

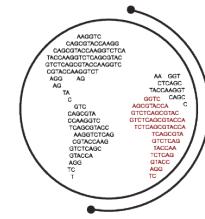


# Stan applications in Human Genetics

|

Manuel A. Rivas

Department of Biomedical Data Science  
Stanford University  
[rivaslab.stanford.edu](http://rivaslab.stanford.edu)



RIVASLAB

# Rare diseases run in families

---

If you have **cystic fibrosis**, what is the risk to:

**Your neighbor?**

**0.04%**

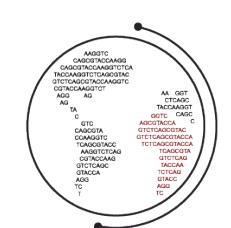
**Your sibling?**

**25%**

**Your identical twin?**

**~99%**

Variation in DNA influences disease risk



**RIVASLAB**

# Common diseases run in families

---

If you have **type 2 diabetes**, what is the risk to:

**Your neighbor?**

**~10%**

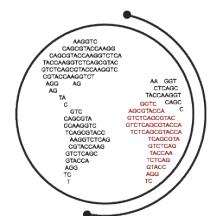
**Your sibling?**

**~30%**

**Your identical twin?**

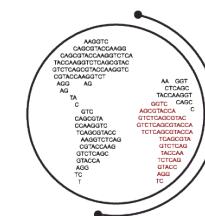
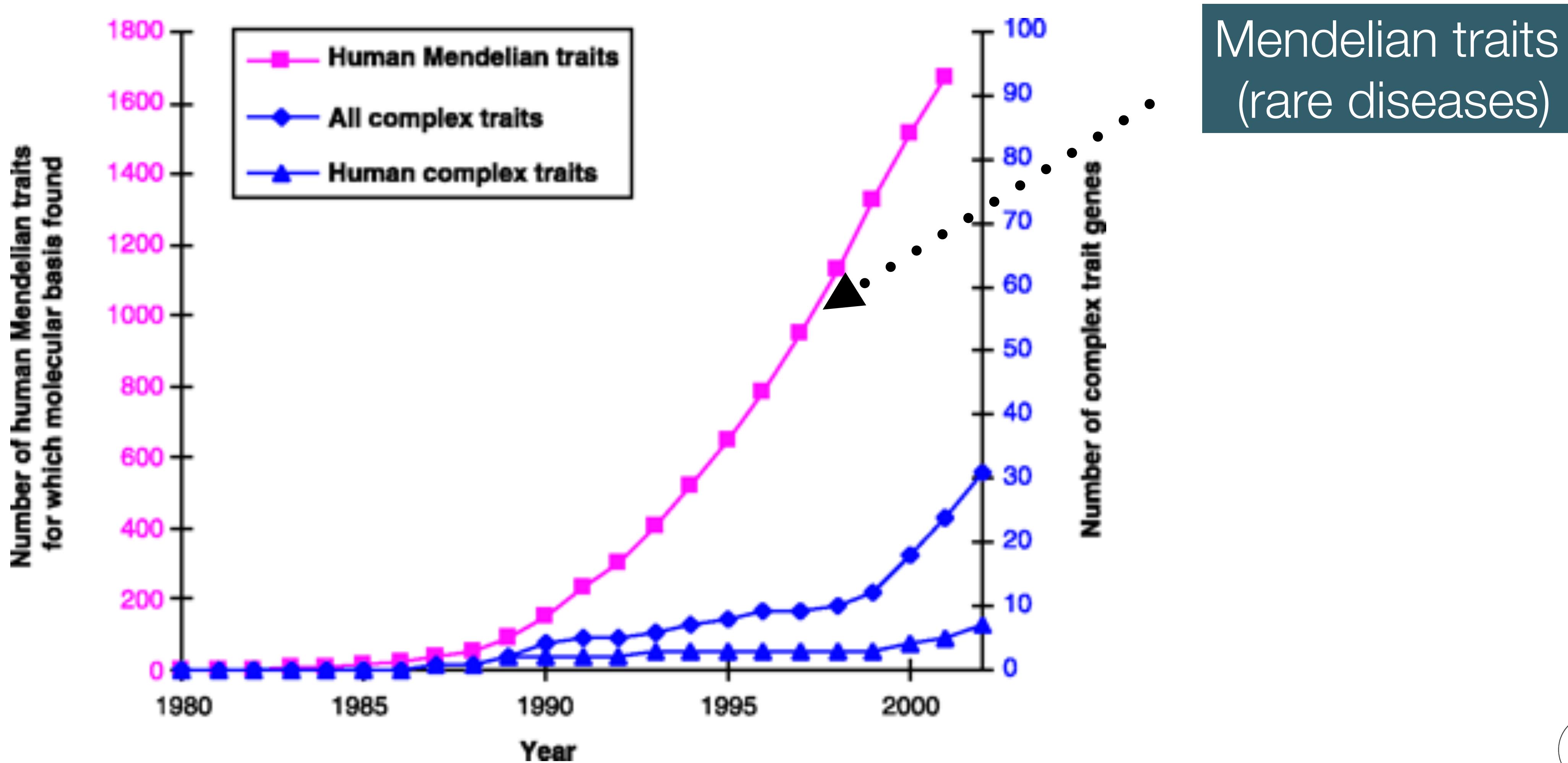
**>50%**

Variation in DNA influences disease risk

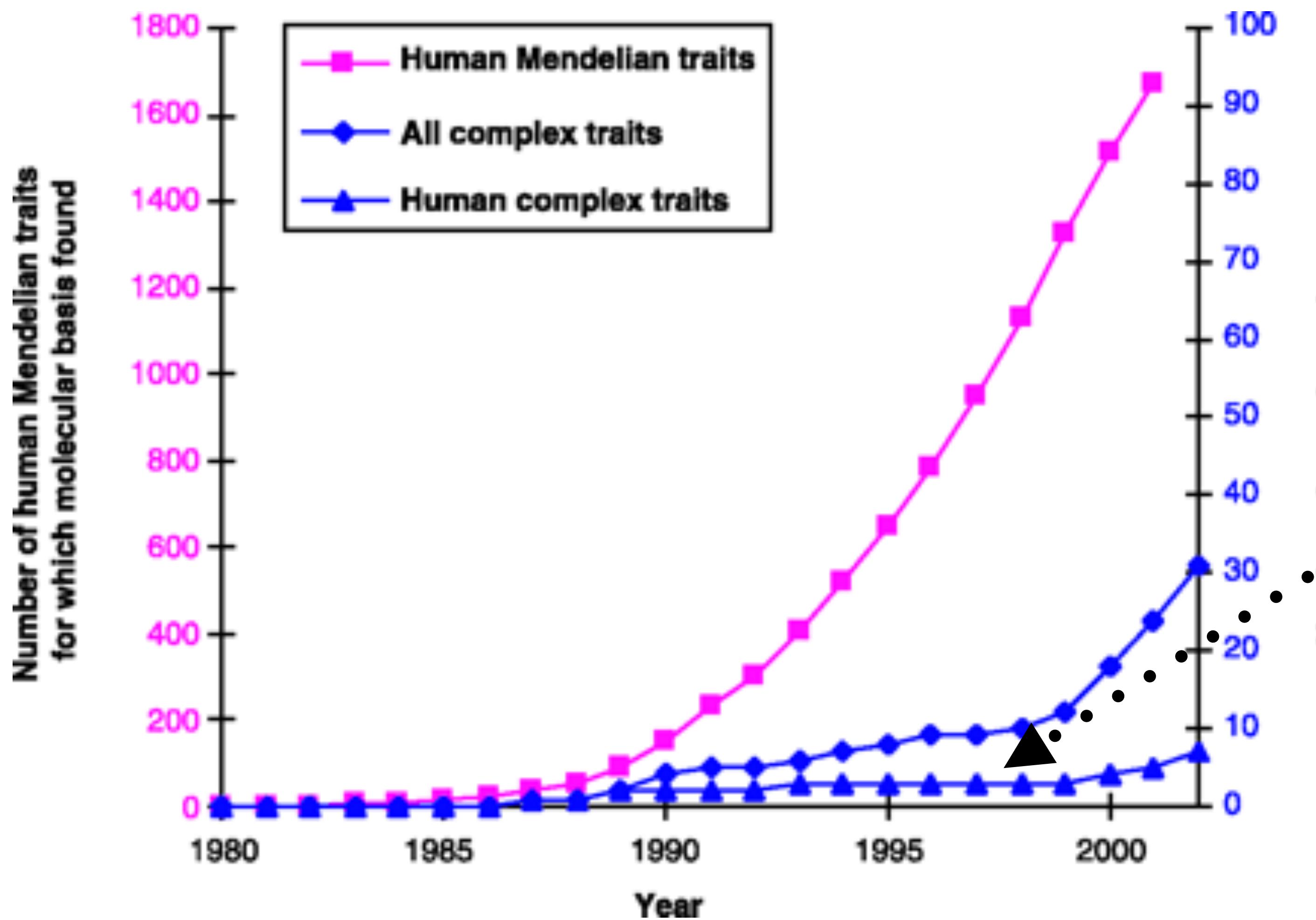


RIVASLAB

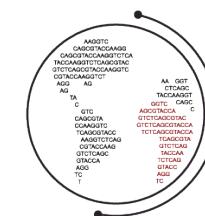
# Dark ages of human genetics of complex traits



# Dark ages of human genetics of complex traits

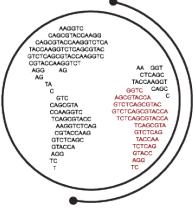


Complex traits  
(common diseases)

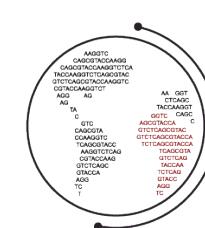
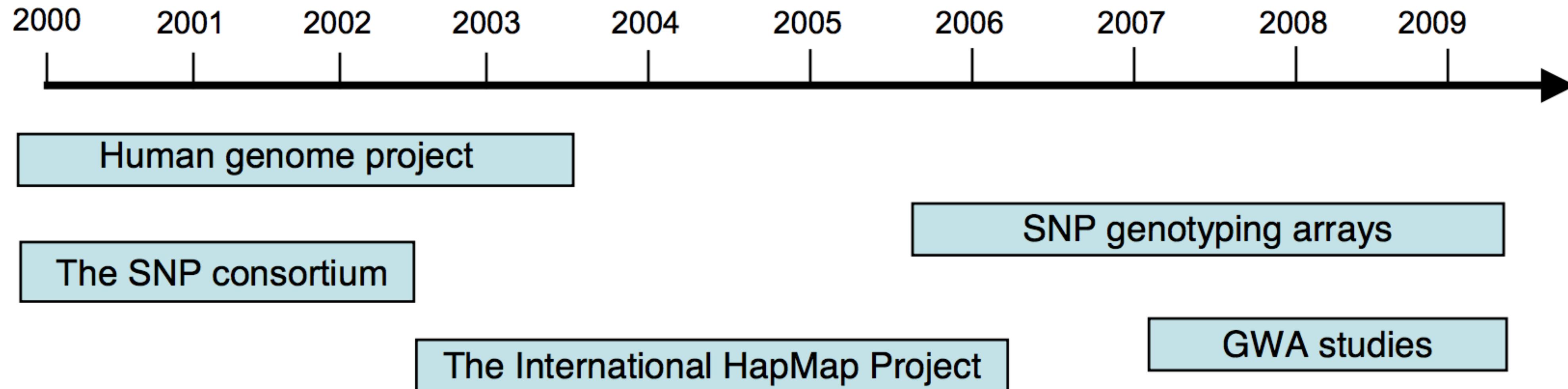


# Human Genome Variation

- Hailed as “Breakthrough Of The Year” in 2007

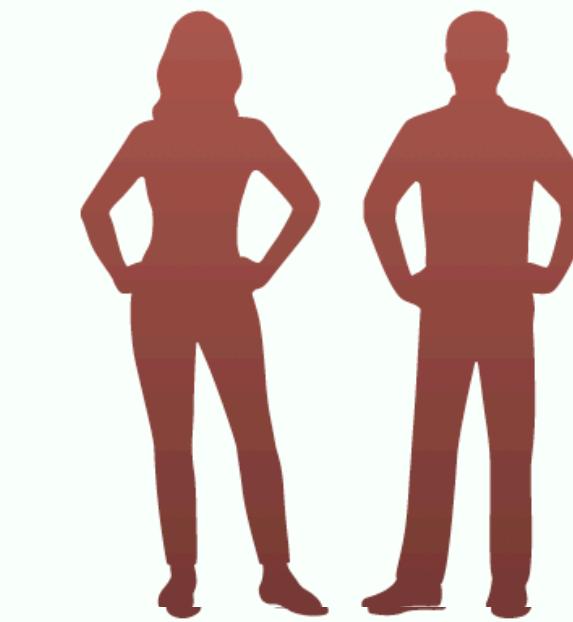
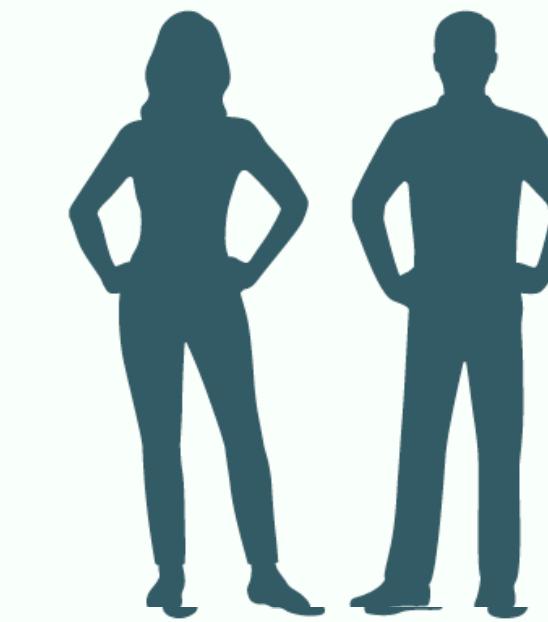


Genome wide association studies have been very successful\*



# Genetic association study

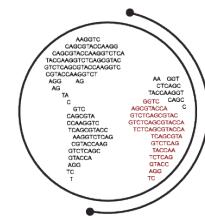
For a given variant..



A



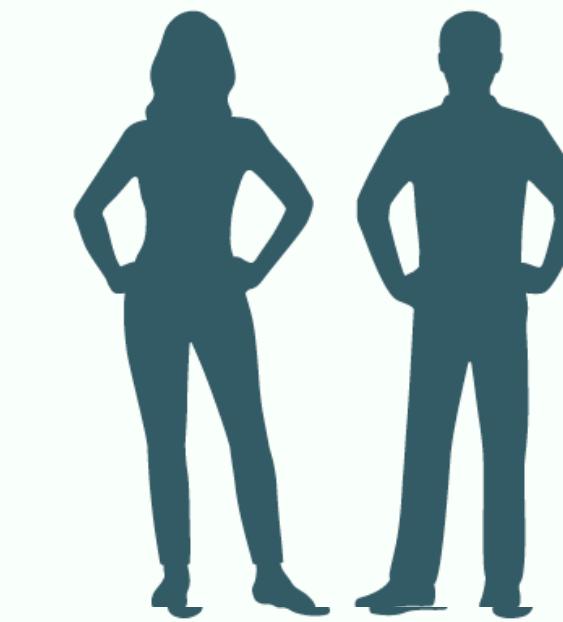
C



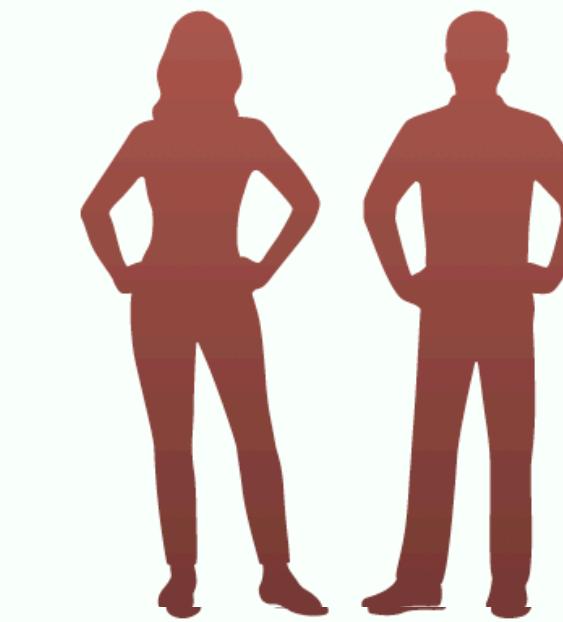
RIVASLAB

# Genetic association study

For a given variant..



Control



Patients



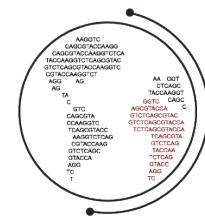
A

80%



C

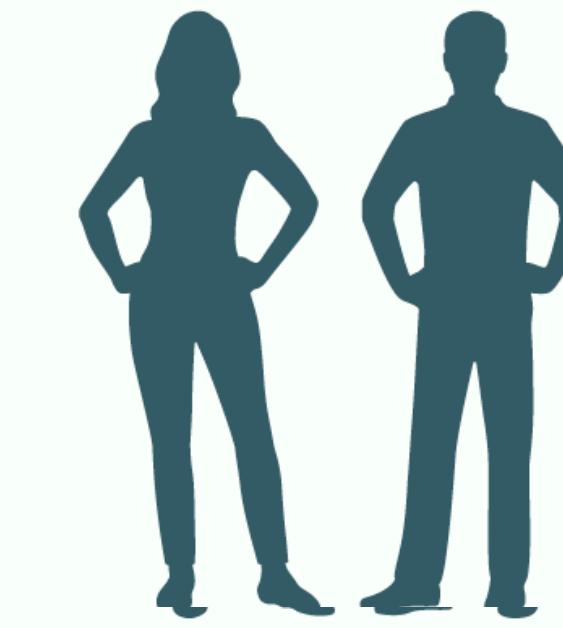
20%



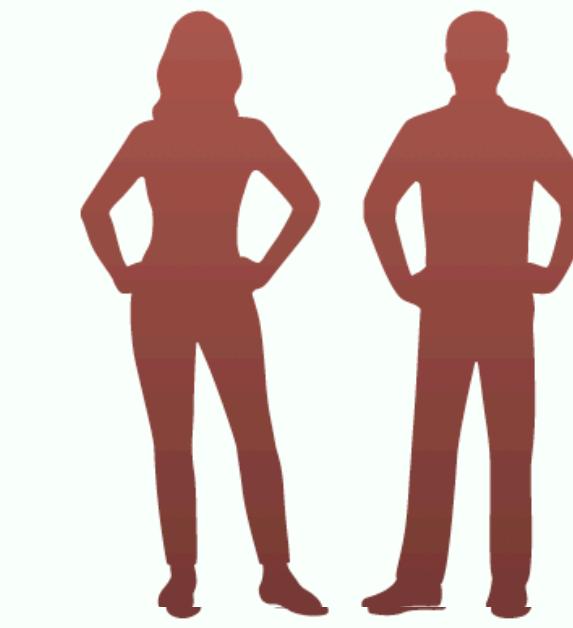
RIVASLAB

# Genetic association study

For a given variant..



Control



Patients



A

80%

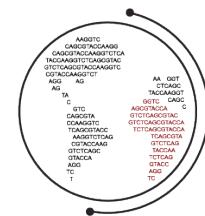
10%



C

20%

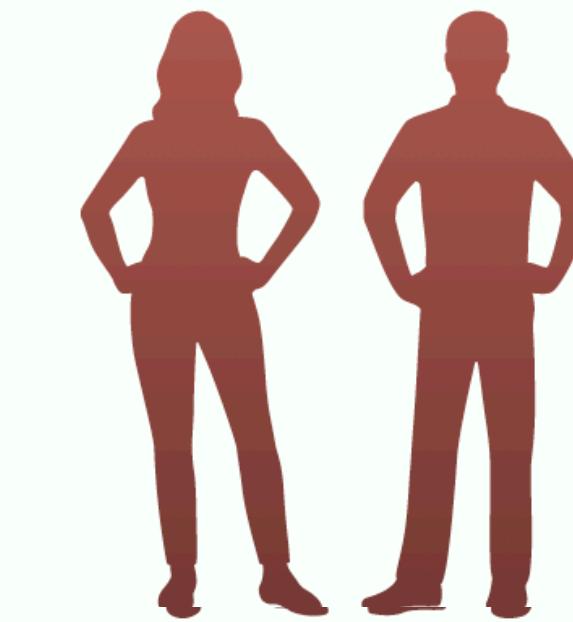
90%



RIVASLAB

# Genetic association study

For a given variant..



Control

Patients



Variation on DNA is associated  
with the trait of interest (disease)

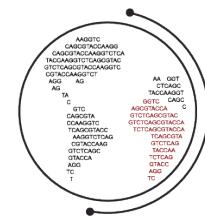
10%



C

20%

90%

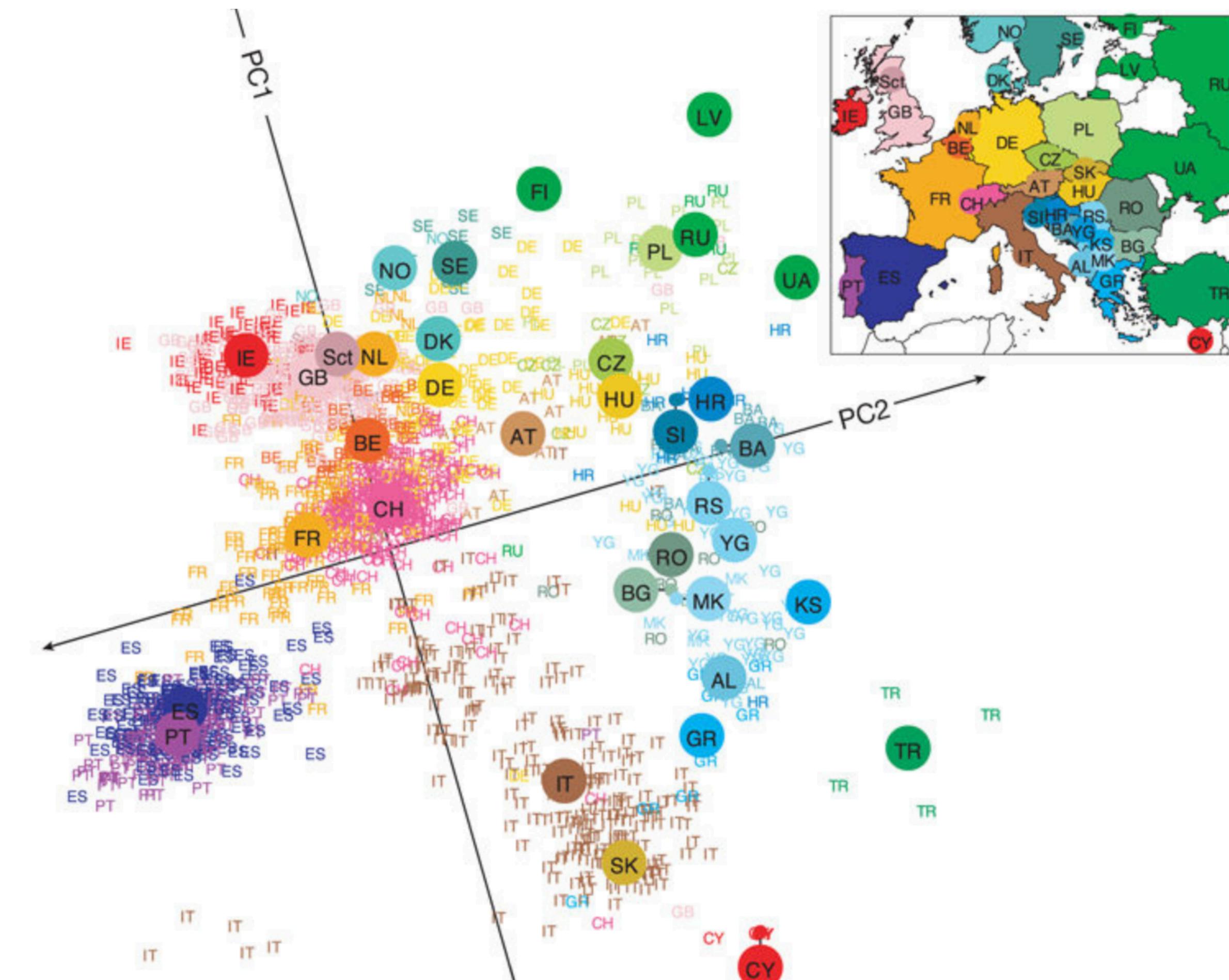


RIVASLAB

# Typical GWAS analysis

# Principal components to adjust for population structure

# Logistic regression analysis with covariates

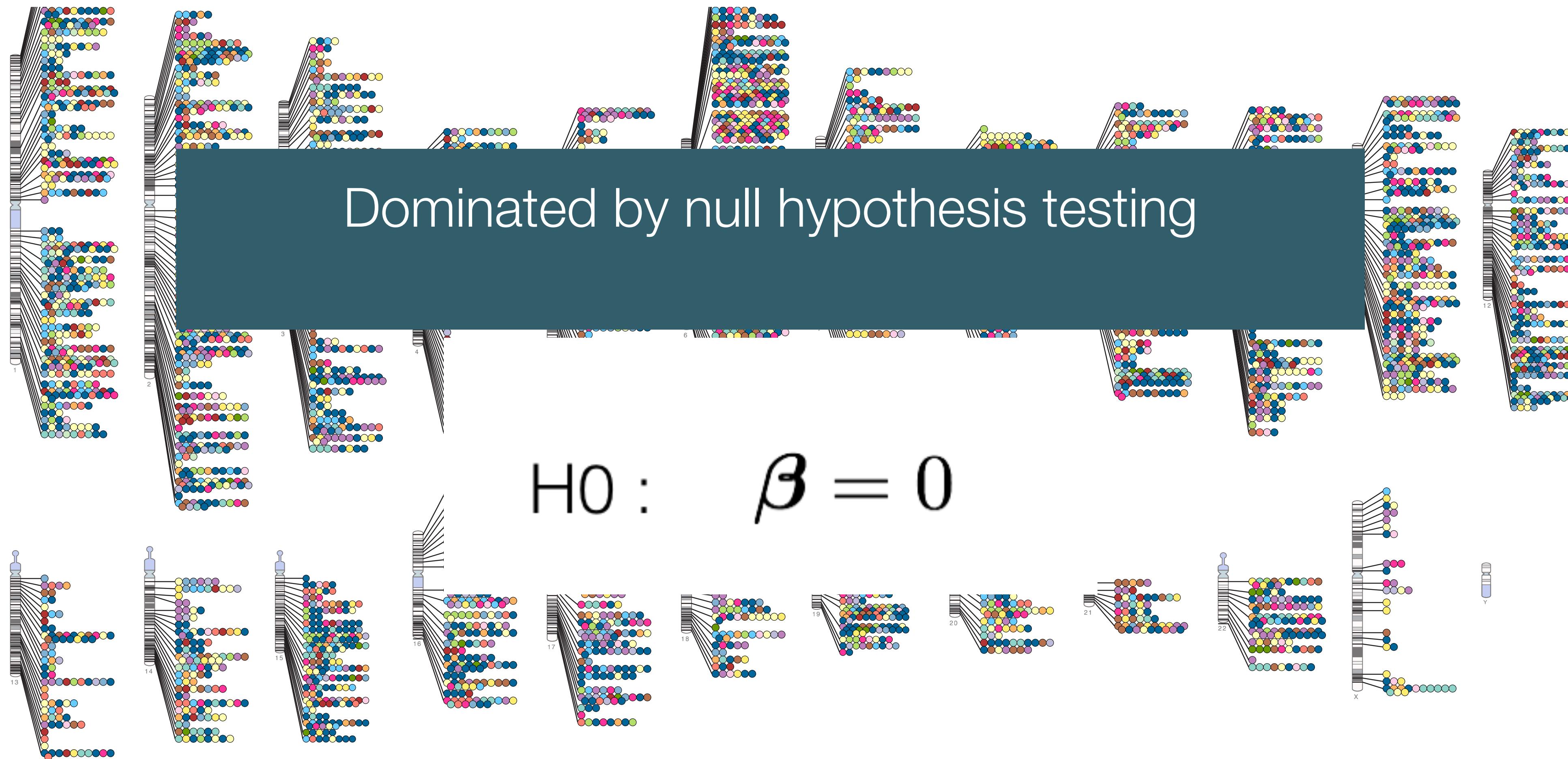


# Genome wide association studies have been very successful\*



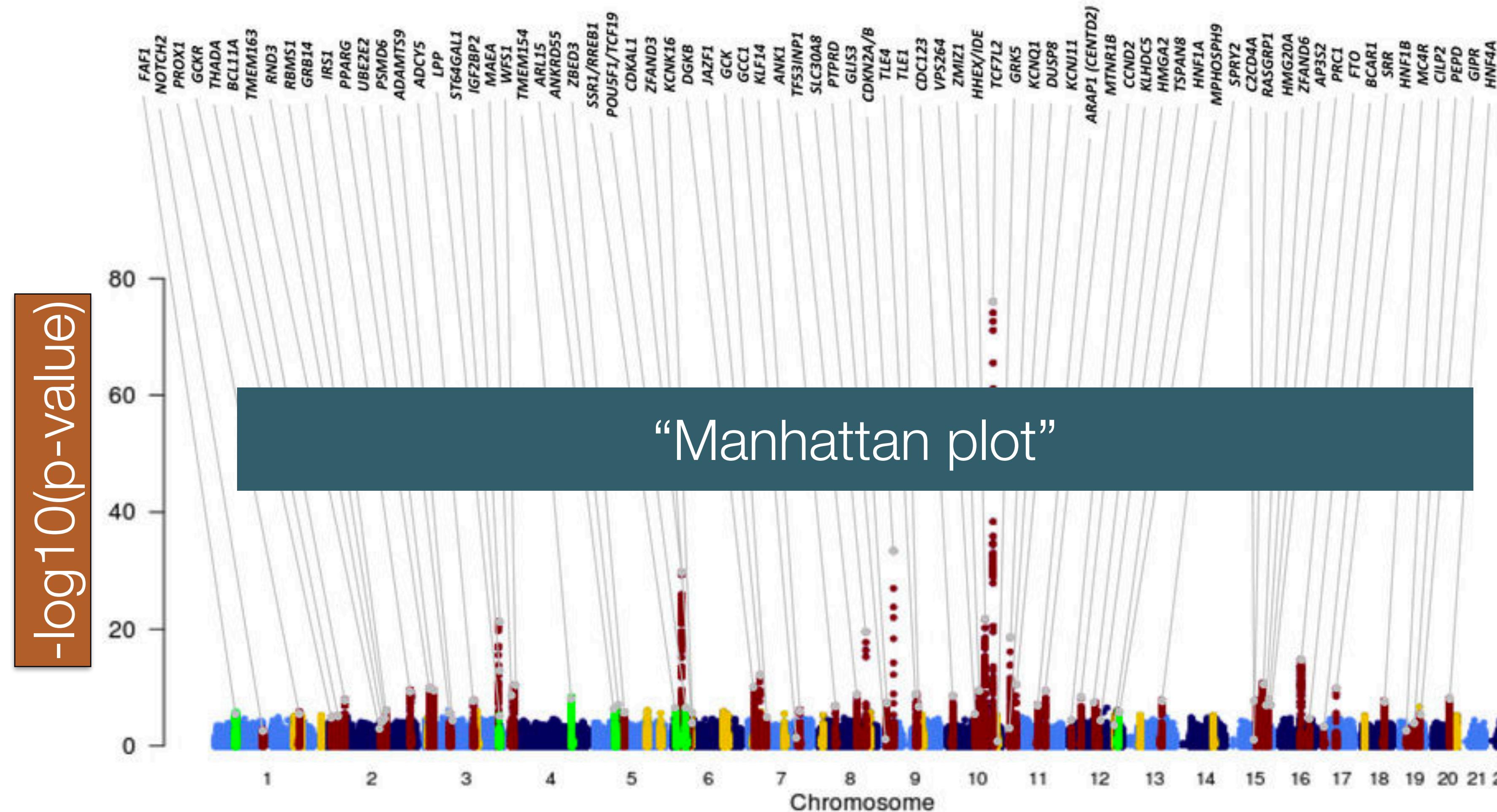
Over 1000 genetic associations

Genome wide association studies have been very successful\*

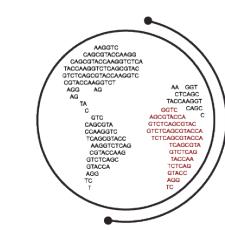
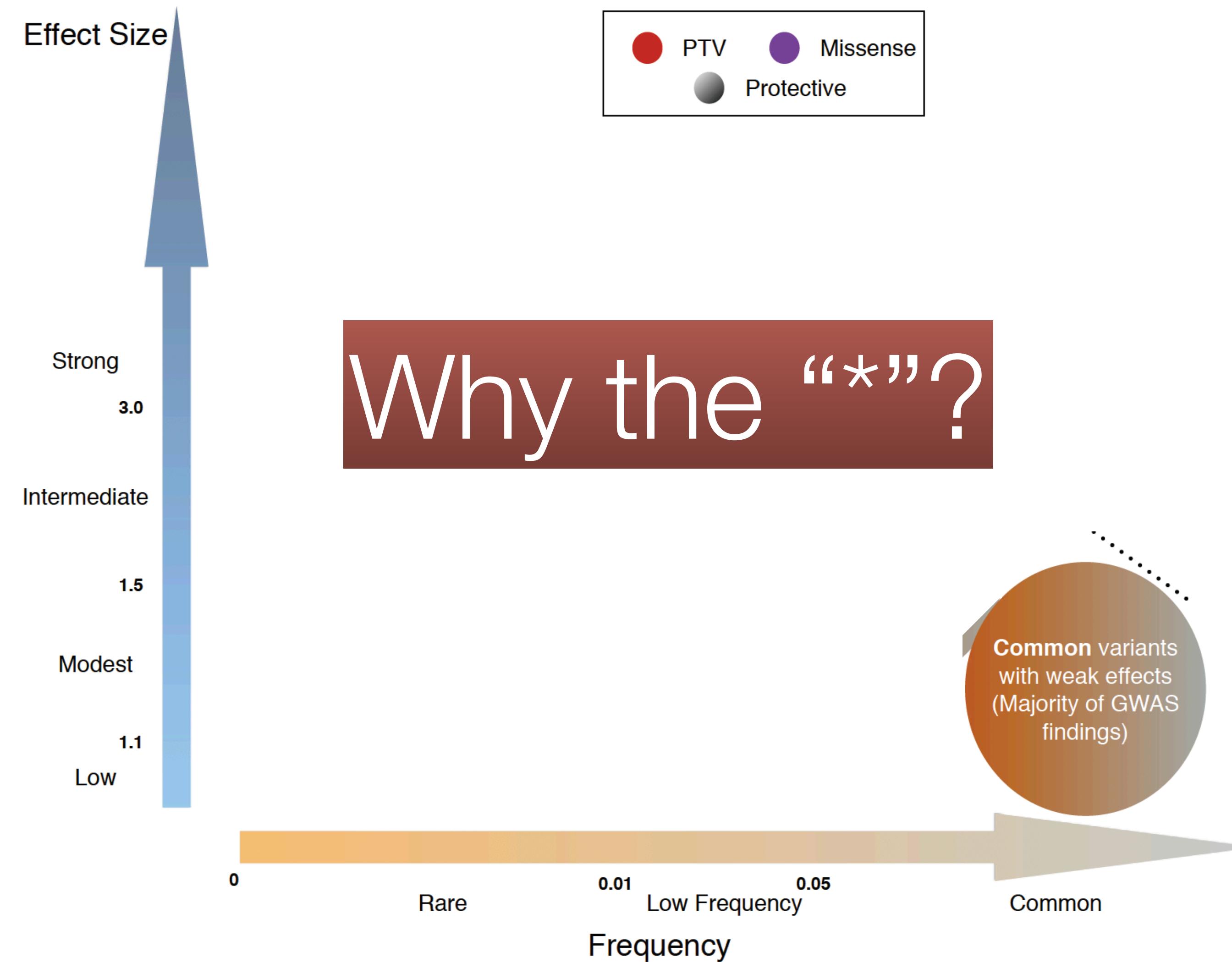


i.e. the effect of the genetic variant is 0.

# Genome wide association studies have been very successful\*

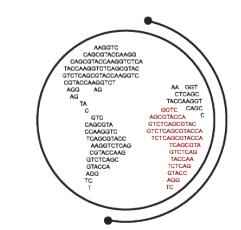
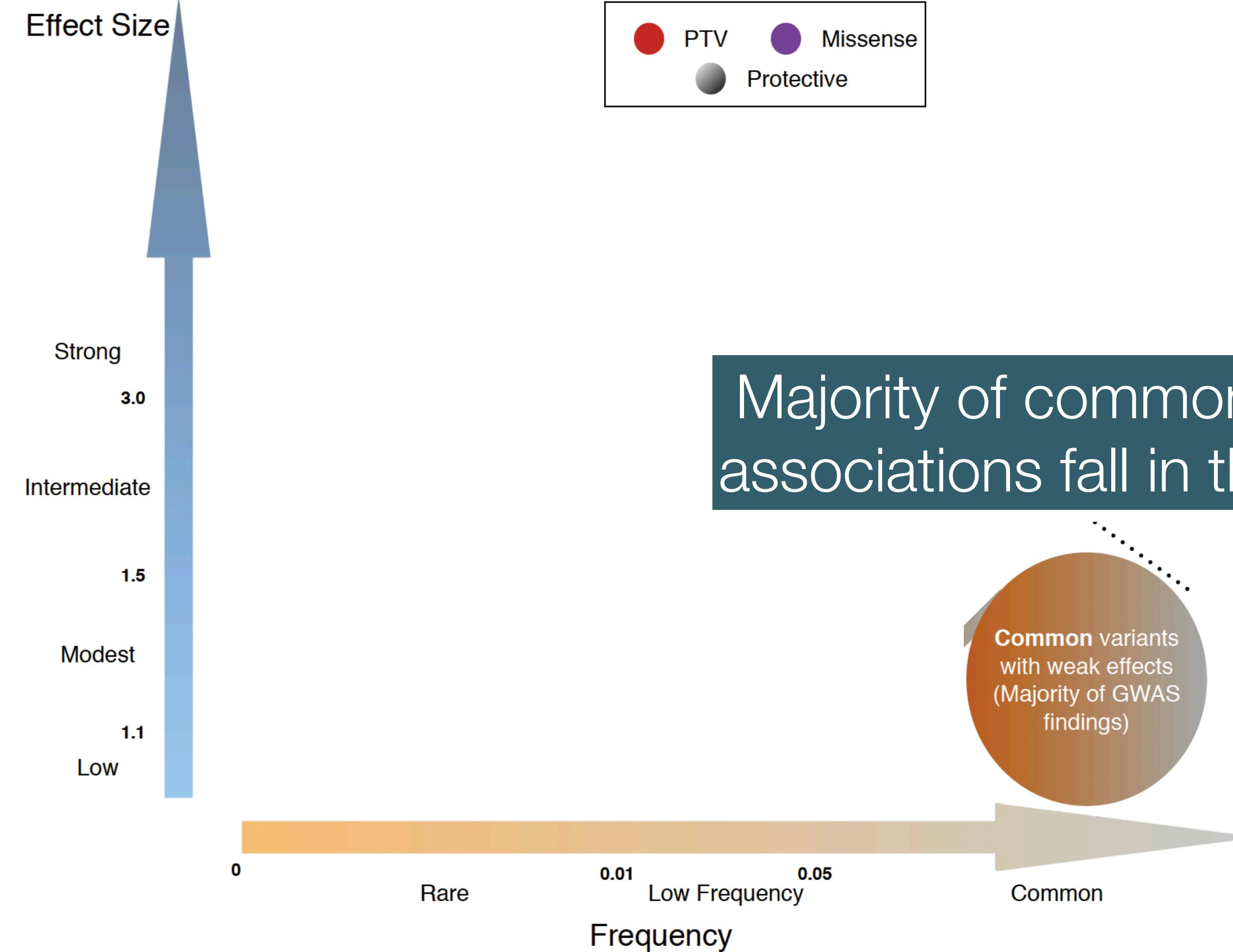


# Genome wide association studies have been very successful\*



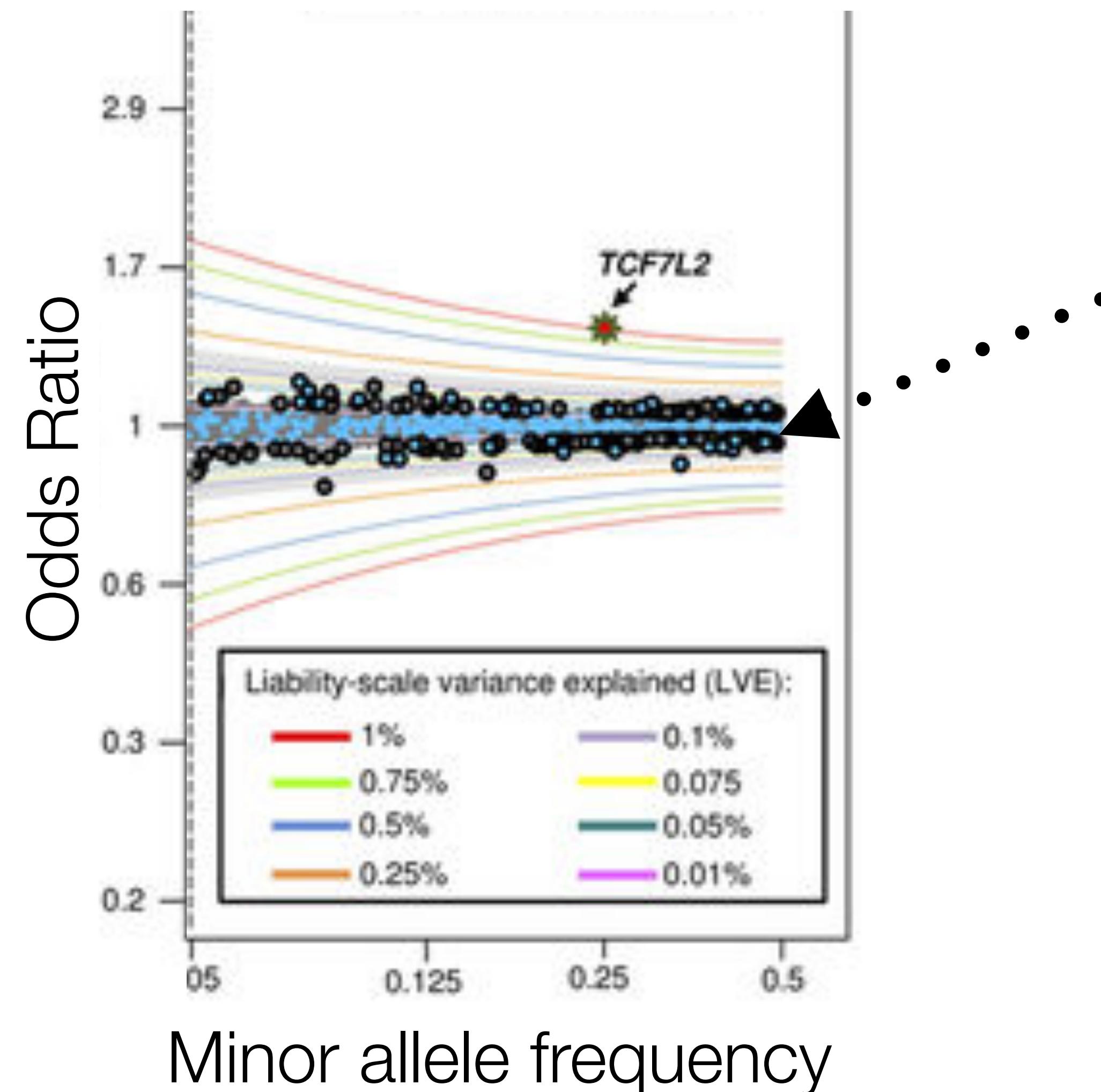
# Genome wide association studies have been very successful\*

Why the “\*”?

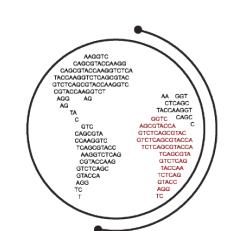


In the context of **type 2 diabetes** all common variant associations were tiny

Why the “\*”?

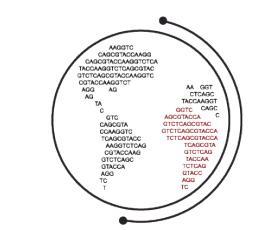
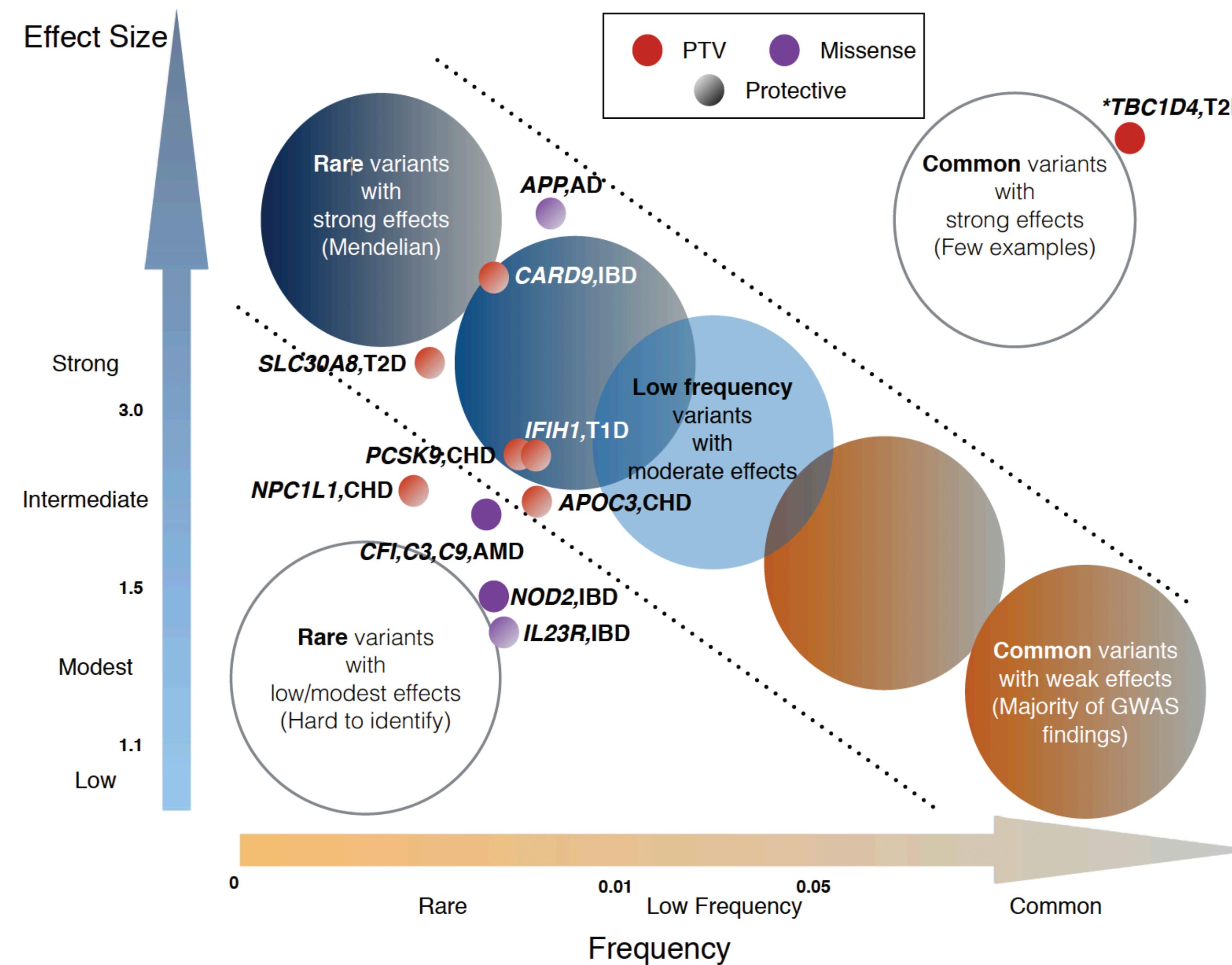


Tiny effect sizes for all associated variants

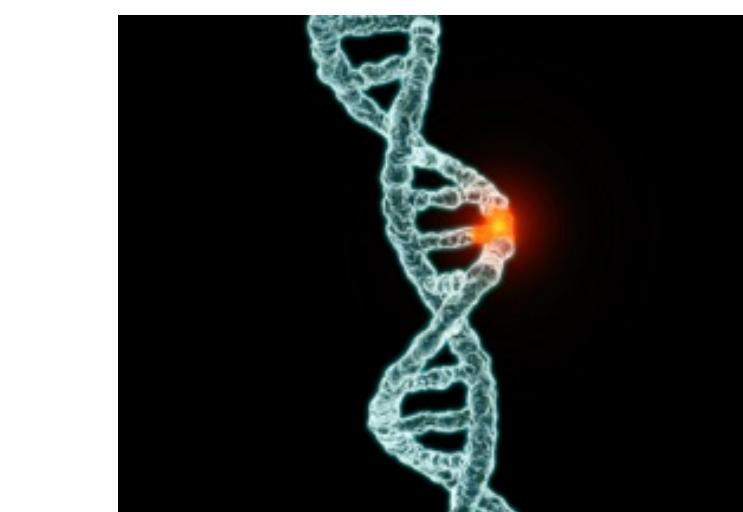


RIVAS LAB

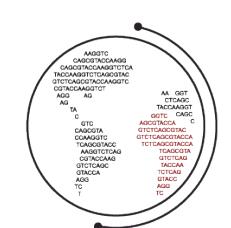
# Additional signals started emerging from rare variant studies



# “Experiments of nature” that protect can guide selection of drug targets



Lower  
risk for  
disease



RIVASLAB

# Examples of protective mutations

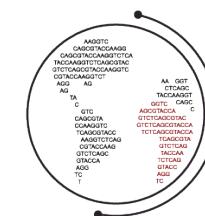
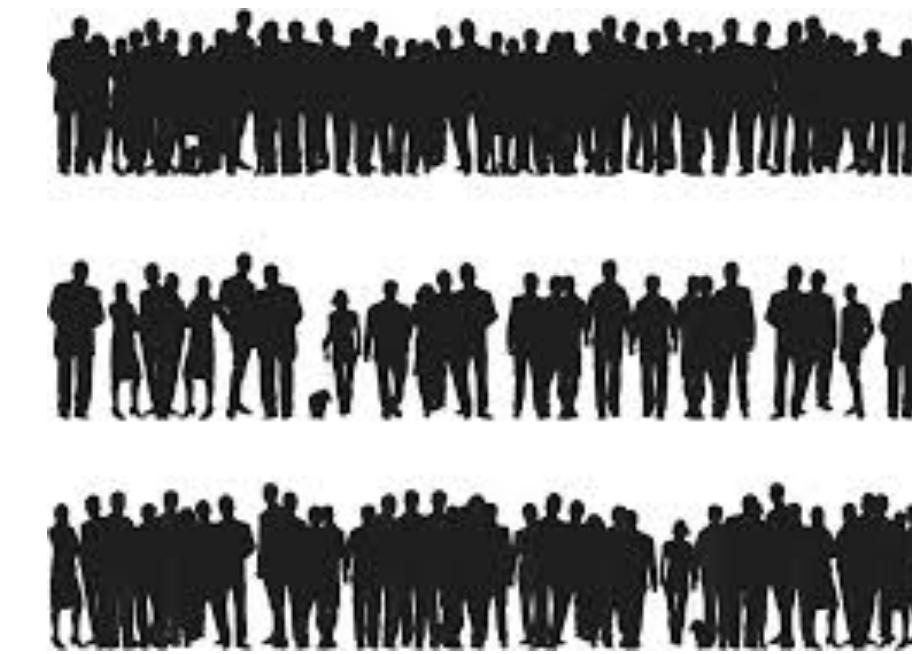
PCSK9 for LDL and MI

Nav 1.7 for pain

CARD9 for Crohn's disease and ulcerative colitis

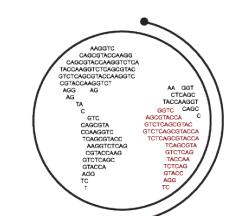
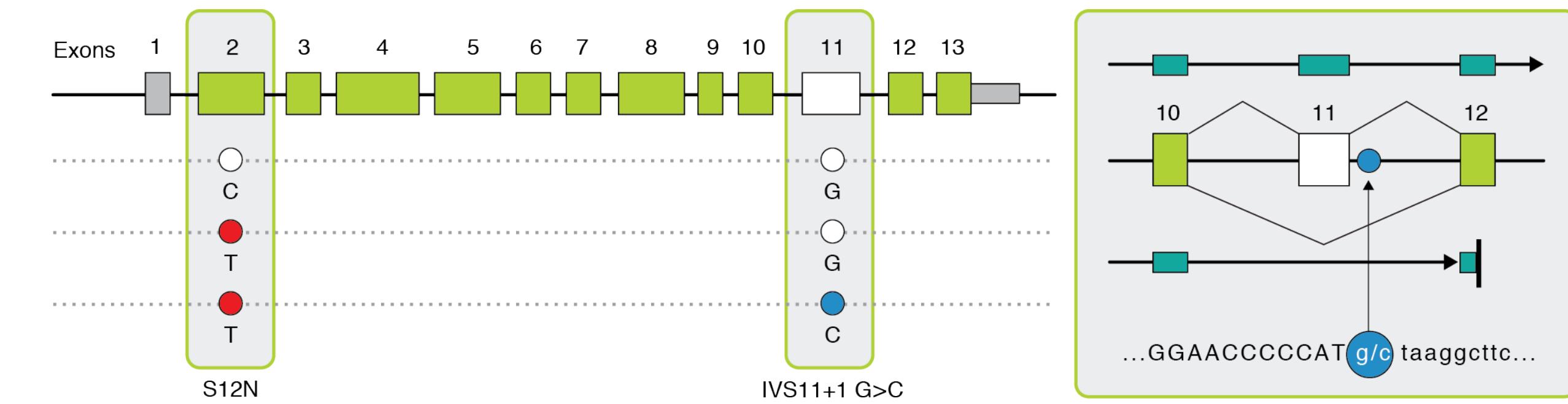
RNF186 for ulcerative colitis

CCR5 for HIV



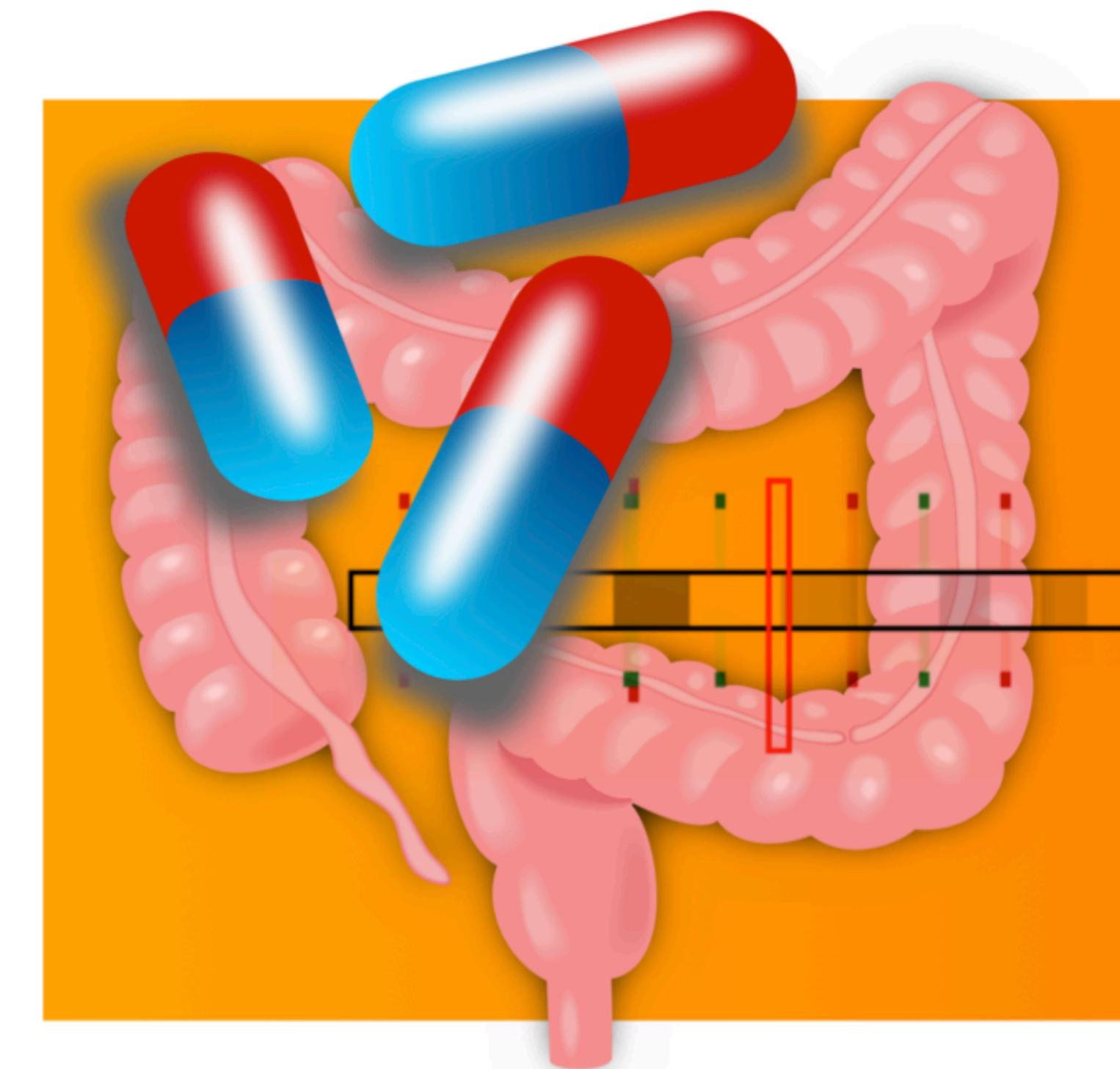
# Rare, strong acting alleles provide interpretation of the GWAS results

- Splice variant in **CARD9** cause premature truncating of protein and **strongly protects** against the development of Crohn's disease and ulcerative colitis ( $p < 10^{-16}$ ).

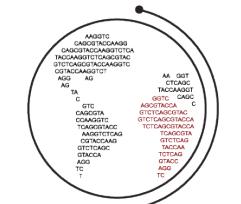


# Rare, strong acting alleles provide interpretation of the GWAS results

- Protective loss-of-function variant in *RNF186* found to confer protection against ulcerative colitis (3-fold protective effect).
- Protective genetic variants reveal process that is:
  - **safe** (naturally occurs in healthy adults)
  - **effective** (proven to reduce risk of disease).



Found in Iceland and Finland



RIVASLAB

Rivas et al., *Nature Communications* 2016.

---

## Loss-of-function mutations in *SLC30A8* protect against type 2 diabetes

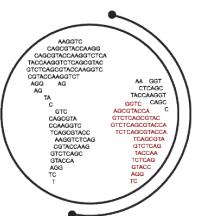
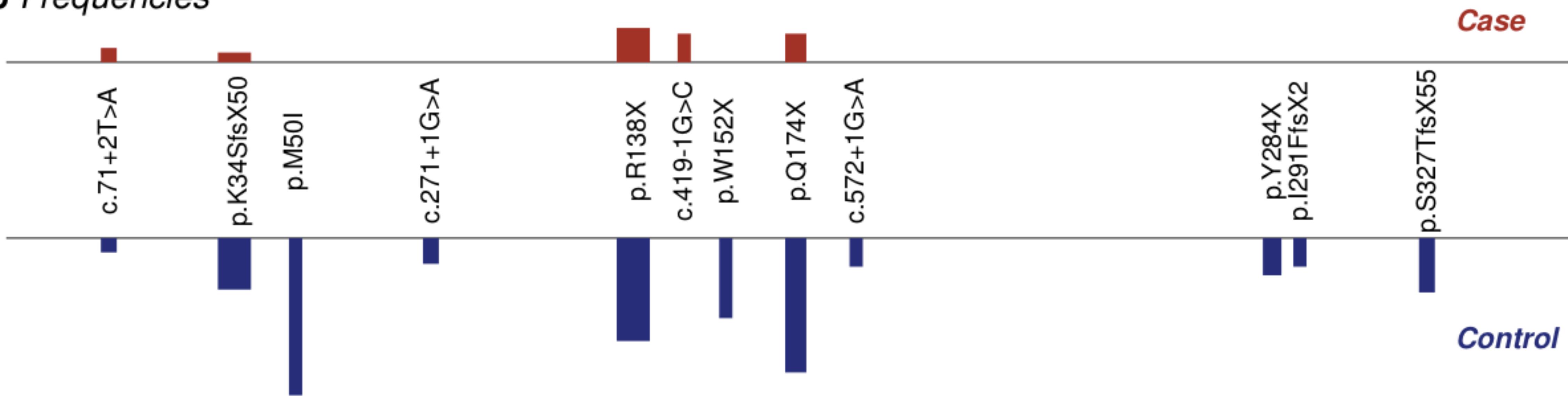
Jason Flannick<sup>1–3</sup>, Gudmar Thorleifsson<sup>4</sup>, Nicola L Beer<sup>1,5</sup>, Suzanne B R Jacobs<sup>1</sup>, Niels Grarup<sup>6</sup>, Noël P Burtt<sup>1</sup>, Anubha Mahajan<sup>7</sup>, Christian Fuchsberger<sup>8</sup>, Gil Atzmon<sup>9,10</sup>, Rafn Benediktsson<sup>11</sup>, John Blangero<sup>12</sup>, Don W Bowden<sup>13–16</sup>, Ivan Brandslund<sup>17,18</sup>, Julia Brosnan<sup>19</sup>, Frank Burslem<sup>20</sup>, John Chambers<sup>21–23</sup>, Yoon Shin Cho<sup>24</sup>, Cramer Christensen<sup>25</sup>, Desirée A Douglas<sup>26</sup>, Ravindranath Duggirala<sup>12</sup>, Zachary Dymek<sup>1</sup>, Yossi Farjoun<sup>1</sup>, Timothy Fennell<sup>1</sup>, Pierre Fontanillas<sup>1</sup>, Tom Forsén<sup>27,28</sup>, Stacey Gabriel<sup>1</sup>, Benjamin Glaser<sup>29,30</sup>, Daniel F Gudbjartsson<sup>4</sup>, Craig Hanis<sup>31</sup>, Torben Hansen<sup>6,32</sup>, Astradur B Hreidarsson<sup>11</sup>, Kristian Hveem<sup>33</sup>, Erik Ingelsson<sup>7,34</sup>, Bo Isomaa<sup>35,36</sup>, Stefan Johansson<sup>37–39</sup>, Torben Jørgensen<sup>40–42</sup>, Marit Eika Jørgensen<sup>43</sup>, Sekar Kathiresan<sup>1,44–46</sup>, Augustine Kong<sup>4</sup>, Jaspal Kooner<sup>22,23,47</sup>, Jasmina Kravic<sup>48</sup>, Markku Laakso<sup>49</sup>, Jong-Young Lee<sup>50</sup>, Lars Lind<sup>51</sup>, Cecilia M Lindgren<sup>1,7</sup>, Allan Linneberg<sup>40,41,52</sup>, Gisli Masson<sup>4</sup>, Thomas Meitinger<sup>53</sup>, Karen L Mohlke<sup>54</sup>, Anders Molven<sup>37,55,56</sup>, Andrew P Morris<sup>7,57</sup>, Shobha Potluri<sup>58</sup>, Rainer Rauramaa<sup>59,60</sup>, Rasmus Ribel-Madsen<sup>6</sup>, Ann-Marie Richard<sup>19</sup>, Tim Rolph<sup>19</sup>, Veikko Salomaa<sup>61</sup>, Ayellet V Segrè<sup>1,2</sup>, Hanna Skärstrand<sup>26</sup>, Valgerdur Steinthorsdottir<sup>4</sup>, Heather M Stringham<sup>8</sup>, Patrick Sulem<sup>4</sup>, E Shyong Tai<sup>62–64</sup>, Yik Ying Teo<sup>62,65–68</sup>, Tanya Teslovich<sup>8</sup>, Unnur Thorsteinsdottir<sup>4,69</sup>, Jeff K Trimmer<sup>19</sup>, Tiinamaija Tuomi<sup>27,35</sup>, Jaakko Tuomilehto<sup>70–72</sup>, Fariba Vaziri-Sani<sup>26</sup>, Benjamin F Voight<sup>1,73,74</sup>, James G Wilson<sup>75</sup>, Michael Boehnke<sup>8</sup>, Mark I McCarthy<sup>5,7,76</sup>, Pål R Njølstad<sup>1,37,77</sup>, Oluf Pedersen<sup>6</sup>, Go-T2D Consortium<sup>78</sup>, T2D-GENES Consortium<sup>78</sup>, Leif Groop<sup>48,79</sup>, David R Cox<sup>58</sup>, Kari Stefansson<sup>4,69</sup> & David Altshuler<sup>1–3,44,45,80,81</sup>

In aggregate, **12 loss of function** mutations associated with **65% decreased risk** of T2D

**a Variants**



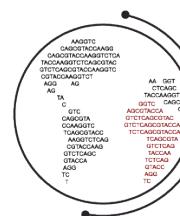
**b Frequencies**



# Why the need for **Stan** in Human Genetics?



*Stan*

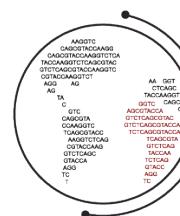


RIVASLAB

**Parameters became more exciting (and informative)...**

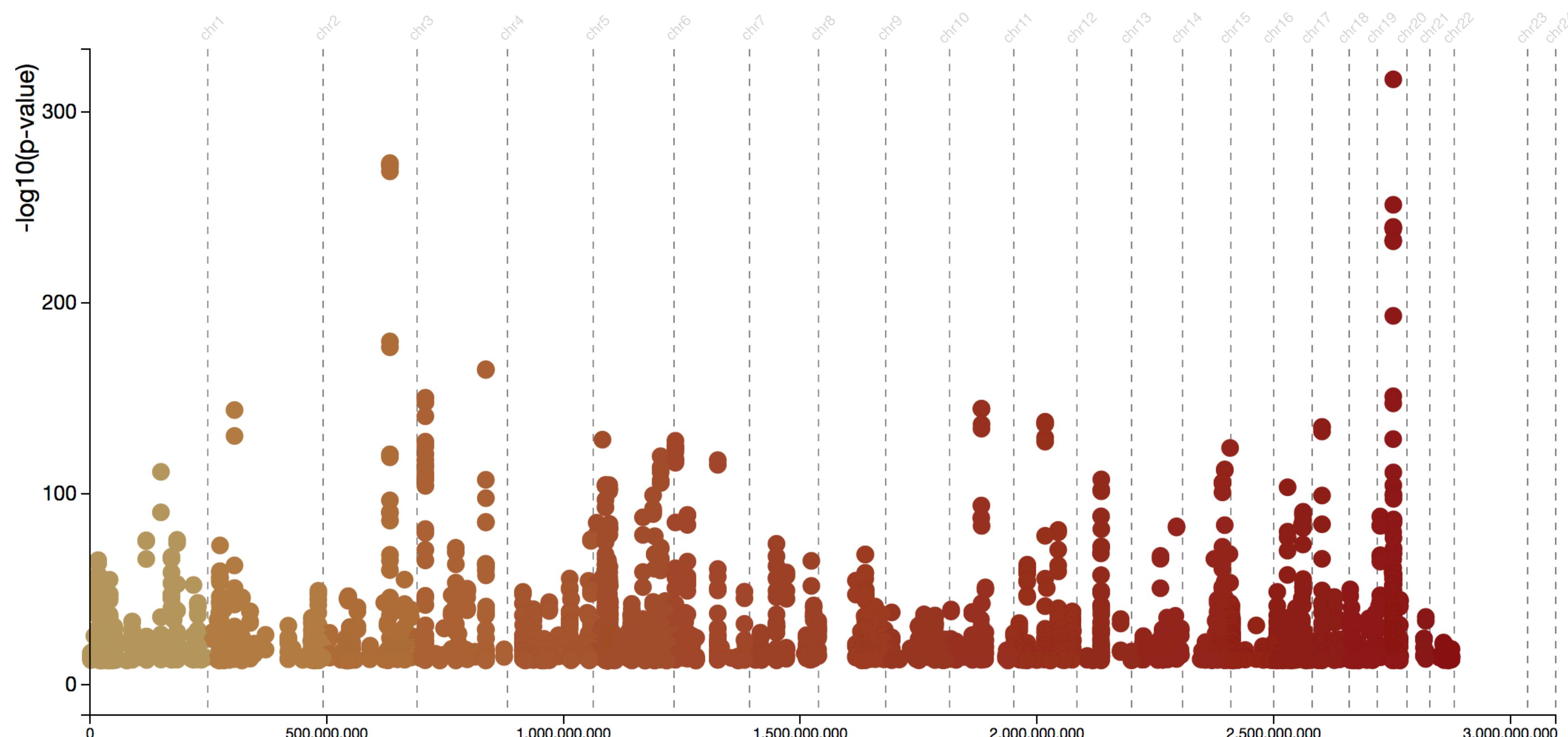
Location, scale, correlation, effect size

p-value is not informative especially...



RIVASLAB

when you realize the entire genome is “associated” to disease

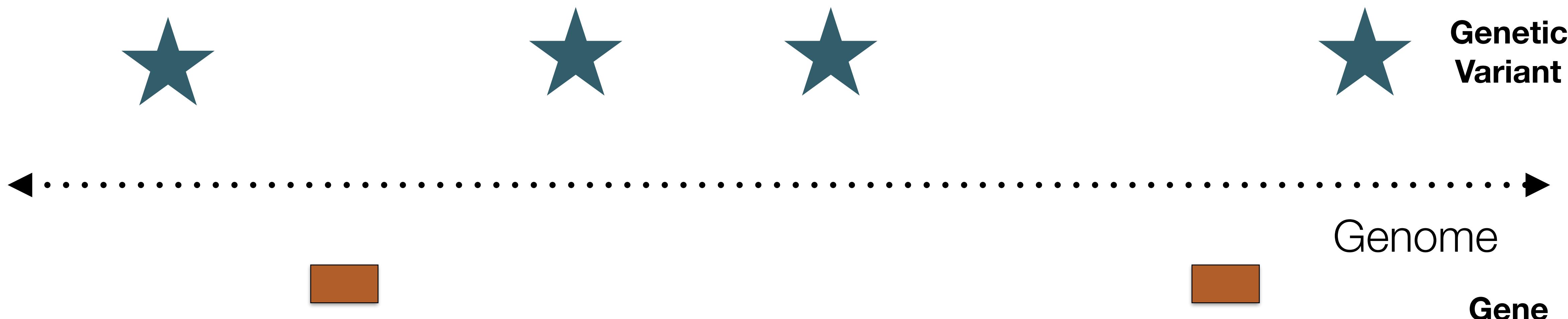


Note: only variants with a p-value less than 0.001 are included in the manhattan plot.



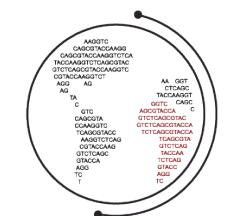
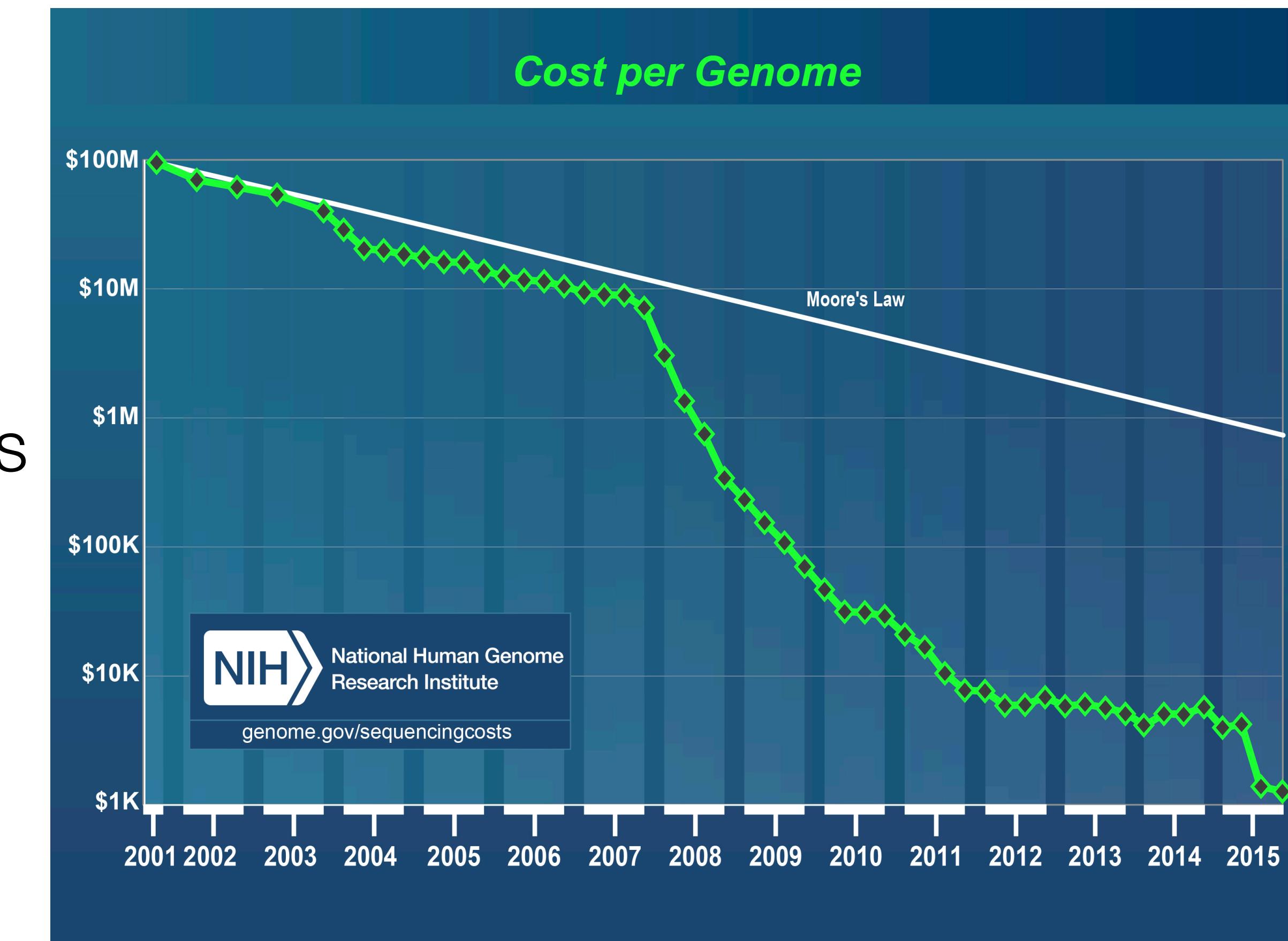
RIVASLAB

# Age of common variants

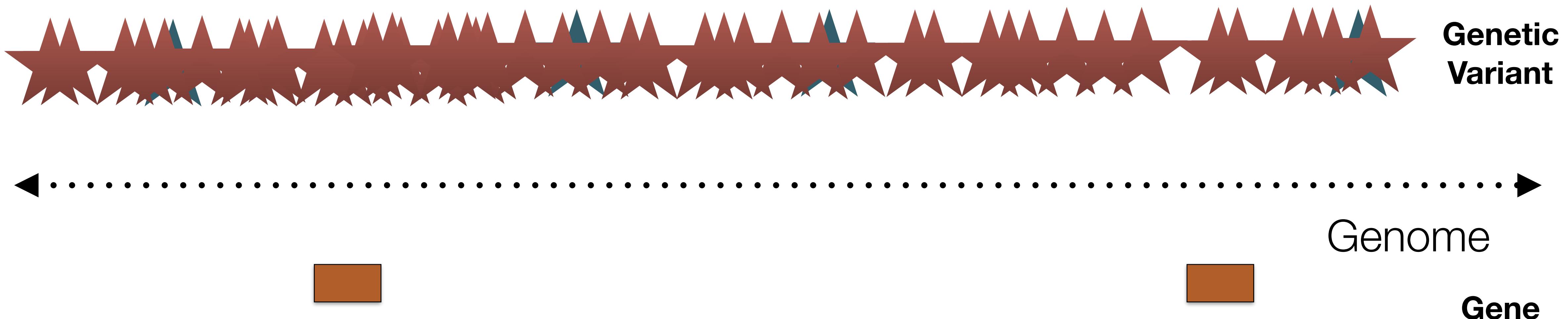


# Technologies transforming biomedicine

Cost of sequencing has plummeted over the past 15 years



# Age of rare variants



# Age of rare variants

## Gene: NOD2



**NOD2** nucleotide-binding oligomerization domain containing 2

Transcripts ▾

**Number of variants**

1988 (Including filtered: 2164)

[16:50727514-50766988](#) ↗

[UCSC Browser](#)

[GeneCards](#)

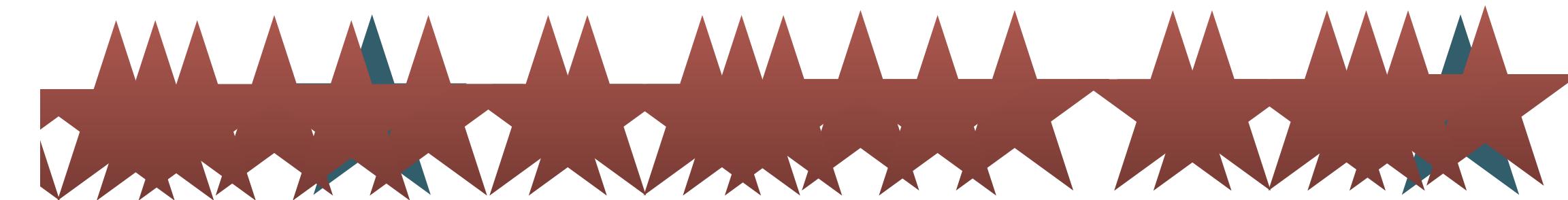
[NOD2](#) ↗

[605956](#) ↗

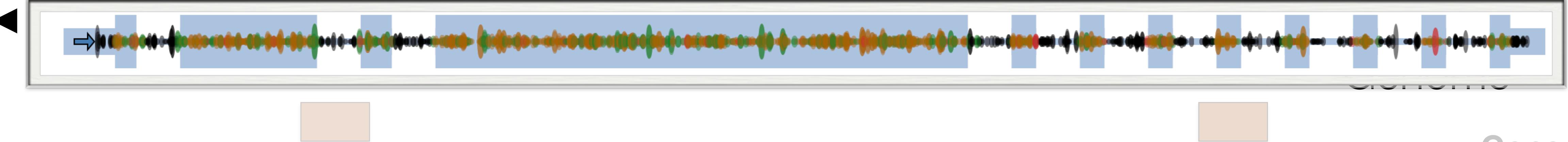
[OMIM](#)

[Other](#)

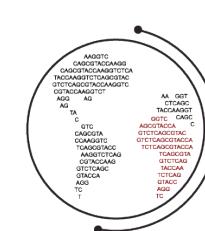
External References ▾



**Genetic Variant**



<https://ibd.broadinstitute.org/gene/ENSG00000167207>



RIVASLAB

# Age of rare variants

## Gene: NOD2



<https://ibd.broadinstitute.org/gene/ENSG00000167207>

# Age of rare variants

Gene: NOD2

NOD2 nucleotide-binding oligomerization domain containing 2

Number of variants 1988 (Including filtered: 2164)

Transcripts ▾

UCSC Browser 16:50727514-5076

GeneCards NOD2

OMIM 605956

Other External References

**Genetic variants have predicted functional consequences**

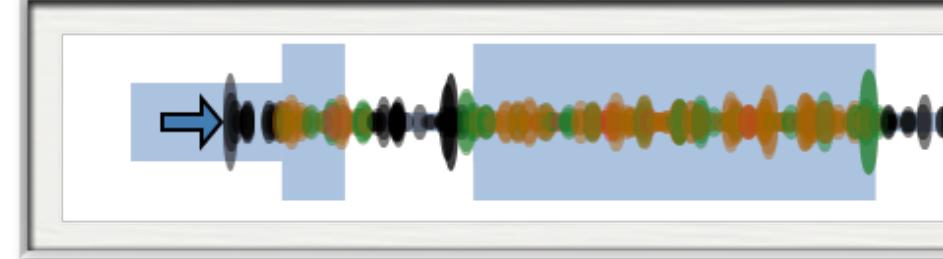
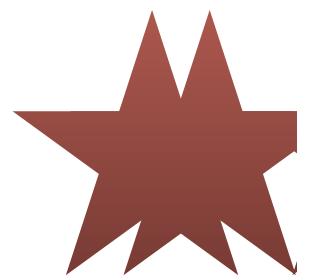
- Silent (no effect)
- Missense (mild effect)
- Loss of function (strong effect)

Gene

<https://ibd.broadinstitute.org/gene/ENSG00000167207>

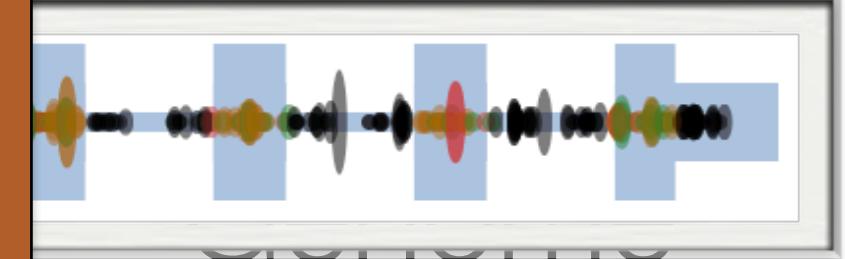
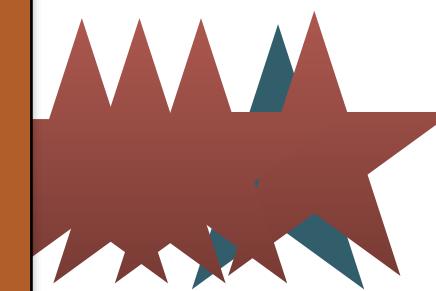
# Age of rare variants

## Gene: NOD2



**Data comes from multiple populations**

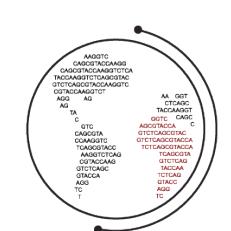
- Ashkenazi Jewish
- Finland
- French Canadian
- Non-Finnish European
- African American
- Latinos



**Genetic Variant**

**Gene**

<https://ibd.broadinstitute.org/gene/ENSG00000167207>



**RIVASLAB**

# Data from genetic study

Crohn's Disease				
BETA	SE	Annotation	Gene	Population
1.2	0.03	lof	NOD2	AJ
1.5	0.9	mis	IL23R	AJ
-1.1	0.006	syn	NOD2	NFE
...	...	...	...	...

Single population

~50 genetic variants

# Modeling of rare variants in *NOD2*, *IL23R*, and *CARD9*

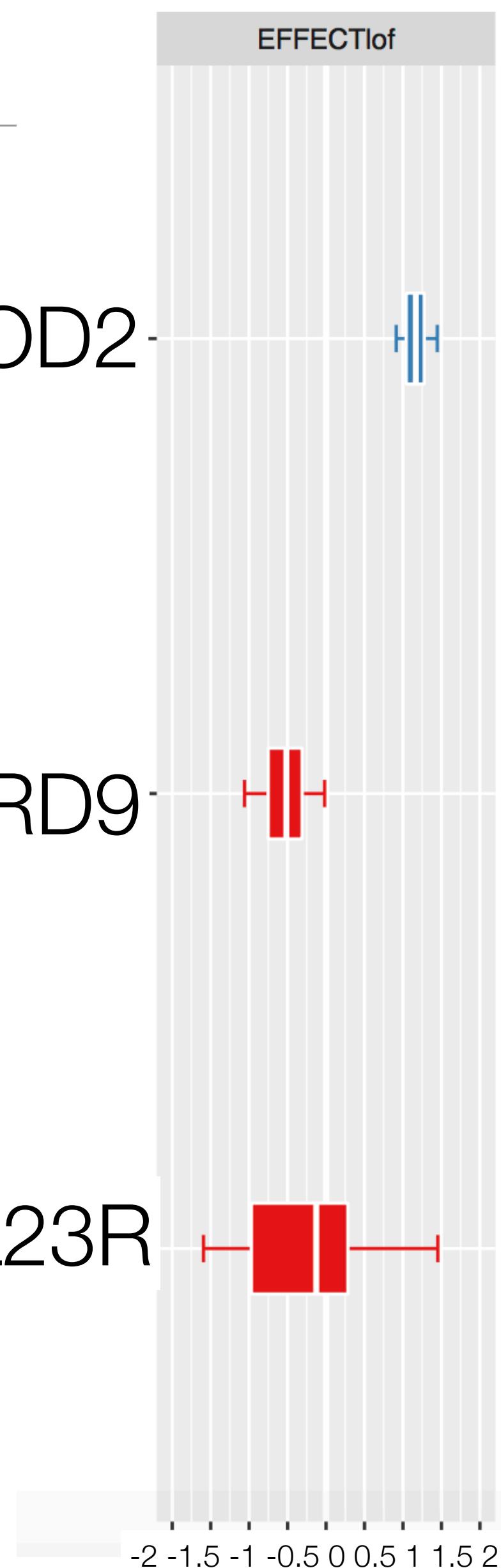
**brms** facilitates multilevel modeling with “meta-analysis” type data

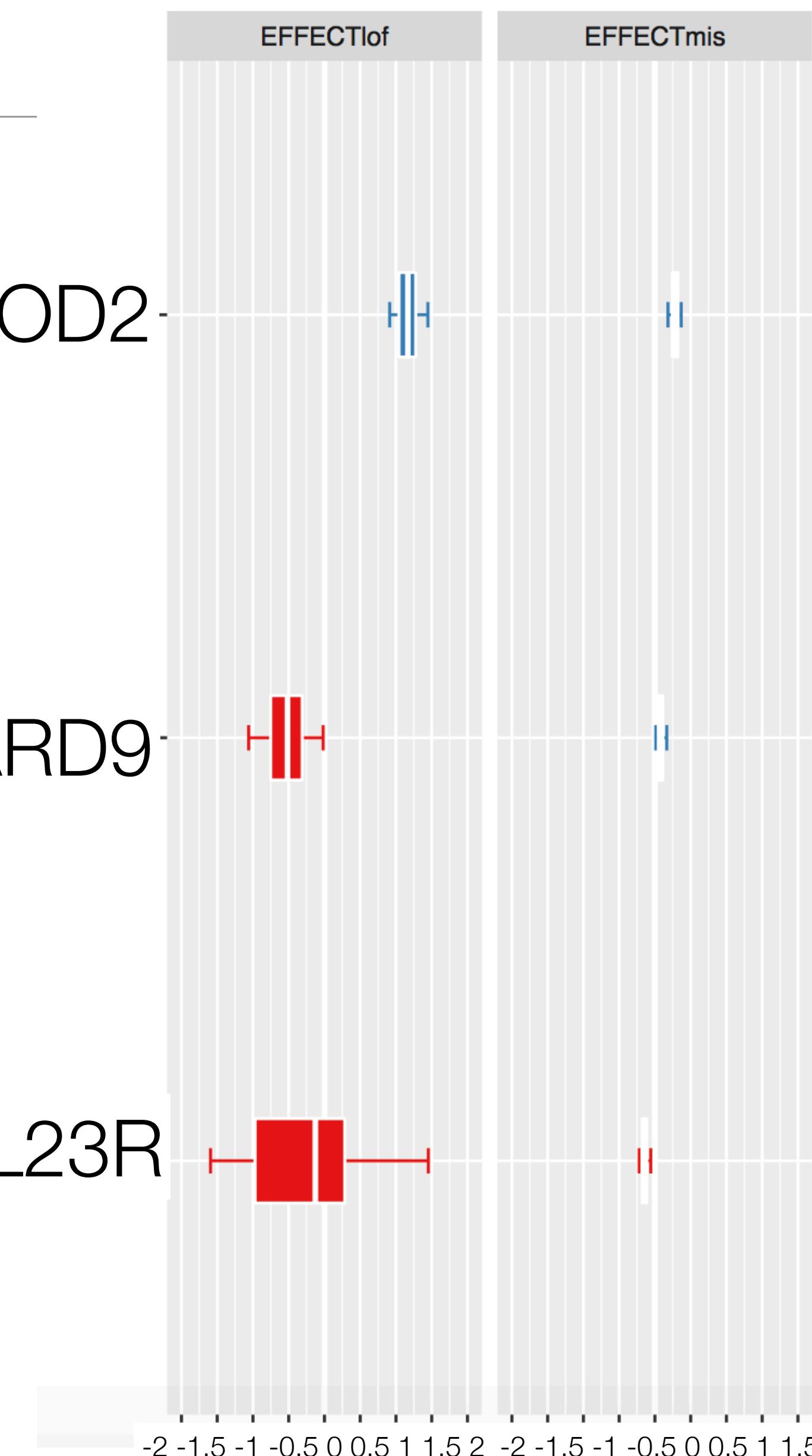
$\hat{B} | se(SE) \sim \text{Annotation} + \text{Gene} + \text{Annotation} | \text{Gene}$ , prior = gaussian, prior = “horseshoe(3)”, algorithm = “meanfield, fullrank, sampling”)

Population level effects      assumes to be the same across population

Group level effects      varying across grouping variables

Single population      ~50 genetic variants





Group-Level Effects:  
~GENE (Number of levels 3)

	Estimate	Est.Error	l-95%	u-95%	CI	Eff.Sample
sd(EFFECTlof)	1.38	0.14	1.11	1.70	1000	
sd(EFFECTmis)	0.20	0.02	0.16	0.25	875	
sd(EFFECTsplice)	0.38	0.14	0.18	0.71	1000	
sd(EFFECTsyn)	0.08	0.02	0.05	0.12	1000	

# Population cohorts enables the study of **shared genetics** with other diseases and phenotypes

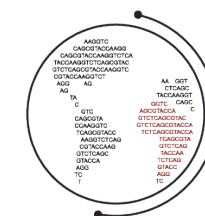
First large-scale population biobank cohort widely available to researchers is **UK Biobank** with approximately **500,000 participants**

- Prospective cohort study
- 500k individuals (337k European British)
- Large number of phenotypes + Array genotypes
  - Hospital in-patient record
  - Verbal questionnaire
  - Cancer and death registry
  - Imaging (MRI) derived features
  - Biomarker measurements
  - Lifestyle measurements



**STUDY COHORT**

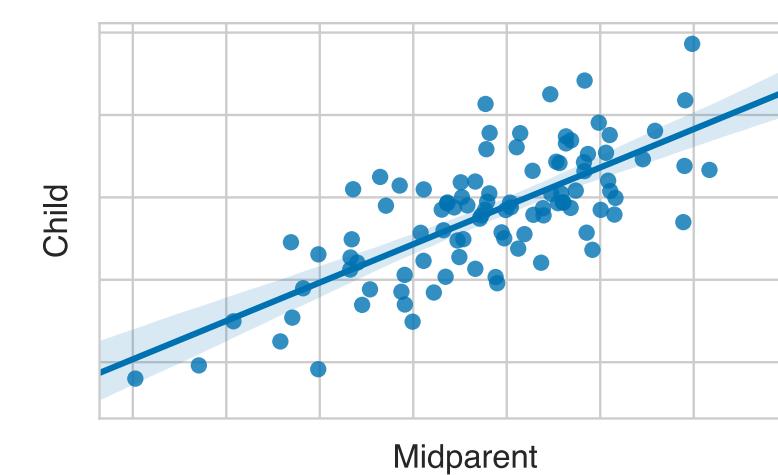
- UK Biobank  
337,208 individuals  
Hospital in-patient record  
Verbal questionnaire data  
Cancer and death registry



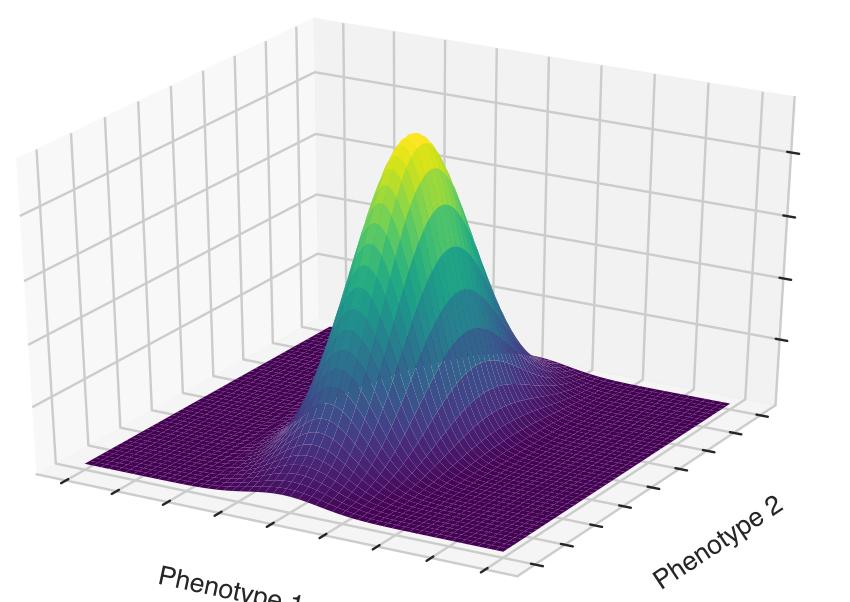
# Bayesian Mixture Model to estimate genetic parameters

Estimate genetic parameters using GWAS summary statistics

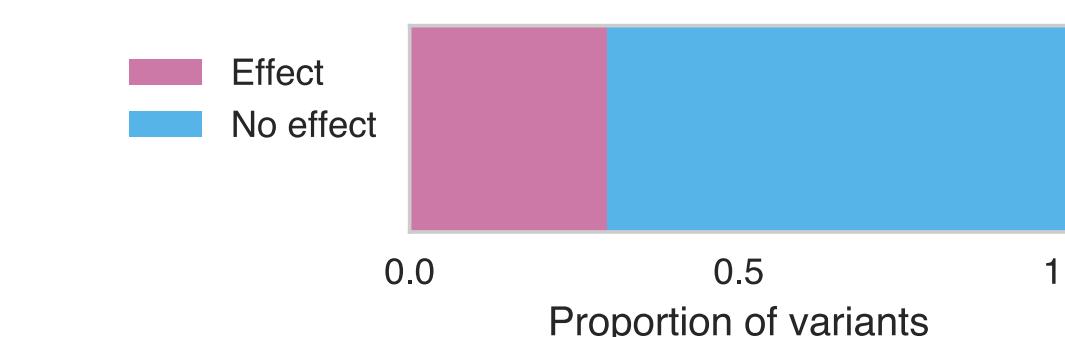
- $h^2$ : heritability



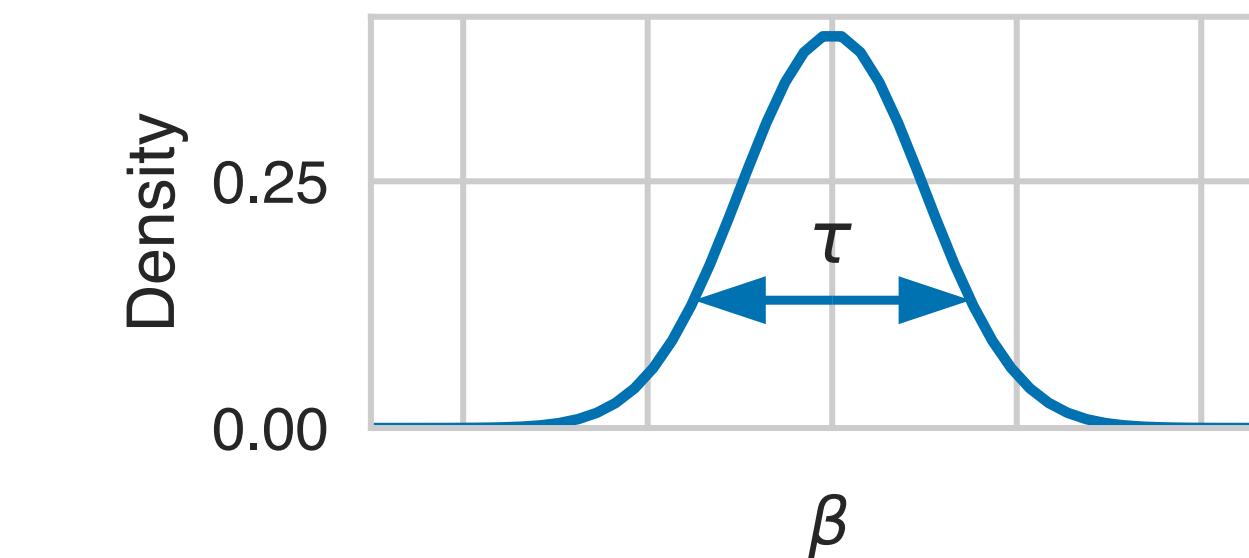
- $\Omega$ : genetic correlation



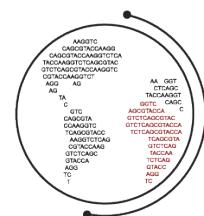
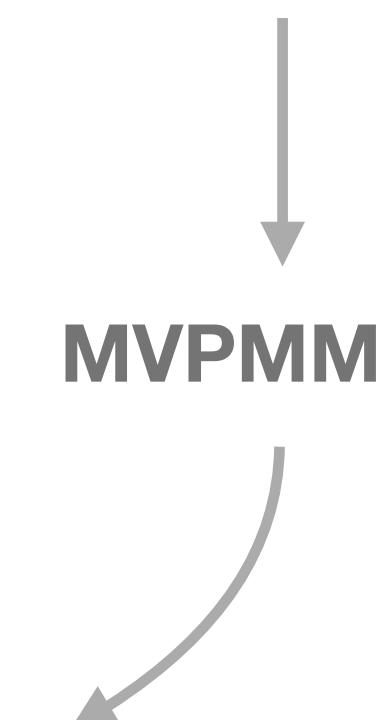
- $\pi$ : membership/polygenicity



- $\tau$ : spread/scale of effects



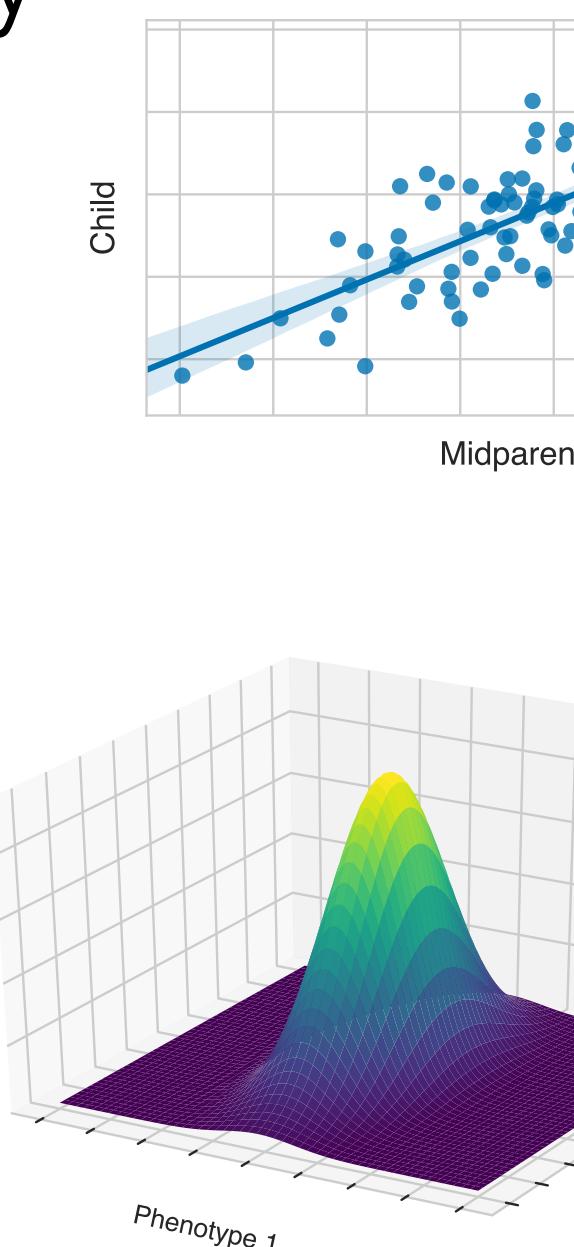
Phenotype 1		Phenotype 2	
BETA	SE	BETA	SE
1.2	0.03	1.15	0.04
1.5	0.9	1.3	0.1
-1.1	0.006	-1.2	0.05
...	...	...	...



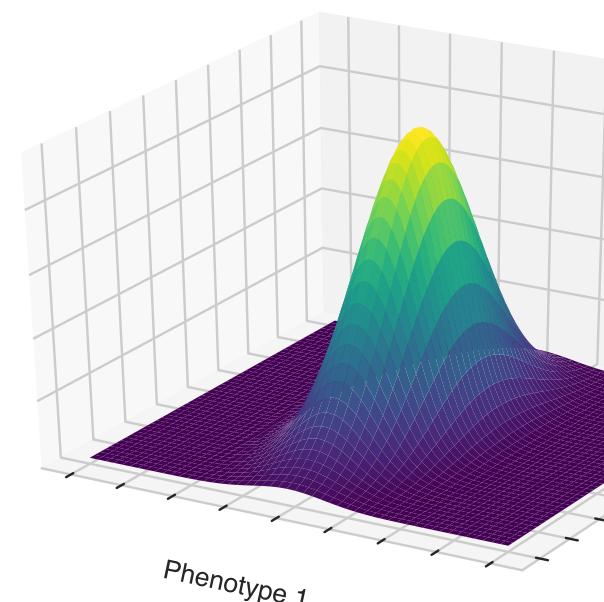
# Bayesian Mixture Model to estimate genetic parameters

Estimate genetic parameters using GWAS summary statistics

- $h^2$ : heritability



- $\Omega$ : genetic correlation

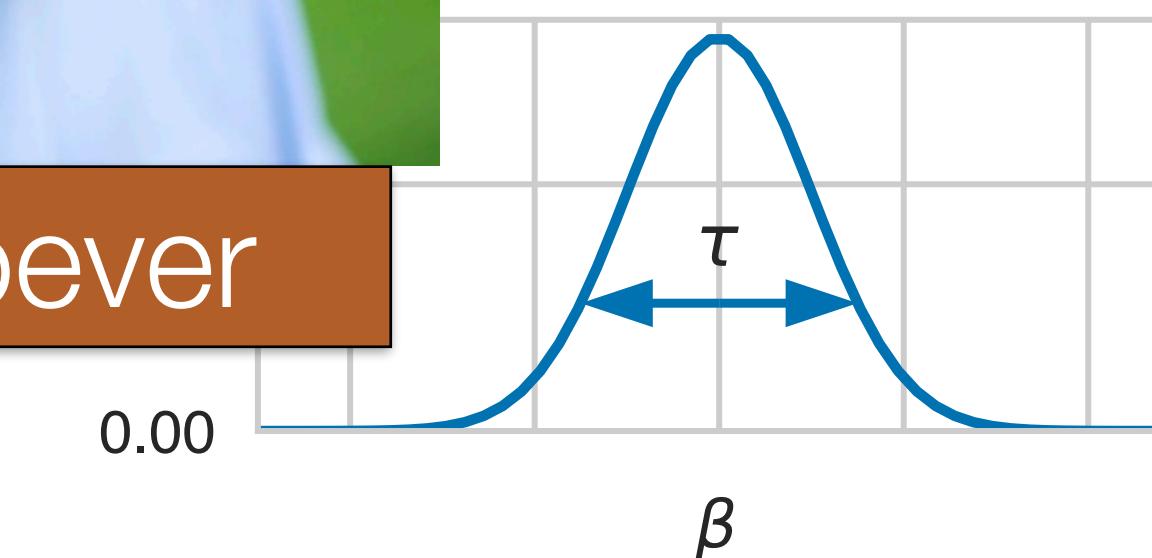


Chris DeBoever

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
...	...	...	...

Genetic parameter estimates across diseases, biomarkers, and lifestyle measures in the UK Biobank

d/scale of effects



# Multivariate Polygenic Mixture Model (MVPMM)

- 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

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

Correlated errors

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

Non-null component

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

Correlated genetic effects

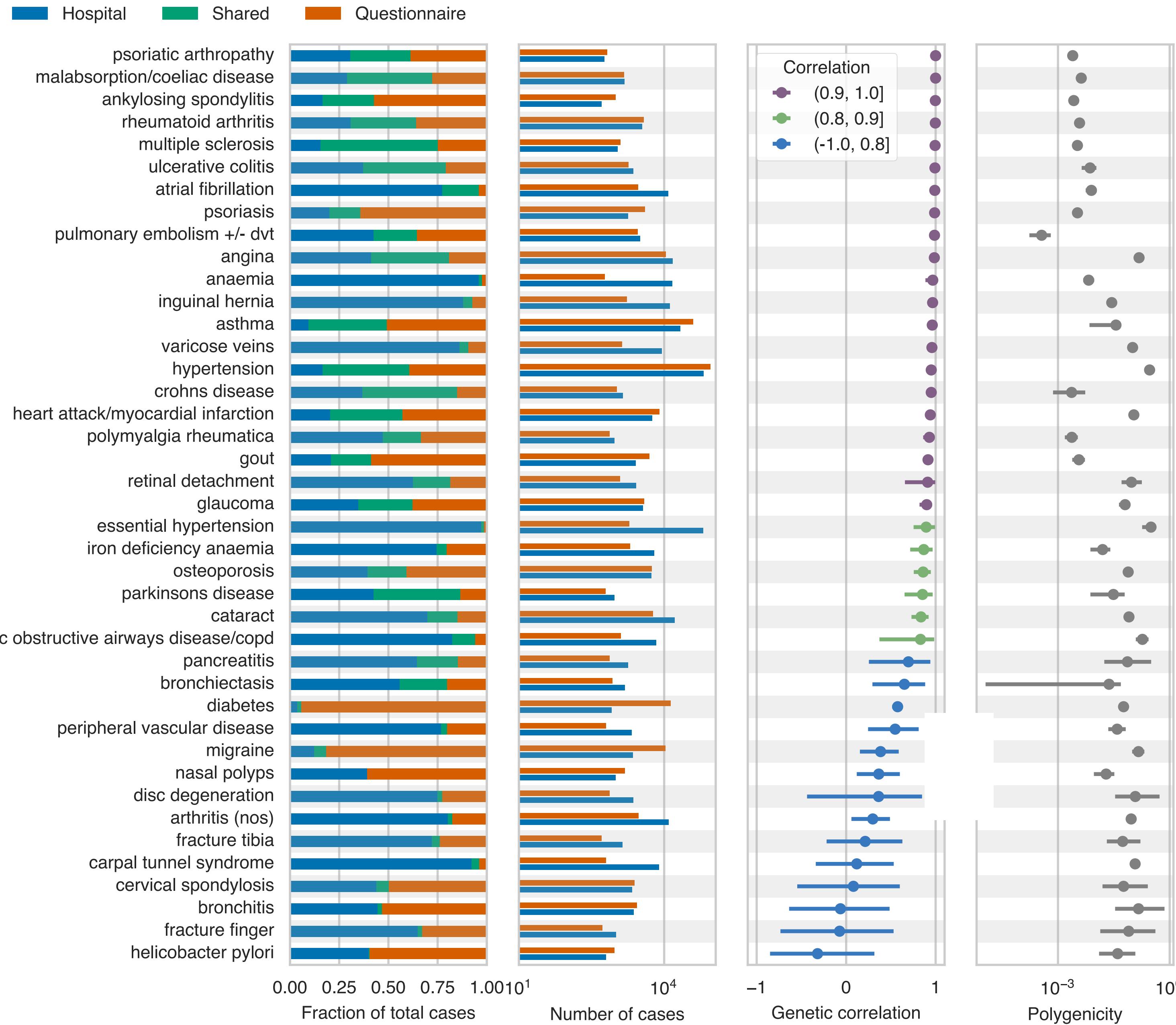
$$\Sigma_{\Omega} = \text{diag}(\tau) \cdot \Omega \cdot \text{diag}(\tau)$$

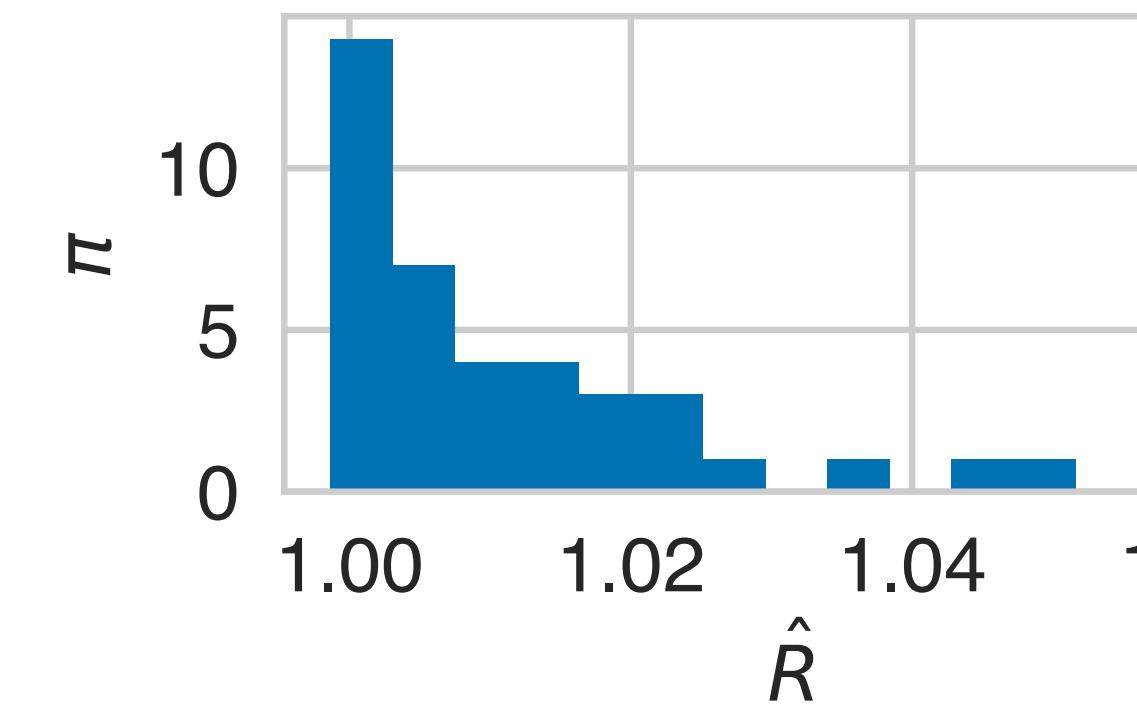
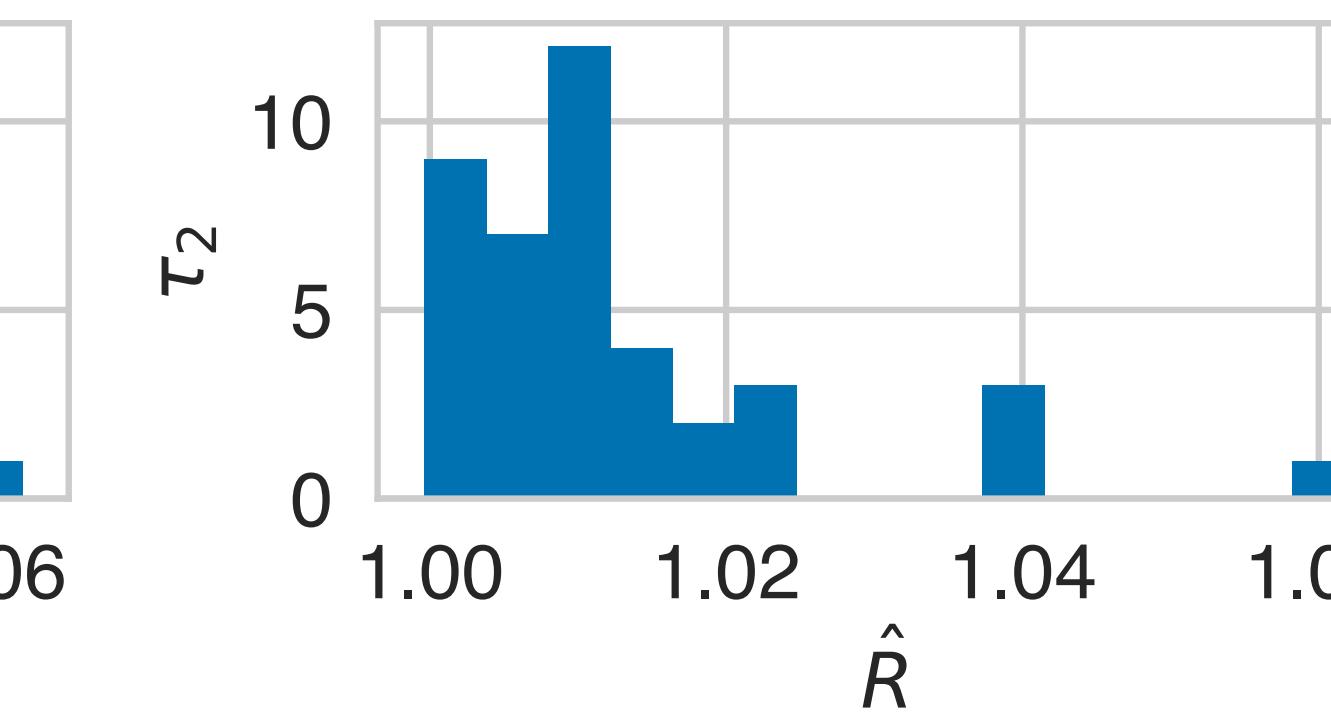
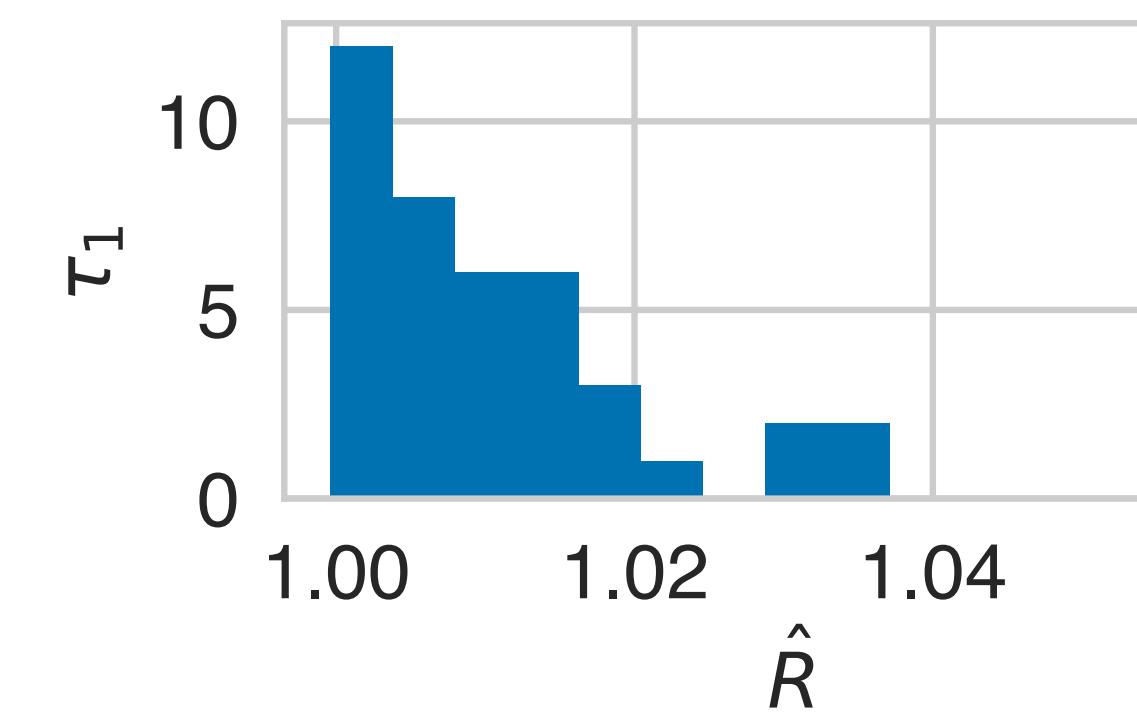
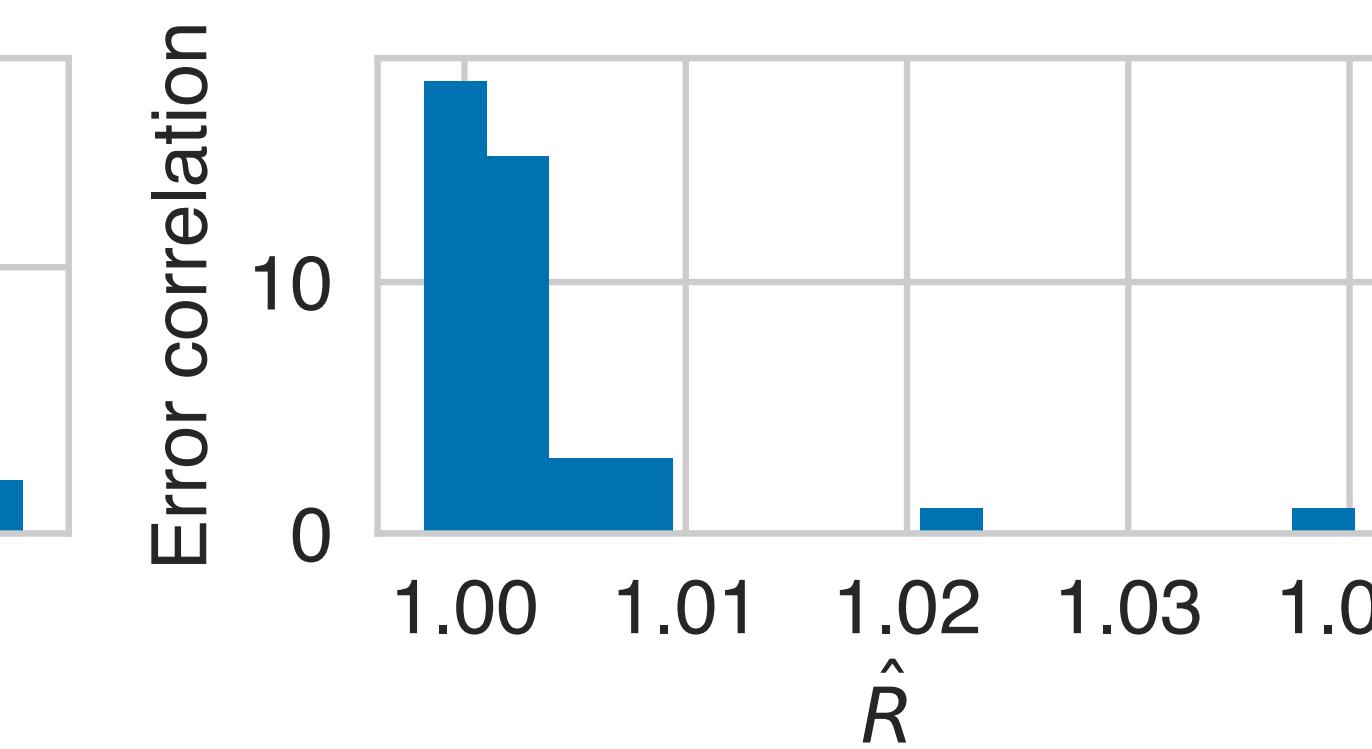
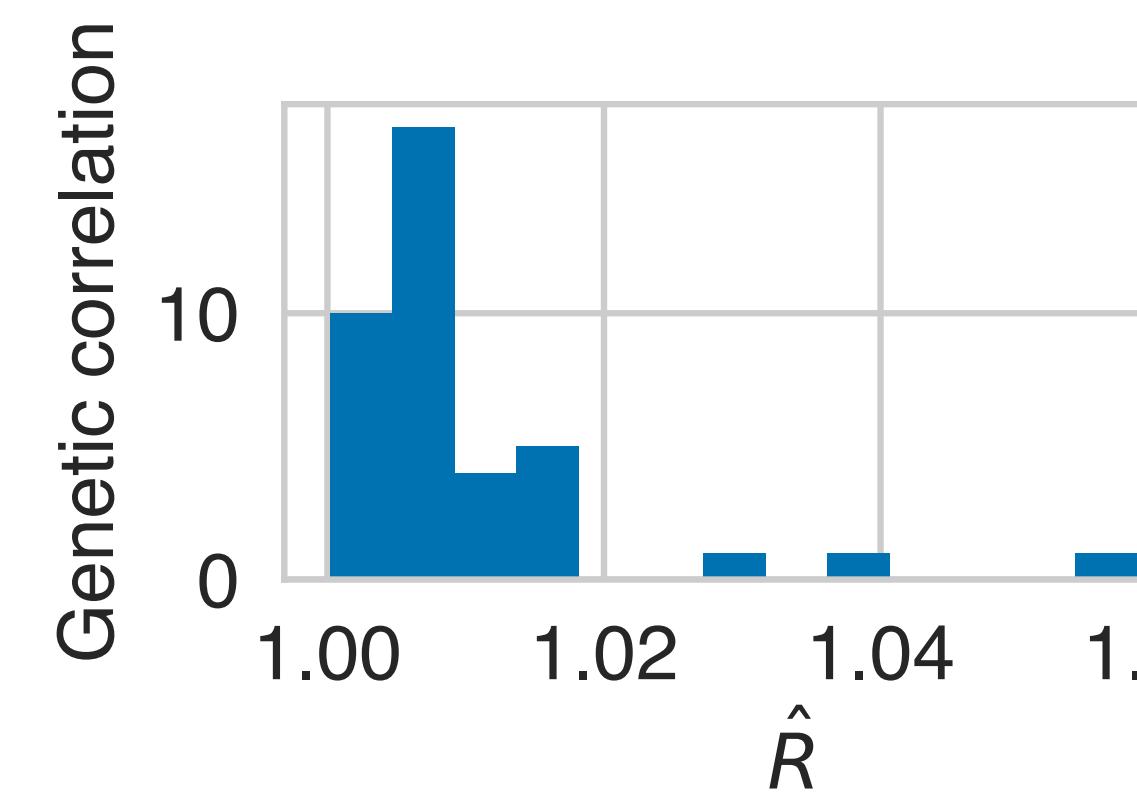
$\pi$  : mixture proportion

# Hospital record vs. verbal questionnaire digital phenotyping

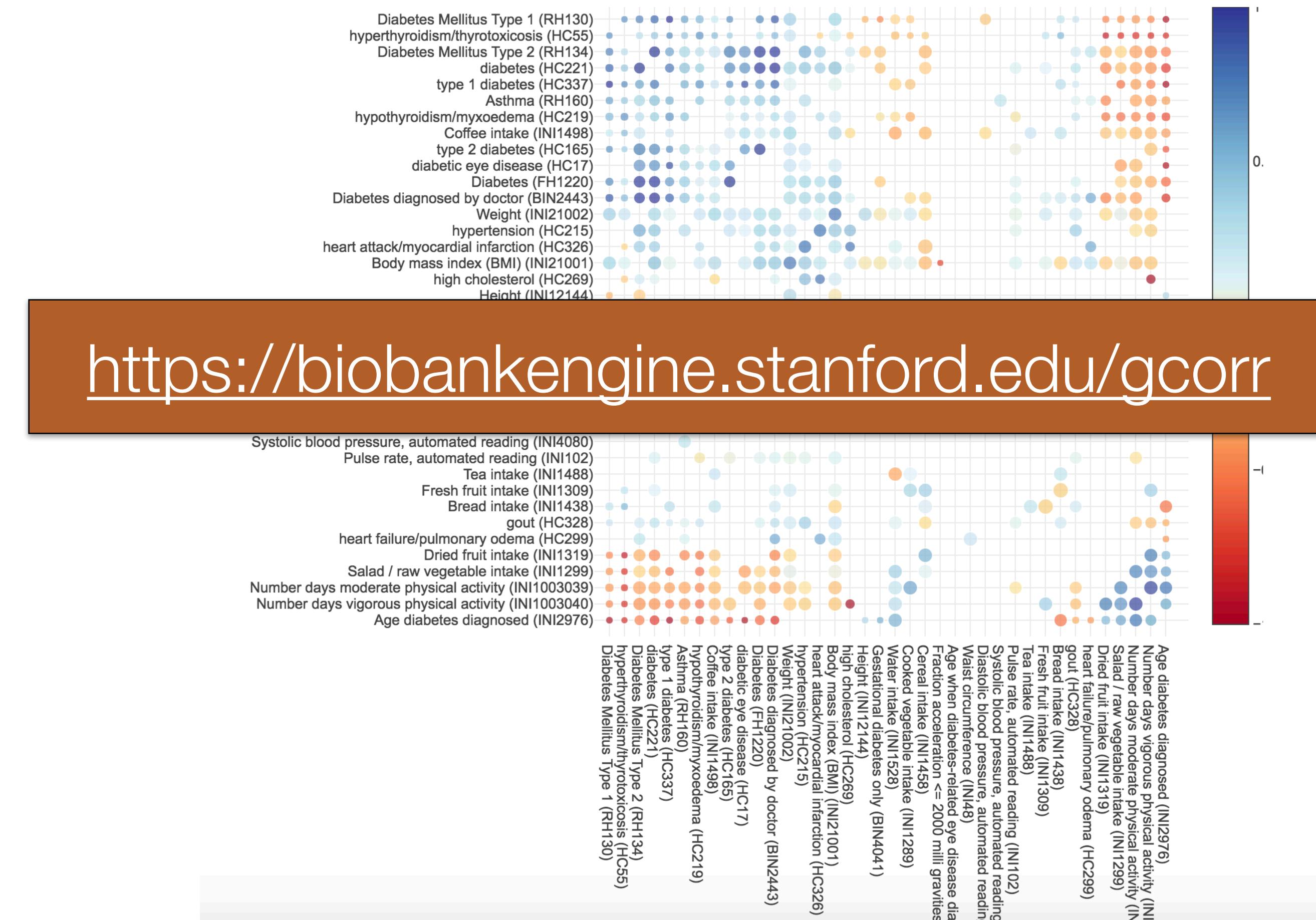
- We have defined binary disease traits using a combination of hospital records and verbal questionnaire answers
  - Question: to what extent do these two phenotyping methods agree?
  - Approach
    - Run GWAS using only cases from hospital records or only cases from questionnaire results
    - Calculate genetic correlation between GWAS results for the “same” traits

# Stanford University

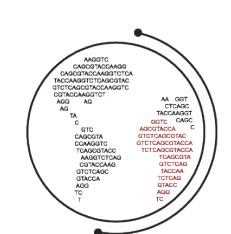




# Estimate sharing of genetic effects across diseases



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



# Bayesian Mixture Models for estimating sex-related contributions to variance

Previous studies have identified sex differences in heritability  
(Rawlik et al. 2017, Ober et al. 2008)

We developed a new approach to estimate sex-related contributions to variance from GWAS summary statistics using a Bayesian Mixture Model (BMM), implemented in Stan

Applied to quantitative traits in UK Biobank

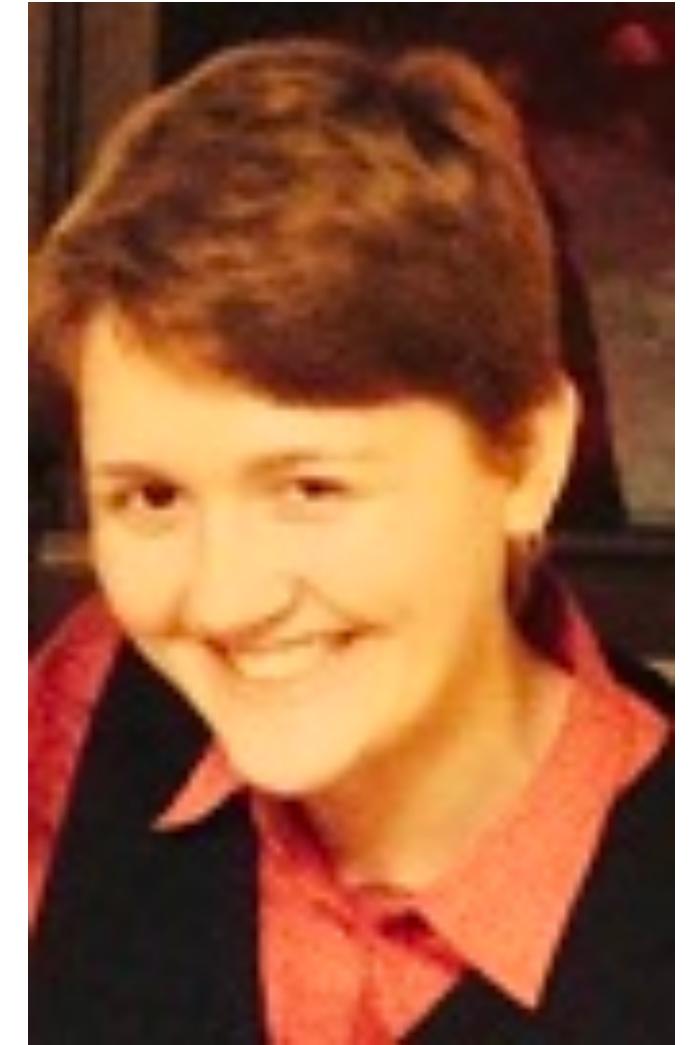


# Bayesian Mixture Models for estimating sex-related contributions to variance

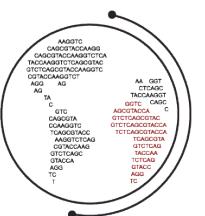
---

**Estimated component fractions and sex-specific variance-covariance matrix, and genetic correlation**

**Results similar to literature values**

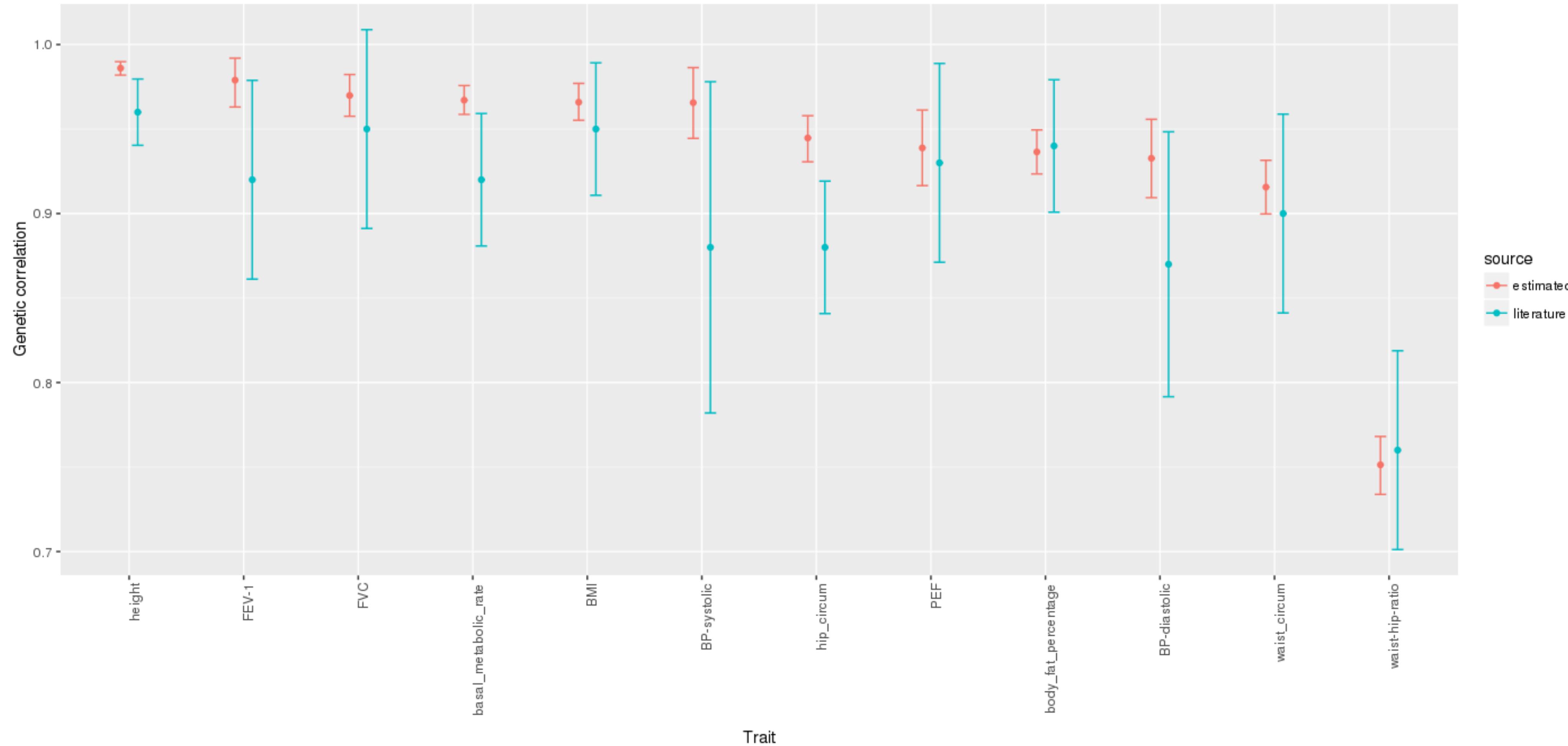


Emily  
Flynn

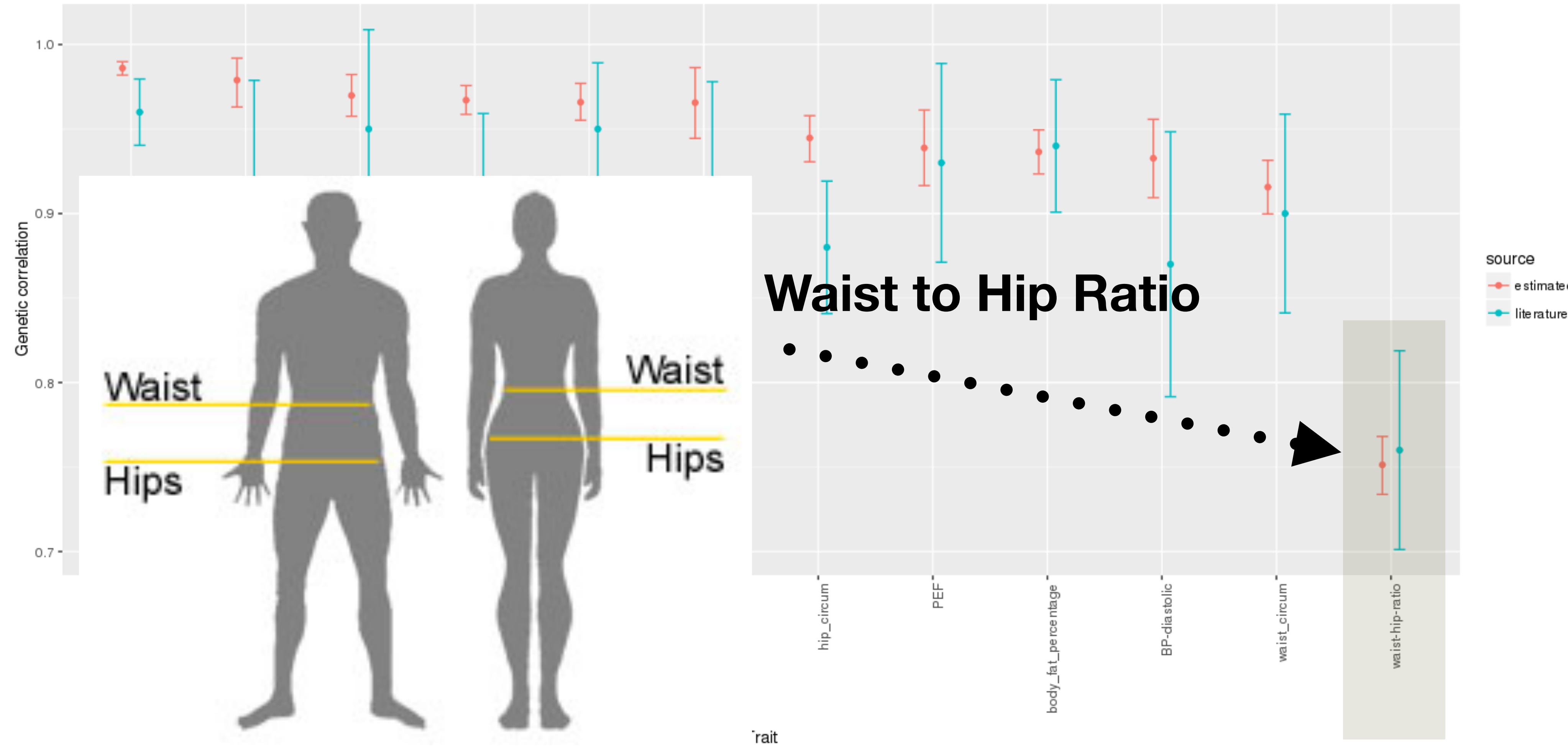


RIVASLAB

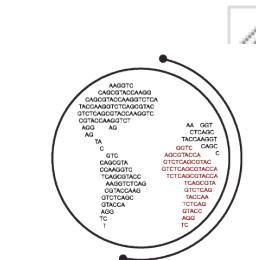
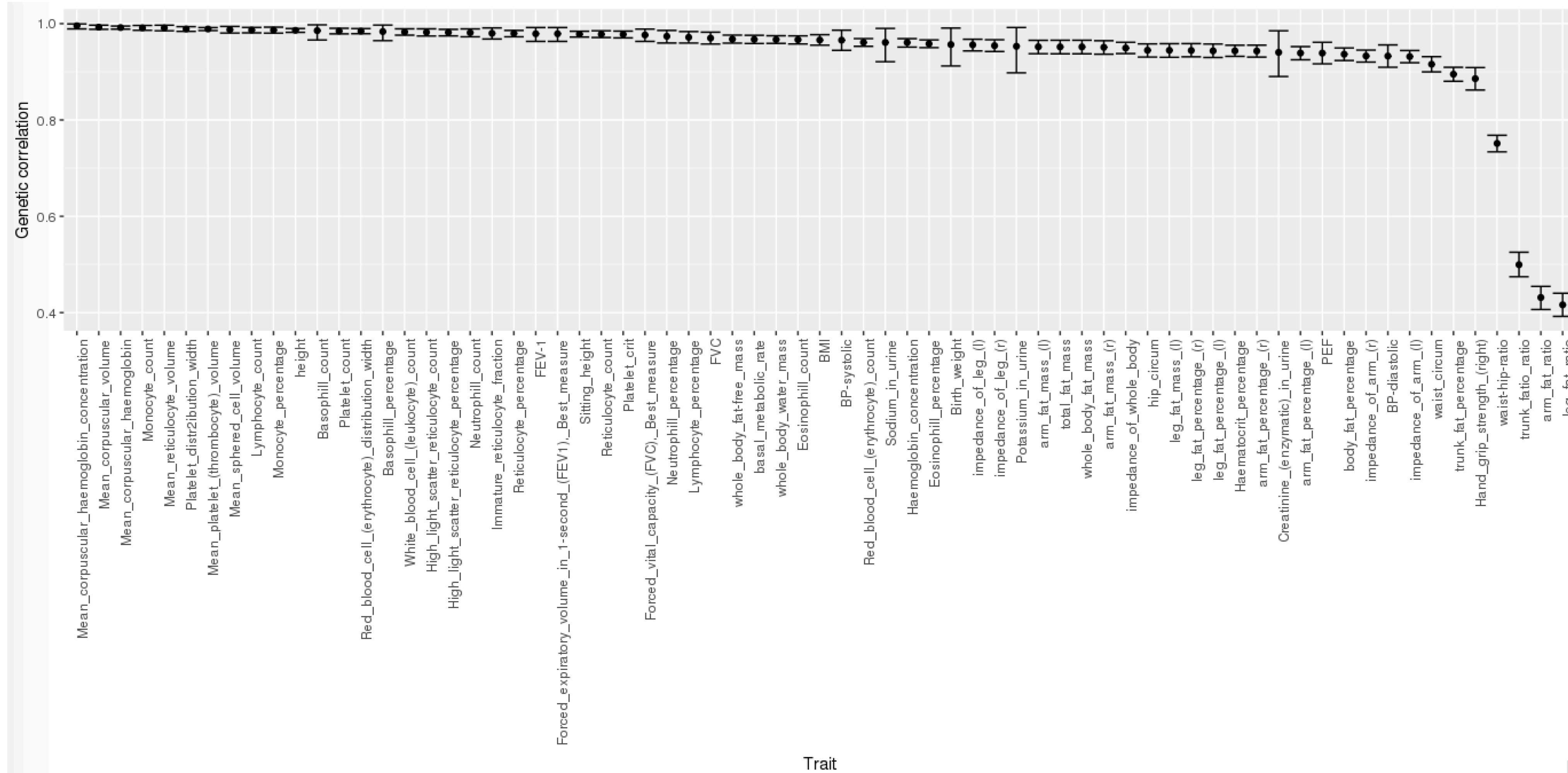
# Genetic correlation estimates are similar to the literature



# Genetic architecture of WHR differs across men and women

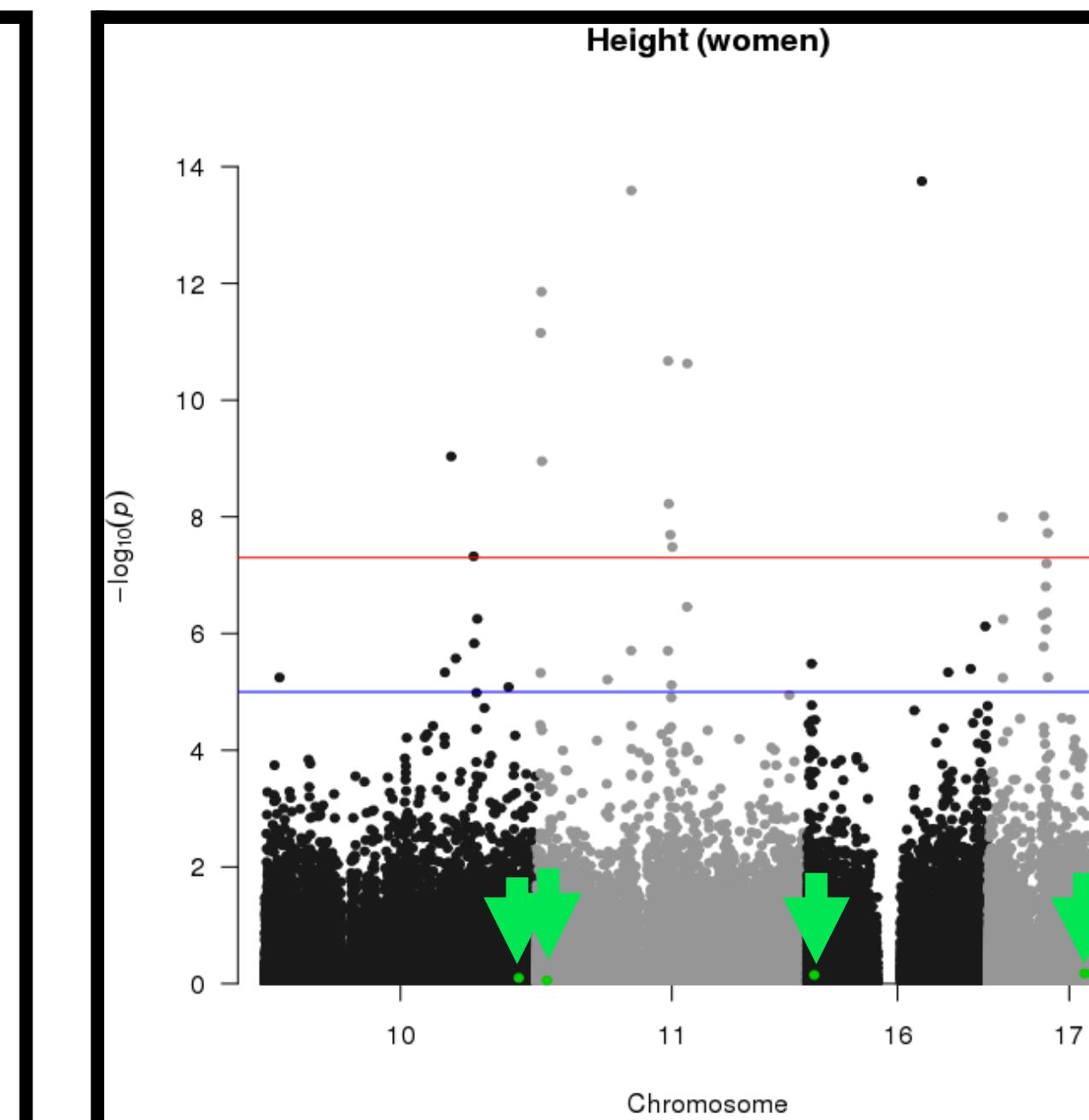
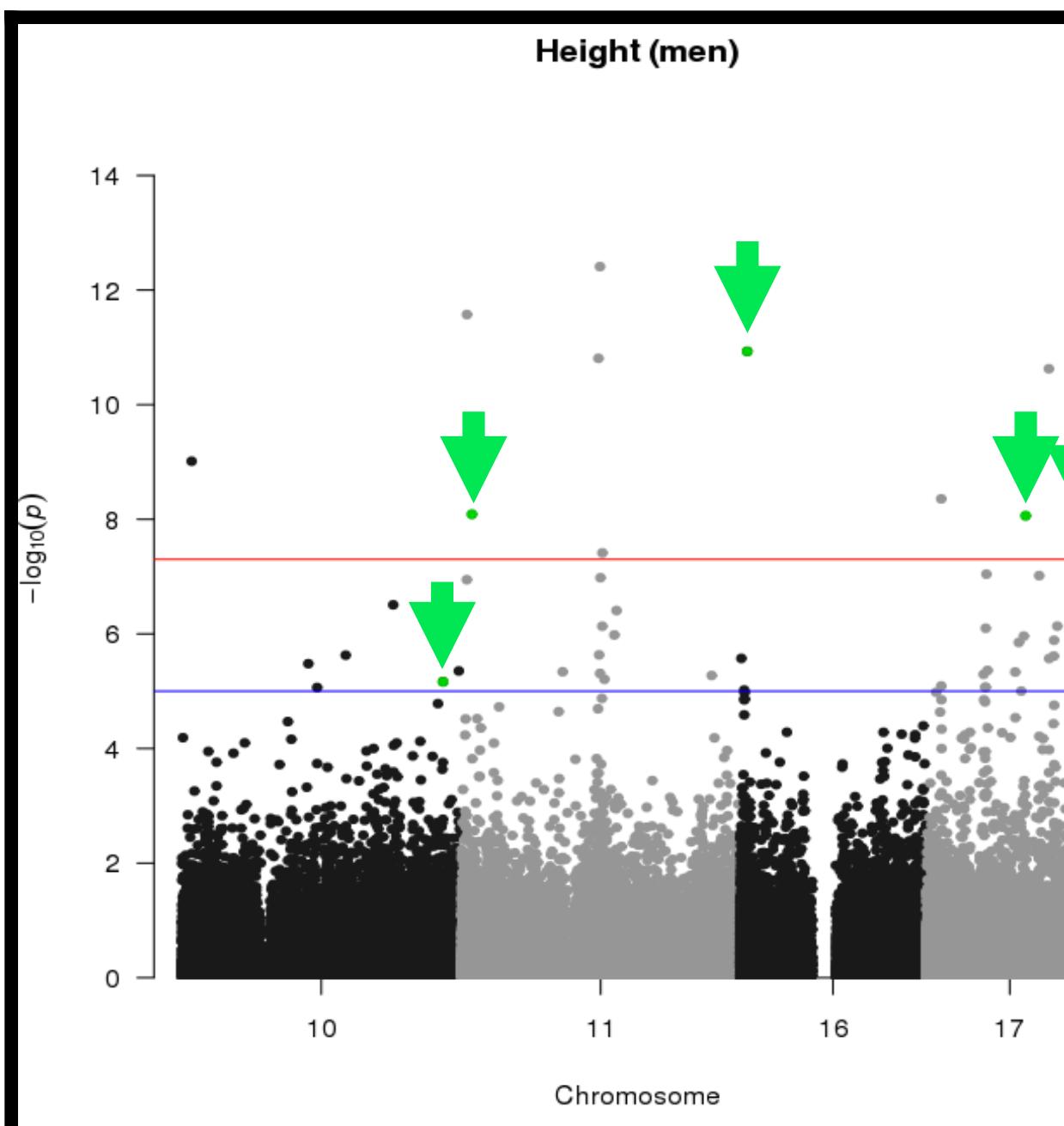


# Sex affects many genetic correlation across many traits



# Identification of sex-specific effects

- Mixture model with four components: (1) null/no effect, (2) female-specific effect, (3) male-specific effect, and (4) effects in both sexes
- Estimated component fractions and parameters, assigned SNPs based on posterior probability
- Identified five male-specific SNPs in height

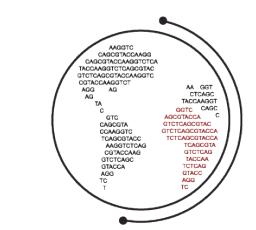


Male-specific SNPs for height

SNP	posterior	gene	B_f	B_m	p_f	p_m
Affx-58847809	1.000	CTBP2	0.0296	-0.476	0.795	6.84E-06
Affx-35292772	0.960	HBB	0.0465	-2.14	0.881	8.17E-09
Affx-80252875	1.000	NLRC3	-0.128	-2.52	0.710	1.17E-11
Affx-80298599	1.000	SGCA	0.121	-1.57	0.668	8.68E-09
Affx-80286332	1.000	DNAI2	-0.123	-1.96	0.692	1.62E-08

# Identification of sex-specific effects

	Female-spec.	Male-specific	Literature genes <sup>1,2</sup>
Waist-hip ratio	20	0	VEGFA, COBLL1, ADAMTS9
Leg fat ratio	28	0	ADAMTS3, ID4
Arm fat ratio	11	4	SLC12A2
Trunk fat ratio	44	0	PCSK5, ADAMTS14, ADAMTS17



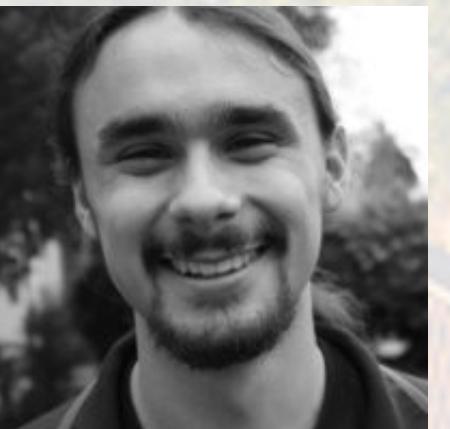
# Conclusions

---

Ample opportunities in linking



**STAN** has been a powerful statistical modeling language for our applications



**Oliver Bear Don't  
Walk IV**

Natural Language  
Processing (NLP)  
Clinical Note Tagger



**Chris DeBoever**

Gene discovery,  
statistical methods  
development



**Yosuke Tanigawa**

Pathway identification,  
Risk modeling, and  
methods development



**Matthew Aguirre**

Ancestry inference,  
Global Biobank Engine



**Johanne Justesen**  
Disease progression  
Temporal data



**Julia Olivieri**  
Inference of HLA  
allelotypes to human  
diseases

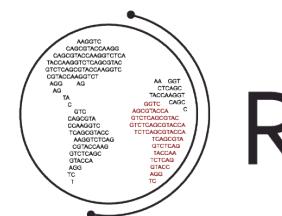


**Guhan Venkataraman**  
Integrating genomic  
data with clinical notes



**Anna Cichonska**  
Machine Learning,  
Drug-protein-genetics  
target repositioning and  
interaction inference

## The Lab



**RIVASLAB**