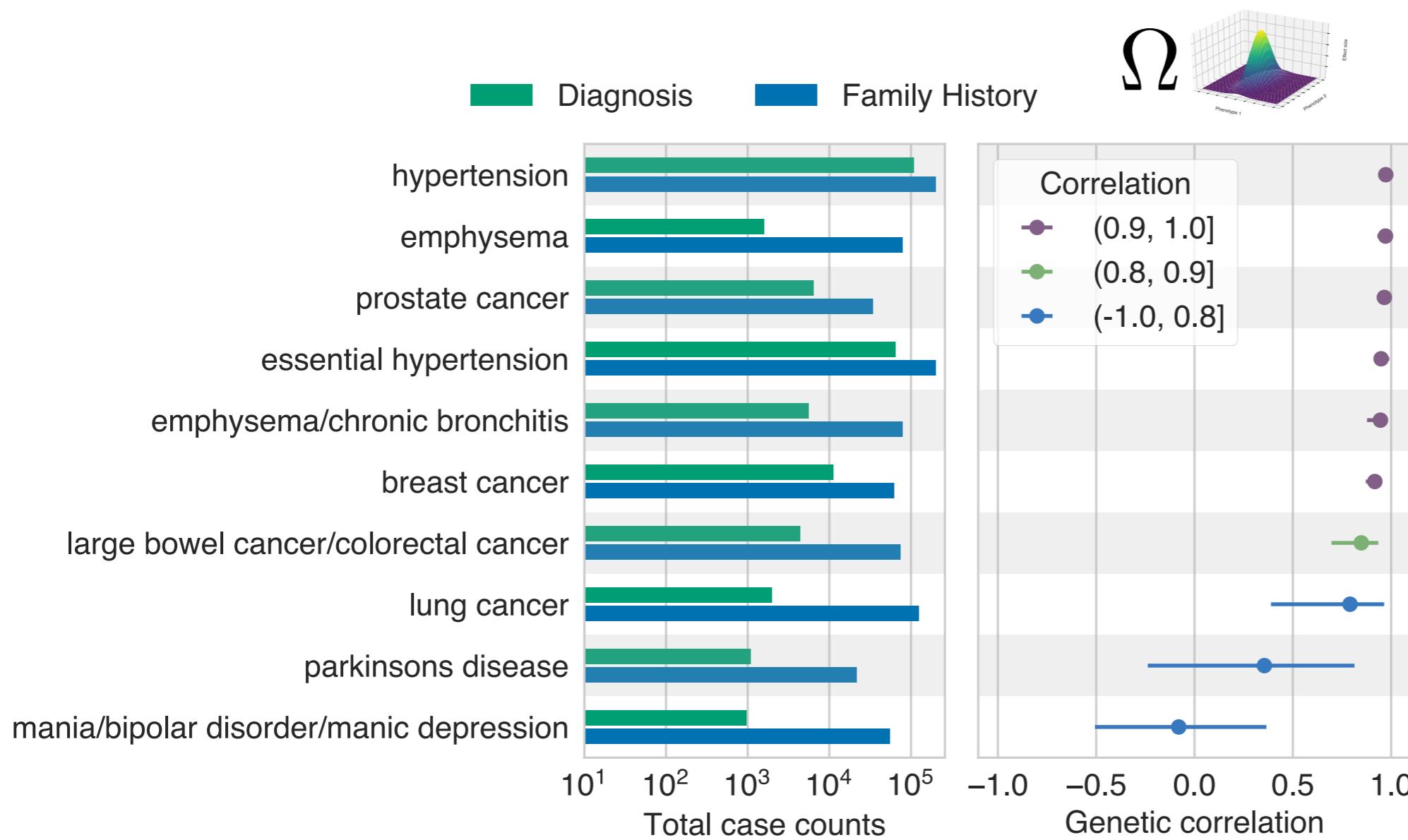
Family history captures disease genetics

- Family history of disease from questionnaire vs. diagnosis



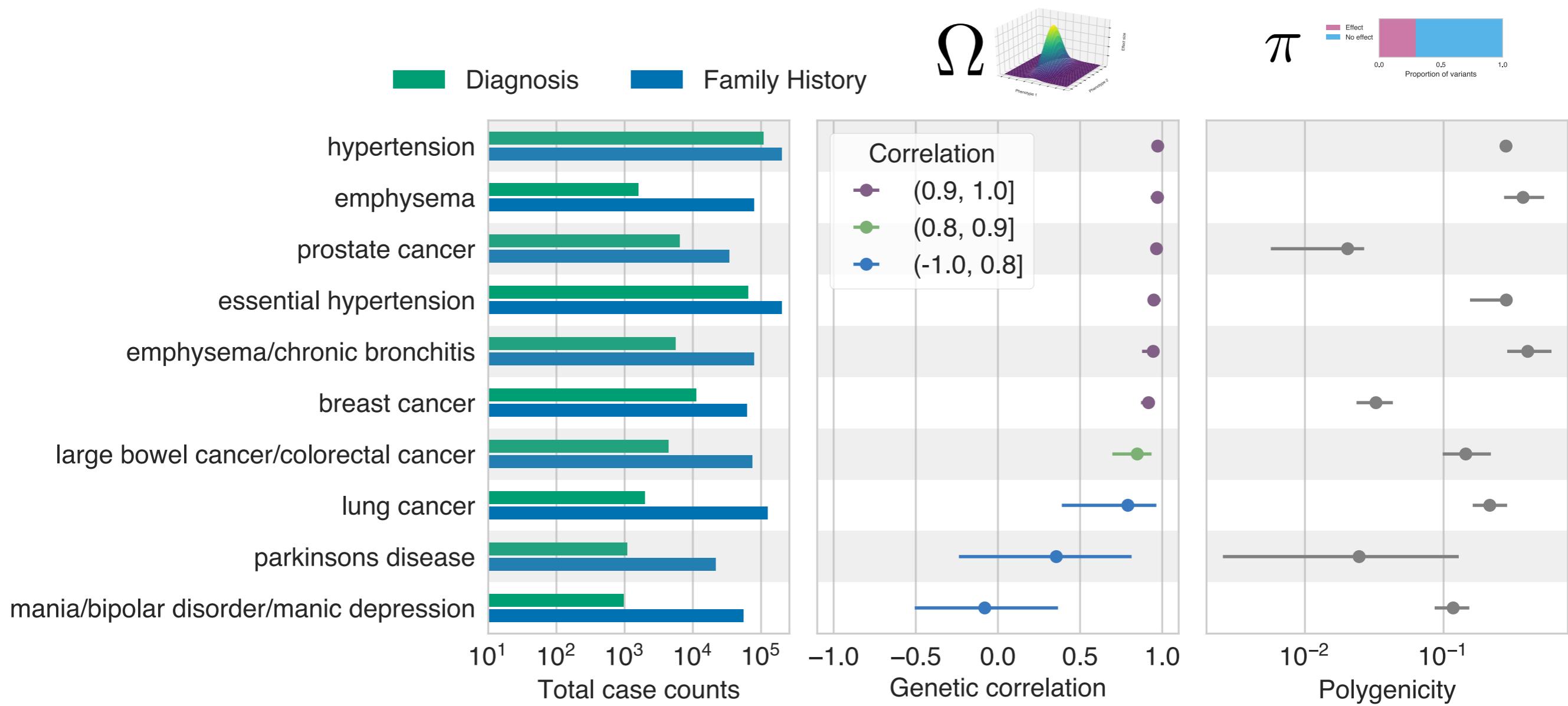
Family history captures disease genetics

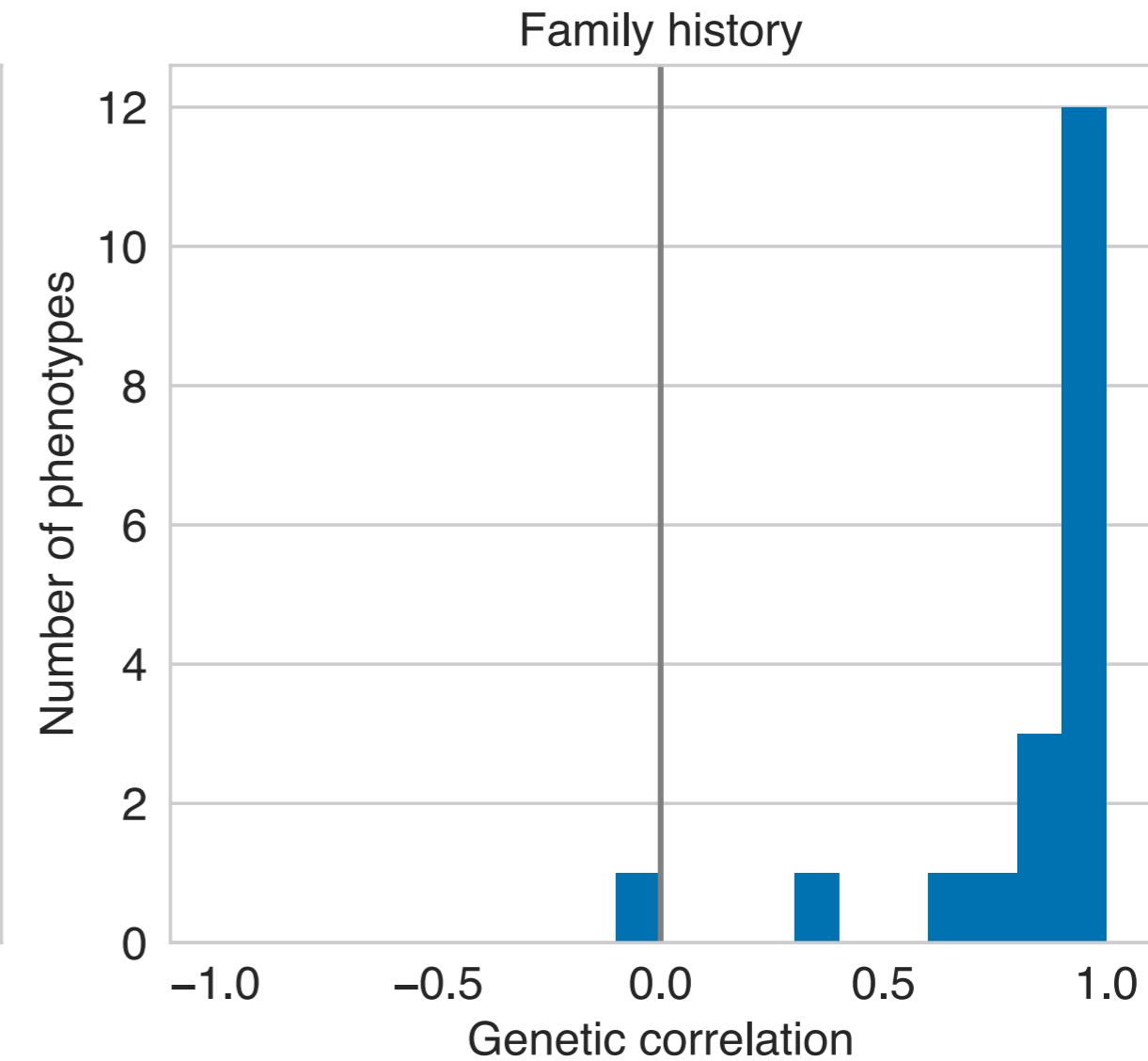
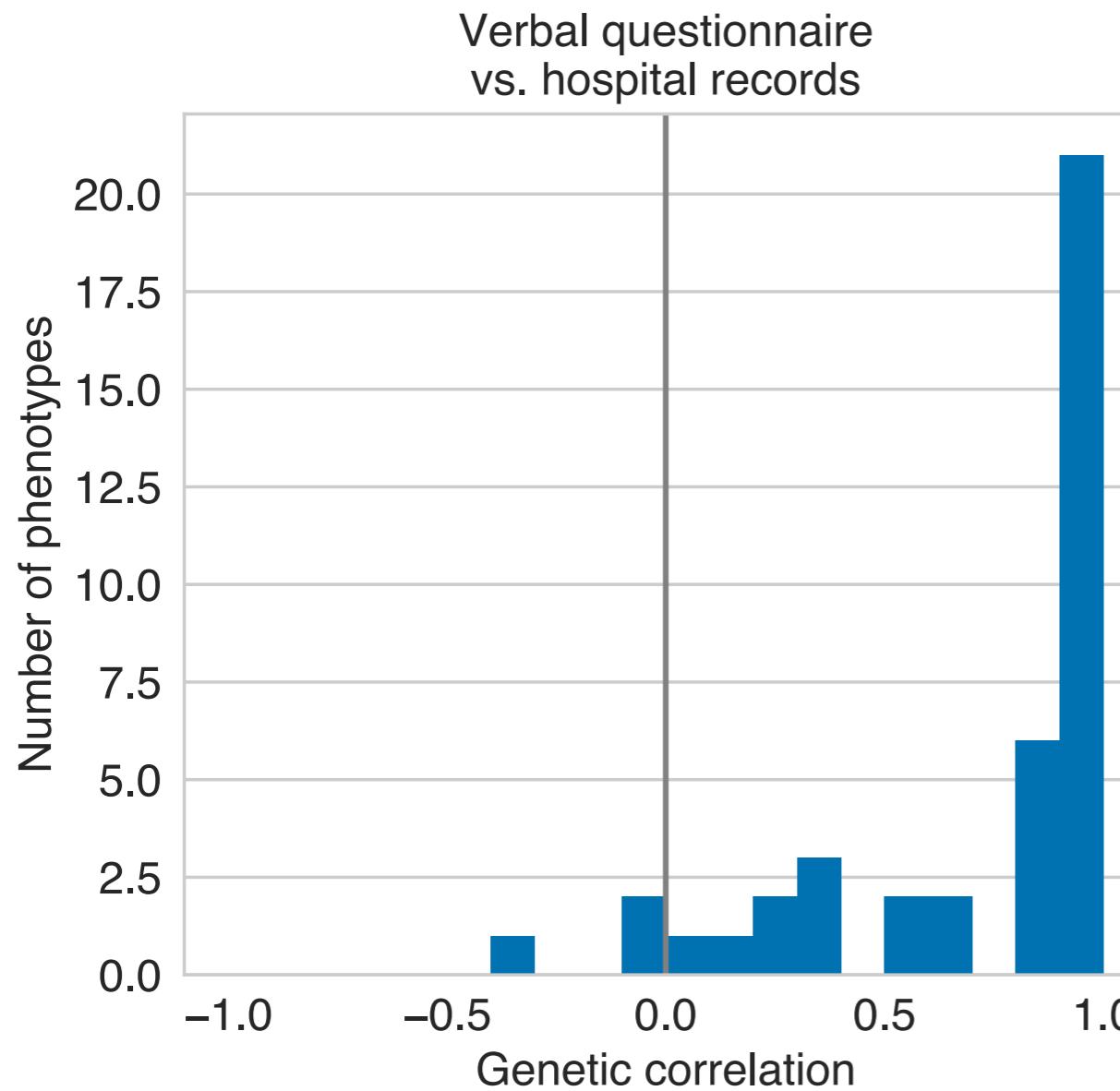
- Family history of disease from questionnaire vs. diagnosis



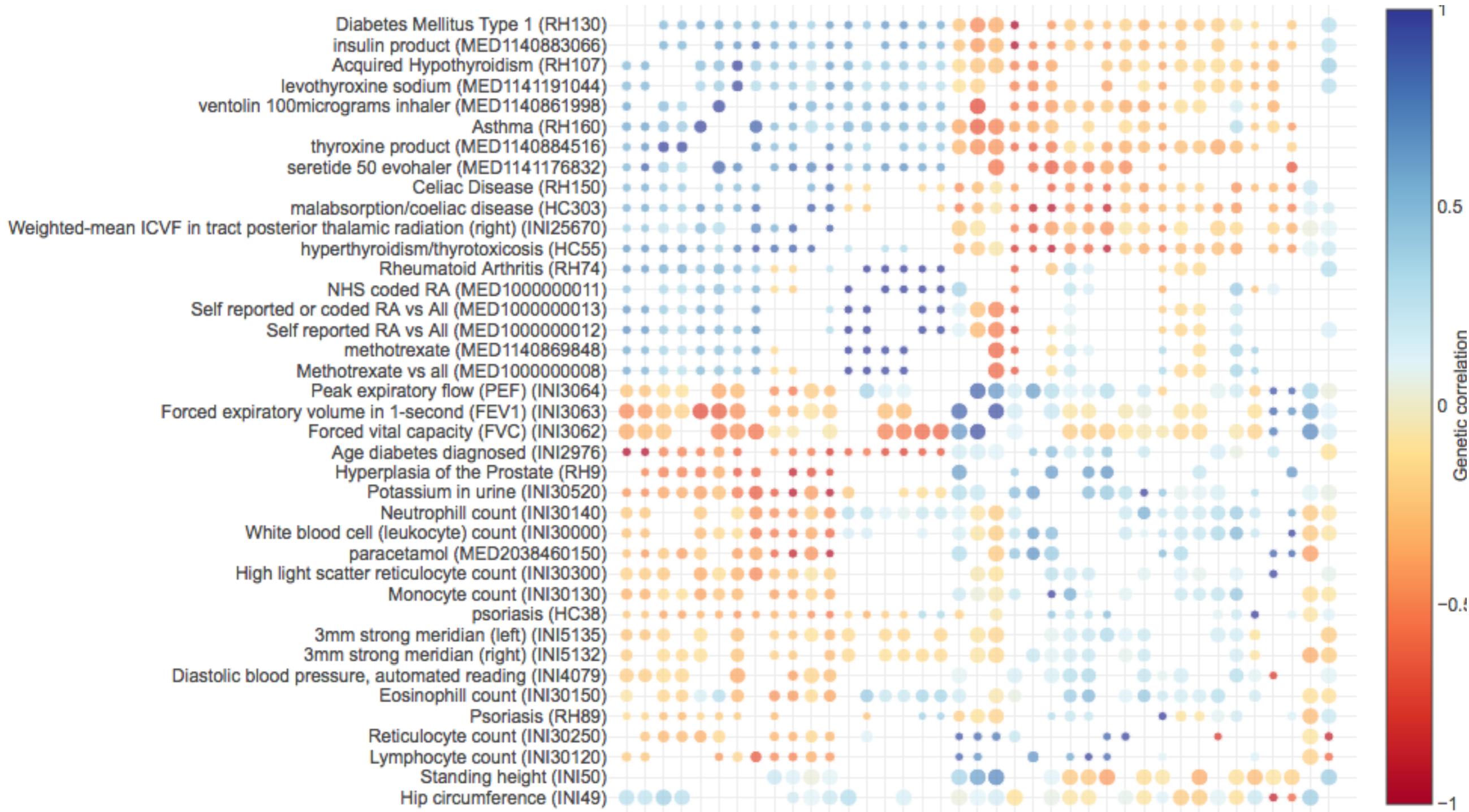
Family history captures disease genetics

- Family history of disease from questionnaire vs. diagnosis



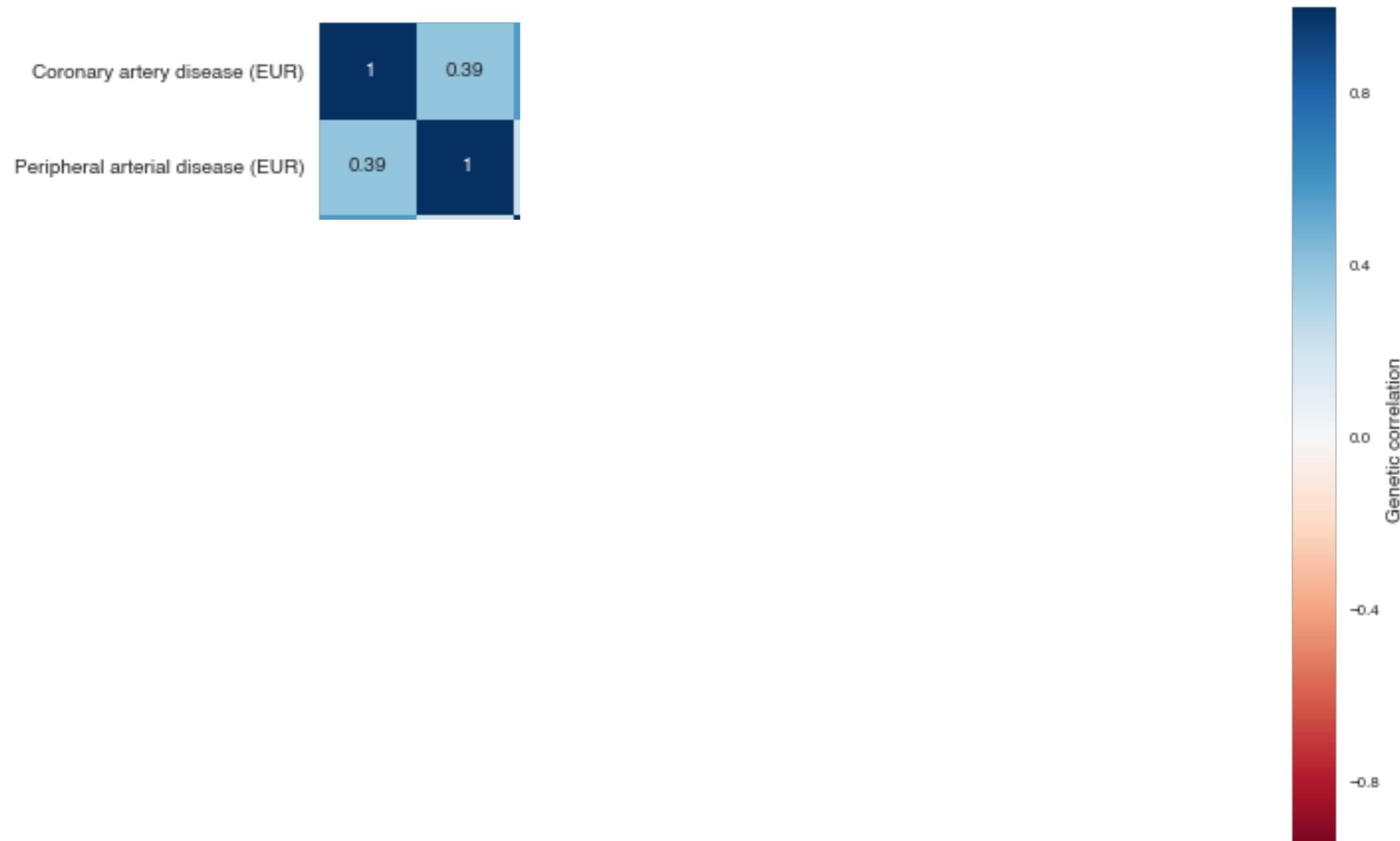


- Compared 41 phenotypes for hospital records vs. verbal questionnaire
- Compared 19 family history phenotypes
- Good agreement for most phenotypes



<https://biobankengine.stanford.edu/gcorr>

Bringing in multiple studies



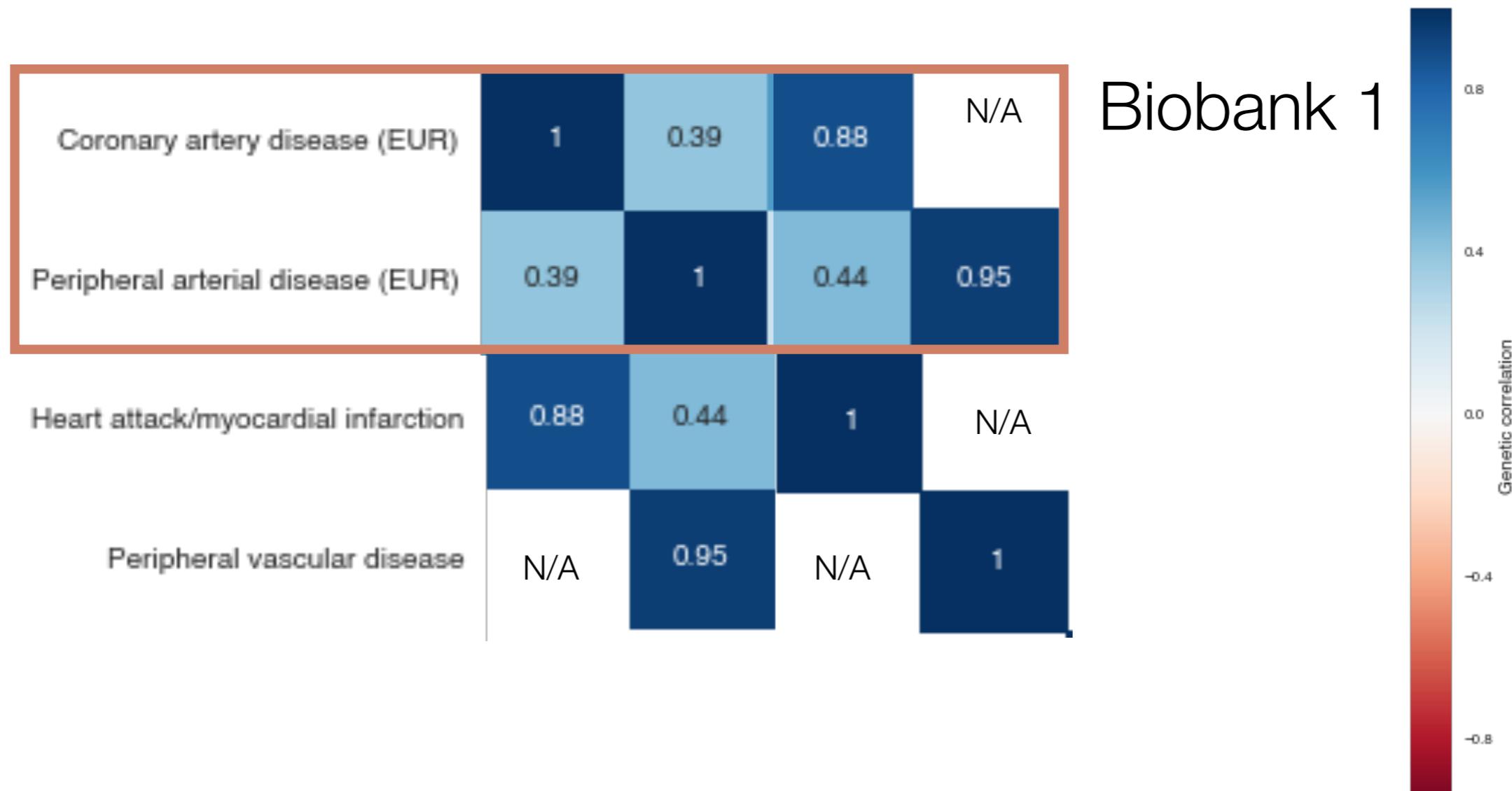
From one biobank to the next...

Bringing in multiple studies



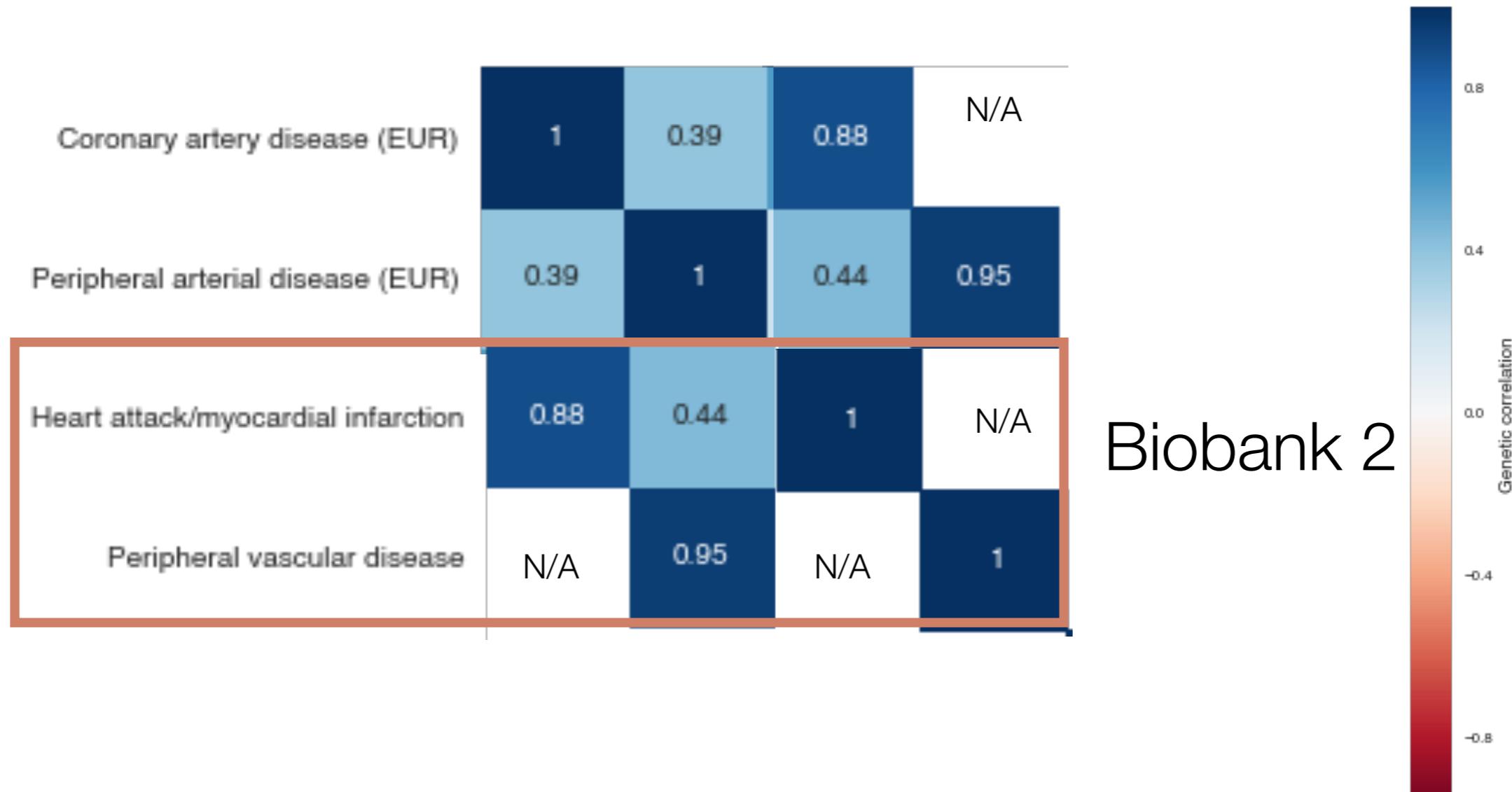
From one biobank to the next...

Bringing in multiple studies



From one biobank to the next...

Bringing in multiple studies



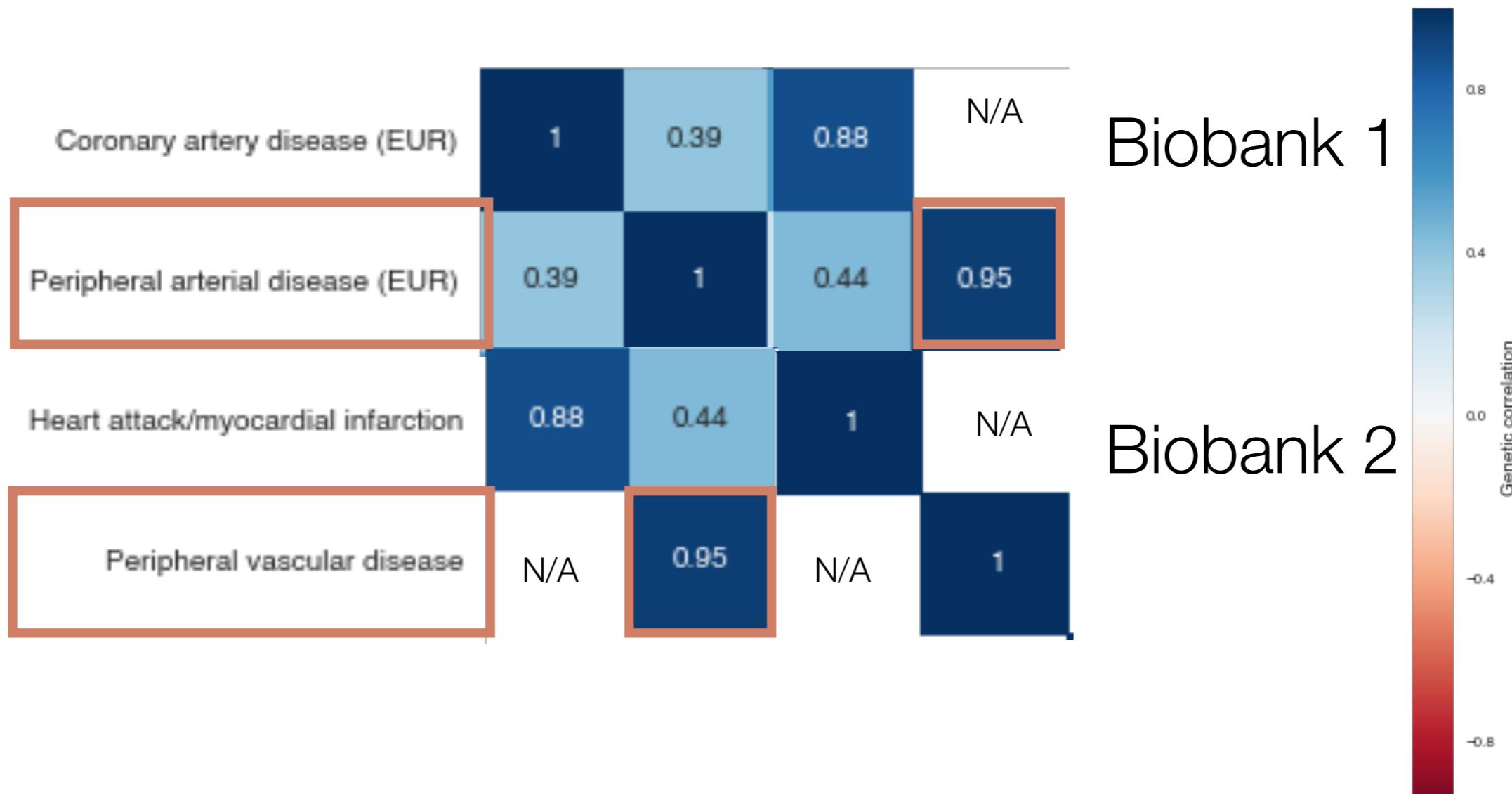
From one biobank to the next...

Coronary artery disease & heart attack/MI



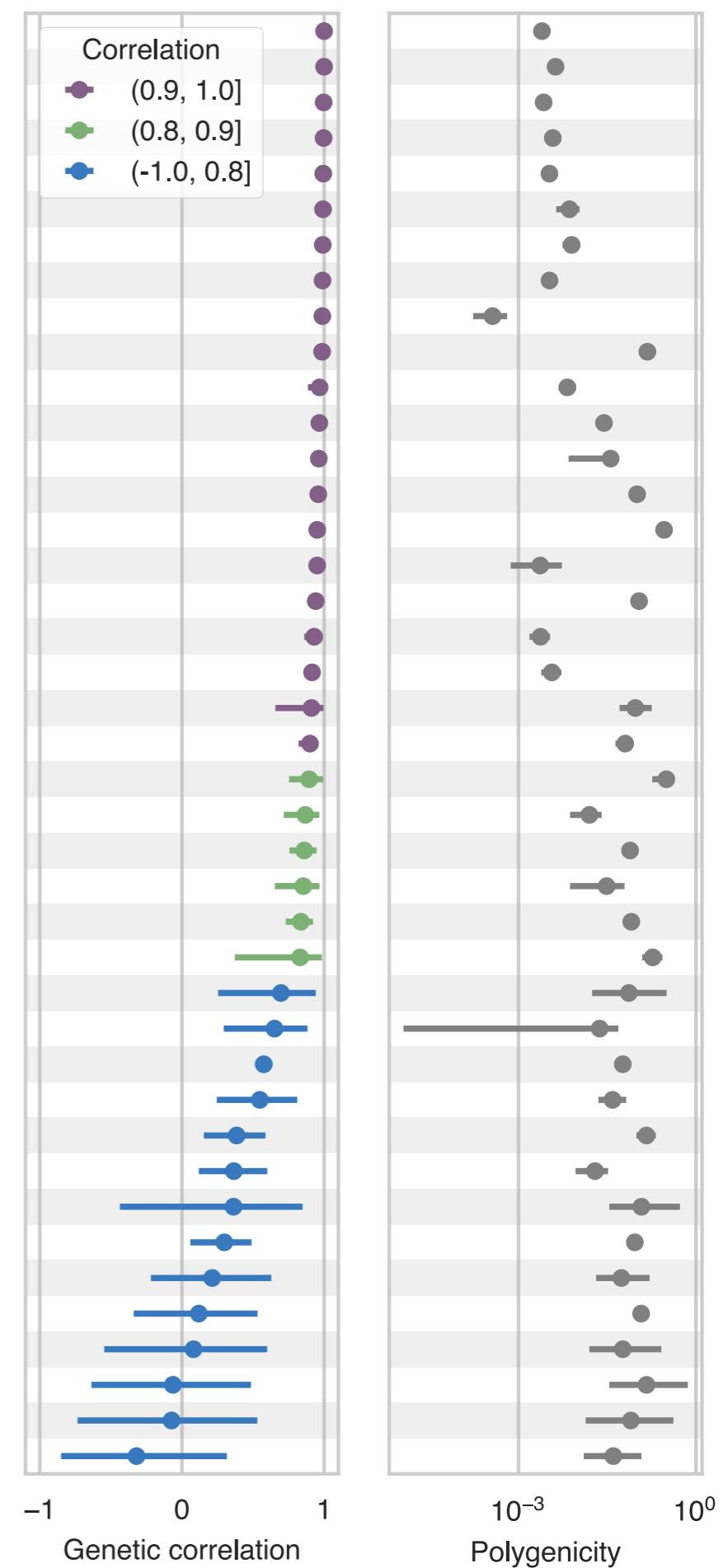
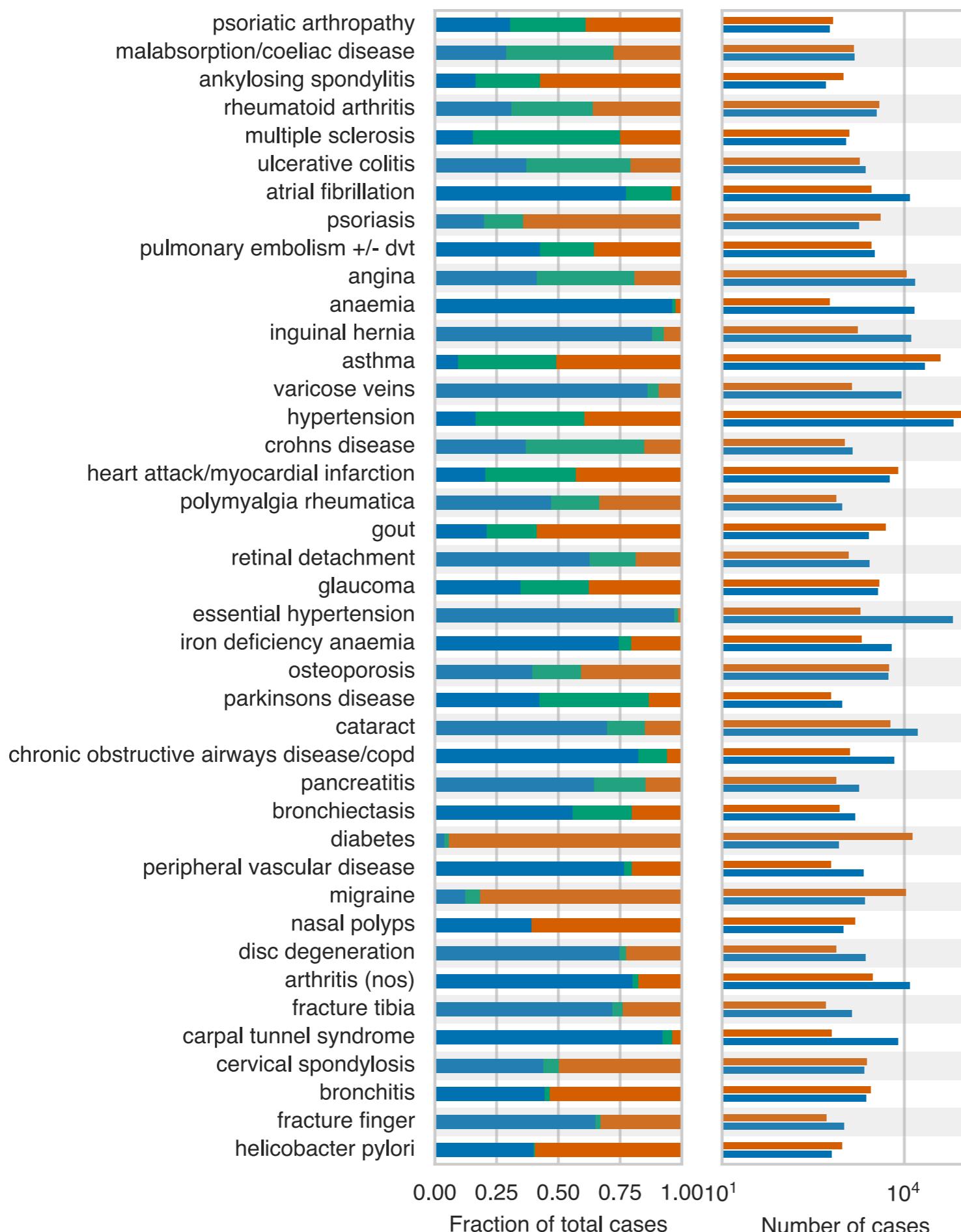
From one biobank to the next...

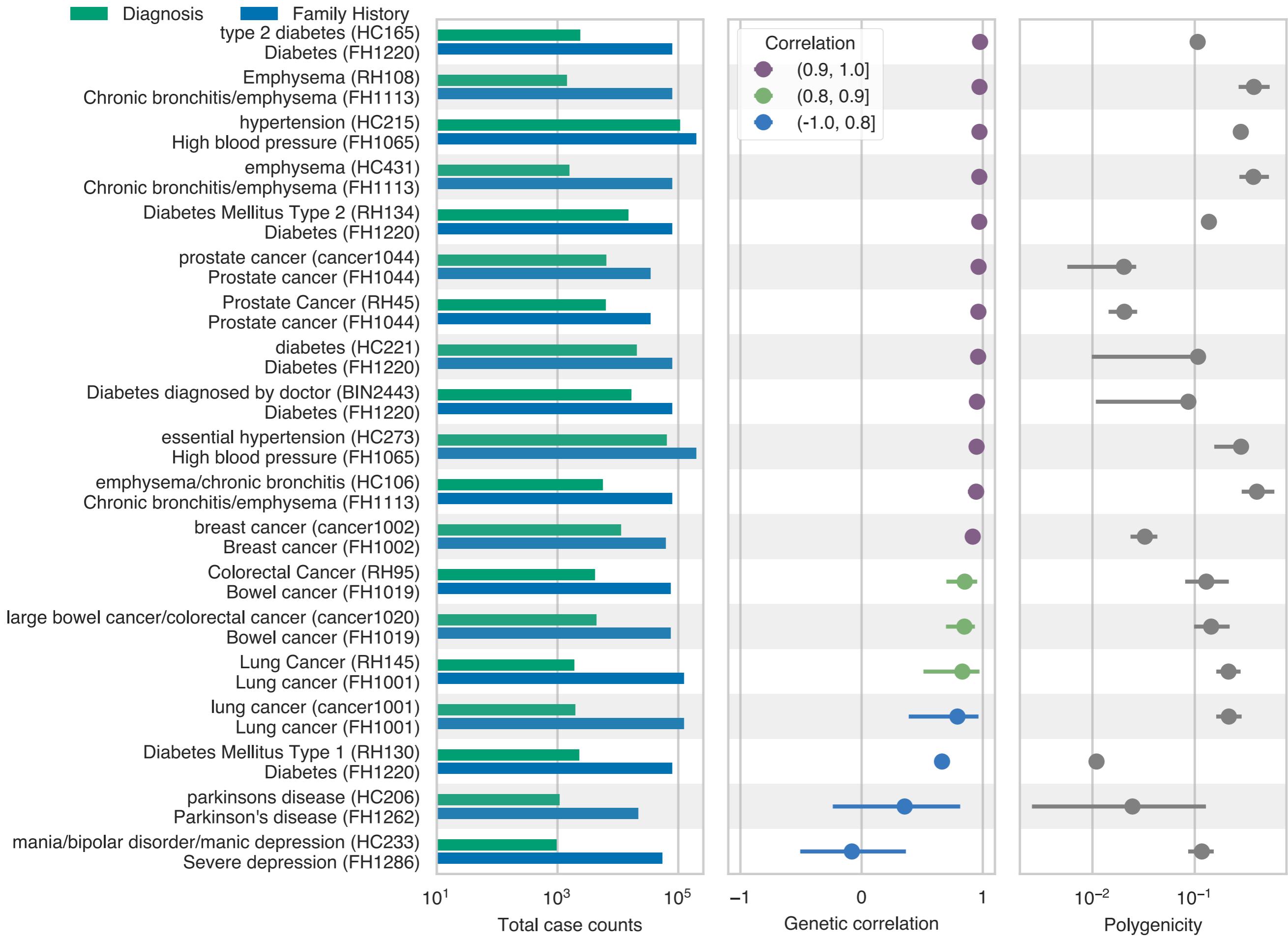
Peripheral arterial disease & Peripheral vascular disease



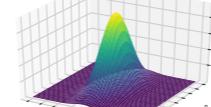
From one biobank to the next...

Hospital Shared Questionnaire

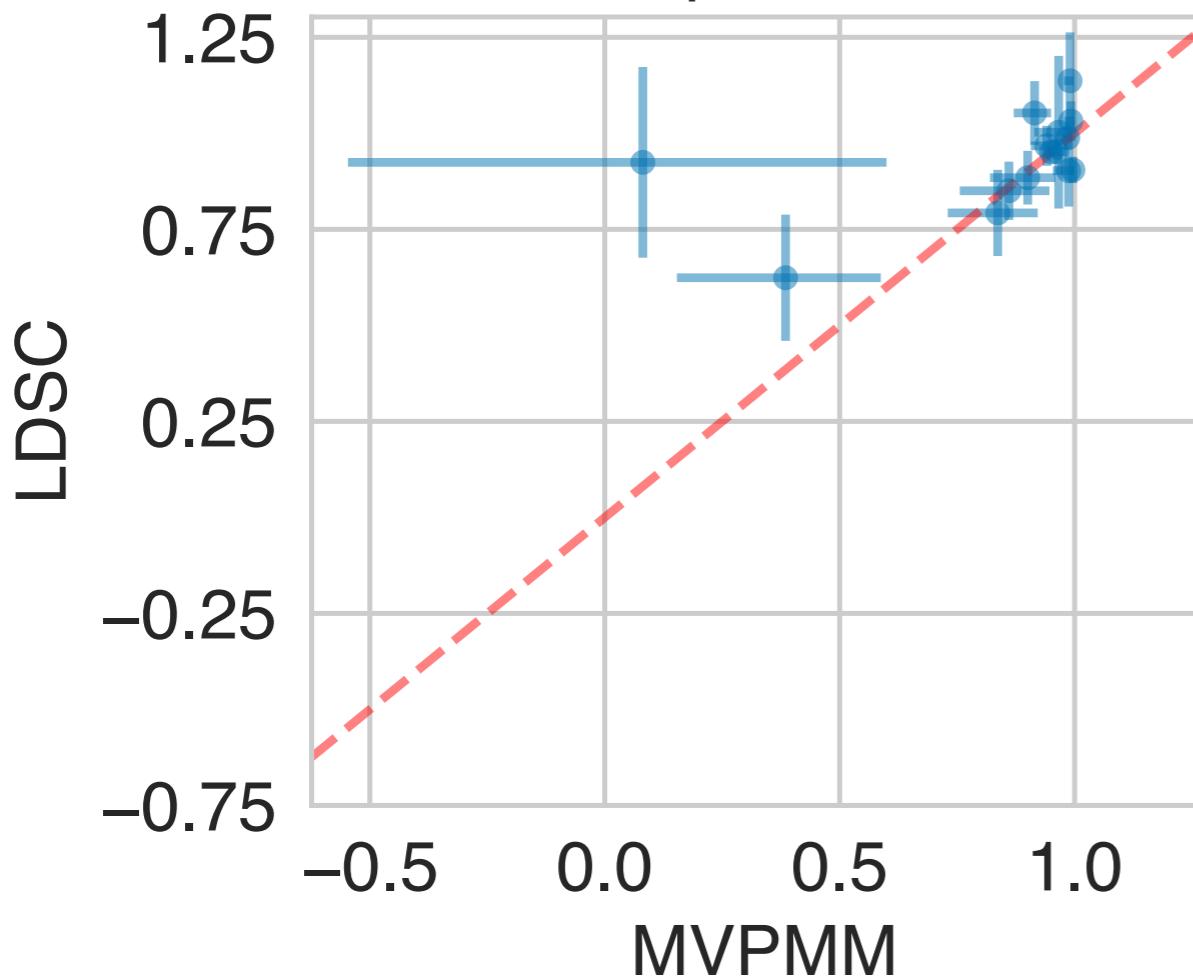




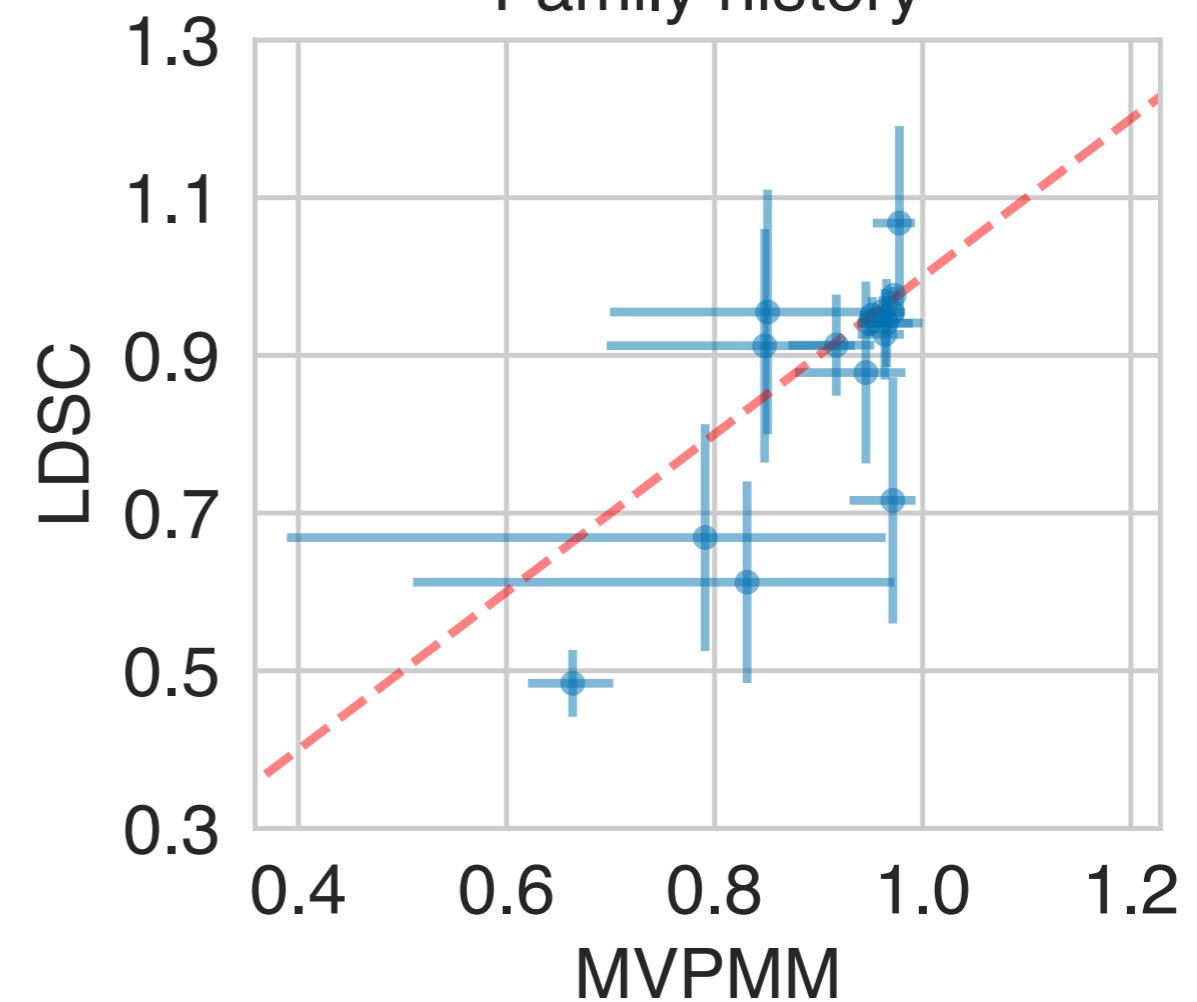
Comparison to LD Score Regression

$$\Omega$$


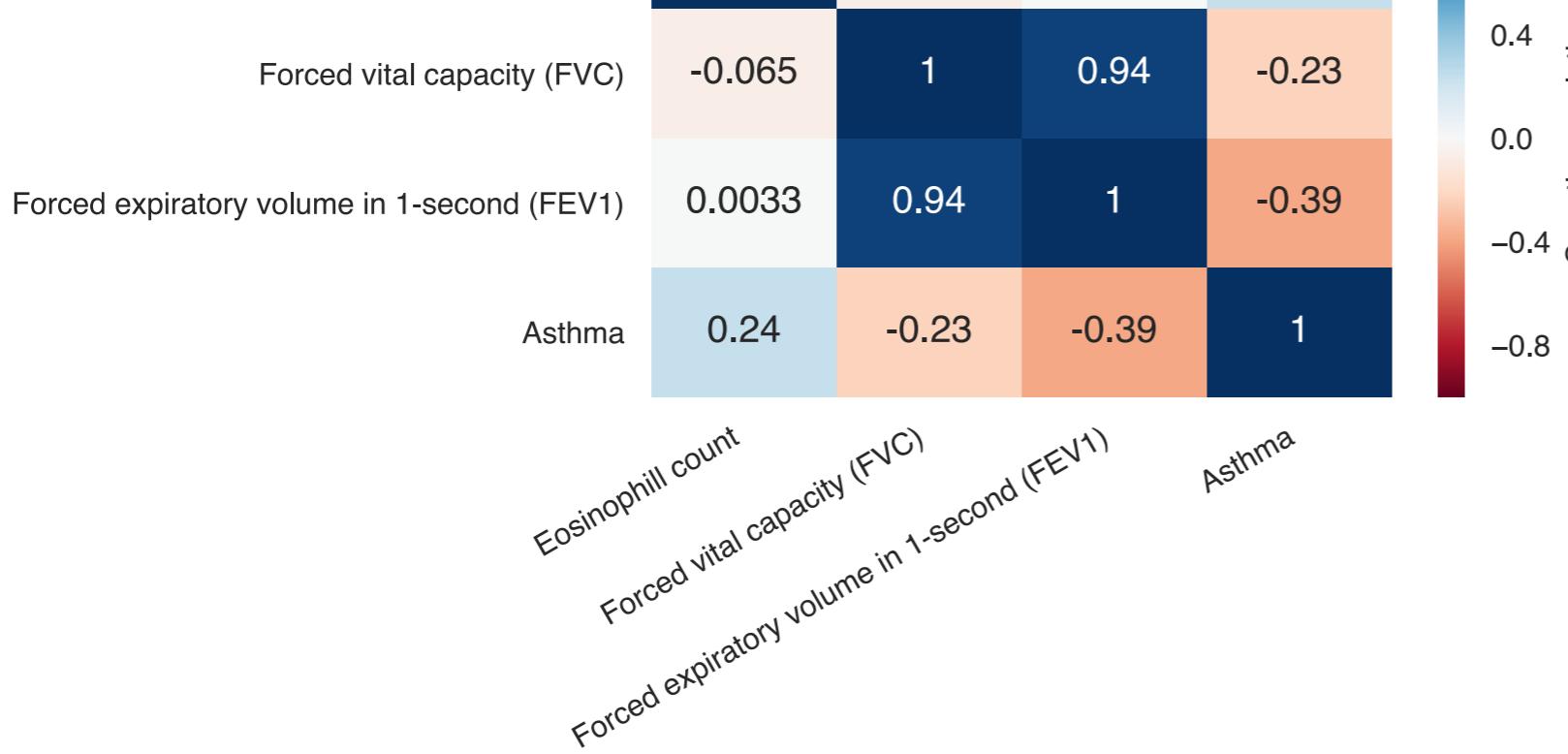
Verbal questionnaire
vs. hospital records



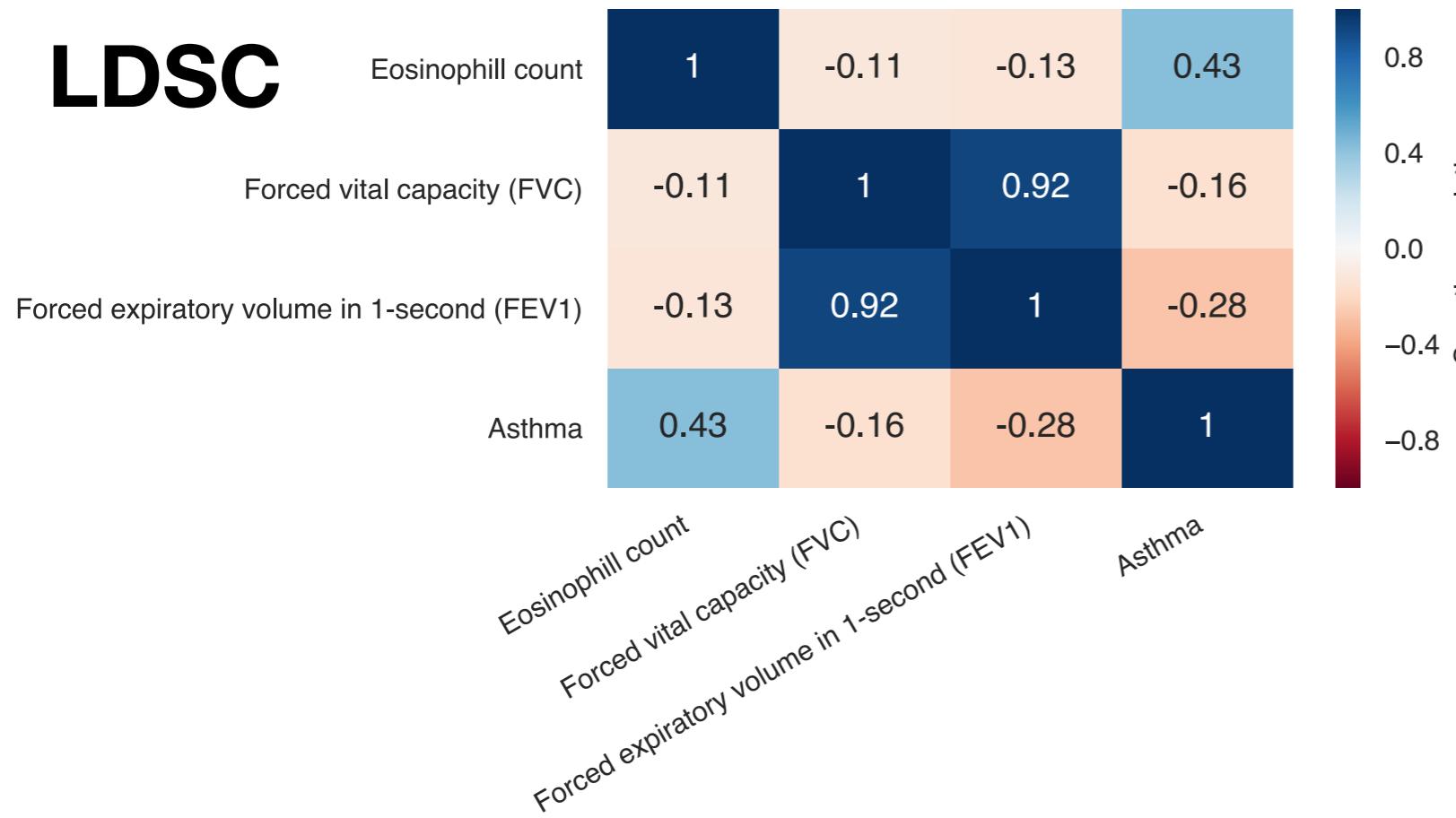
Family history

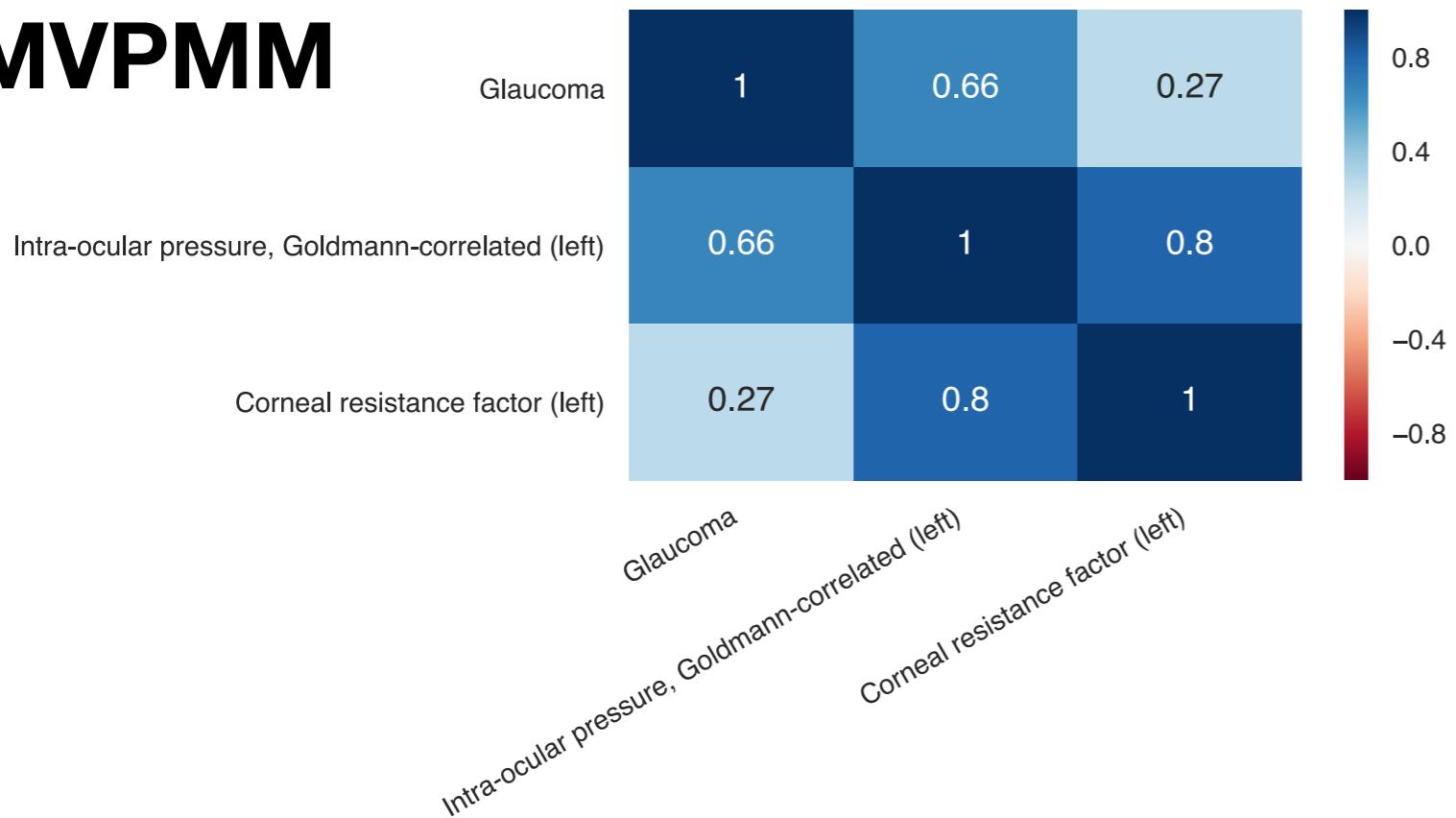


RIVASLAB

MVPMM**Genetic correlation estimates: asthma**

- Genetic correlation estimates from MVPMM and LD score-regression

LDSC

MVPMM

Genetic correlation estimates: glaucoma

- Genetic correlation estimates from MVPMM and LD score-regression

LDSC