

# WORKSHOP POLYGENIC SCORES

**Palle Duun Rohde** [[palledr@hst.aau.dk](mailto:palledr@hst.aau.dk)]

Associate Professor & Research Group Leader

Genomic Medicine

Department of Health Science and Technology

Aalborg University

<https://pdrohde.github.io/>

# THE PURPOSE OF TODAY

- ❖ Give an introduction to polygenic scores (PGS)
- ❖ Provide an introduction to complex trait genetics
  - Monogenic vs multifactorial aetiology
- ❖ How we can utilize genomic data to elucidate molecular genetic aetiology underlying complex traits
  - Genome-wide association studies (GWAS)
- ❖ Stratify a population/cohort by their inherent genetic load towards common complex diseases
  - Polygenic scores (PGS)
- ❖ Identify future projects of common interests



# AGENDA

<b>08:00 – 08:30</b>	Welcome and common introductions
<b>08:30 – 09:10</b>	Session 1: Introduction to Polygenic Scores (PGS)
<b>09:10 – 09:20</b>	Break
<b>09:20 – 10:00</b>	Session 2: Data Sources and Computational Methods
<b>10:00 – 10:10</b>	Break
<b>10:10 – 10:40</b>	Session 3: Evaluating and Interpreting Polygenic Scores
<b>10:40 – 11:00</b>	Break
<b>11:00 – 11:45</b>	Session 4: Advanced Applications and Future Directions
<b>11:45 – 12:30</b>	Lunch and short walk
<b>12:30 – 15:30</b>	Identification of 2-3 projects of common interest
<b>15:30 – 16:00</b>	Next steps and thank you for today

# AGENDA

<b>08:00 – 08:30</b>	Welcome and common introductions
08:30 – 09:10	Session 1: Introduction to Polygenic Scores (PGS)
09:10 – 09:20	Break
09:20 – 10:00	Session 2: Data Sources and Computational Methods
10:00 – 10:10	Break
10:10 – 10:40	Session 3: Evaluating and Interpreting Polygenic Scores
10:40 – 11:00	Break
11:00 – 11:45	Session 4: Advanced Applications and Future Directions
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today

# AGENDA

08:00 – 08:30	Welcome and common introductions
<b>08:30 – 09:10</b>	<b>Session 1: Introduction to Polygenic Scores (PGS)</b>
09:10 – 09:20	Break
09:20 – 10:00	Session 2: Data Sources and Computational Methods
10:00 – 10:10	Break
10:10 – 10:40	Session 3: Evaluating and Interpreting Polygenic Scores
10:40 – 11:00	Break
11:00 – 11:45	Session 4: Advanced Applications and Future Directions
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today

# SESSION 1

- Precision Medicine?
- Complex traits?
- Genetic variants as chilies
- What is a polygenic score?
- What is needed to compute a polygenic score?
- Why so many different methods?



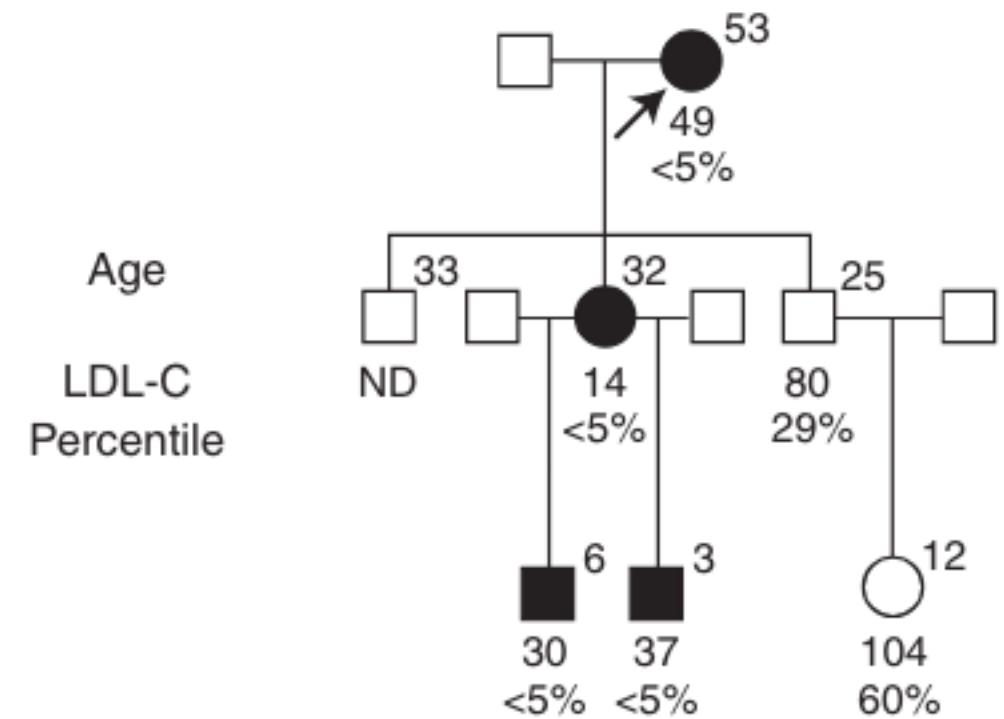
# WHY STUDY THE HUMAN GENOME

- ❖ A story about the gene *PCSK9*.



# A GENETIC VARIANT IN *PCSK9* IS ASSOCIATED WITH CHOLESTROL LEVELS

- LDL is a strong risk factor for heart disease.
- Carriers of nonsense mutation (filled symbols, heterozygous) within the *PCSK9* gene display markedly decreased LDL-C levels.
- Carriers are apparently healthy.
- **Why is this a very interesting observation?**
  - *PCSK9* could be a new drug target for lowering LDL cholesterol.

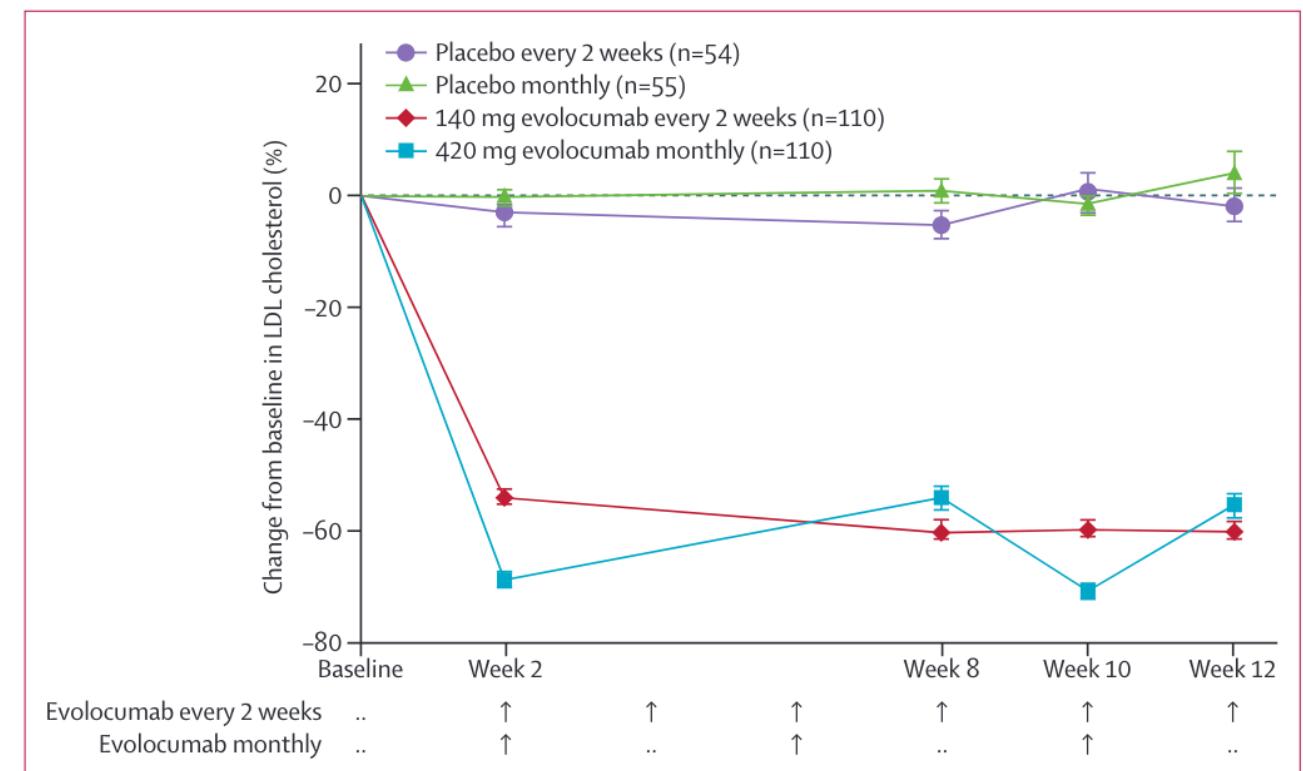


# PCSK9 INHIBITION WITH EVOLOCUMAB ON LDL CHOLESTEROL

Familial hypercholesterolaemia (FH) is a monogenic disease affecting 1 in 250 people and increases the likelihood of having coronary heart disease at a younger age.

FH increase LDL-C in blood → increase heart disease risk.

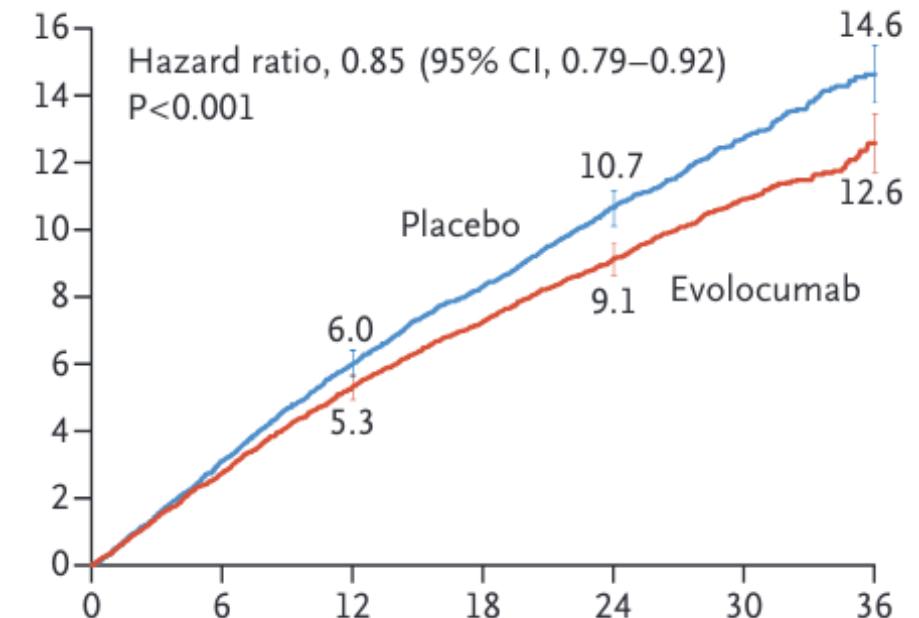
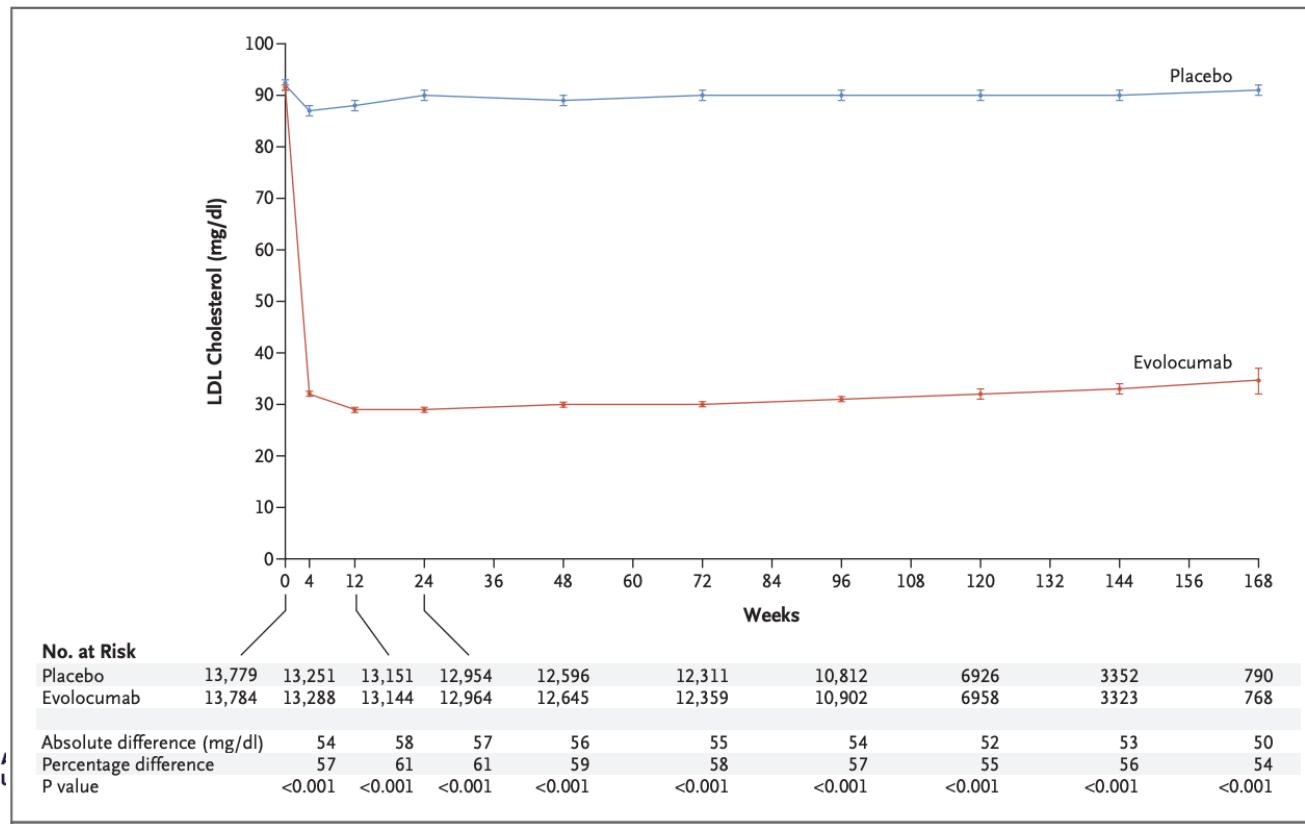
In patients that are heterozygous for FH gene mutation evolocumab reduced LDL-C.



# PCSK9 INHIBITION WITH EVOLOCUMAB ON LDL CHOLESTEROL

In patients with cardiovascular disease.

At 48 weeks reduction of 59% LDL-C with evolocumab.



Reduction in primary end points in evolocumab group compared to placebo.



WE LOOK AT VARIATION  
NATURE CREATED

Nature as a big genetic laboratory

# GENOMIC-INFORMED MEDICINE

Personalised medicine  
vs  
Precision medicine  
vs  
Individualised medicine  
  
Uhh?!

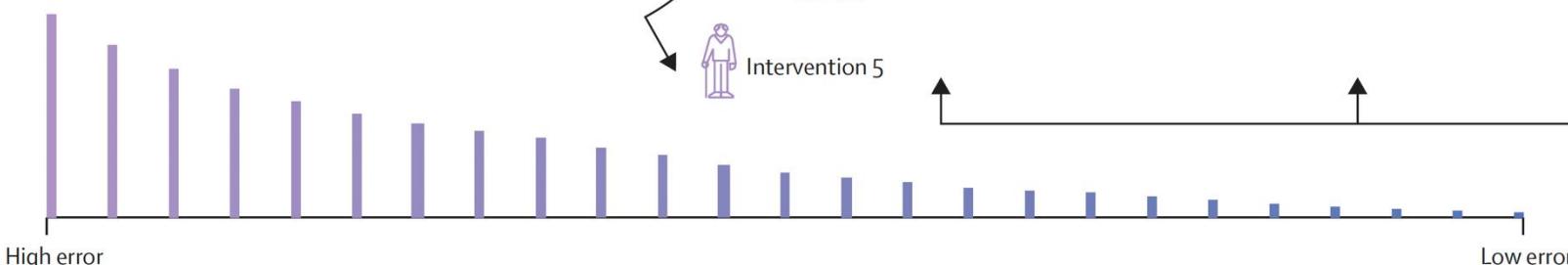
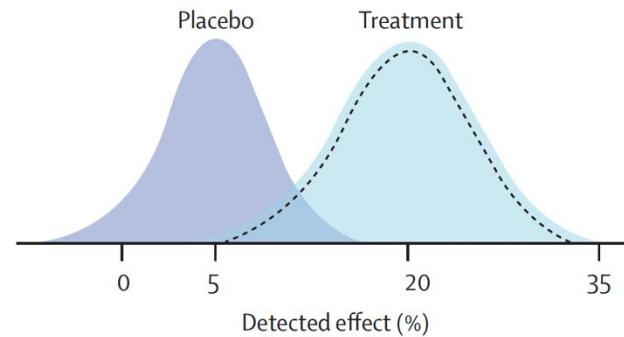


# IMPLEMENTATION OF PRECISION MEDICINE

EPPOS [evidence-based precision personalised objective subjective]

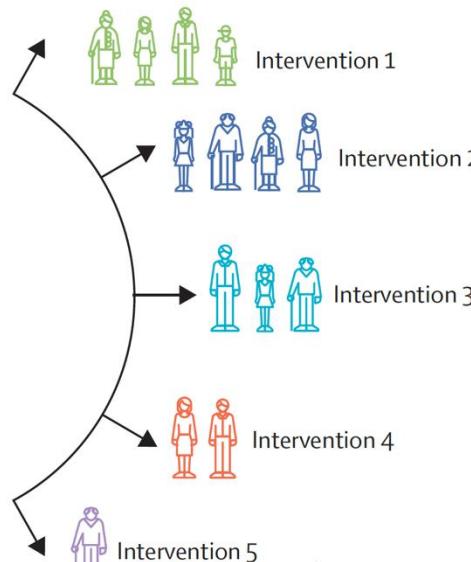
## Evidence-based Medicine

- (1) Contemporary evidence-based medicine  
Estimate average risk or response using epidemiological and clinical trial cohorts



## Precision Medicine

- (2) Probability scoring and stratification  
Maximise response and minimise risk using subclassification



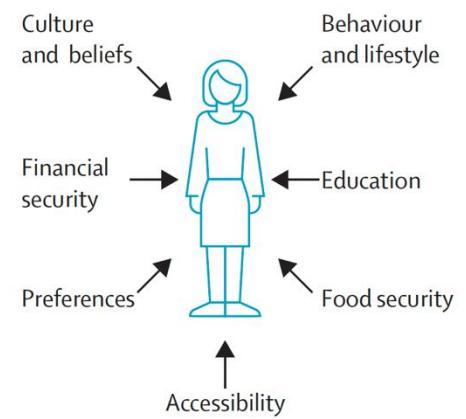
## Personalised Medicine

- (3) Personalisation (objective)  
Monitor response to optimise dose, timing, and delivery



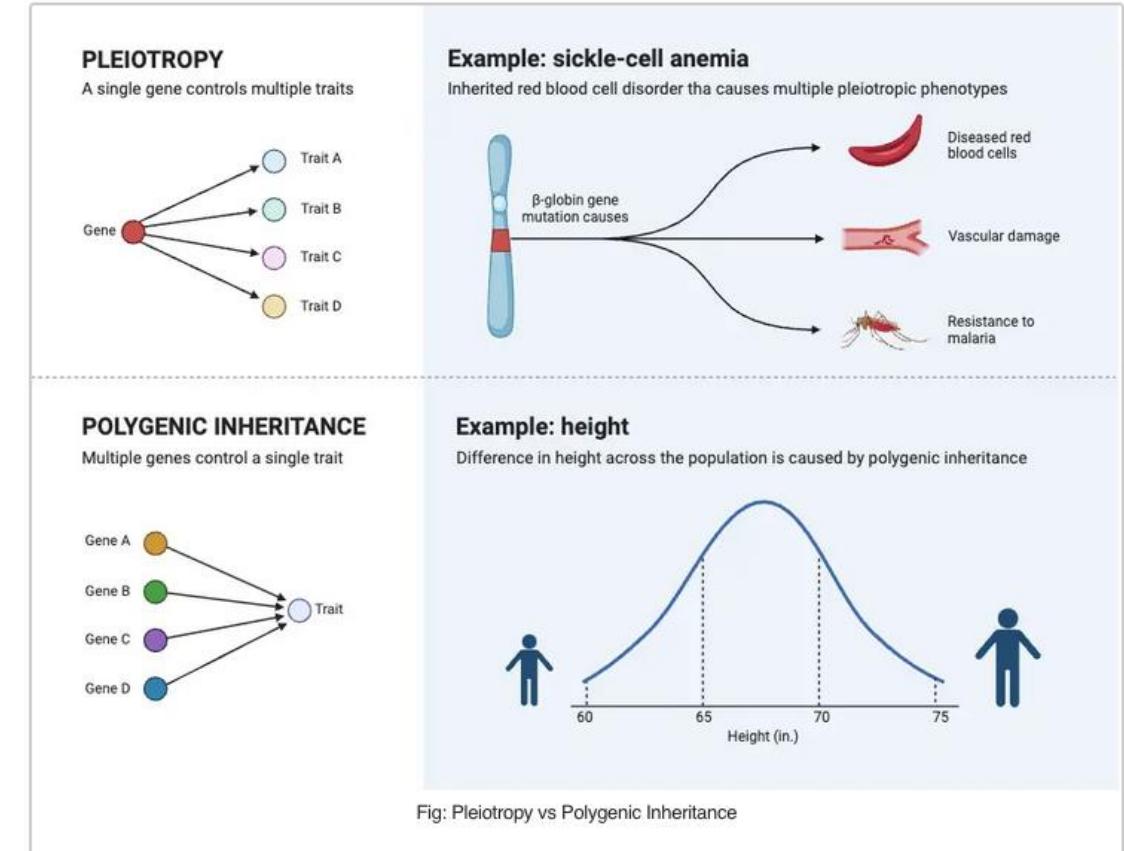
## Individualised Medicine

- (4) Personalisation (subjective)  
Adapt intervention to fit the person's needs, capabilities, and preferences



# DIFFERENT MODE OF INHERITANCES

- ❖ Monogenic (single gene variant)
- ❖ Polygenic (many gene variants)
- ❖ Multifactorial (many gene variants plus environment exposures)



# QUANTITATIVE TRAITS IN DIFFERENT SHAPES

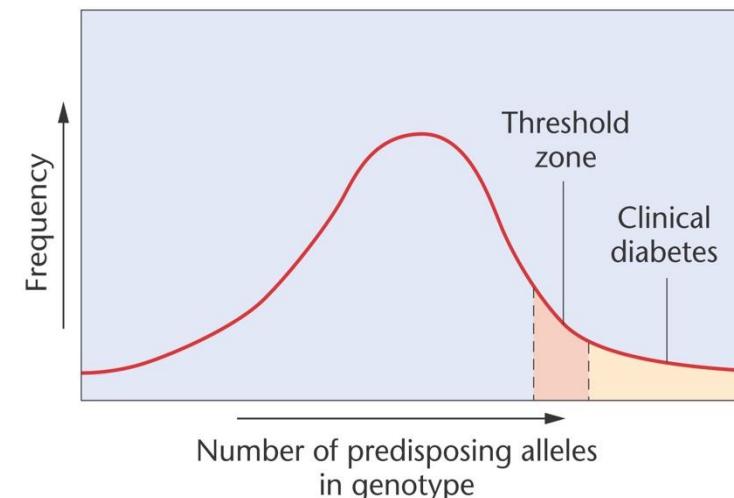
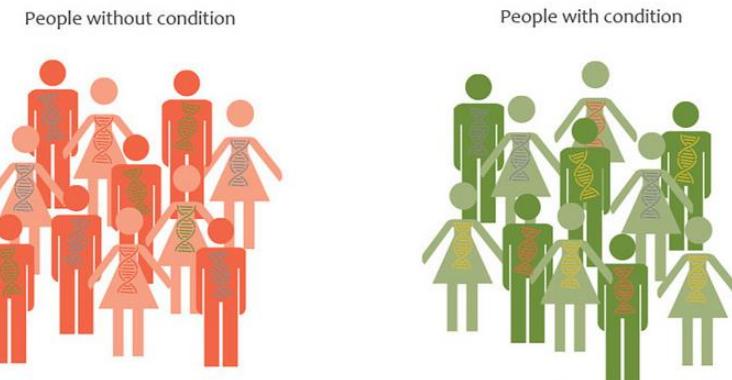
Continuous variation



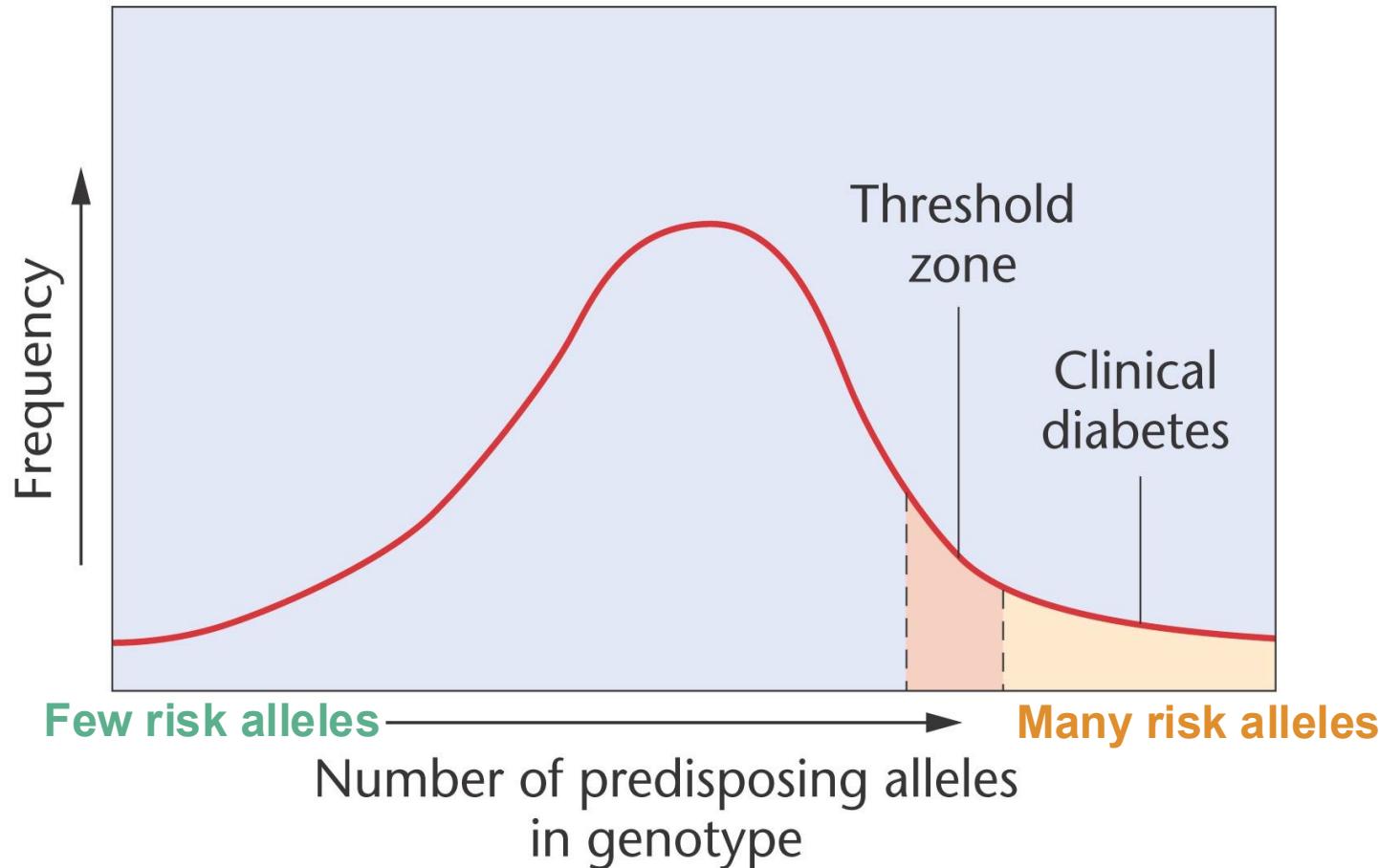
Categorical variation



Threshold traits



# LIABILITY (THRESHOLD) MODEL



## Liability model

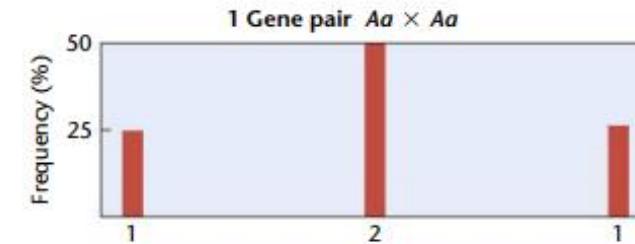
Only individuals with a liability over a certain threshold will become affected

The **sum** of many genetic variants with **small effect/risk**.

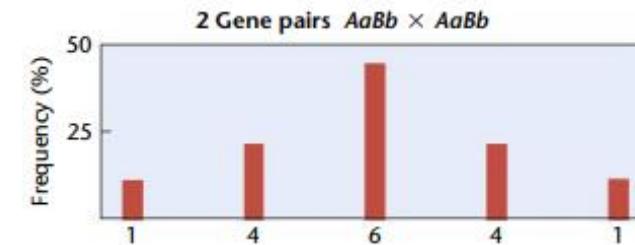
Each locus follow Mendelian inheritance pattern, although the trait does not

# POLYGENIC TRAITS

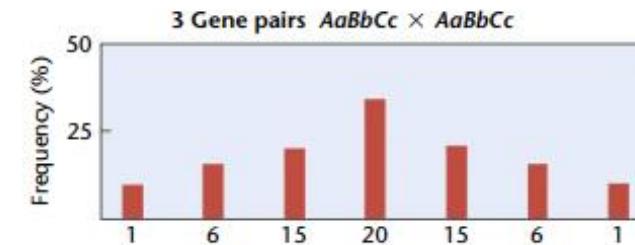
A simple ‘Mendelian’ explanation for why polygenic traits follows a normal distribution



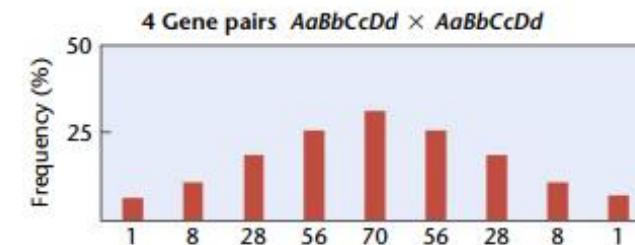
3 classes



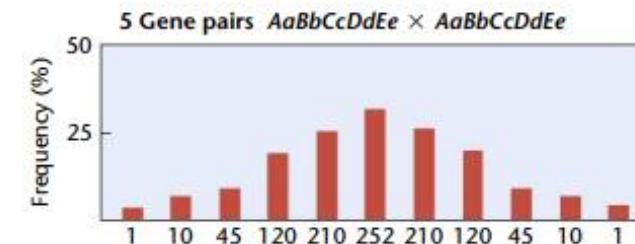
5 classes



7 classes



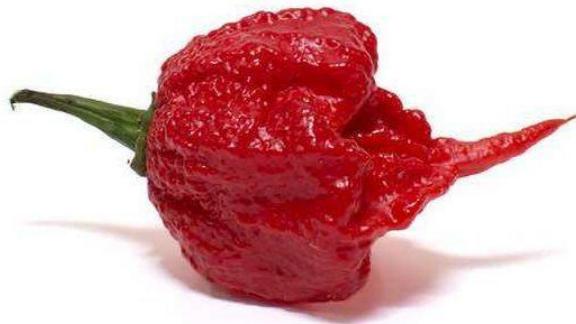
9 classes



11 classes

# GENETIC VARIATION IS LIKE CHILI

Carolina reaper



Strong effect on phenotype

Habanero Lemon



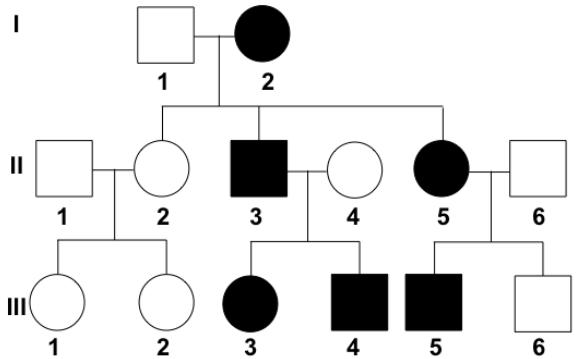
Moderate effect on phenotype

Bell pepper



Weak effect on phenotype  
(or none at all)

# GENETIC VARIATION IS LIKE CHILI



Mutation

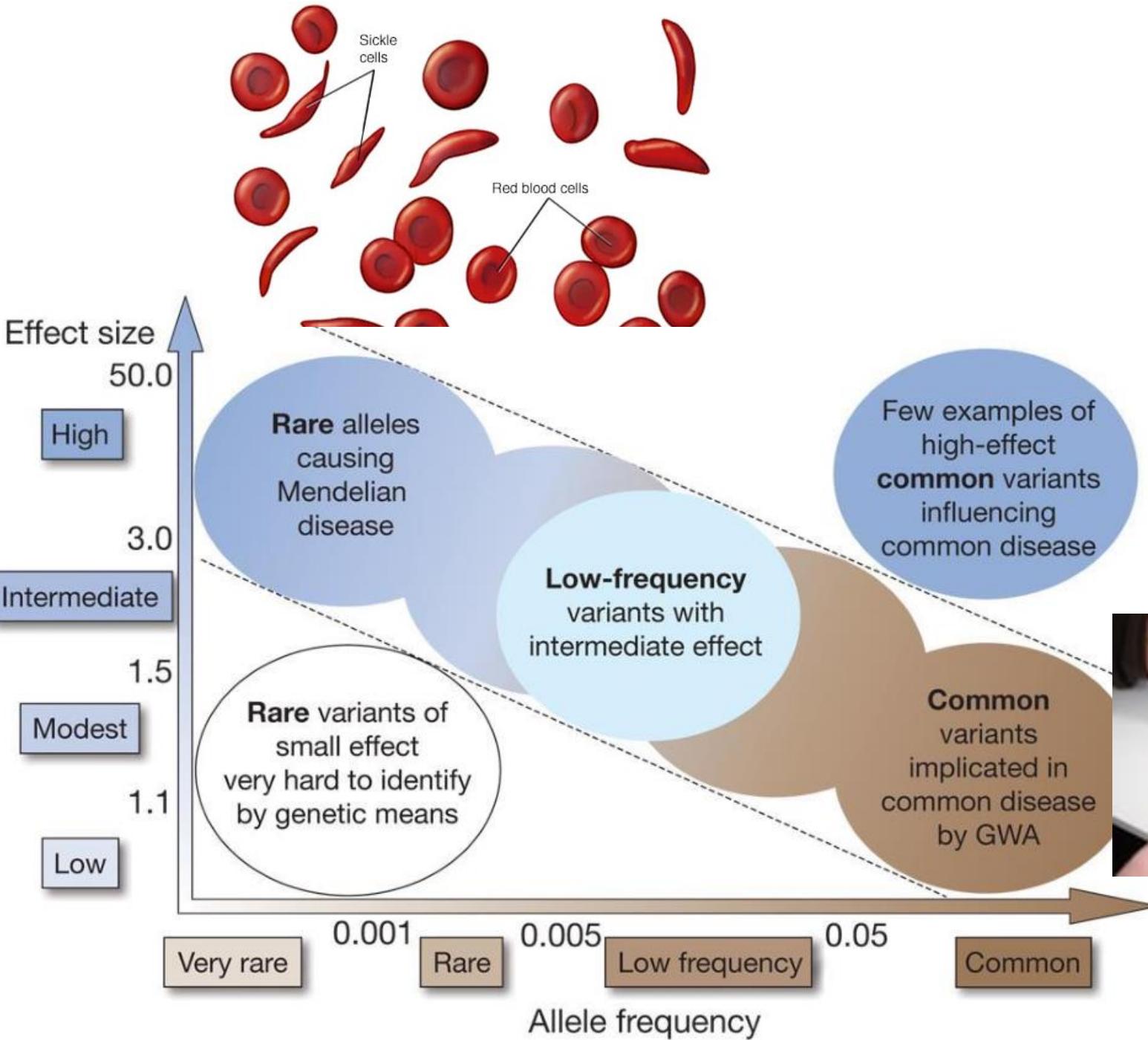


Each genetic variant is **both** necessary and sufficient



Each genetic variant is **neither** necessary nor sufficient





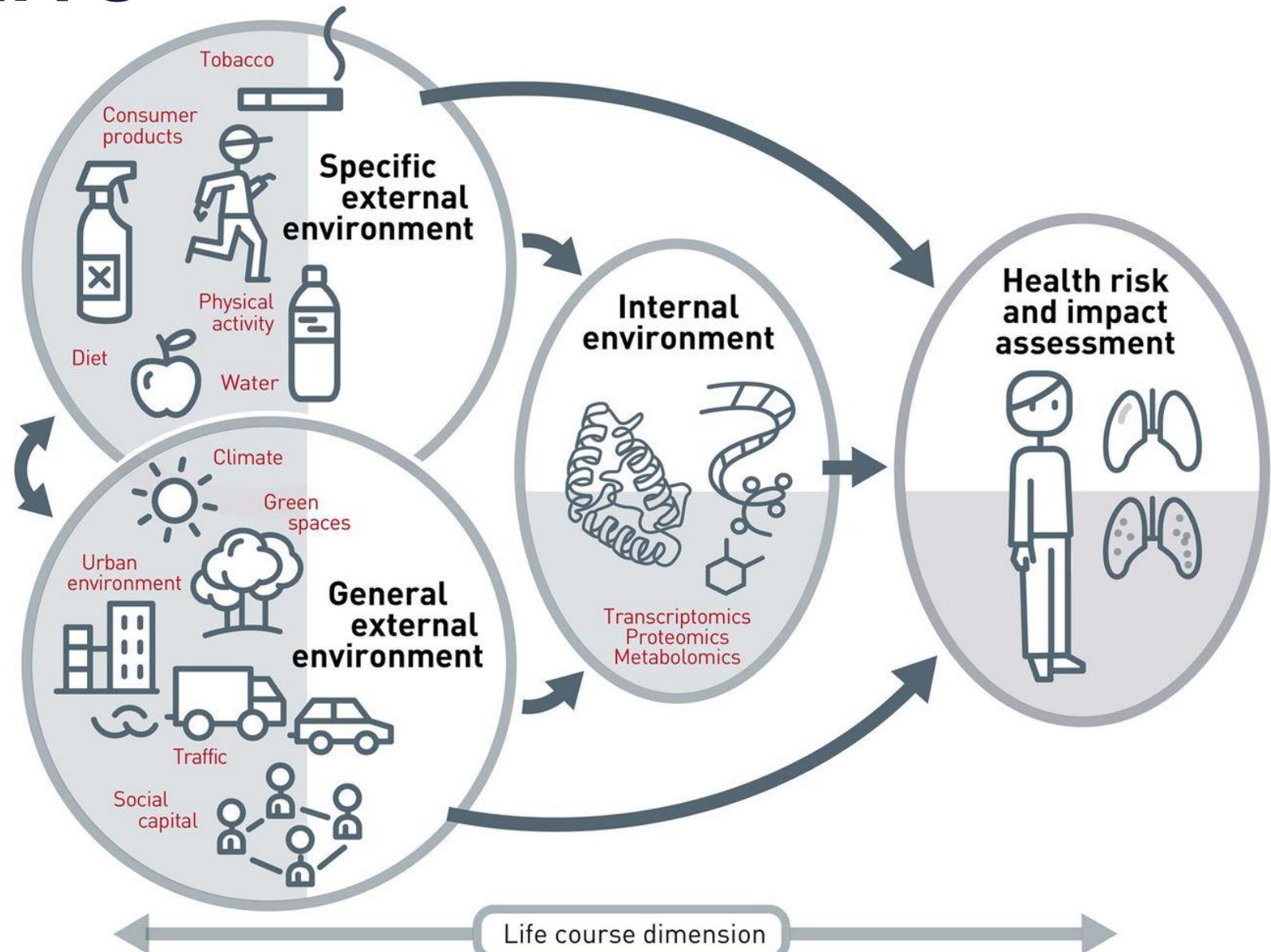
# COMPLEX TRAITS are complex...

Complex traits have a **genetic component** and an **environmental component**

The relative genetic contribution is called **heritability**.

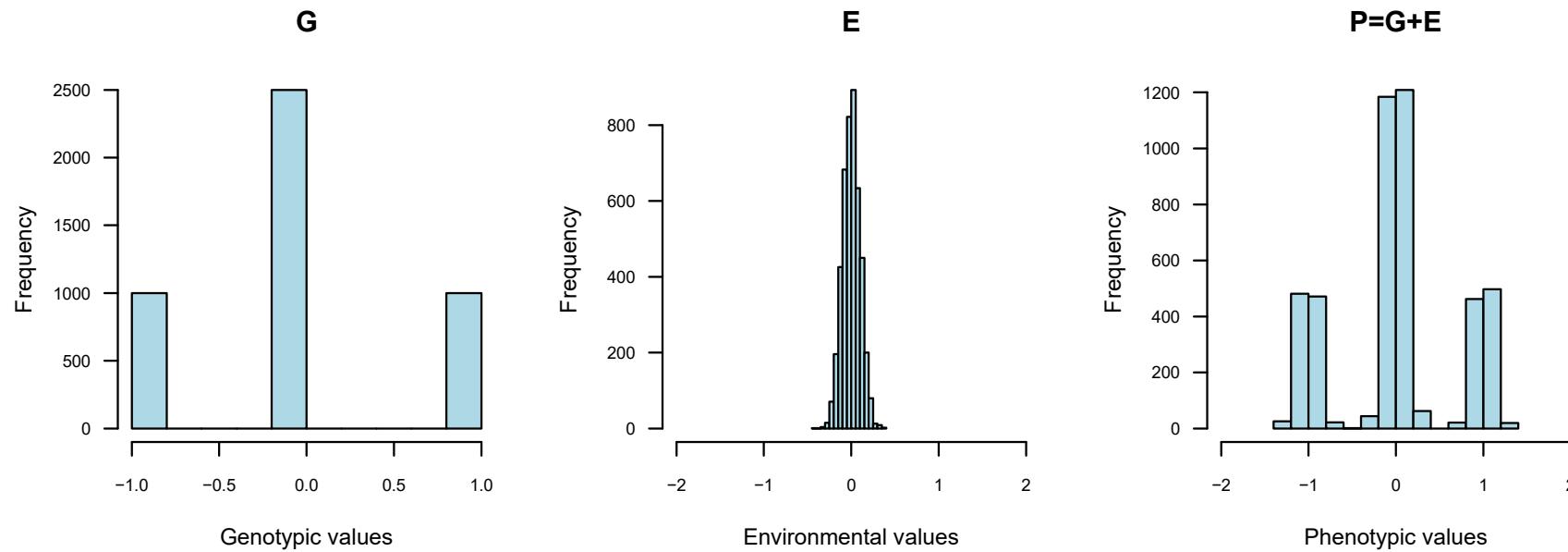
$$V_G / (V_G + V_E)$$

Thus, set a limits for the genetic predictor



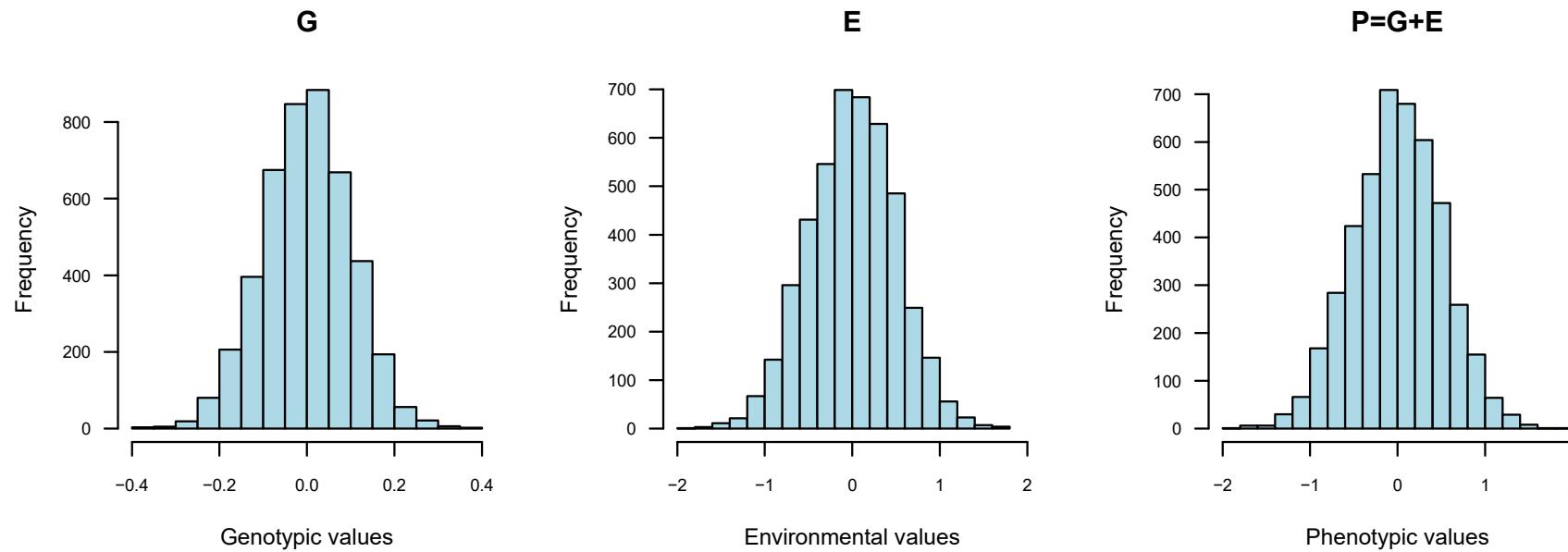
# MULTIFACTORIAL TRAITS

Environmental variance (non-genetic factors) blurs phenotypic classes



# MULTIFACTORIAL TRAITS

Genotypic and environmental variance creates infinite many phenotypic classes



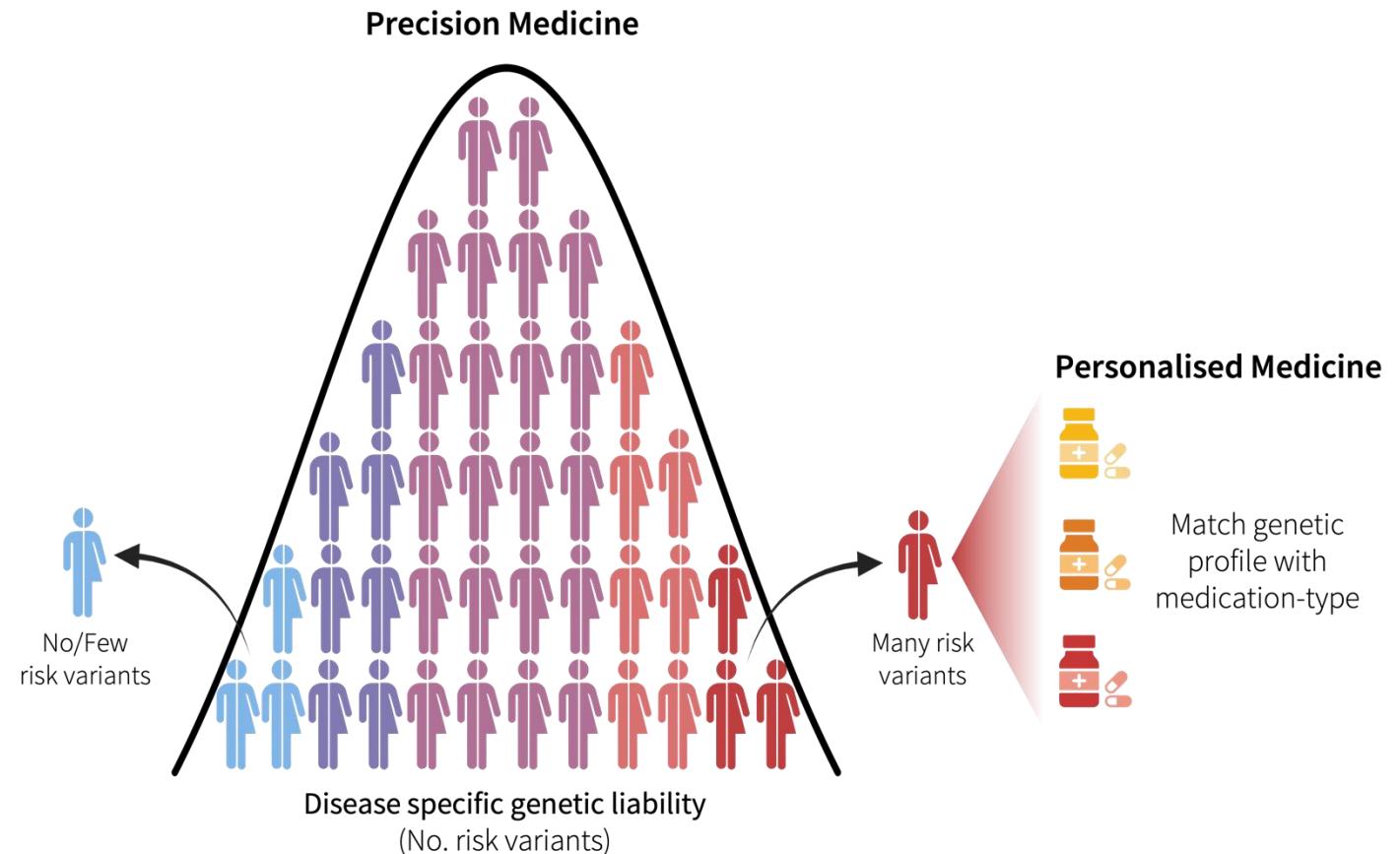
# PREDICTING DISEASE RISK FROM GENETIC DATA

A “polygenic score” is one way by which people can learn about their risk of developing a disease, based on the total number of changes (i.e., SNPs) related to the disease (NHI)

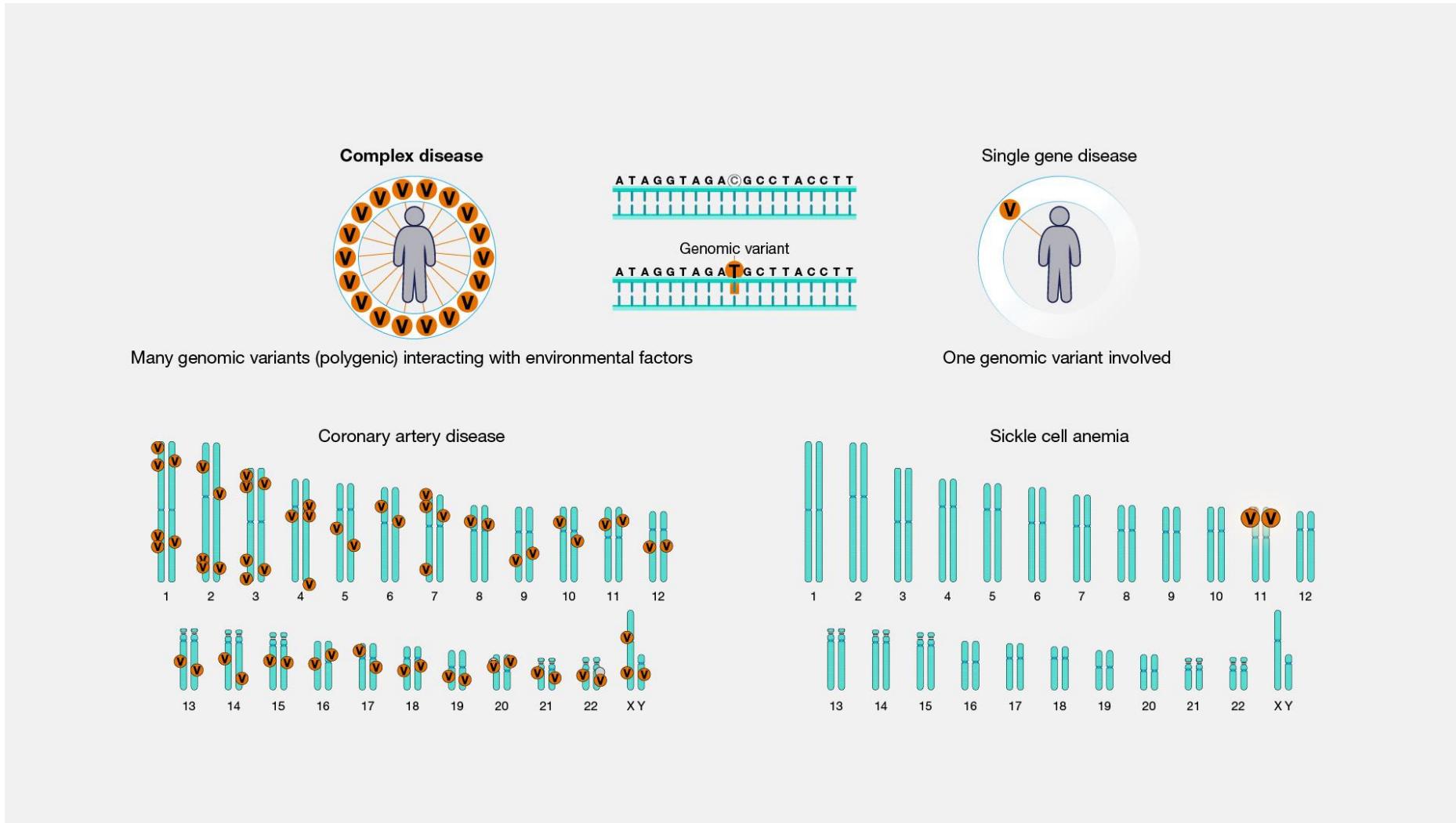


# DIFFERENT NAMES BUT THE SAME

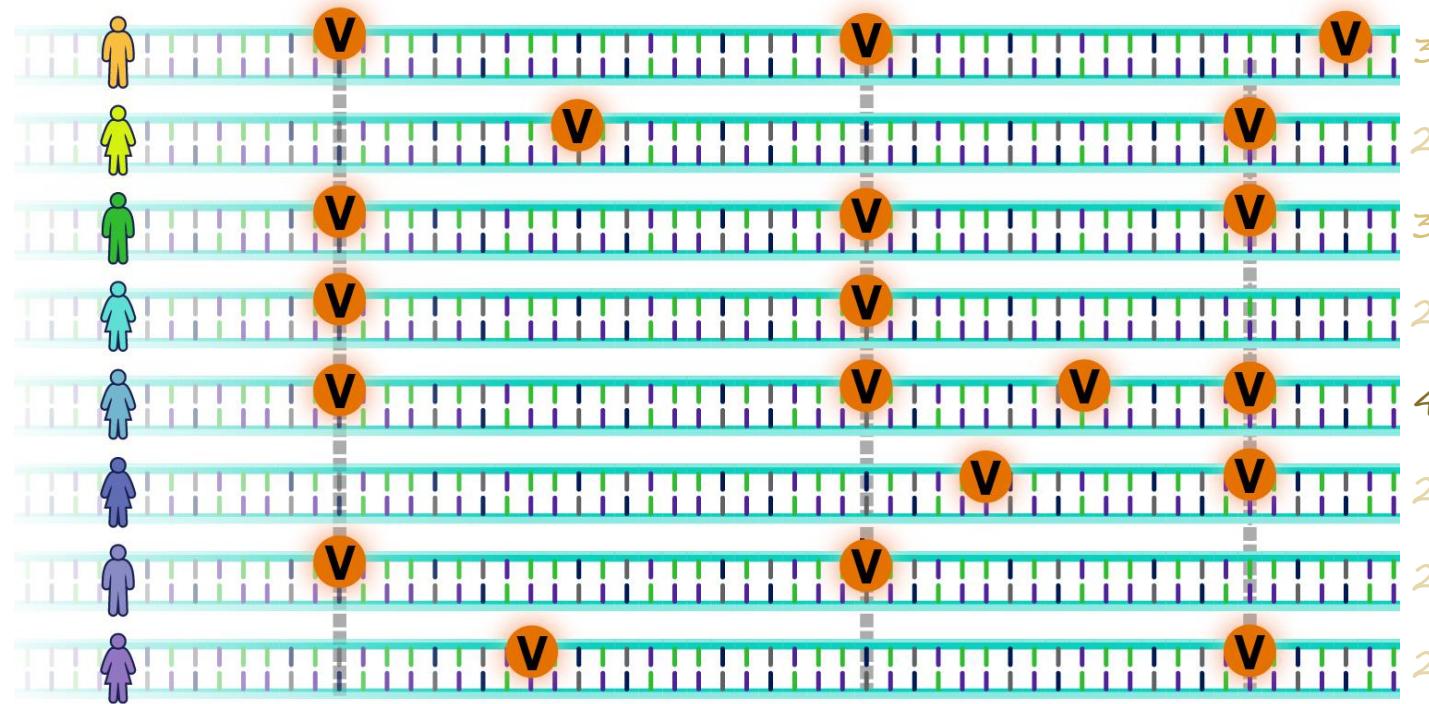
- Polygenic risk score (PRS)
- Polygenic score (PGS)
- Genetic score (GS)
- Genetic risk score (GRS)
- Genetic value
- Genetic liability
- Breeding value
- ...



# THE INHERENT DISEASE RISK



# THE INHERENT DISEASE RISK



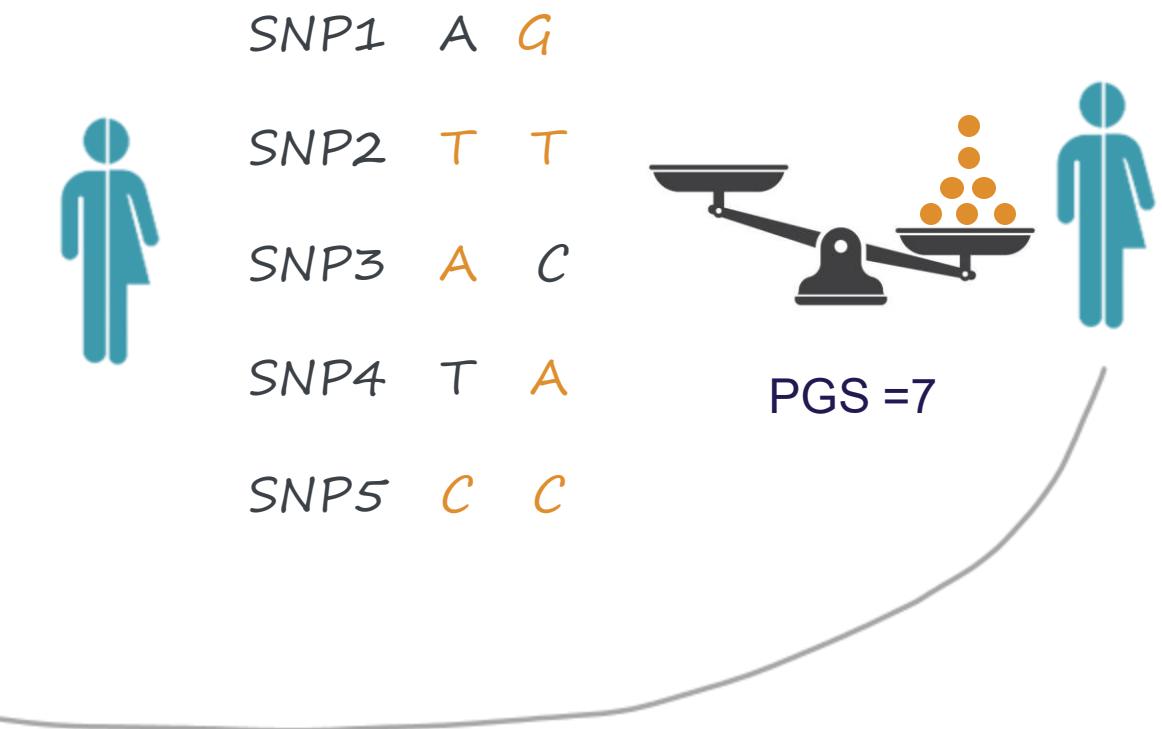
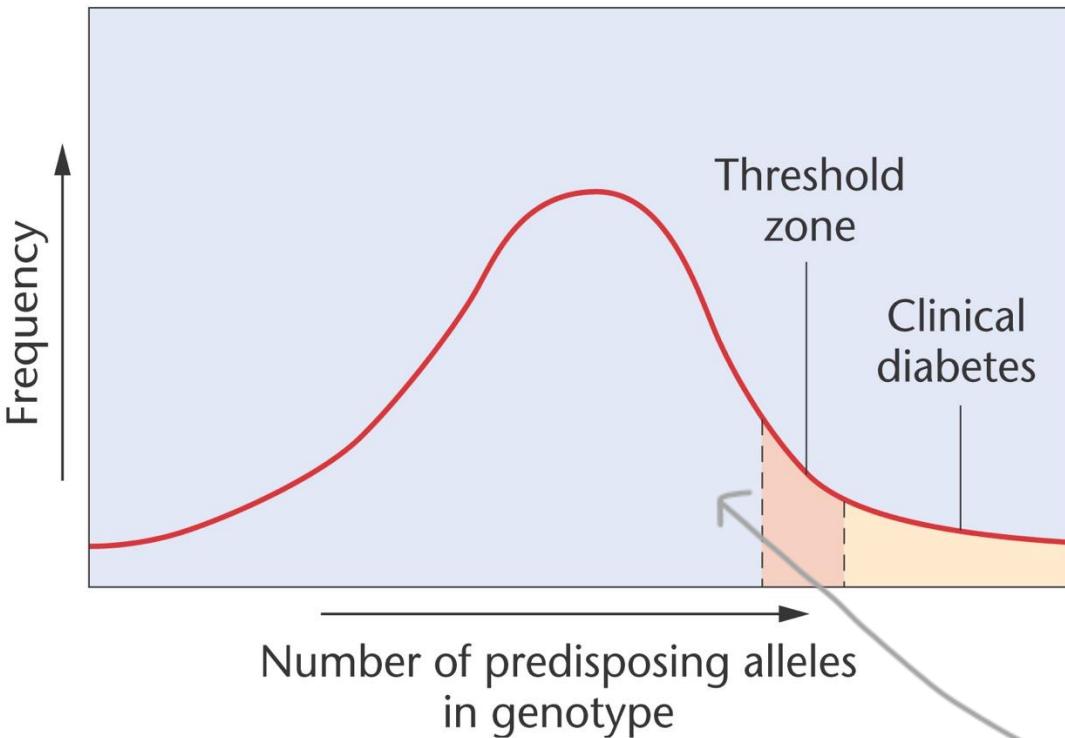
Variation in the  
number of risk  
variants



Polygenic score  
(PGS)

# WHAT IS A PGS?

## A RELATIVE RISK

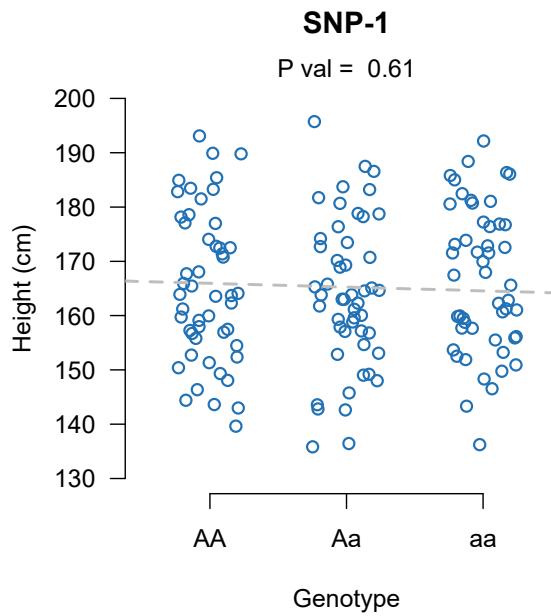


How do we find the 'orange' alleles?



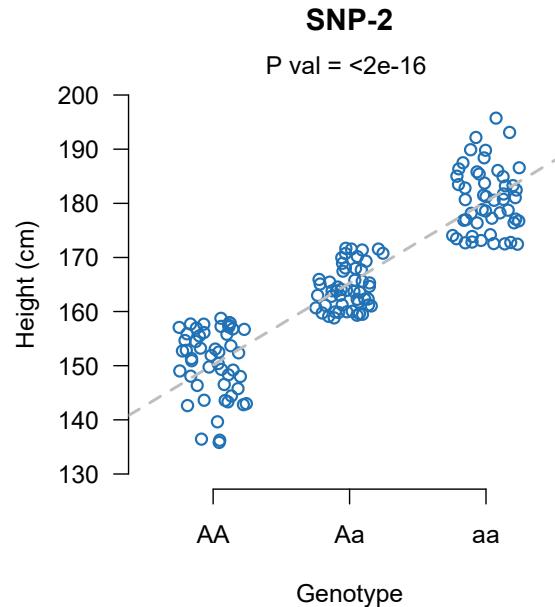
# IDENTIFY RISK VARIANTS

## GENOME-WIDE ASSOCIATION STUDY (GWAS)

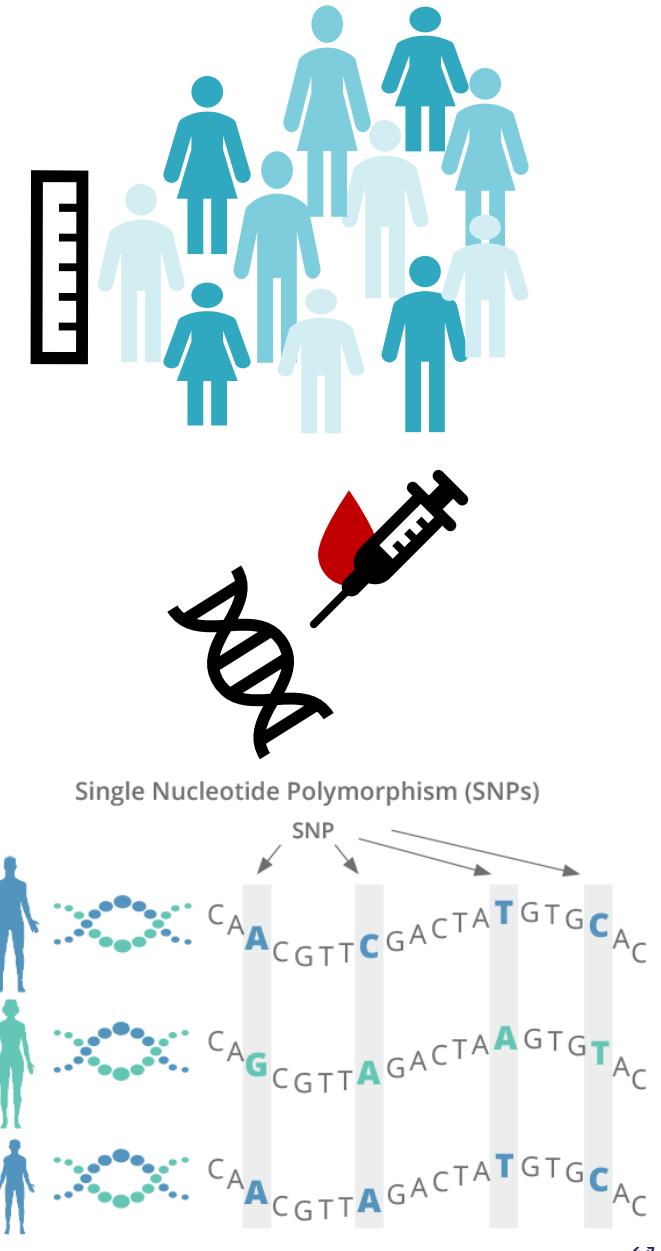


No association

Systematic hypothesis-free scanning of all  
common genetic variants



Association



# WHAT IS A PGS?

*“A PGS combines information from large numbers of markers across the genome (hundreds to millions) to give a single numerical score for an individual’s risk for developing a specific disease on the basis of the DNA variants they have inherited.“*

$b$  is the slope (effect size) from regression

*The effect size of the SNP – obtain from the GWAS*

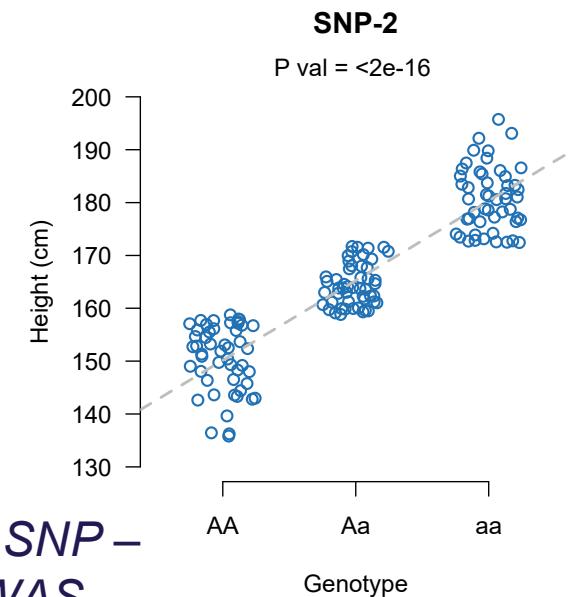
$$PGS = \sum X_i b_i$$

*The genotype of the individual for SNP  $i$  (0, 1, 2 – counting the number of the alternative allele)*

AA = 0

Aa = 1

aa = 2



$$PGS = \sum X_i b_i$$

# HOW TO COMPUTE A (simple) PGS?

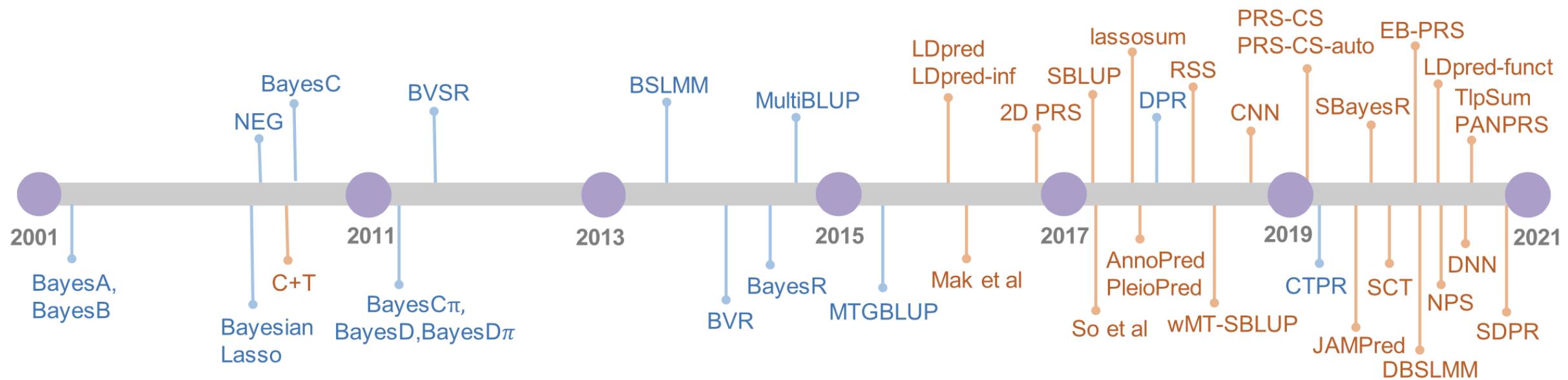
SNPs	Adams Genotypes	Ref allele	Alt allele	X	b	Xb
SNP-1	TC	T	C	1	0.04	0.04
SNP-2	GG	G	T	0	0.02	0.00
SNP-3	CC	A	C	2	0.05	0.10
SNP-4	TG	T	G	1	0.02	0.02
SNP-5	AA	A	G	0	0.06	0.00



$PGS = 0.16$

# A LARGE PALETTE OF PGS METHODS

$$PGS = \sum X_i b_i$$



# WHY DIFFERENT PGS SCORING METHODS?

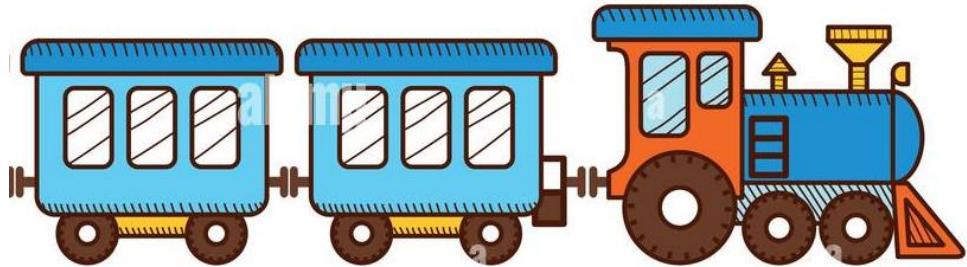
Complex traits have different underlying genetic architectures

- ❖ some are influenced by <100 genetic loci
- ❖ some are influenced by >1000 genetic loci
- ❖ some loci have very small effects
- ❖ some loci have moderate effects
- ❖ correlation structure among loci (linkage disequilibrium)

The different scoring algorithms attempts to account for this.



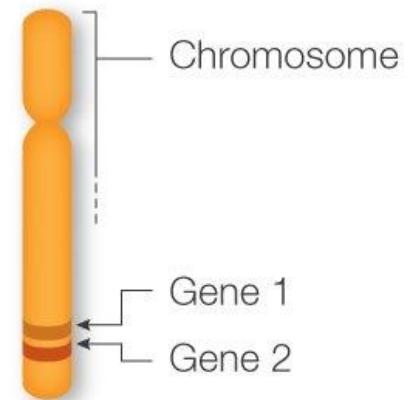
# GENETIC LINKAGE



Linked train wagons



Linked prisoners

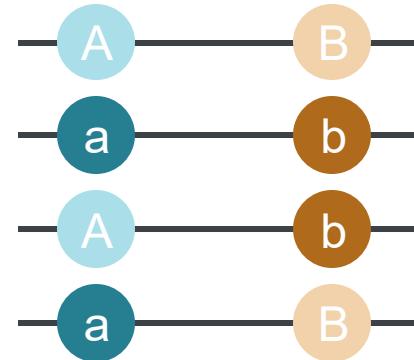


Linked loci  
(Physical proximity)

# WHAT IS LD?



Which gametes can  
be produced?



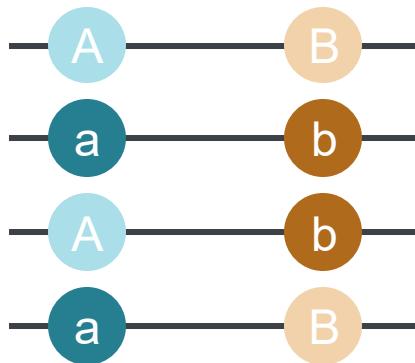
What are the  
frequencies of the  
alleles?

$$P(A), P(a)  
P(B), P(b)$$

What are the  
frequencies of the  
haplotypes?

$$P(AB), P(Ab)  
P(aB), P(ab)$$

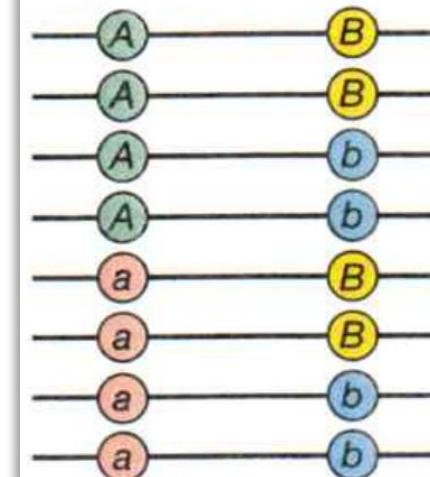
# WHAT IS LD?



If there is random relationship among alleles at the two loci then the frequency of the haplotypes will be the product of the frequencies of the two alleles:

$$P(AB) = P(A) \times P(B)$$

(a) Linkage equilibrium



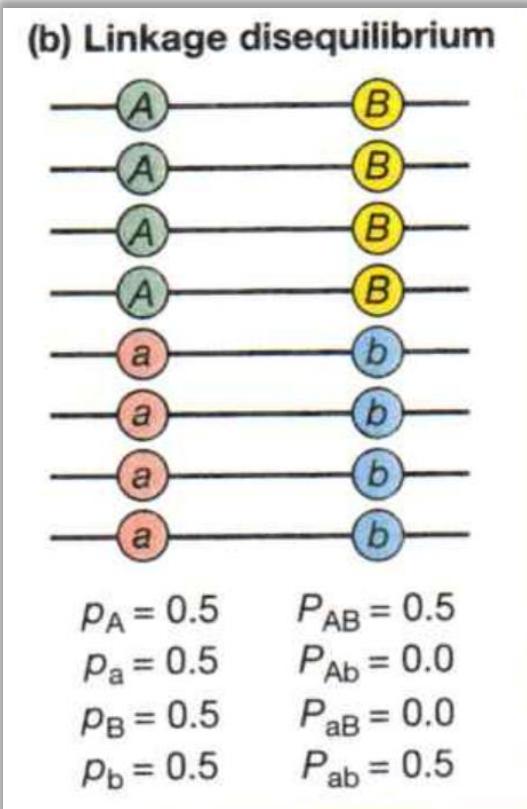
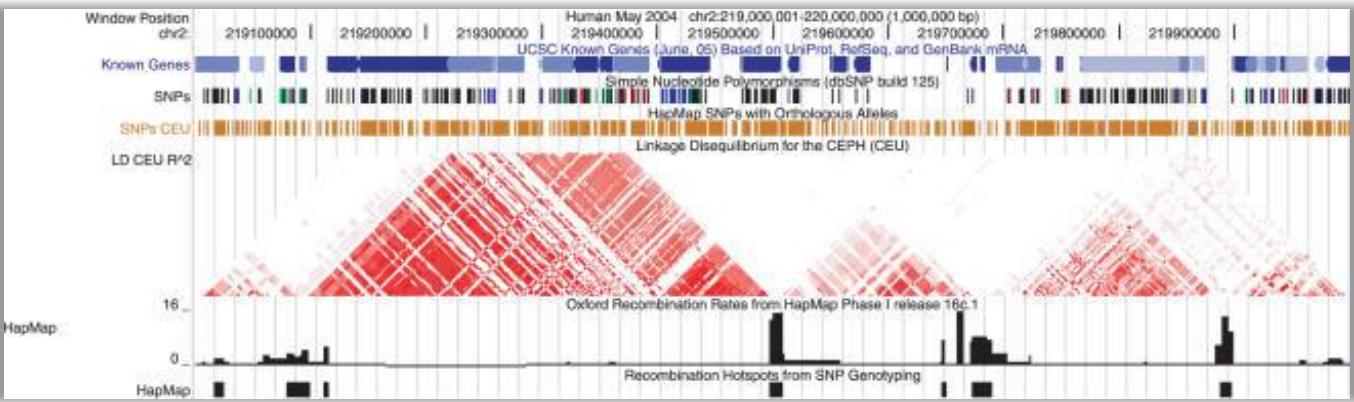
$$p_A = 0.5 \quad P_{AB} = 0.25$$

$$p_a = 0.5 \quad P_{Ab} = 0.25$$

$$p_B = 0.5 \quad P_{aB} = 0.25$$

$$p_b = 0.5 \quad P_{ab} = 0.25$$

# WHAT IS LD?

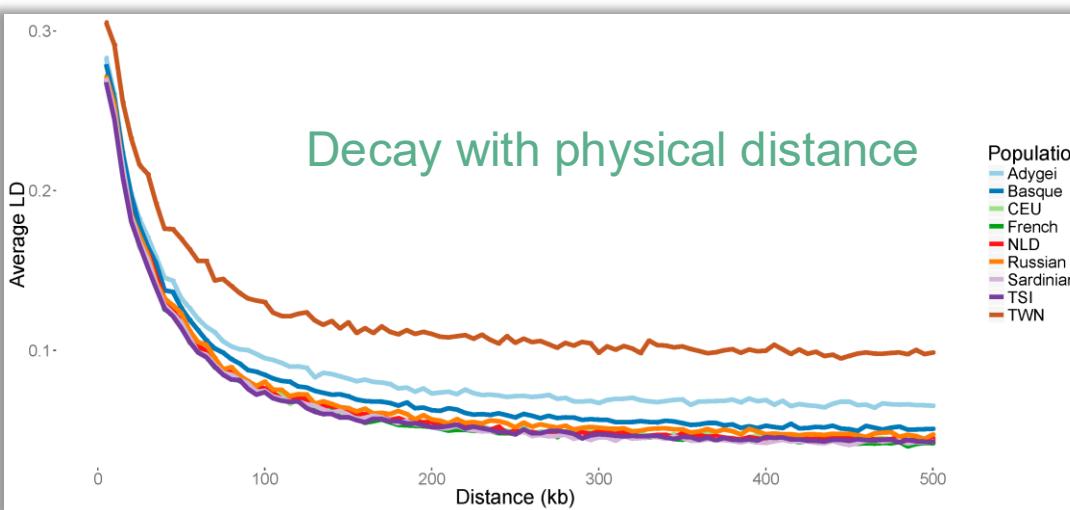


When the association between alleles at two loci is **non-random** they are said to be in **linkage disequilibrium**

The degree of LD can be measured in several ways – the simplest one is:

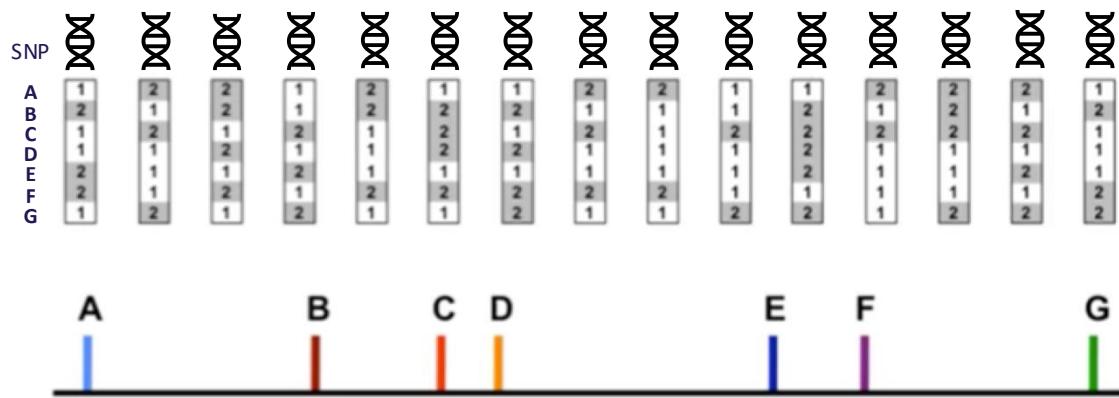
$$D = P_{AB} - P_A P_B$$

If  $D=0$ , no LD, if  $D>0$  LD

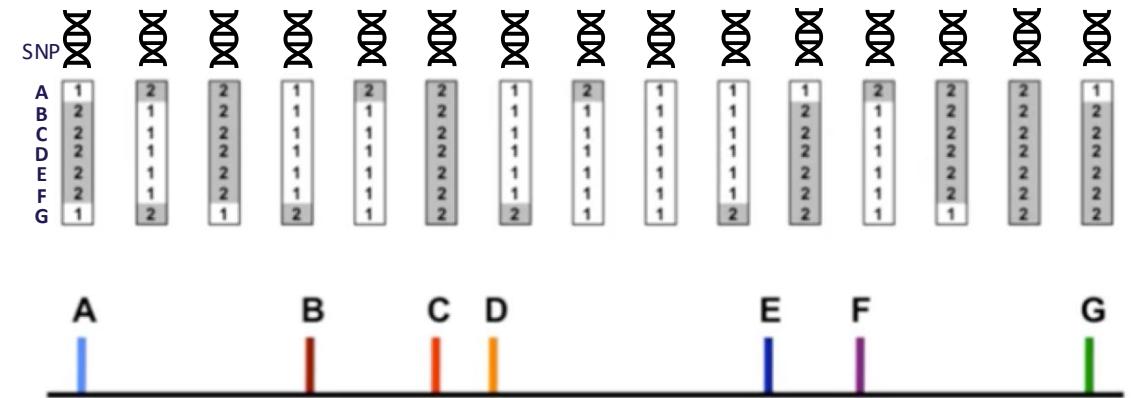


# LD AND GENE MAPPING

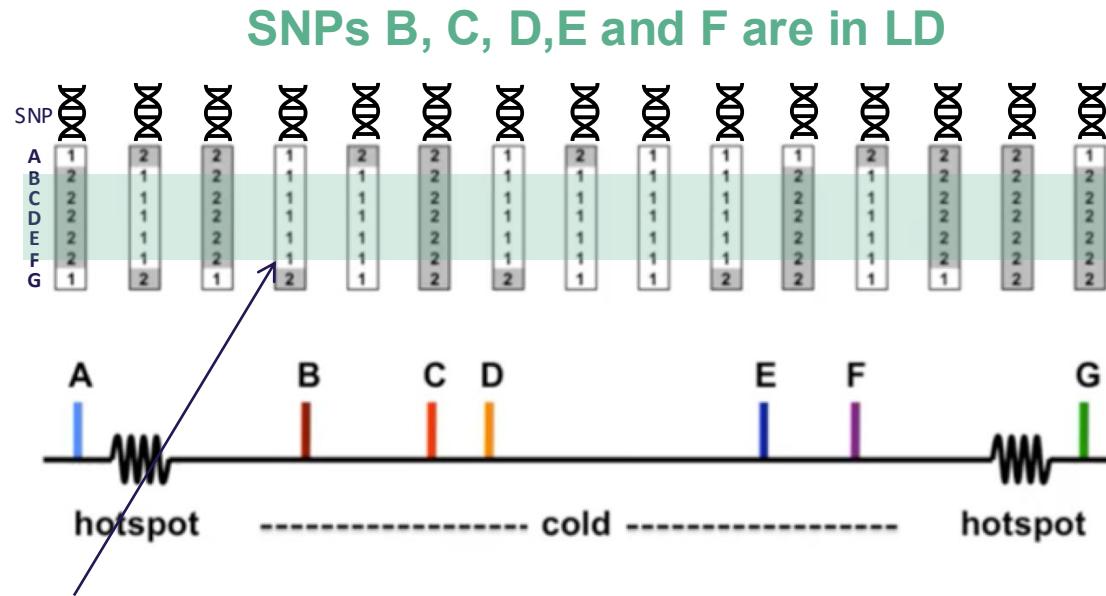
Linkage equilibrium – *random association*



Linkage disequilibrium – *non-random association*

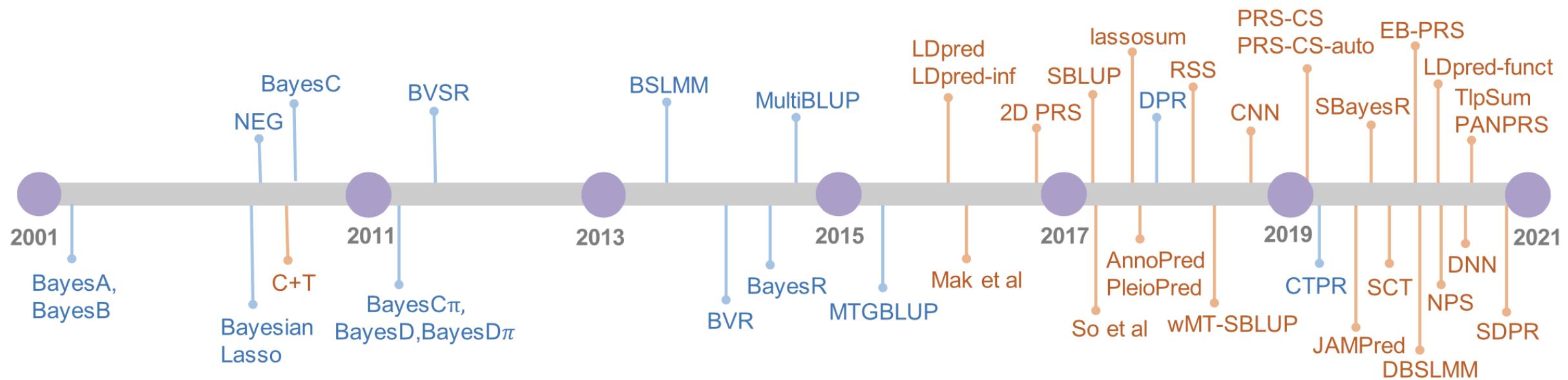


# LD AND GENE MAPPING



If you have allele 1 here, I know  
what you are at the remaining sites  
in this haploblok

# A LARGE PALETTE OF PGS METHODS



# SESSION 1

- Precision Medicine?
- Complex traits?
- Genetic variants as chilies
- What is a polygenic score?
- What is needed to compute a polygenic score?
- Why so many different methods?





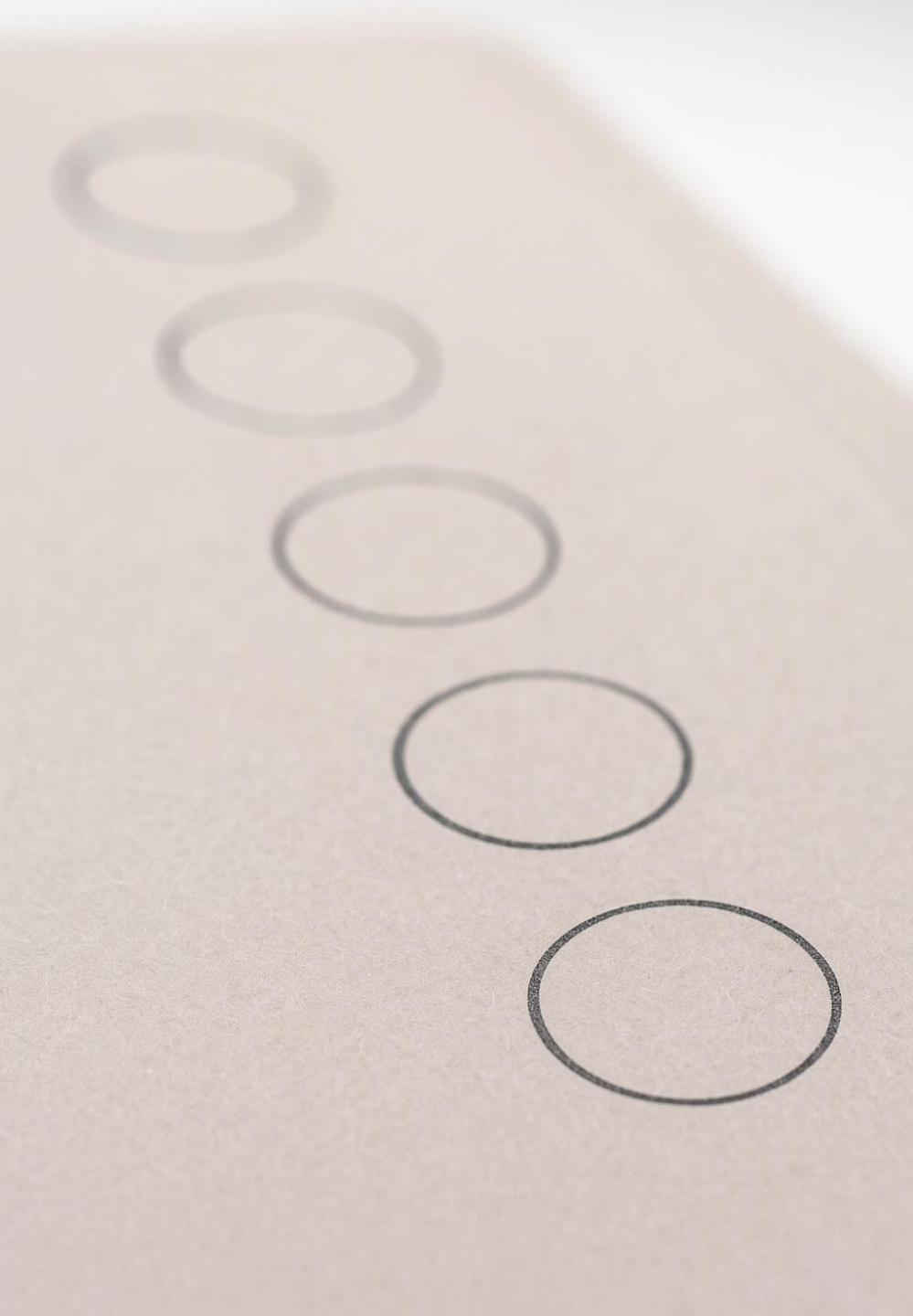
**BREAK**

# AGENDA

08:00 – 08:30	Welcome and common introductions
08:30 – 09:10	Session 1: Introduction to Polygenic Scores (PGS)
09:10 – 09:20	Break
<b>09:20 – 10:00</b>	<b>Session 2: Data Sources and Computational Methods</b>
10:00 – 10:10	Break
10:10 – 10:40	Session 3: Evaluating and Interpreting Polygenic Scores
10:40 – 11:00	Break
11:00 – 11:45	Session 4: Advanced Applications and Future Directions
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today

# SESSION 2

- The first polygenic score
- What you need is...
- Commonly used scoring algorithms
- Workflow
- Current challenges with PGS?

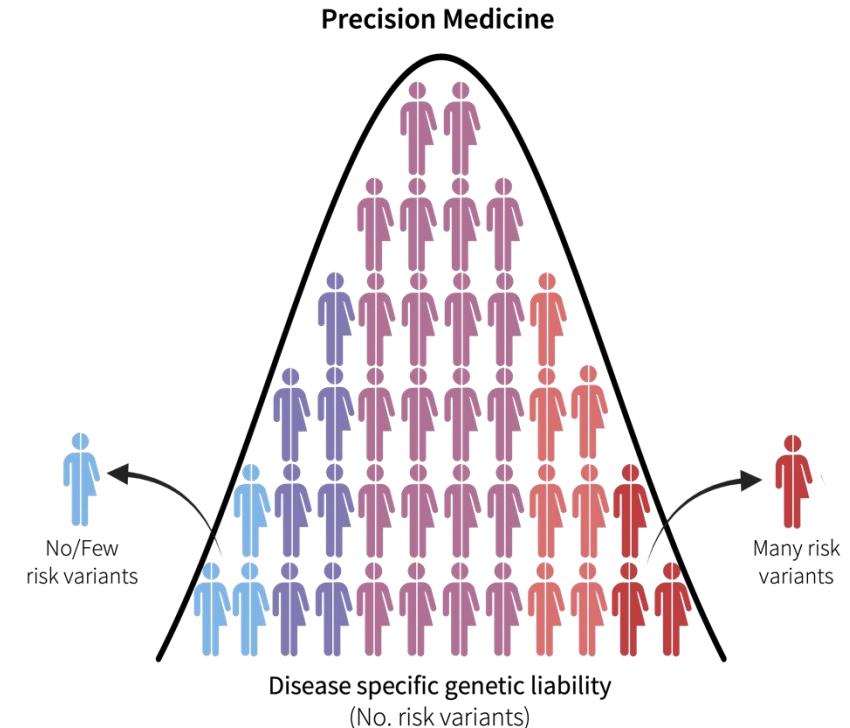


# THE FIRST PGS

Idea originating back to animal and plant breeding - finding the 'best' animals that should establish the next generation.

In 2007, Wray proposed to 'predict' human disease traits from SNP data.

Wray, Goddard & Visscher (2007) Genome Research, 17:1520–1528

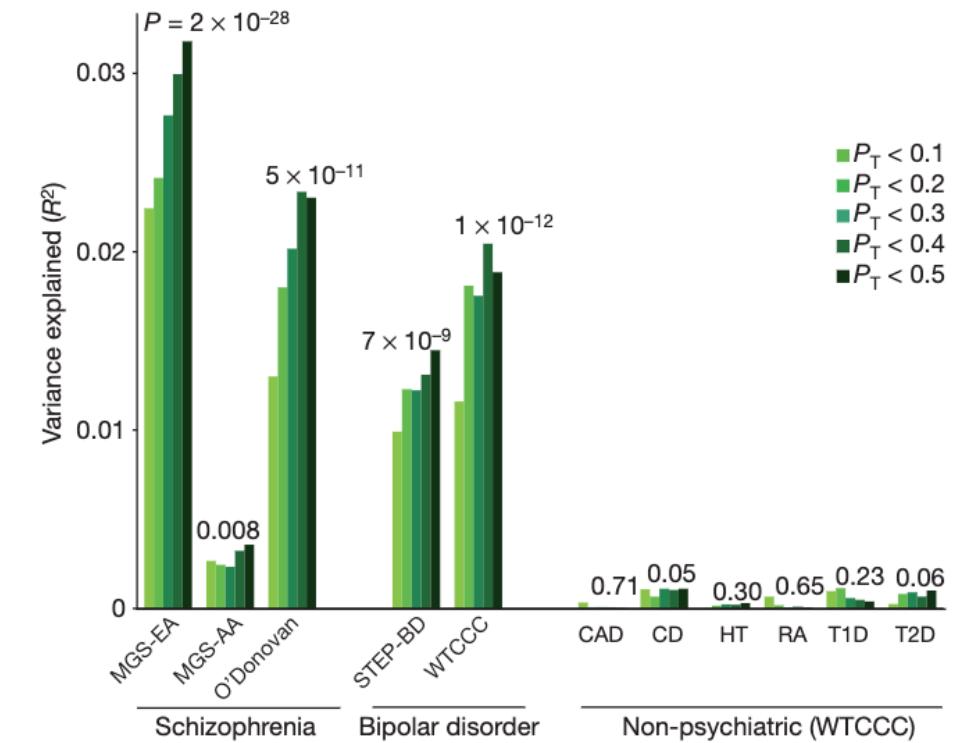


# THE FIRST PGS

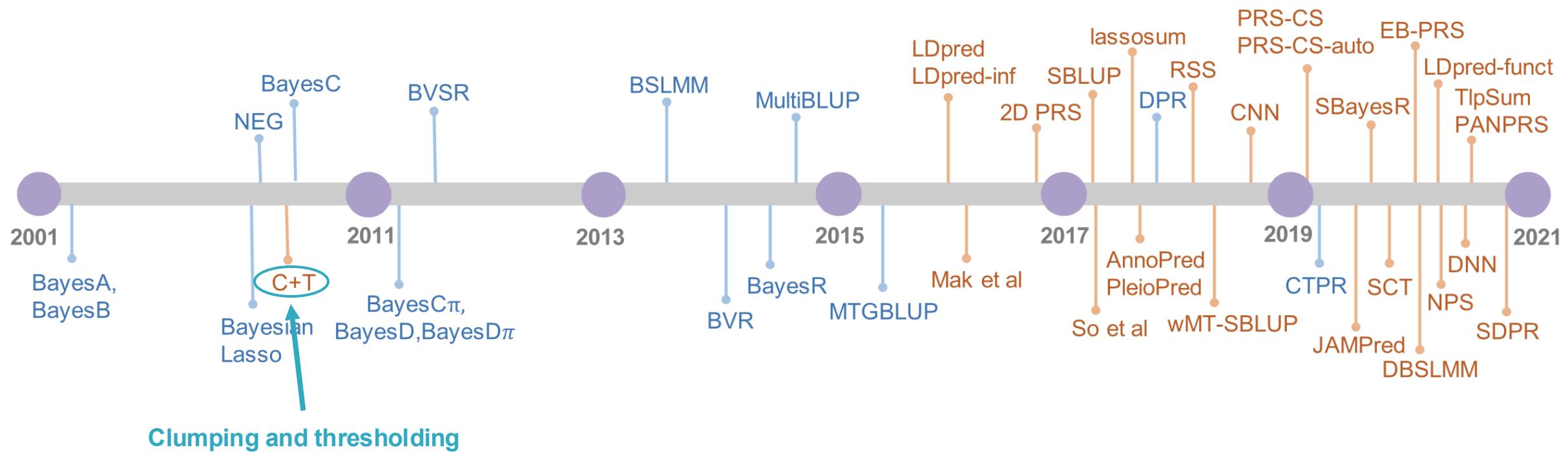
In 2009, Purcell constructed the first polygenic score for a human disease:

*'We summarized variation across nominally associated loci into quantitative scores'*

The International Schizophrenia Consortium (2009) Nature, 460



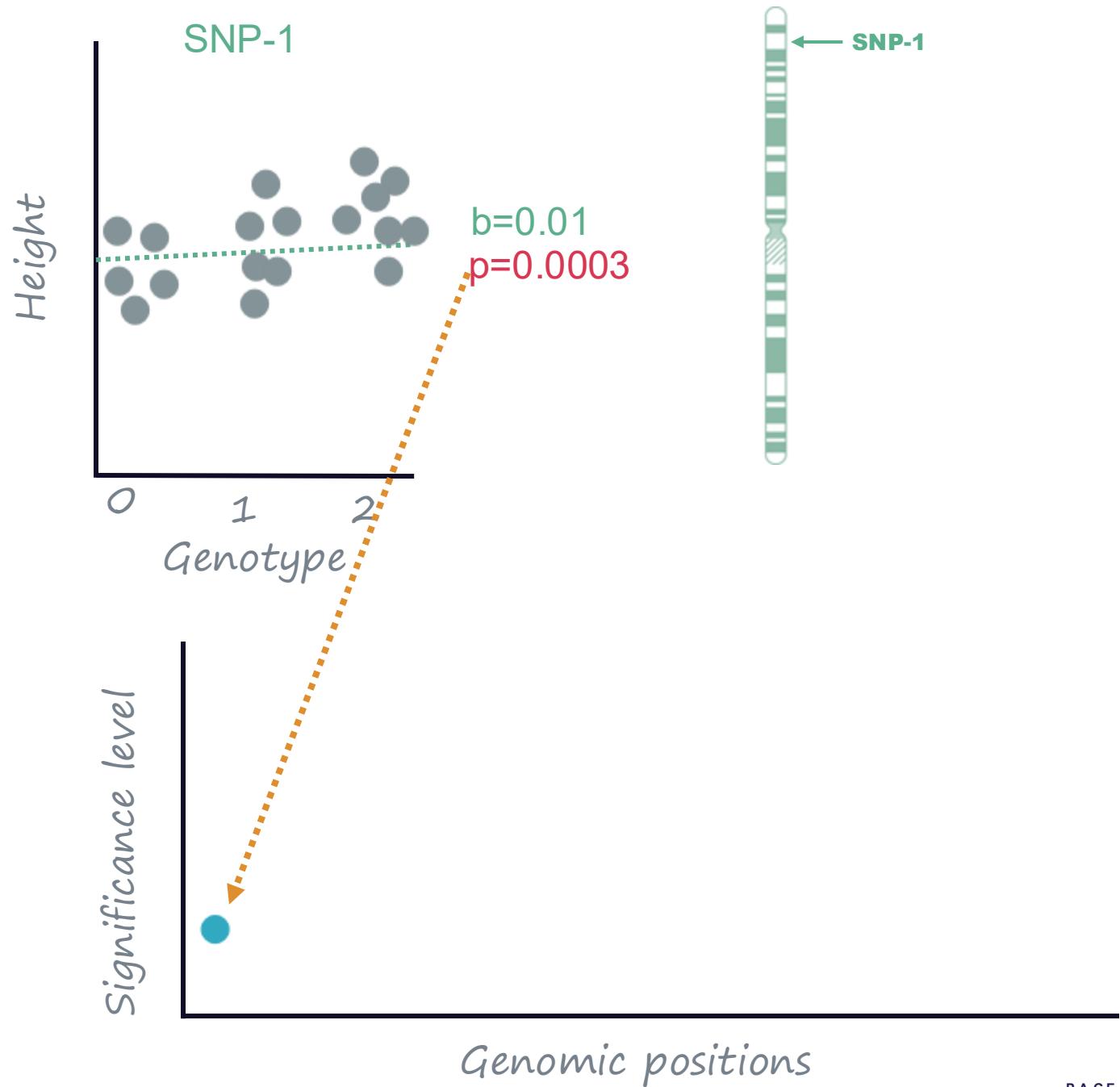
# A LARGE PALETTE OF PGS METHODS



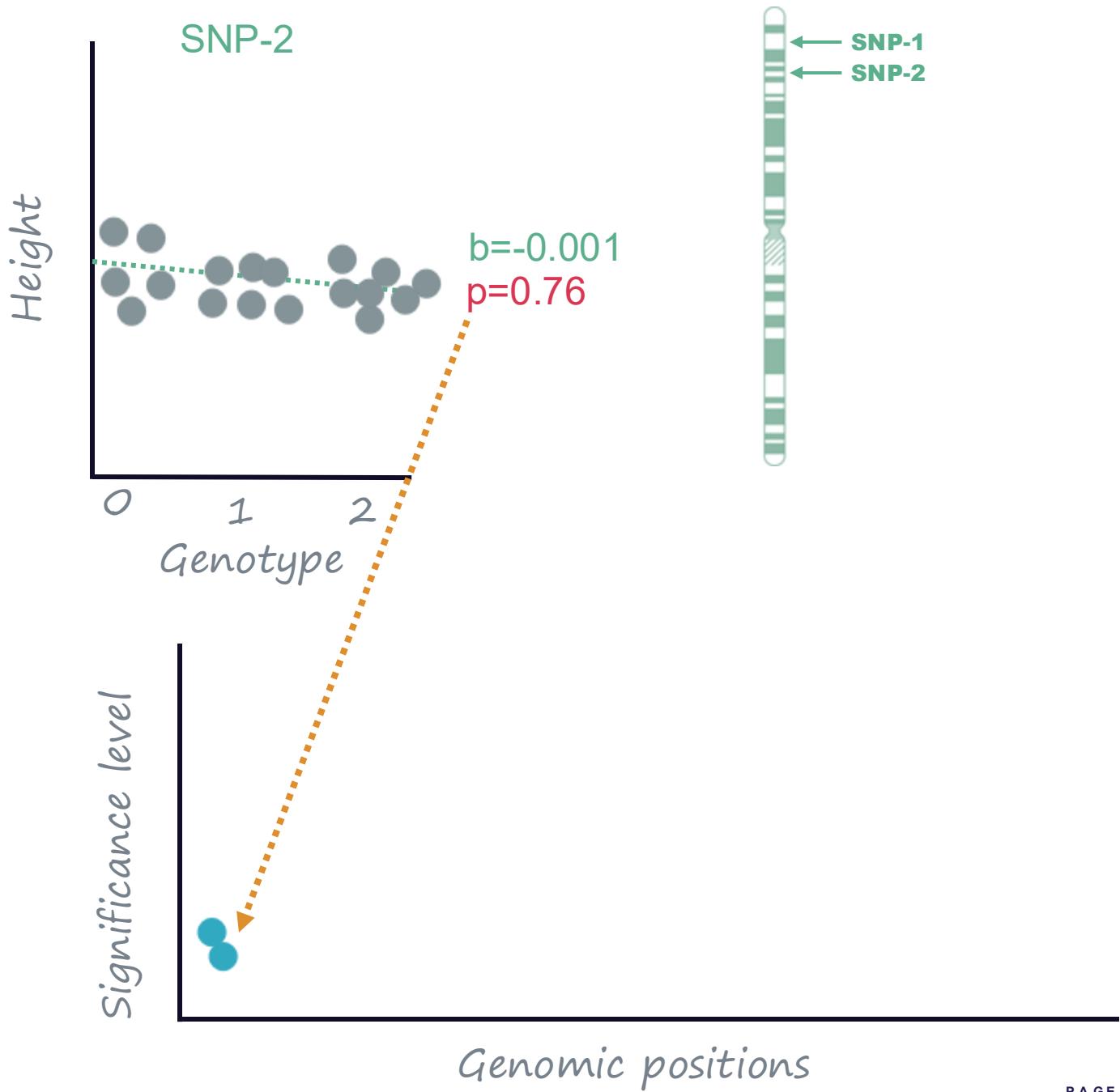
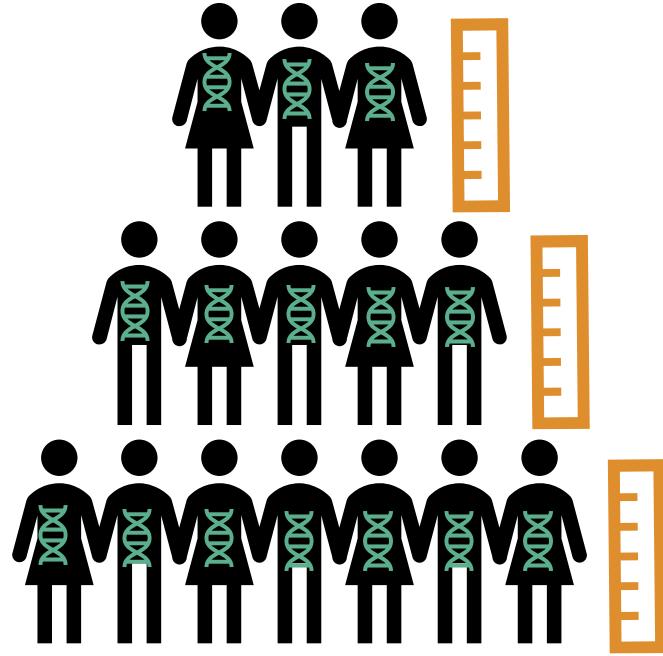
# GWAS RECAP



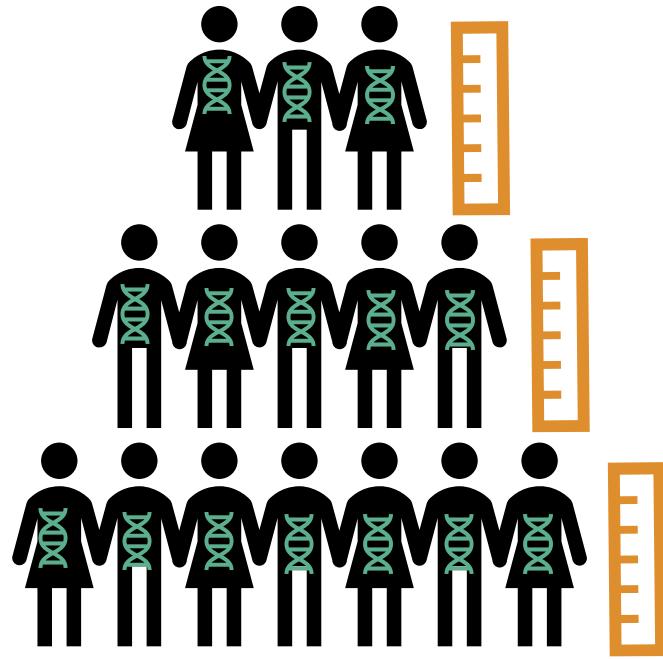
Which SNPs associate with height?



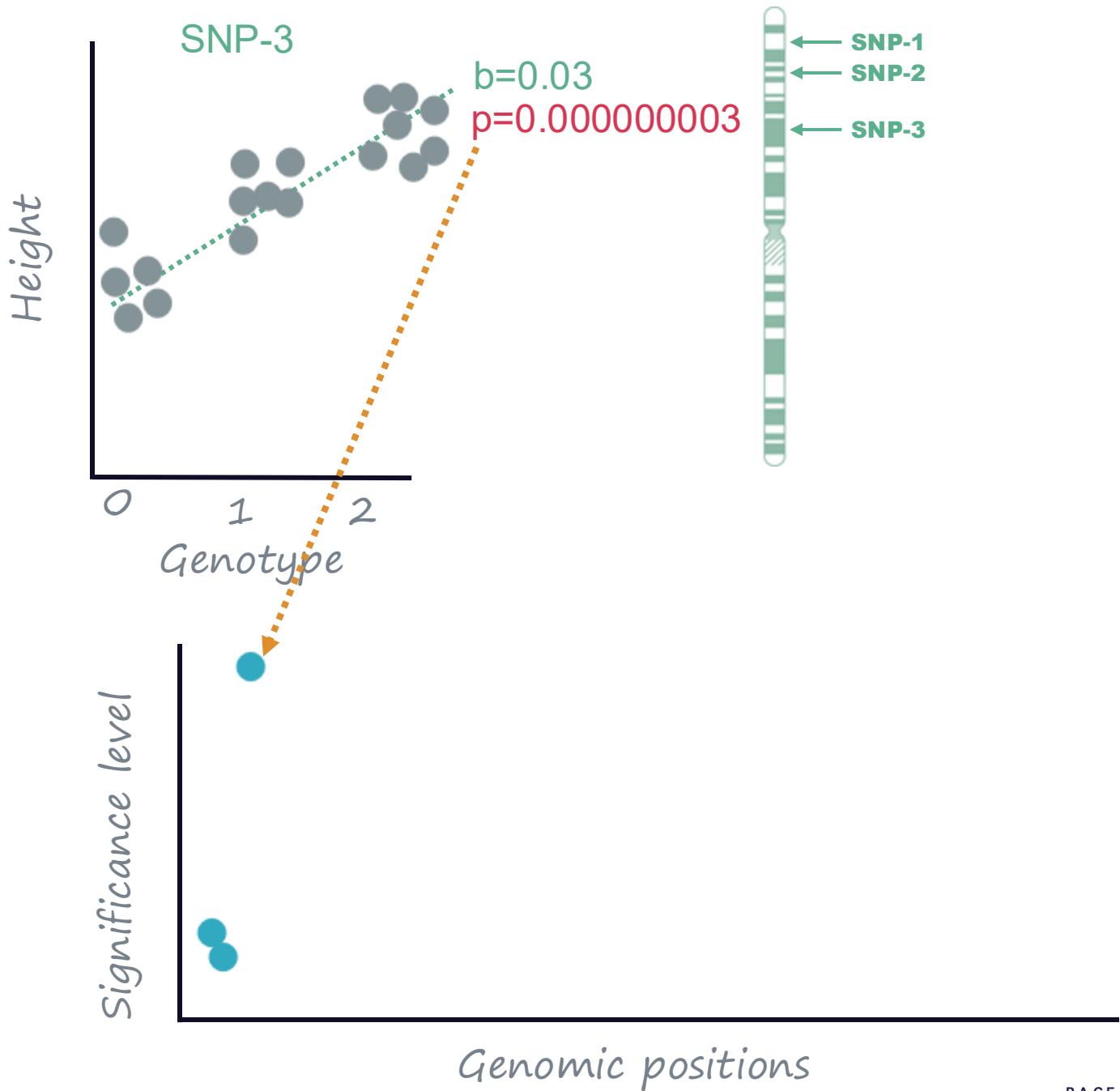
# GWAS RECAP



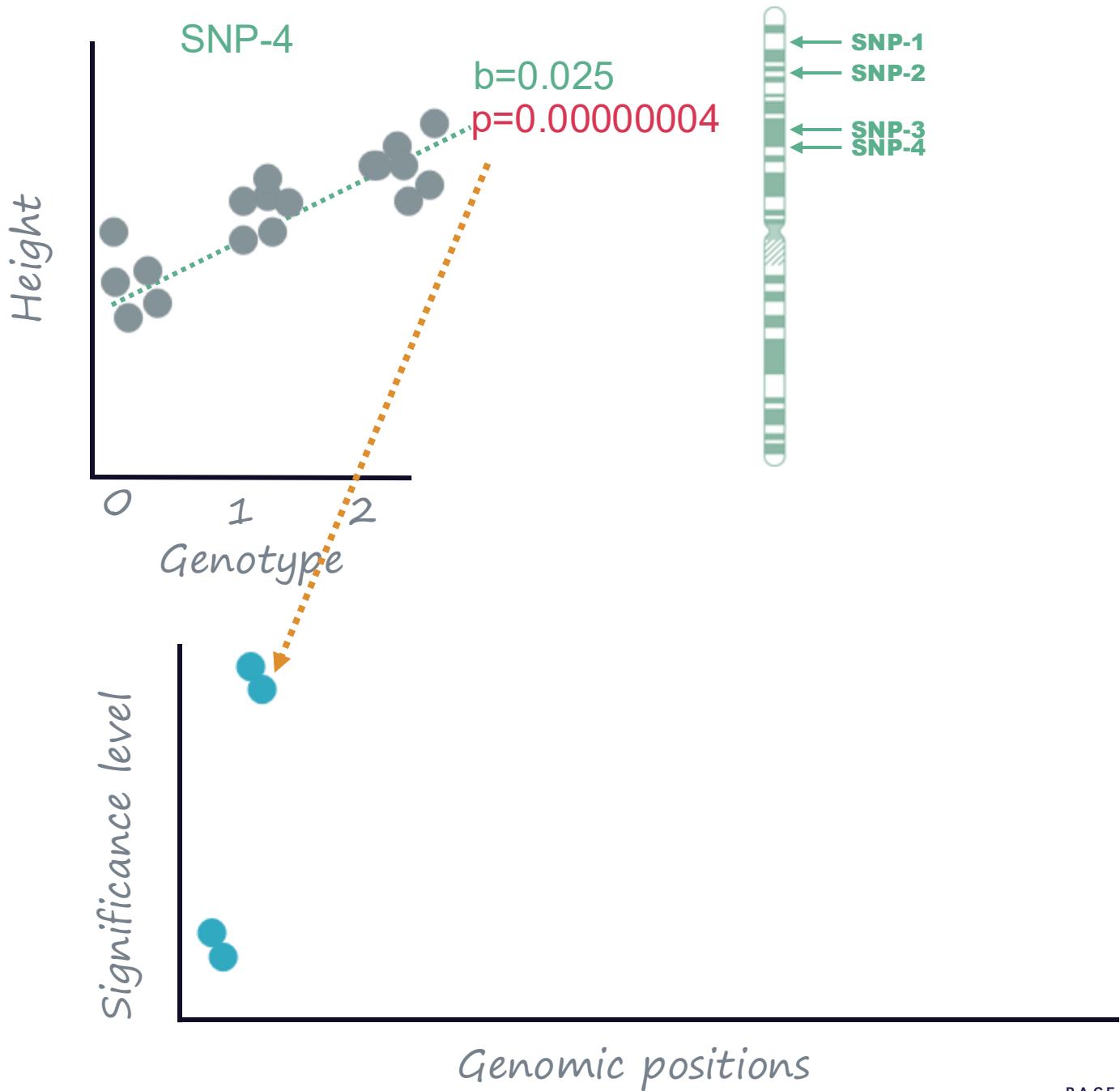
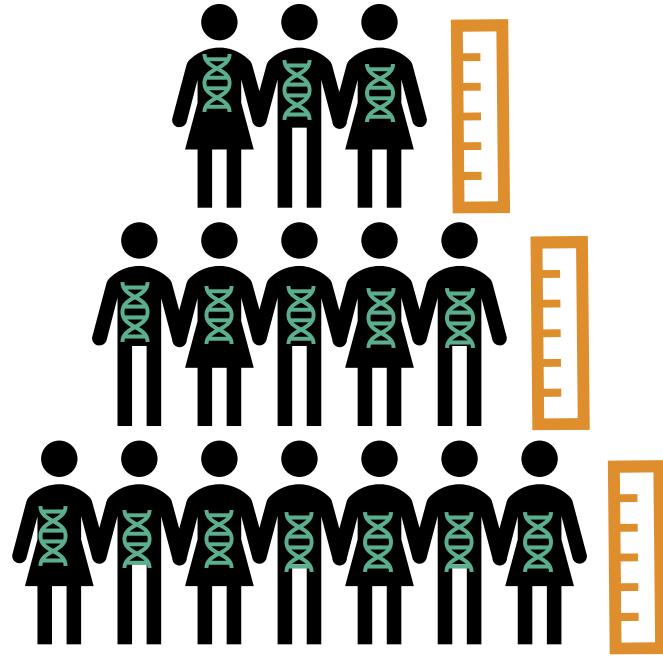
# GWAS RECAP



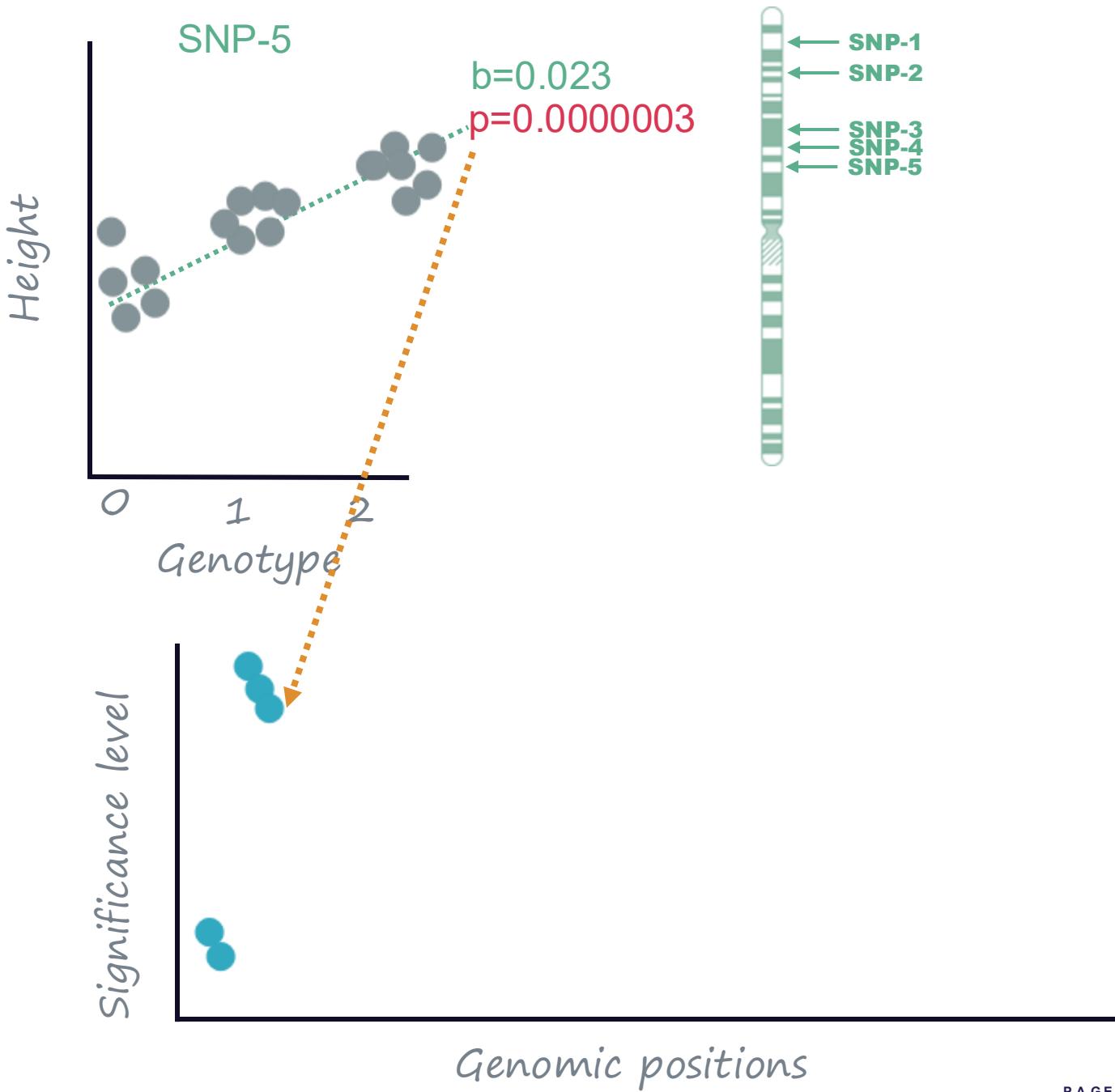
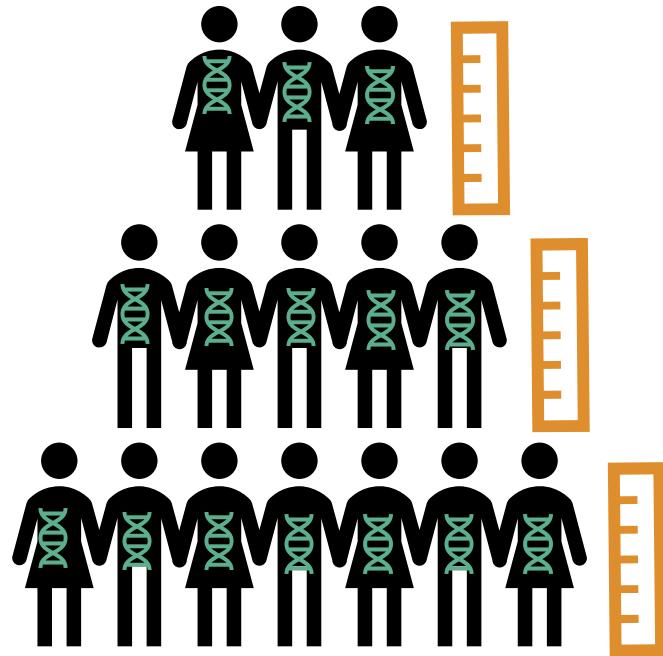
Which SNPs associate with height?



# GWAS RECAP



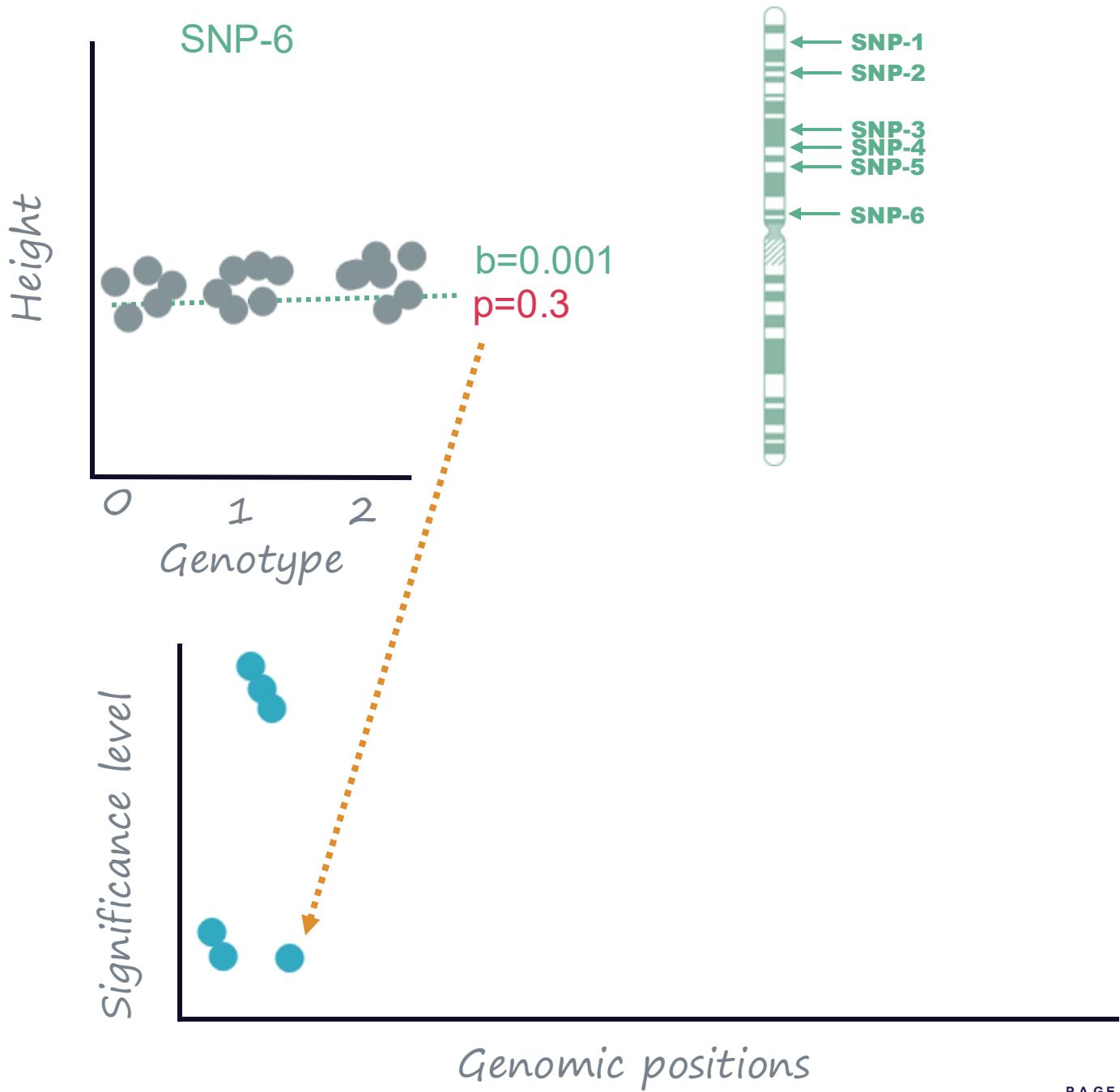
# GWAS RECAP



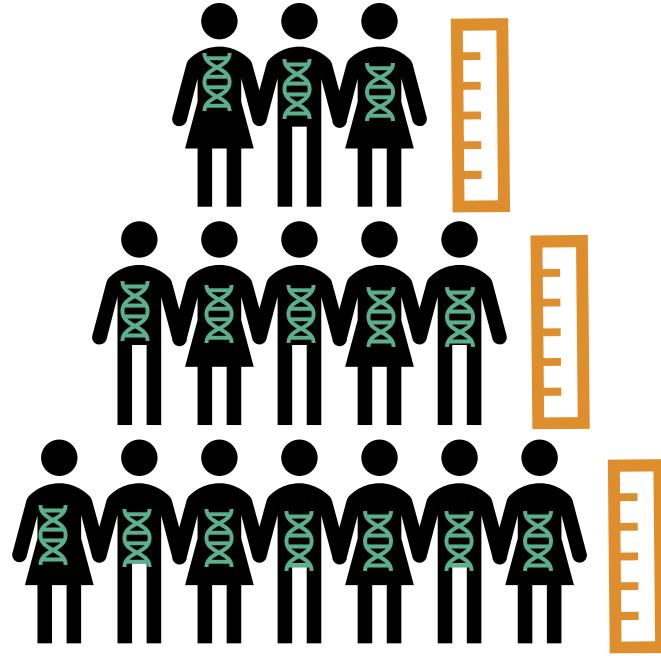
# GWAS RECAP



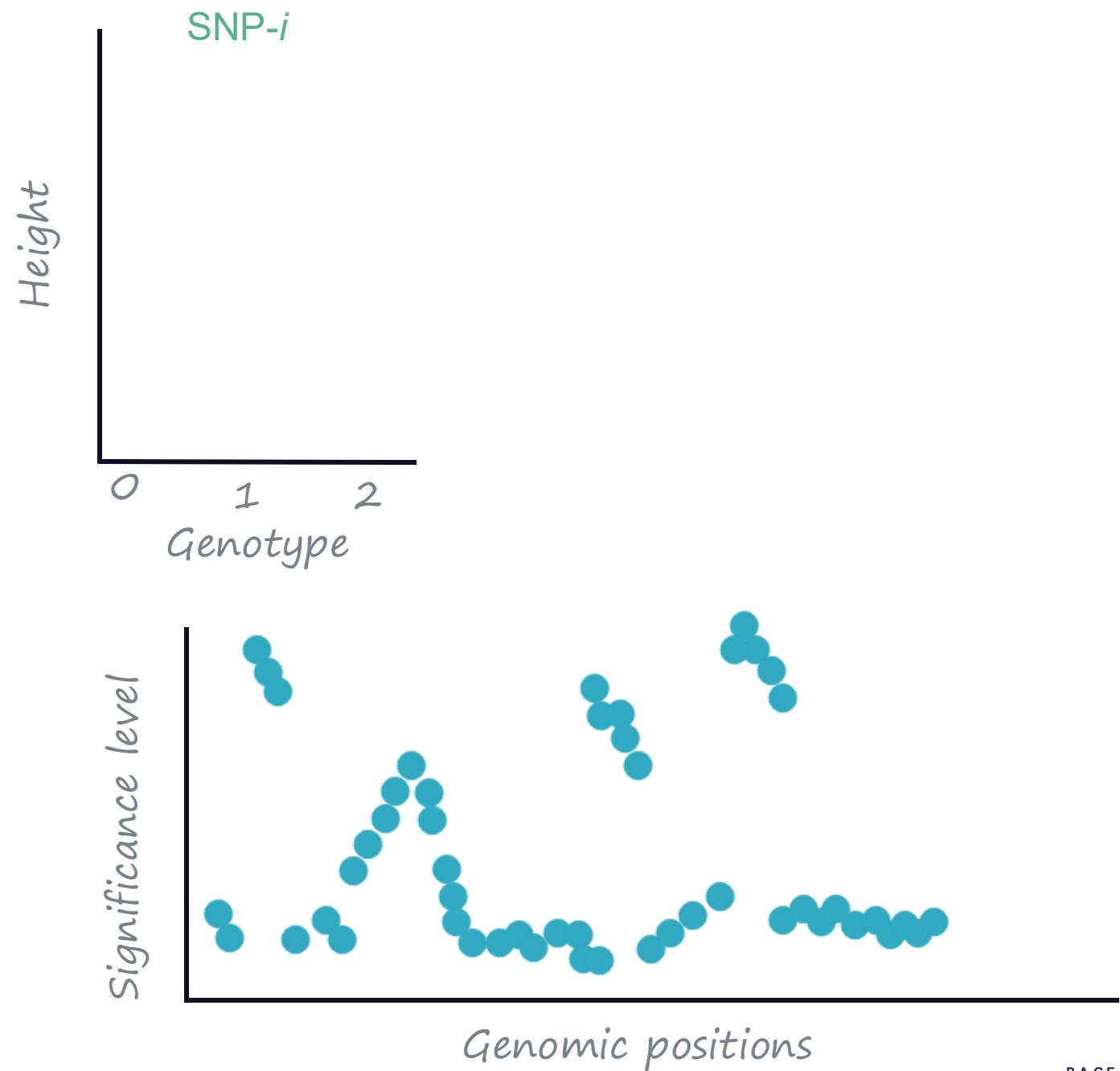
**Which SNPs associate with height?**



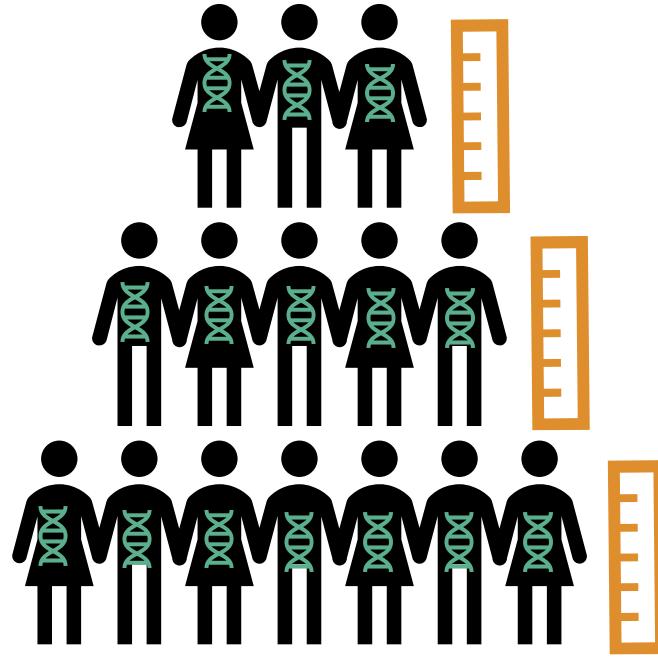
# GWAS RECAP



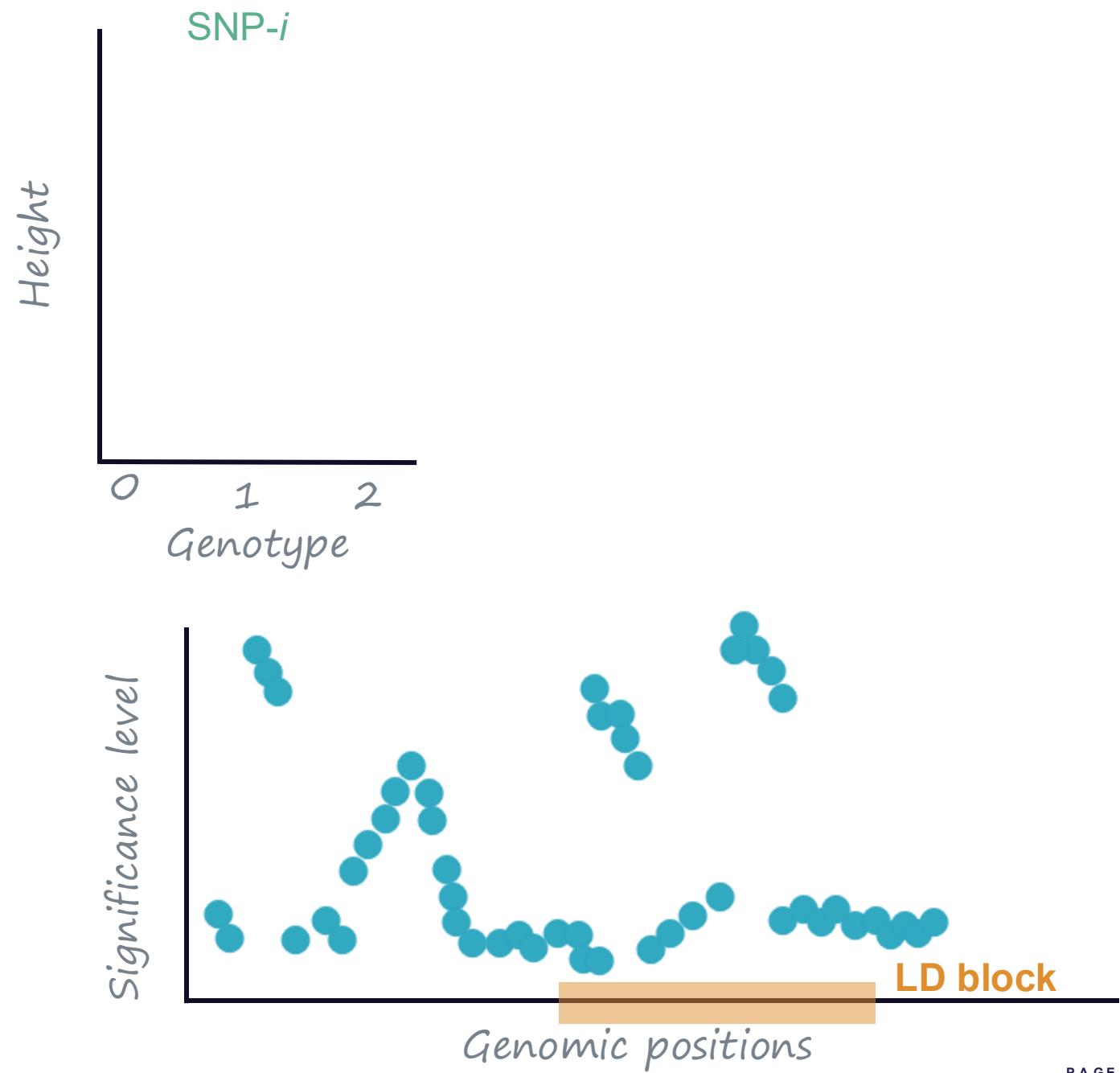
**Which SNPs associate with height?**



# GWAS RECAP



**Which SNPs associate with height?**



# CLUMPING AND THRESHOLDING (C+T)

0: Set LD (=0.8) and  $P$  values (0.01)

SNP	b	p
1	0.21	0.005
2	0.22	0.0048
3	0.25	0.0003
4	0.1	0.04
5	0.05	0.15
6	0.02	0.49
7	0.03	0.87
8	0.12	0.003
9	0.14	0.0034
10	0.18	0.0004
11	0.21	0.00003
12	0.12	0.15
13	0.14	0.12
14	0.03	0.84
15	0.02	0.32

1: Sort by  $P$ -value

SNP	b	p
11	0.21	0.00003
3	0.25	0.0003
10	0.18	0.0004
8	0.12	0.003
9	0.14	0.0034
2	0.22	0.0048
1	0.21	0.005
4	0.1	0.04
13	0.14	0.12
5	0.05	0.15
12	0.12	0.15
15	0.02	0.32
6	0.02	0.49
14	0.03	0.84
7	0.03	0.87

2: Compute LD and select variants based of thresholds

SNP	b	p	r <sup>2</sup>
11	0.21	0.00003	1st variant in LD-pair
3	0.25	0.0003	0.96
10	0.18	0.0004	0.93
8	0.12	0.003	0.88
9	0.14	0.0034	0.74
2	0.22	0.0048	0.4
1	0.21	0.005	0.03
4	0.1	0.04	0.04
13	0.14	0.12	0.05
5	0.05	0.15	0.03
12	0.12	0.15	0.04
15	0.02	0.32	0.01
6	0.02	0.49	0.01
14	0.03	0.84	0.01
7	0.03	0.87	0.01

Have LD>r<sup>2</sup> – ignore those

# CLUMPING AND THRESHOLDING (C+T)

0: Set LD (=0.8) and  $P$  values (0.01)

SNP	b	p
1	0.21	0.005
2	0.22	0.0048
3	0.25	0.0003
4	0.1	0.04
5	0.05	0.15
6	0.02	0.49
7	0.03	0.87
8	0.12	0.003
9	0.14	0.0034
10	0.18	0.0004
11	0.21	0.00003
12	0.12	0.15
13	0.14	0.12
14	0.03	0.84
15	0.02	0.32

1: Sort by P-value

SNP	b	p
11	0.21	0.00003
3	0.25	0.0003
10	0.18	0.0004
8	0.12	0.003
9	0.14	0.0034
2	0.22	0.0048
1	0.21	0.005
4	0.1	0.04
13	0.14	0.12
5	0.05	0.15
12	0.12	0.15
15	0.02	0.32
6	0.02	0.49
14	0.03	0.84
7	0.03	0.87

2: Compute LD and select variants based of thresholds

SNP	b	p	r <sup>2</sup>
11	0.21	0.00003	
3	0.25	0.0003	
10	0.18	0.0004	
8	0.12	0.003	
9	0.14	0.0034	
2	0.22	0.0048	0.98
1	0.21	0.005	0.96
4	0.1	0.04	0.96
13	0.14	0.12	0.52
5	0.05	0.15	0.34
12	0.12	0.15	0.10
15	0.02	0.32	0.04
6	0.02	0.49	0.01
14	0.03	0.84	0.01
7	0.03	0.87	0.01

1st variant in LD-pair

Have LD>r<sup>2</sup> – ignore those

# CLUMPING AND THRESHOLDING (C+T)

0: Set LD (=0.8) and  $P$  values (0.01)

SNP	b	p
1	0.21	0.005
2	0.22	0.0048
3	0.25	0.0003
4	0.1	0.04
5	0.05	0.15
6	0.02	0.49
7	0.03	0.87
8	0.12	0.003
9	0.14	0.0034
10	0.18	0.0004
11	0.21	0.00003
12	0.12	0.15
13	0.14	0.12
14	0.03	0.84
15	0.02	0.32

1: Sort by P-value

SNP	b	p
11	0.21	0.00003
3	0.25	0.0003
10	0.18	0.0004
8	0.12	0.003
9	0.14	0.0034
2	0.22	0.0048
1	0.21	0.005
4	0.1	0.04
13	0.14	0.12
5	0.05	0.15
12	0.12	0.15
15	0.02	0.32
6	0.02	0.49
14	0.03	0.84
7	0.03	0.87

2: Compute LD and select variants based of thresholds

SNP	b	p	r <sup>2</sup>
11	0.21	0.00003	
3	0.25	0.0003	
10	0.18	0.0004	
8	0.12	0.003	
9	0.14	0.0034	
2	0.22	0.0048	
1	0.21	0.005	
4	0.1	0.04	
13	0.14	0.12	1st variant in LD-pair
5	0.05	0.15	0.86
12	0.12	0.15	0.82
15	0.02	0.32	0.81
6	0.02	0.49	0.85
14	0.03	0.84	0.85
7	0.03	0.87	0.81

Have LD>r<sup>2</sup> – ignore those

# CLUMPING AND THRESHOLDING (C+T)

0: Set LD (=0.8) and  $P$  values (0.01)

SNP	b	p
1	0.21	0.005
2	0.22	0.0048
3	0.25	0.0003
4	0.1	0.04
5	0.05	0.15
6	0.02	0.49
7	0.03	0.87
8	0.12	0.003
9	0.14	0.0034
10	0.18	0.0004
11	0.21	0.00003
12	0.12	0.15
13	0.14	0.12
14	0.03	0.84
15	0.02	0.32

1: Sort by P-value

SNP	b	p
11	0.21	0.00003
3	0.25	0.0003
10	0.18	0.0004
8	0.12	0.003
9	0.14	0.0034
2	0.22	0.0048
1	0.21	0.005
4	0.1	0.04
13	0.14	0.12
5	0.05	0.15
12	0.12	0.15
15	0.02	0.32
6	0.02	0.49
14	0.03	0.84
7	0.03	0.87

2: Compute LD and select variants based on LD

SNP	b	p	r <sup>2</sup>
11	0.21	0.00003	←
3	0.25	0.0003	
10	0.18	0.0004	
8	0.12	0.003	
9	0.14	0.0034	←
2	0.22	0.0048	
1	0.21	0.005	
4	0.1	0.04	
13	0.14	0.12	
5	0.05	0.15	
12	0.12	0.15	
15	0.02	0.32	
6	0.02	0.49	
14	0.03	0.84	
7	0.03	0.87	

3: Compute PGS based on effect sizes (b) and  $P$ -values

$$PGS = \sum X_i b_i$$

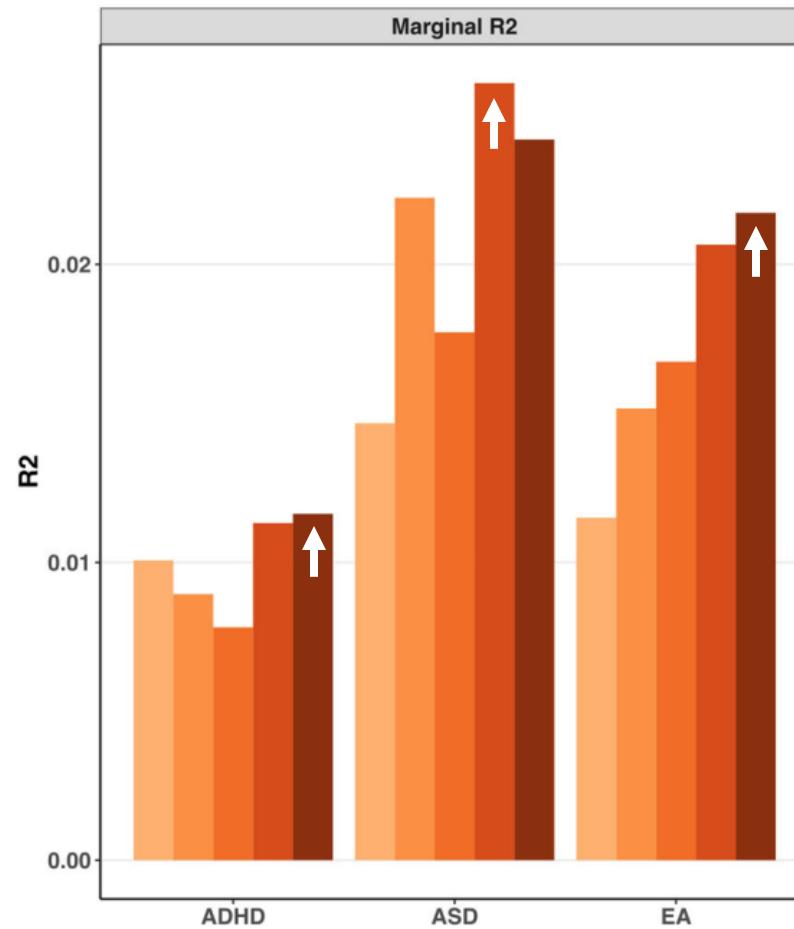
$$= X_{11} \times 0.21 + X_9 \times 0.14$$

# CLUMPING AND THRESHOLDING (C+T)

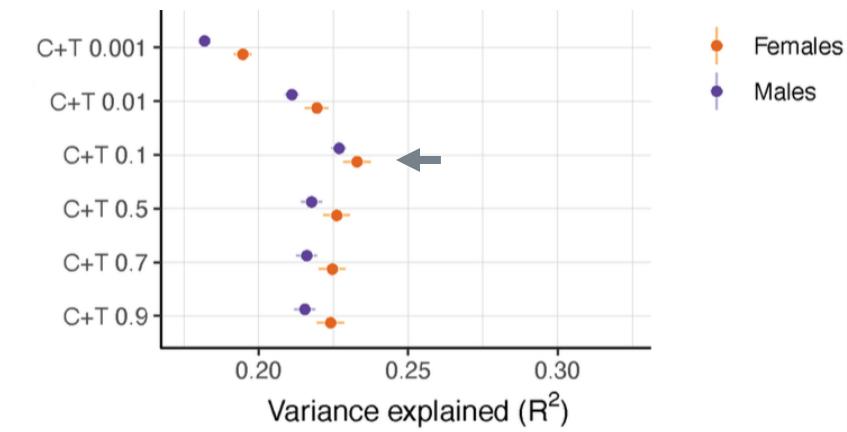
Repeat for other *P*-value cutoffs (and LD values)

How does the PGS associate with the disease

$$y_{trait} = PGS + \varepsilon$$



Finding optimal r<sup>2</sup> and *P*-cutoff and apply in second cohort



# CLUMPING AND THRESHOLDING (C+T)

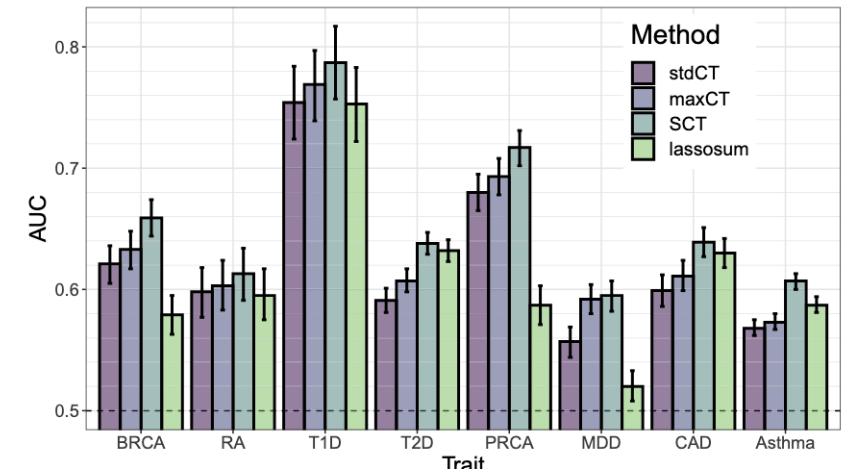
## Clumping and Thresholding (C+T) and Stacked C+T (SCT)

We compute C+T scores for each chromosome separately and for several parameters:

- Threshold on imputation INFO<sub>T</sub> within {0.3, 0.6, 0.9, 0.95}.
- Squared correlation threshold of clumping  $r_c^2$  within {0.01, 0.05, 0.1, 0.2, 0.5, 0.8, 0.95}.
- Base size of clumping window within {50, 100, 200, 500}. The window size  $w_c$  is then computed as the base size divided by  $r_c^2$ . For example, for  $r_c^2 = 0.2$ , we test values of  $w_c$  within {250, 500, 1000, 2500} (in kb). This is motivated by the fact that linkage disequilibrium is inversely proportional to genetic distance between variants.<sup>11</sup>
- A sequence of 50 thresholds on p values between the least and the most significant p values, equally spaced on a log-log scale.

Thus, for individual  $i$ , chromosome  $k$ , and the four hyper-parameters INFO<sub>T</sub>,  $r_c^2$ ,  $w_c$ , and  $p_T$ , we compute C+T predictions

$$V_i^{(k)}(\text{INFO}_T, r_c^2, w_c, p_T) = \sum_{\substack{j \in S_{\text{clumping}}(k, \text{INFO}_T, r_c^2, w_c) \\ p_j < p_T}} \hat{\beta}_j \cdot G_{i,j},$$



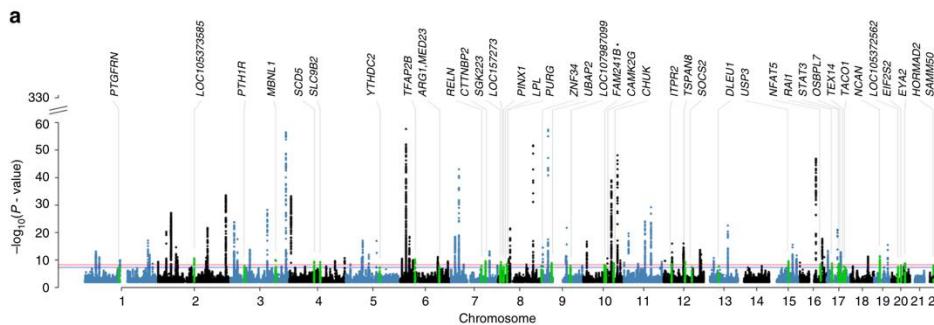
**Table 2. Optimal Choices of C+T Parameters**

Trait	$w_c$	$r_c^2$	INFO <sub>T</sub>	$p_T$
Breast cancer (BRCA)	2,500	0.2	0.95	$2.2 \times 10^{-4}$
Rheumatoid arthritis (RA)	200	0.5	0.95	$7.5 \times 10^{-2}$
Type 1 diabetes (T1D)	10K–50K	0.01	0.90	$2.6 \times 10^{-5}$
Type 2 diabetes (T2D)	625	0.8	0.95	$1.1 \times 10^{-2}$
Prostate cancer (PRCA)	10K–50K	0.01	0.90	$4.2 \times 10^{-6}$
Depression (MDD)	625	0.8	0.95	$1.0 \times 10^{-1}$
Coronary artery disease (CAD)	526	0.95	0.95	$3.5 \times 10^{-2}$
Asthma	2,500	0.2	0.90	$2.2 \times 10^{-4}$

Choice of C+T parameters is based on the maximum AUC in the training set. Hyper-parameters of C+T are the squared correlation threshold  $r_c^2$  and the window size  $w_c$  of clumping, the p value threshold  $p_T$  and the threshold on the quality of imputation INFO<sub>T</sub>. Choosing optimal hyper-parameters for C+T use 63%–90% of the individuals reported in Table 1. Resulting predictions of maxCT in the test set are reported in Figure 2.

# WHAT DO YOU NEED?

**1. A large well-powered GWAS for your trait of interest**

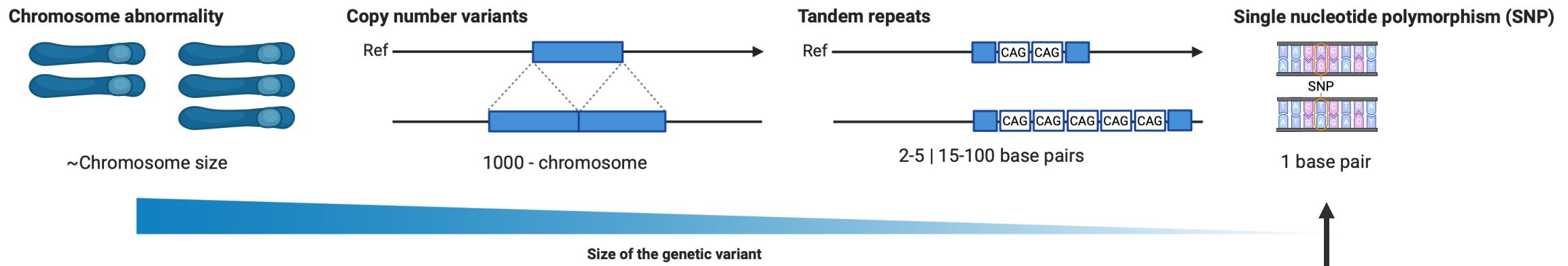


**2. An independent cohort that has been genotyped**



**(3. That some individuals in the cohort has the phenotype)**

# GENETIC DATA / GENETIC VARIATION



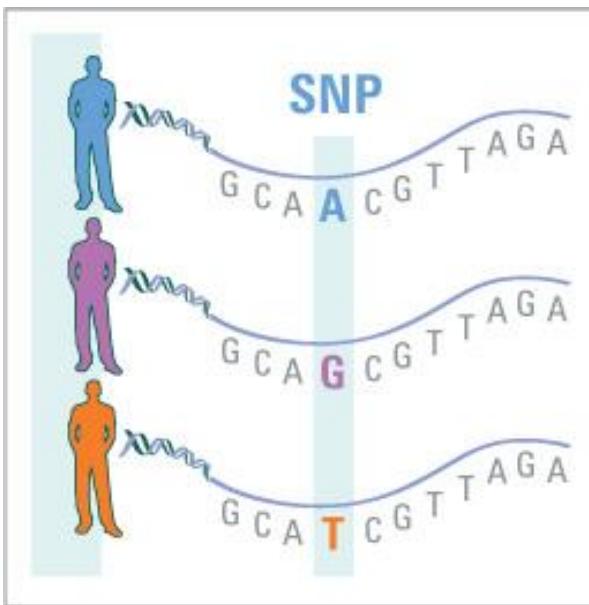
For PGS (GWAS) we are interested in SNP/SNV data

# GENETIC VARIATION

## SINGLE NUCLEOTIDE POLYMORPHISMS (SNPs)

A common change in a single base pair; ~1/1000 bp

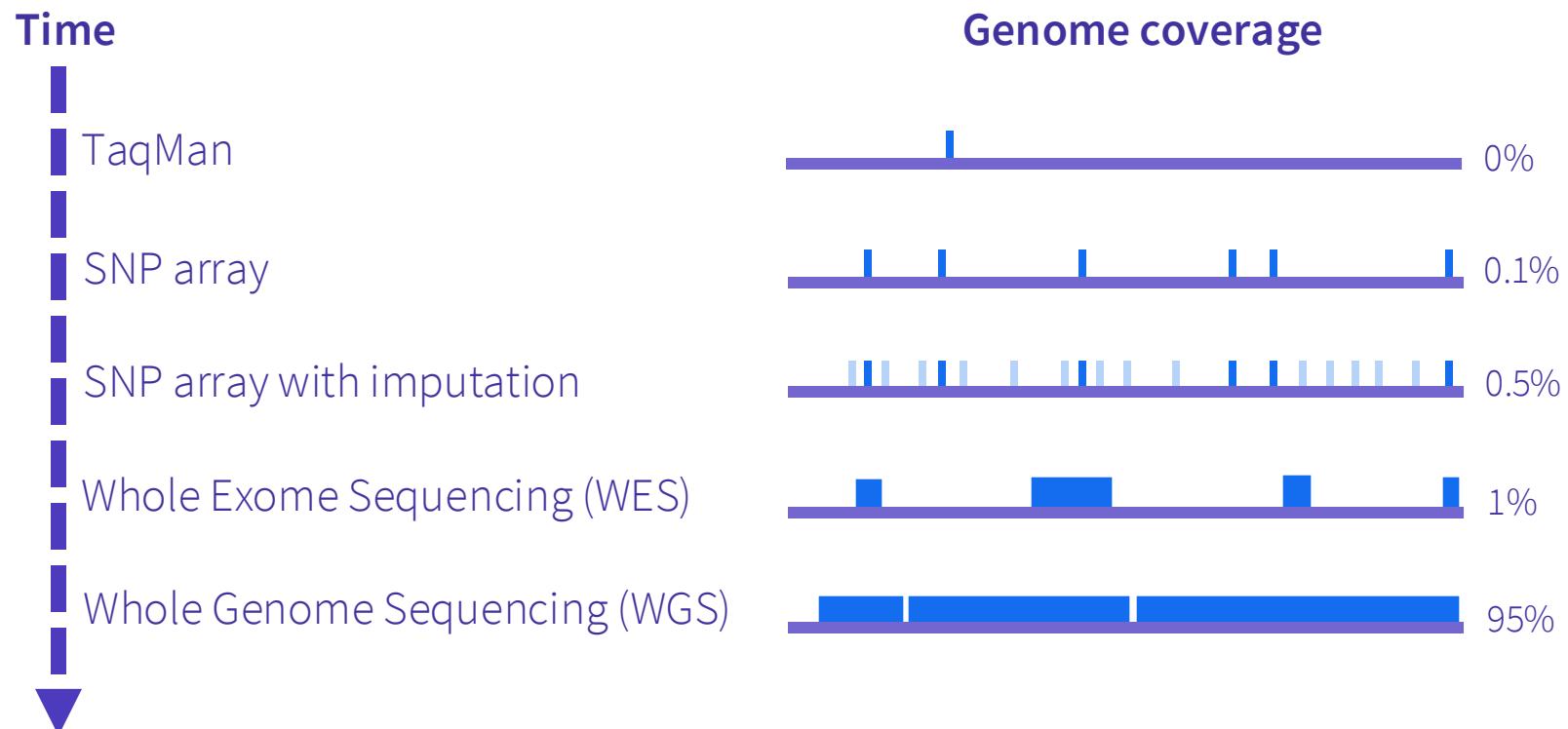
Accounts for ~90% of all variation in the human genome



All (known) SNPs have a unique identifier  
(independent of alleles)

*rsXXX – Ref-SNP cluster ID number*

# GENOMIC COVERAGE



# GENOTYPING VS SEQUENCING

# GENOTYPING VS SEQUENCING

# Genotyping

# GENOTYPING VS SEQUENCING

# Genotyping

WES

# GENOTYPING VS SEQUENCING

# Genotyping

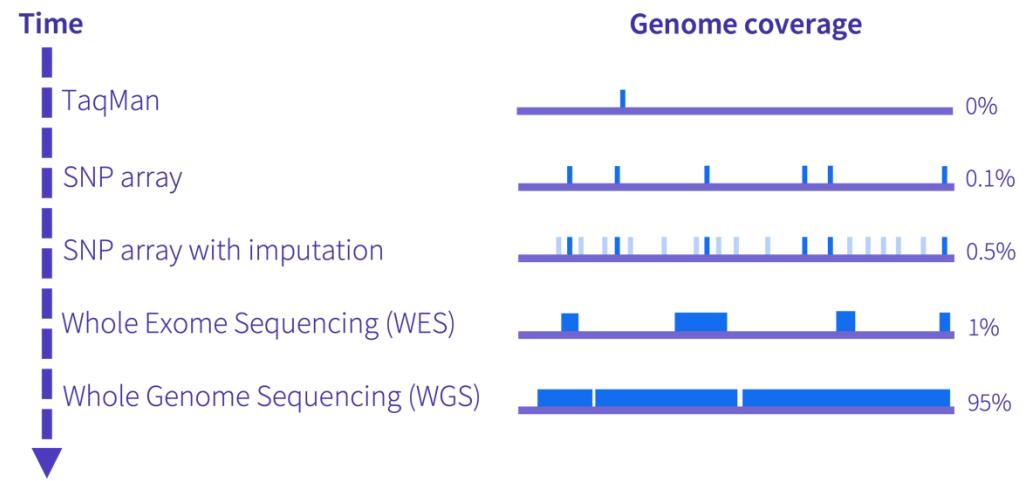
WES

WGS

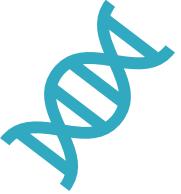
# WHICH TECHNOLOGY?

The choice of technology for detecting single nucleotide polymorphisms (SNPs) depends upon the application.

**For GWAS / PGS we use '*SNP array with imputation*'**



# GENOTYPING



Because of LD you do not have to analyse all 3,000,000,000 variants in the genome.

Typically, we genotype  $\frac{1}{2}$  - 1 million variants

Because of LD we can impute (“guess” what variants are next to the genotyped variant) up to 10 million common genetic variants.

# IMPUTATION USING HAPLOTYPES

The true haplotypes

A	T	C
G	C	A

This individual has inherited a chromosome with alleles A-T-C from one parent, and G-C-A from the other parent

We observe only the genotypes

A/G      T/C      C/A

Genotype data does not carry information about the haplotypes.

We do not know whether A at SNP1 is coming from the same parent as T or C at SNP2

Different haplotypes

A	C	A
A	C	C
A	T	A
A	C	A
G	C	A
G	C	C
G	T	A
G	T	C

**Phasing** = estimate the most likely haplotypes

# IMPUTATION

Genotypes

**Study sample**

.....	A	.....	A	.....	A	.....
.....	G	.....	C	.....	A	.....

**Reference haplotypes**

CGAGATCTCCTTCTTCTGTGC
CGAGATCTCCCGACCTCATGG
CCAAGCTCTTTCTTCTGTGC
CGAAGCTCTTTCTTCTGTGC
CGAGACTCTCCGACCTTATGC
TGGGATCTCCCGACCTCATGG
CGAGATCTCCCGACCTTGTGC
CGAGACTCTTTCTTTGTAC
CGAGACTCTCCGACCTCGTGC
CGAAGCTCTTTCTTCTGTGC

From a sequencing study

**Study sample**

.....	A	.....	A	.....	A	.....
.....	G	.....	C	.....	A	.....

**Reference haplotypes**

CGAGATCTCCTTCTTCTGTGC
CGAGATCTCCCGACCTCATGG
CCAAGCTCTTTCTTCTGTGC
CGAAGCTCTTTCTTCTGTGC
CGAGACTCTCCGACCTTATGC
TGGGATCTCCCGACCTCATGG
CGAGATCTCCCGACCTTGTGC
CGAGACTCTTTCTTTGTAC
CGAGACTCTCCGACCTCGTGC
CGAAGCTCTTTCTTCTGTGC

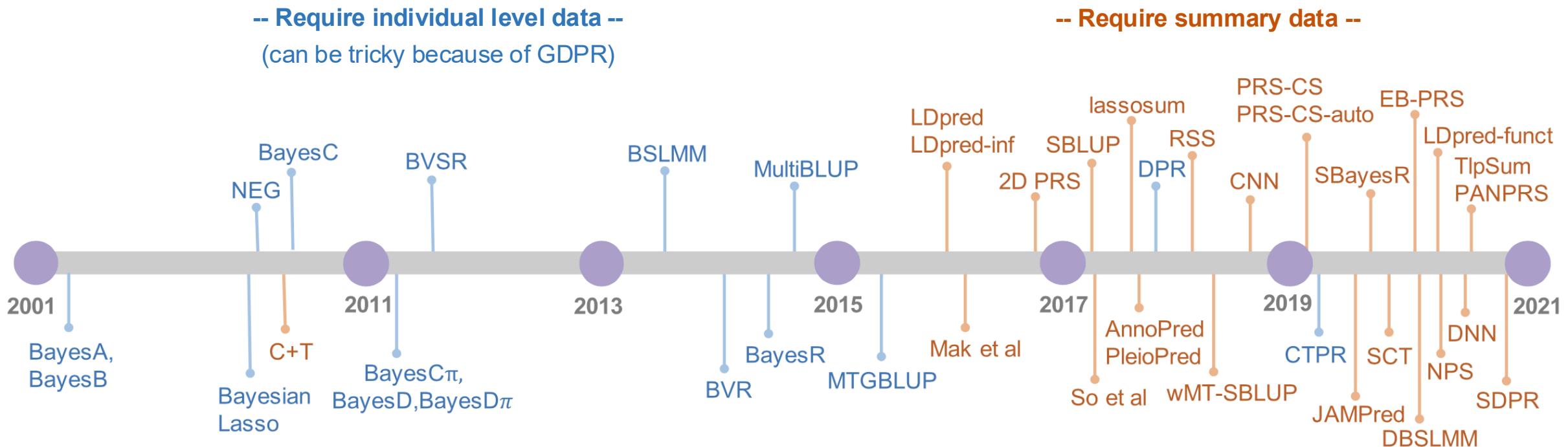
**Study sample**

cgagAtctcccgAcctcAtgg
cgaAGctcttttCtttcAtgg

**Reference haplotypes**

CGGCCCGGGCAATTTTTTTT
CGAGATCTCCCGACCTCATGG
CCAAGCTCTTTCTTCTGTGC
CGAAGCTCTTTCTTCTGTGC
CGAGACTCTCCGACCTTATGC
TGGGATCTCCCGACCTCATGG
CGAGATCTCCCGACCTTGTGC
CGAGACTCTTTCTTTGTAC
CGAGACTCTCCGACCTCGTGC
CGAAGCTCTTTCTTCTGTGC

# A LARGE PALETTE OF PGS METHODS



# COMMONLY USED METHODS

 **(GIGA)<sup>n</sup> SCIENCE**

GigaScience, 8, 2019, 1–6  
doi: 10.1093/gigascience/giz082  
Technical Note

TECHNICAL NOTE

**PRSice-2: Polygenic Risk Score software for biobank-scale data**

Shing Wan Choi  <sup>1,2,\*</sup> and Paul F. O'Reilly  <sup>1,2,\*</sup>

## ARTICLE

### Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores

Bjarni J. Vilhjálmsdóttir <sup>1,2,3,4,\*</sup>, Jian Yang <sup>5,6</sup>, Hilary K. Finucane <sup>1,2,3,7</sup>, Alexander Gusev <sup>1,2,3</sup>, Sara Lindström <sup>1,2</sup>, Stephan Ripke <sup>8,9,10</sup>, Giulio Genovese <sup>3,8,11</sup>, Po-Ru Loh <sup>1,2,3</sup>, Gaurav Bhatia <sup>1,2,3</sup>, Ron Do <sup>12,13</sup>, Tristan Hayeck <sup>1,2,3</sup>, Hong-Hee Won <sup>3,14</sup>, Schizophrenia Working Group of the Psychiatric Genomics Consortium, Discovery, Biology, and Risk of Inherited Variants in Breast Cancer (DRIVE) study, Sekar Kathiresan <sup>3,14</sup>, Michael Pato <sup>15</sup>, Carlos Pato <sup>15</sup>, Rulla Tamimi <sup>1,2,16</sup>, Eli Stahl <sup>3,13,17,18</sup>, Noah Zaitlen <sup>19</sup>, Bogdan Pasaniuc <sup>20</sup>, Gillian Belbin <sup>12,13</sup>, Eimear E. Kenny <sup>12,13,18,21</sup>, Mikkel H. Schierup <sup>4</sup>, Philip De Jager <sup>3,22,23</sup>, Nikolaos A. Patsopoulos <sup>3,22,23</sup>, Steve McCarroll <sup>3,8,11</sup>, Mark Daly <sup>3,8</sup>, Shaun Purcell <sup>3,13,17,18</sup>, Daniel Chasman <sup>22,24</sup>, Benjamin Neale <sup>3,8</sup>, Michael Goddard <sup>25,26</sup>, Peter M. Visscher <sup>5,6</sup>, Peter Kraft <sup>1,2,3,27</sup>, Nick Patterson <sup>3</sup>, and Alkes L. Price <sup>1,2,3,27,\*</sup>

*Bioinformatics*, 36(22–23), 2020, 5424–5431

doi: 10.1093/bioinformatics/btaa1029

Advance Access Publication Date: 16 December 2020

Original Paper

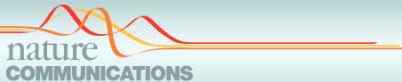


### Genetics and population analysis

### LDpred2: better, faster, stronger

Florian Privé<sup>1,\*</sup>, Julyan Arbel<sup>2</sup> and Bjarni J. Vilhjálmsdóttir<sup>1,3,\*</sup>

<sup>1</sup>National Centre for Register-Based Research, Aarhus University, Aarhus 8210, Denmark, <sup>2</sup>Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, Grenoble 38000, France and <sup>3</sup>Bioinformatics Research Centre, Aarhus University, Aarhus 8000, Denmark



ARTICLE

<https://doi.org/10.1038/s41467-019-109718-5> OPEN

Polygenic prediction via Bayesian regression and continuous shrinkage priors

Tian Ge<sup>1,2,3</sup>, Chia-Yen Chen  <sup>1,2,3,4</sup>, Yang Ni  <sup>5</sup>, Yen-Chen Anne Feng<sup>1,2,3,4</sup> & Jordan W. Smoller<sup>1,2,3</sup>



ARTICLE

<https://doi.org/10.1038/s41467-019-12653-0> OPEN

Improved polygenic prediction by Bayesian multiple regression on summary statistics

Luke R. Lloyd-Jones  <sup>1,9,\*</sup>, Jian Zeng  <sup>1,9,\*</sup>, Julia Sidorenko<sup>1,2</sup>, Loïc Yengo<sup>1</sup>, Gerhard Moser<sup>3,4</sup>, Kathryn E. Kemper<sup>1</sup>, Huanwei Wang  <sup>1</sup>, Zhili Zheng<sup>1</sup>, Reedik Magi<sup>2</sup>, Tõnu Esko<sup>2</sup>, Andres Metspalu<sup>2,5</sup>, Naomi R. Wray  <sup>1,6</sup>, Michael E. Goddard<sup>7</sup>, Jian Yang  <sup>1,8,\*</sup> & Peter M. Visscher  <sup>1,\*</sup>

*Bioinformatics*, 36(8), 2020, 2614–2615  
doi: 10.1093/bioinformatics/btz955  
Advance Access Publication Date: 27 December 2019  
Applications Note

Genetics and population analysis  
**qgg: an R package for large-scale quantitative genetic analyses**

Palle Duun Rohde  <sup>\*†</sup>, Izel Fourie Sørensen<sup>2</sup>, Peter Sørensen<sup>2,\*</sup>

*Bioinformatics*, 2023, 39(11), btad656  
<https://doi.org/10.1093/bioinformatics/btad656>  
Advance Access Publication Date: 26 October 2023  
Applications Note

Genetics and population analysis  
**Expanded utility of the R package, qgg, with applications within genomic medicine**

Palle Duun Rohde  <sup>\*†</sup>, Izel Fourie Sørensen<sup>2</sup>, Peter Sørensen<sup>2,\*</sup>

<sup>1</sup>Genomic Medicine, Department of Health Science and Technology, Aalborg University, 9260 Gistrup, Denmark  
<sup>2</sup>Center for Quantitative Genetics and Genomics, Aarhus University, 8000 Aarhus, Denmark  
\*Corresponding authors. Genomic Medicine, Department of Health Science and Technology, Aalborg University, Selma Langerleffs Vej 249, 9260 Gistrup, Denmark. E-mail: pdr@qgg.aau.dk (P.D.R.); Center for Quantitative Genetics and Genomics, Aarhus University, C. F. Møllers Allé 3, 8000 Aarhus, Denmark. E-mail: psø@qgg.aau.dk (P.S.)  
Associate Editor: Christina Kendziora

# COMMONLY USED METHODS

Different shrinkage methods

**Clumping and thresholding (C+T)**

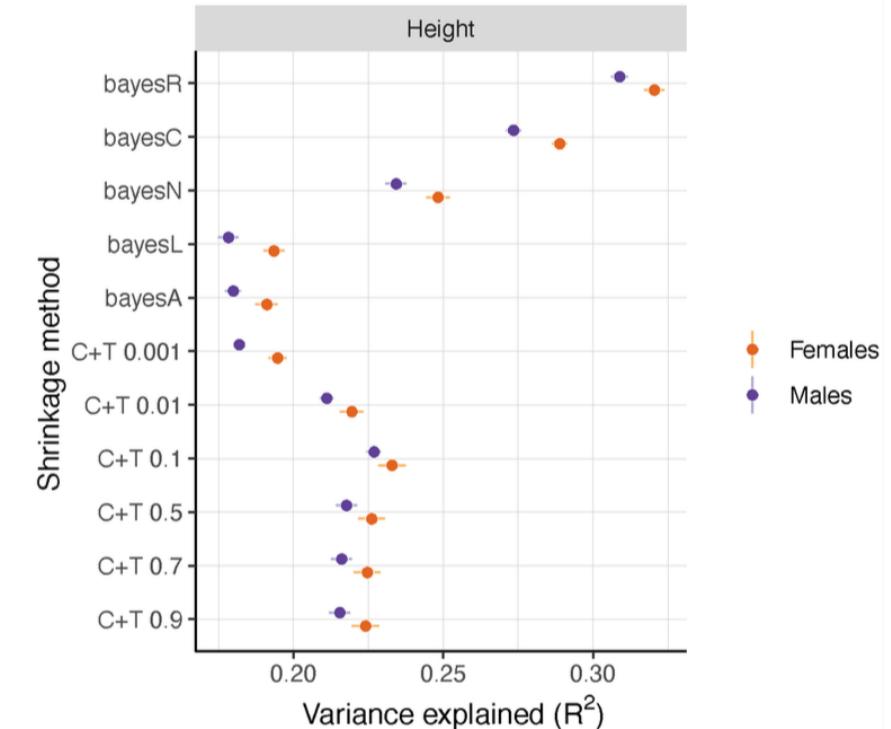
**Bayes-N**;  $\beta \sim N(0, \sigma_\beta^2)$

**Bayes-L**;  $f(\beta_j | \tau_j^2, \sigma_e^2) \sim N(\beta_j | 0, \tau_j^2 \times \sigma_e^2)$

**Bayes-A**;  $\beta_j \sim N(0, \sigma_{\beta_j}^2)$

**Bayes-C**;  $\beta_j \sim N(0, \sigma_{\beta_j}^2)$  with probability  $\pi$ , and  $\beta_j = 0$  with probability  $(1 - \pi)$ , where  $\pi$  is assumed to follow a beta-distribution.

**Bayes-R**;  $\beta_j \sim N(0, \gamma_C \sigma_{\beta_j}^2)$ , where  $C$  defines number of classes (e.g.,  $C=4, \gamma = (0, 0.01, 0.1, 1.0)$ )



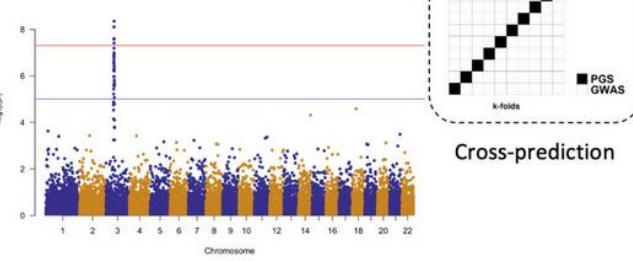
## 1. Discovery

### a. Individual-level data

	SNP 1	SNP 2	SNP 3	y
1	CT	AA	CA	12
2	CT	AA	CA	5
3	TT	TT	CC	10
4	CC	AT	AA	8

$$\mathbf{b} = \frac{\mathbf{x}_j^t \mathbf{y}}{\mathbf{x}_j^t \mathbf{x}_j}$$

### b. GWAS results



Cross-prediction

### c. Summary-level data

	SNP 1	SNP 2	SNP 3	SNP 4
--	-------	-------	-------	-------

Effect allele	C	A	C	T
Weight	0.2	-0.3	0.1	0.2

#### GWAS Repositories



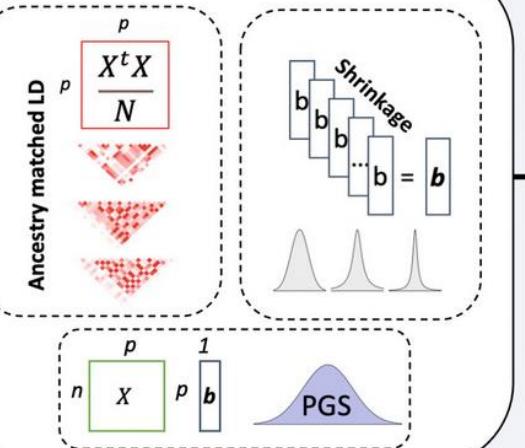
GWAS catalog



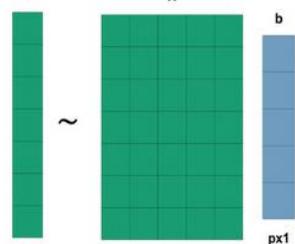
Open GWAS

## 2. Validation

d.

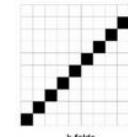


### e. Optimization vs phenotype



f.

#### □ Cross-validation for parameter tuning



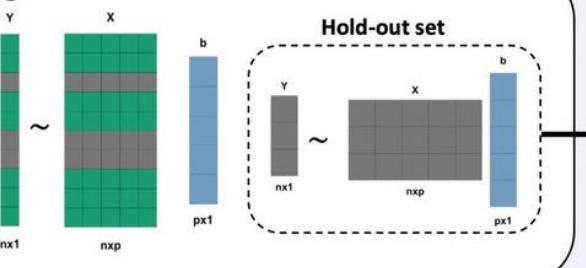
#### □ PGS repositories:

- PGS catalog
- PGI repository

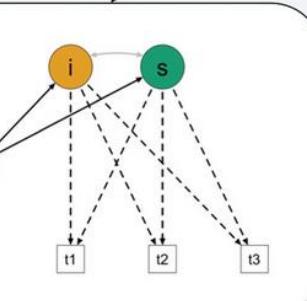
#### □ Single score methods (e.g. SBLUP)

## 3. Target

g.



h.



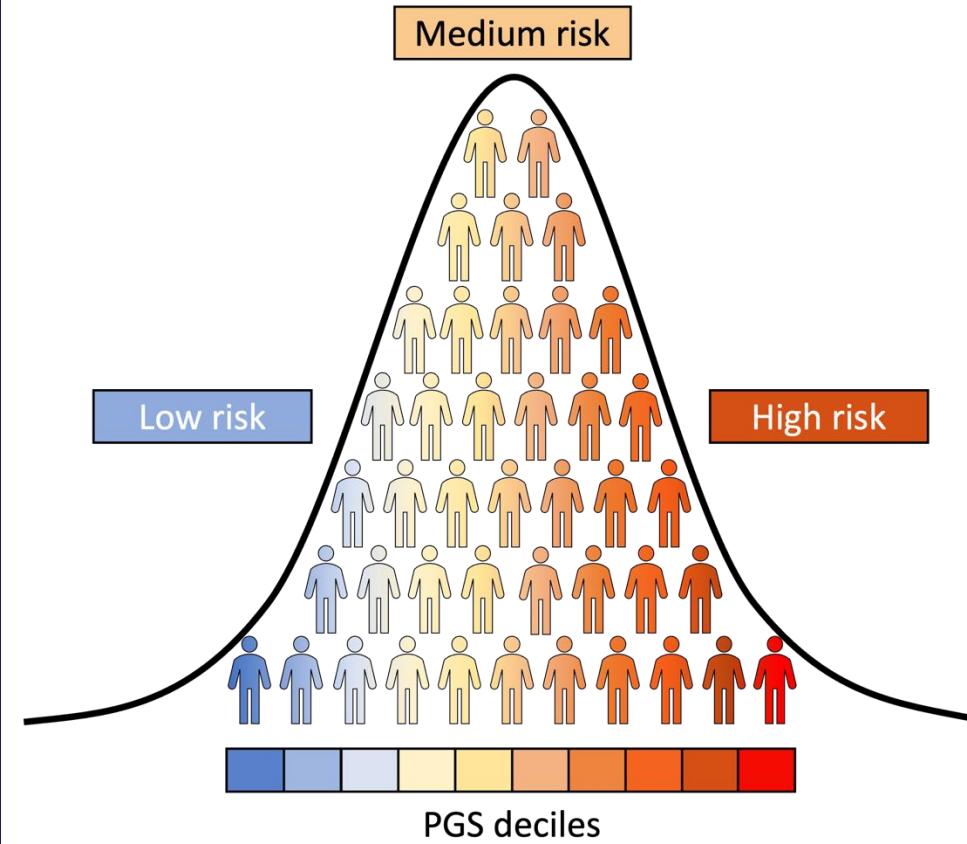
i.

#### □ Optimization of PGS holding-out data for testing

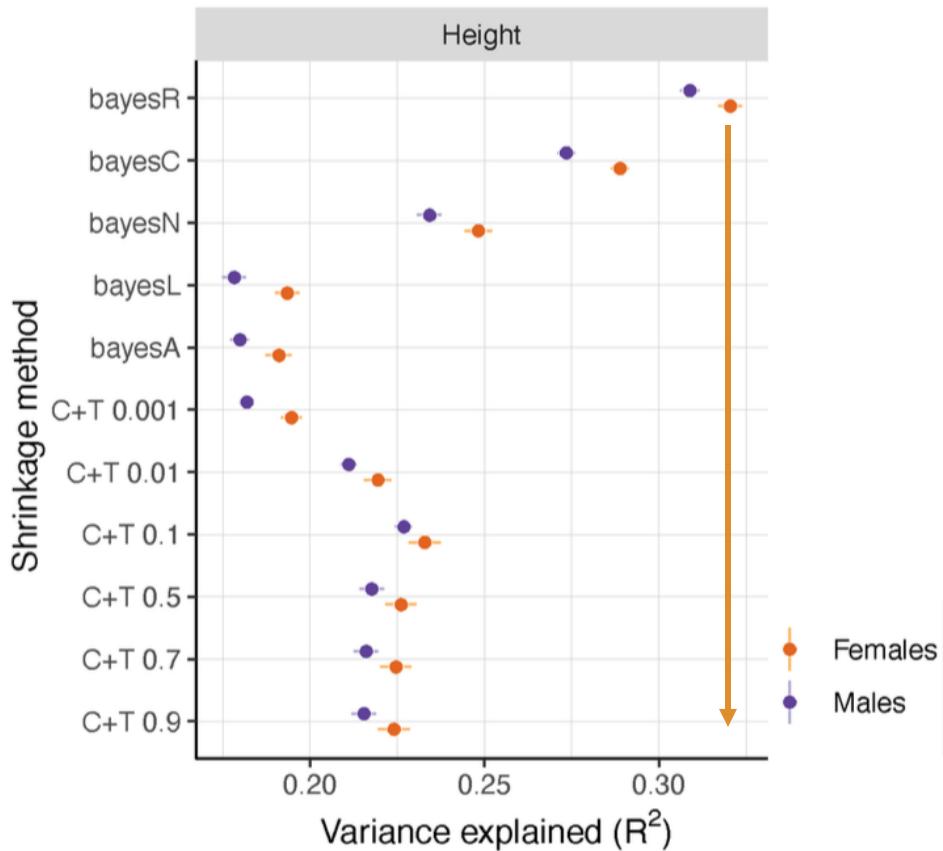
#### □ Split-validation

Allegrini et al (2022) J Child Psych,  
63:1111-1124

# CHALLENGES WITH PGS



# HOW GOOD ARE PGS?



For human height the best PGS has an accuracy of ~40%.

Human height has a heritability ( $h^2$ ) of ~65%.

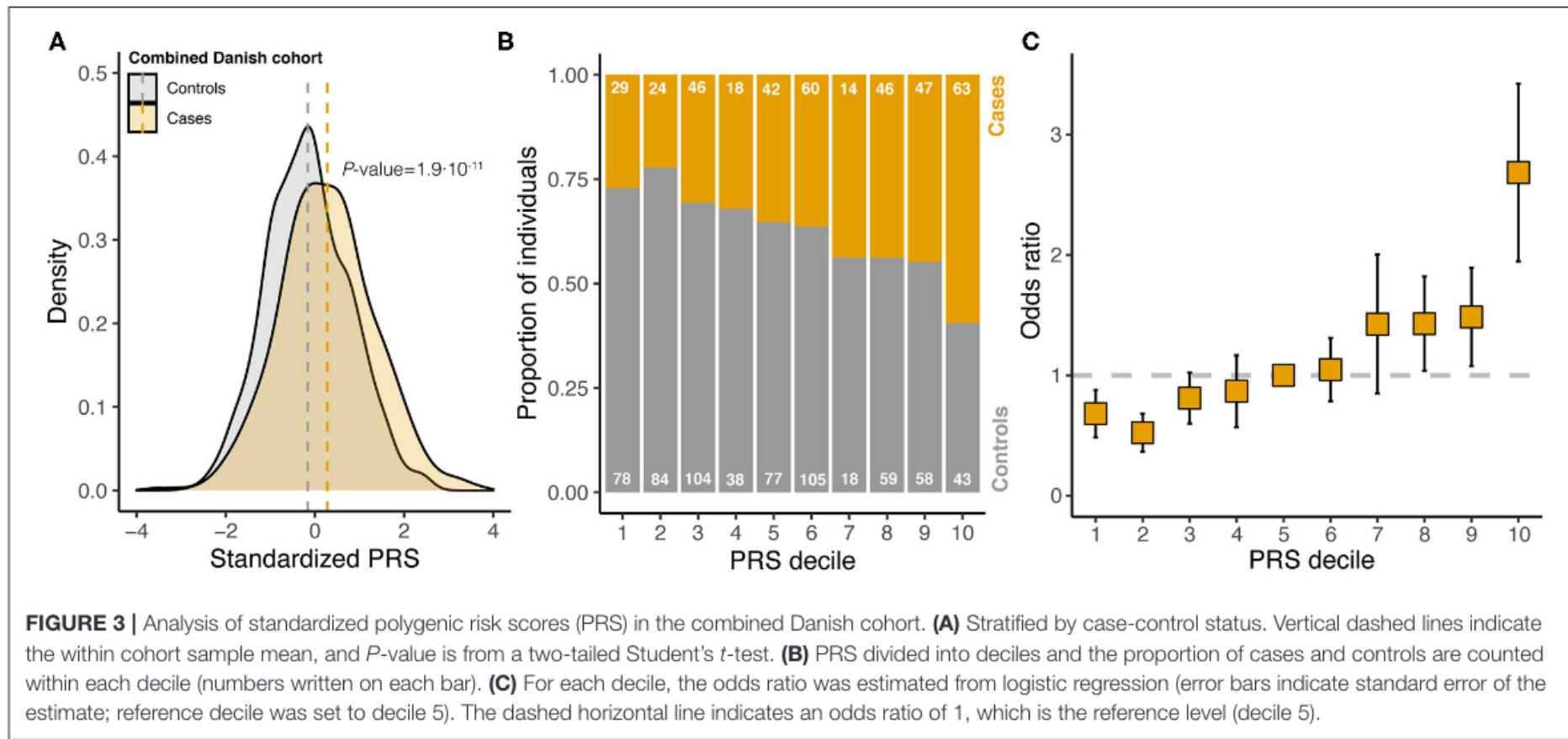
The heritability is the boundary of how accurately we can predict the phenotype, thus, for human height the maximum accuracy is expected to be 65%.

Thus, for human height we currently lack 30% variance explained

**Several reasons/challenges for lack of predictive ability**

# CHALLENGE 1

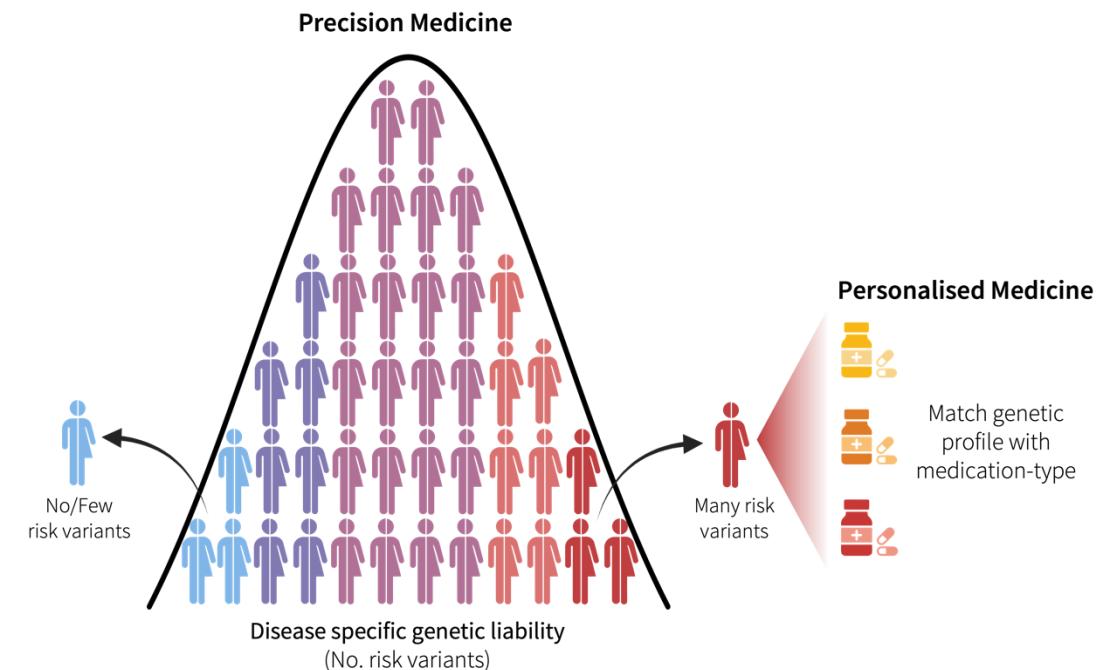
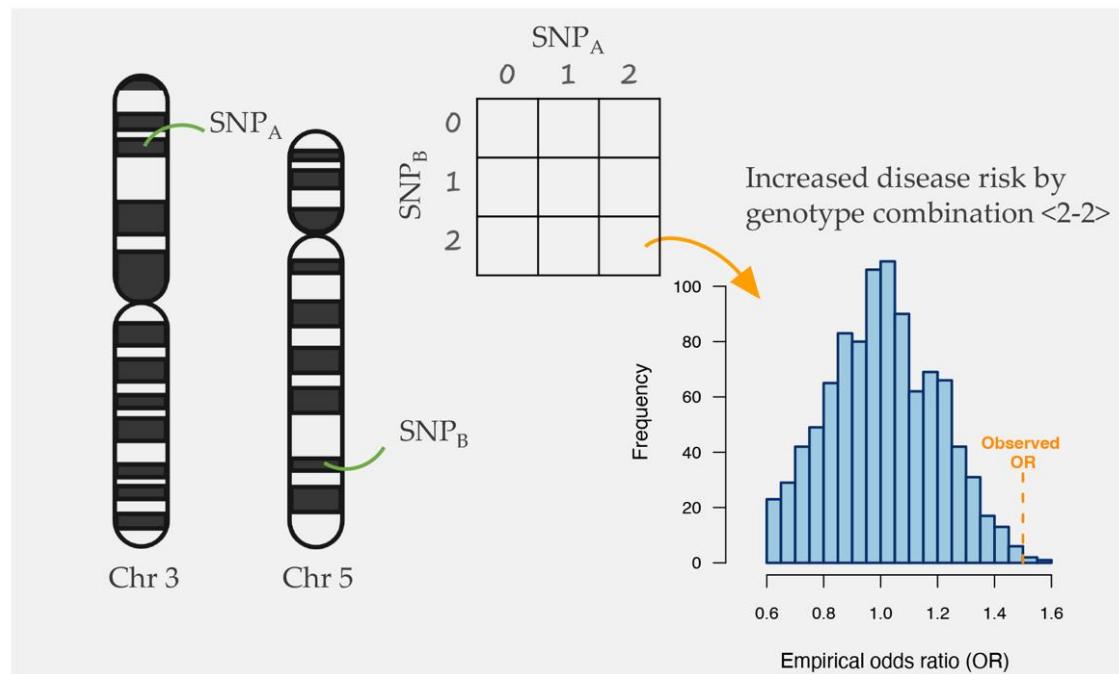
## TOO MANY WITH "AVERAGE" SCORE



$$PGS = \sum X_i b_i$$

# CHALLENGE 2

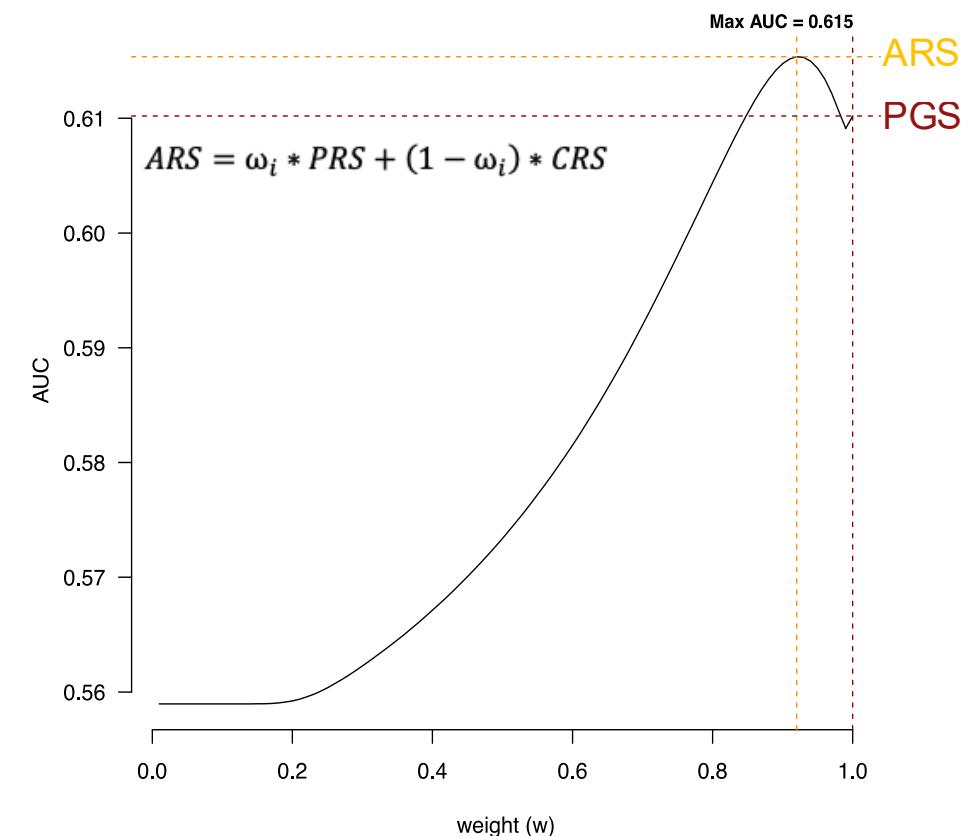
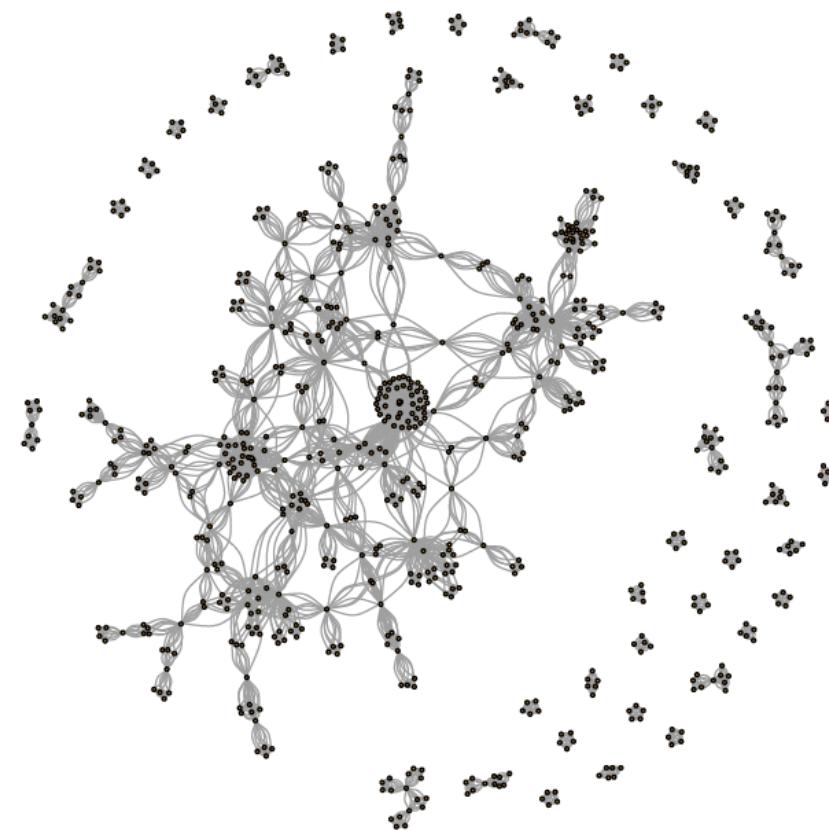
## PGS ARE BASED ON ADDITIVE EFFECTS



$$PGS = \sum X_i b_i$$

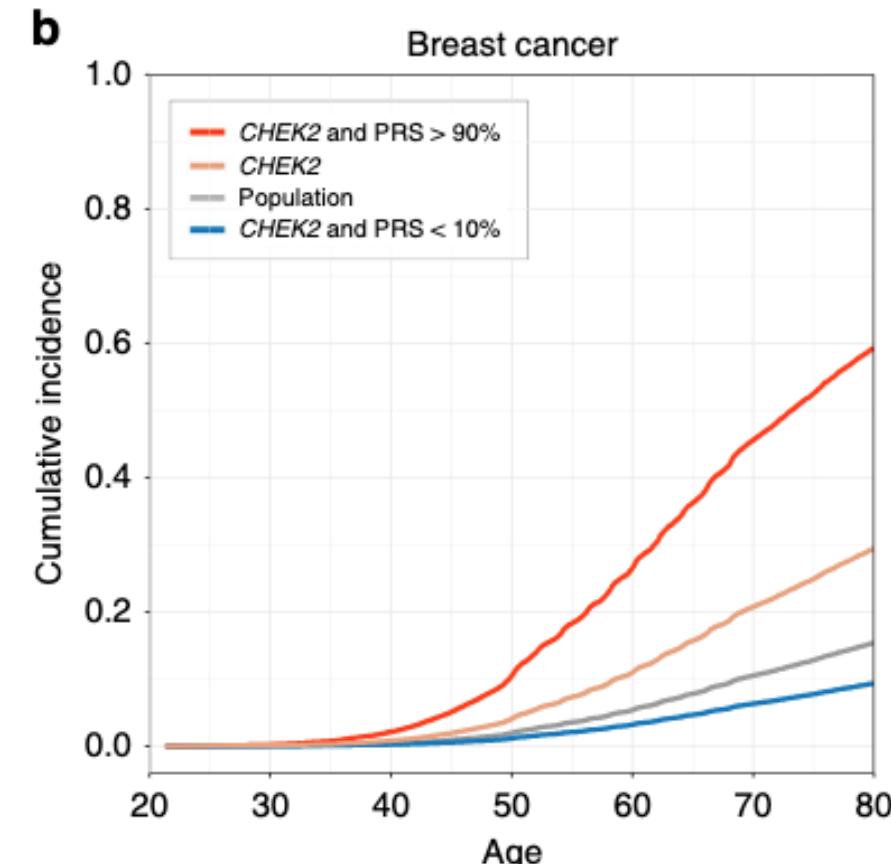
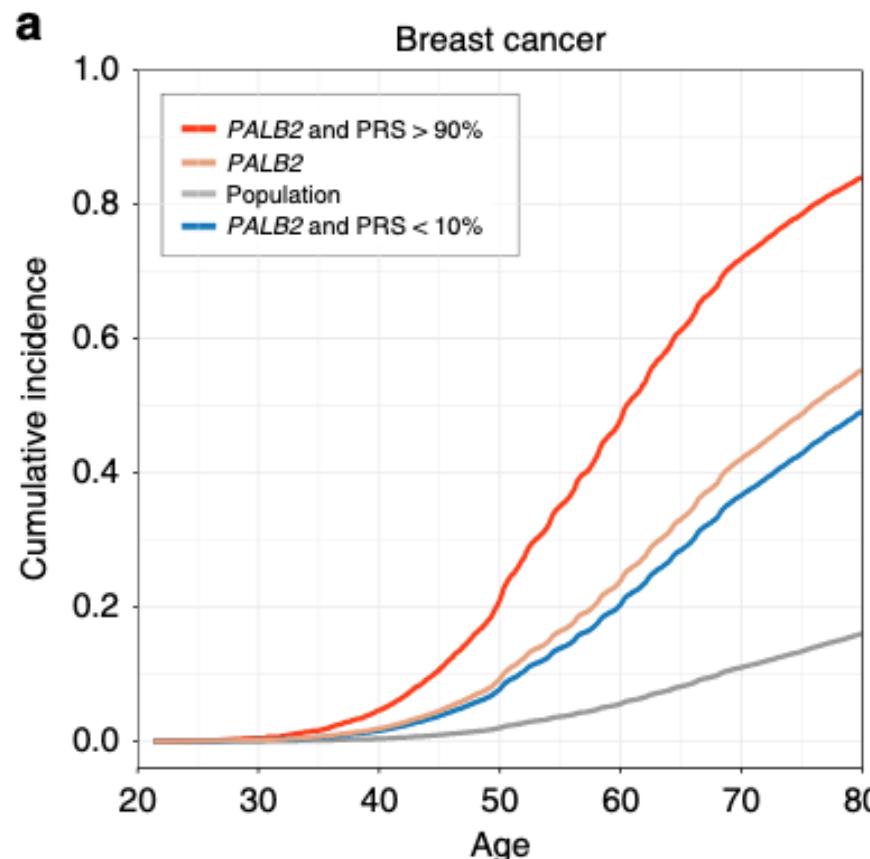
# CHALLENGE 2

PGS ARE BASED ON ADDITIVE EFFECTS



# CHALLENGE 3

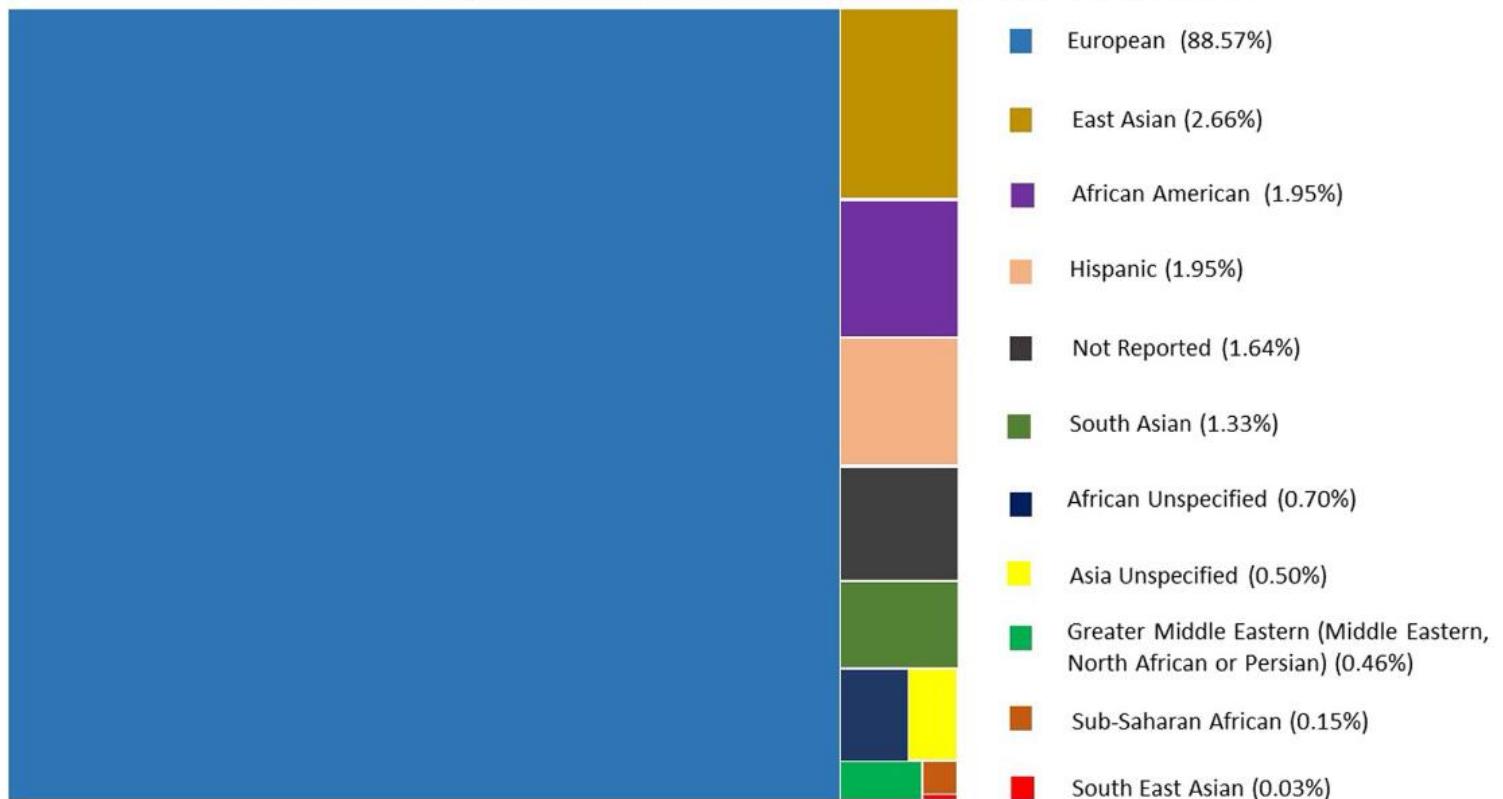
## PGS ARE BASED ON COMMON SNPs



# CHALLENGE 4

## POPULATION SPECIFIC

Broad ancestry categories contributing to development and evaluation of PGS scores



Prediction of T2D

Population	Liability $R^2$	Covariates-adjusted AUC
European	9.2%	0.66
African	2.8%	0.58
Hispanic	8.0%	0.63

# SESSION 2

- The first polygenic score
- What you need is...
- Commonly used scoring algorithms
- Workflow
- Current challenges with PGS?





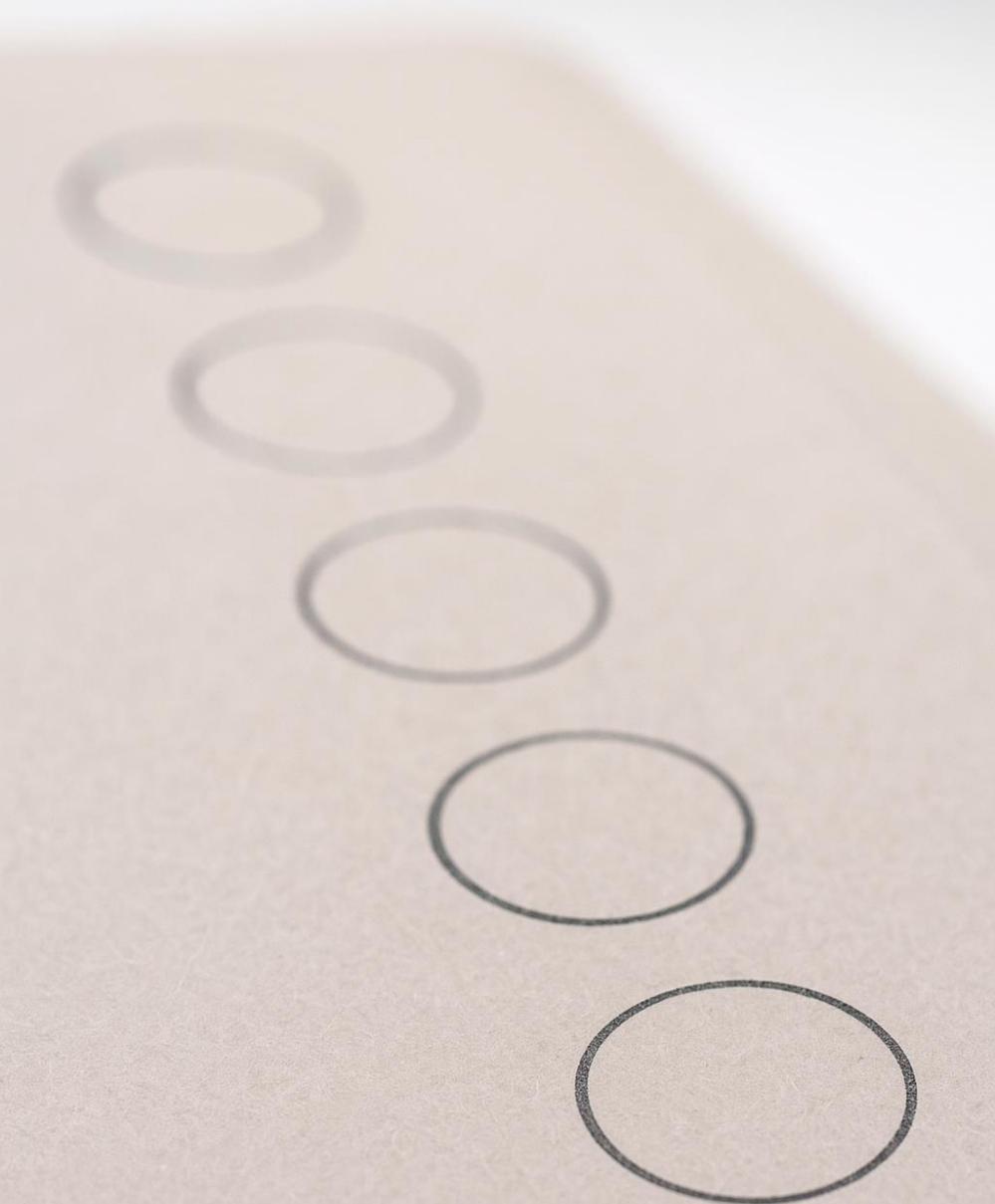
**BREAK**

# AGENDA

08:00 – 08:30	Welcome and common introductions
08:30 – 09:10	Session 1: Introduction to Polygenic Scores (PGS)
09:10 – 09:20	Break
09:20 – 10:00	Session 2: Data Sources and Computational Methods
10:00 – 10:10	Break
<b>10:10 – 10:40</b>	<b>Session 3: Evaluating and Interpreting Polygenic Scores</b>
10:40 – 11:00	Break
11:00 – 11:45	Session 4: Advanced Applications and Future Directions
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today

# SESSION 3

- How to measure 'accuracy'?
- Interpretability and risk communication
- Lack of transferability
- Applications...?



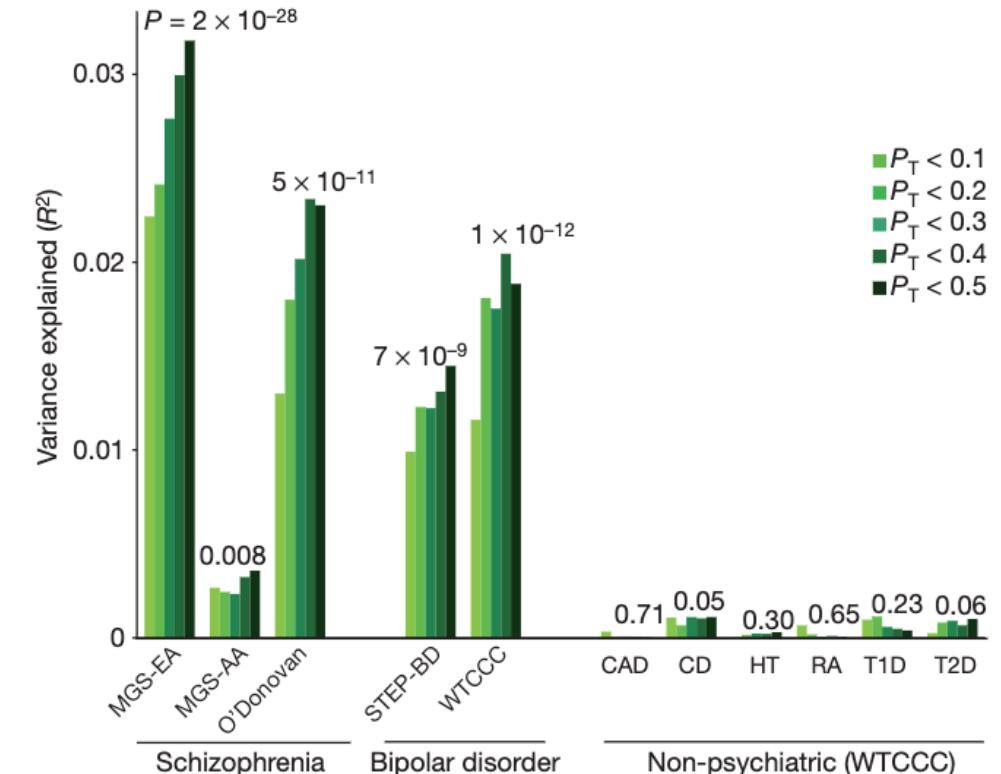
# EVALUATE POLYGENIC PROFILES

$$y = Xb + Zc + e$$

$y$  = phenotype;  $X$  = PGS;  $Z$  = covariates

Compare variance explained from the full model (with  $X$ +covariates) compared to a reduced model (covariates only)

Variance explained ( $R^2$ ) for quantitative traits, and Nagelkerke's  $R^2$  for binary traits (**however, NagR2 is biased with disease prevalence!**)



The International Schizophrenia Consortium (2009) Nature, 460

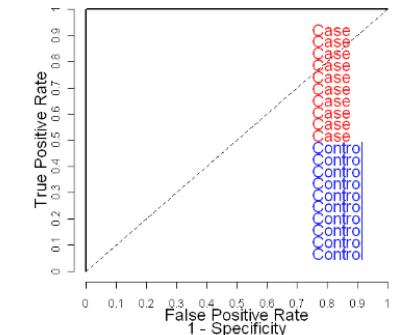
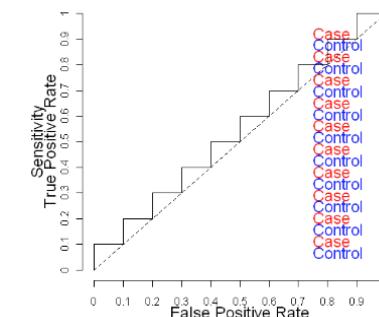
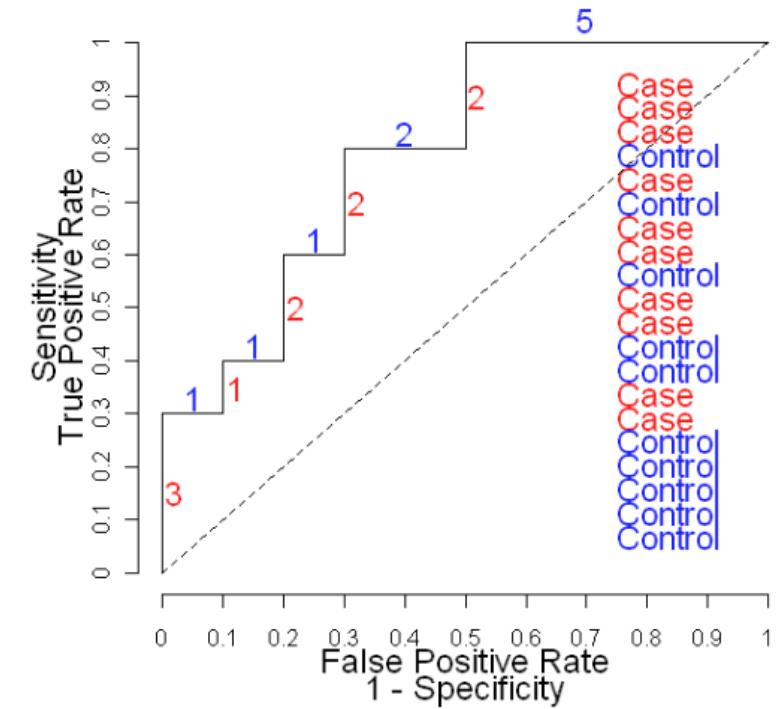
# EVALUATE POLYGENIC PROFILES

## Area Under Receiver Operator Characteristic Curve (AUC)

Well established measure of validity of tests for classifier diseased vs non-diseased individuals

- Nice property – independent to proportion of cases and controls in sample
- Range 0.5 to 1
- 0.5 the score has no predictive value
- **Probability that a randomly selected case has a score higher than a randomly selected control**

Rank individuals on score from highest ranked to lowest



# EVALUATE POLYGENIC PROFILES

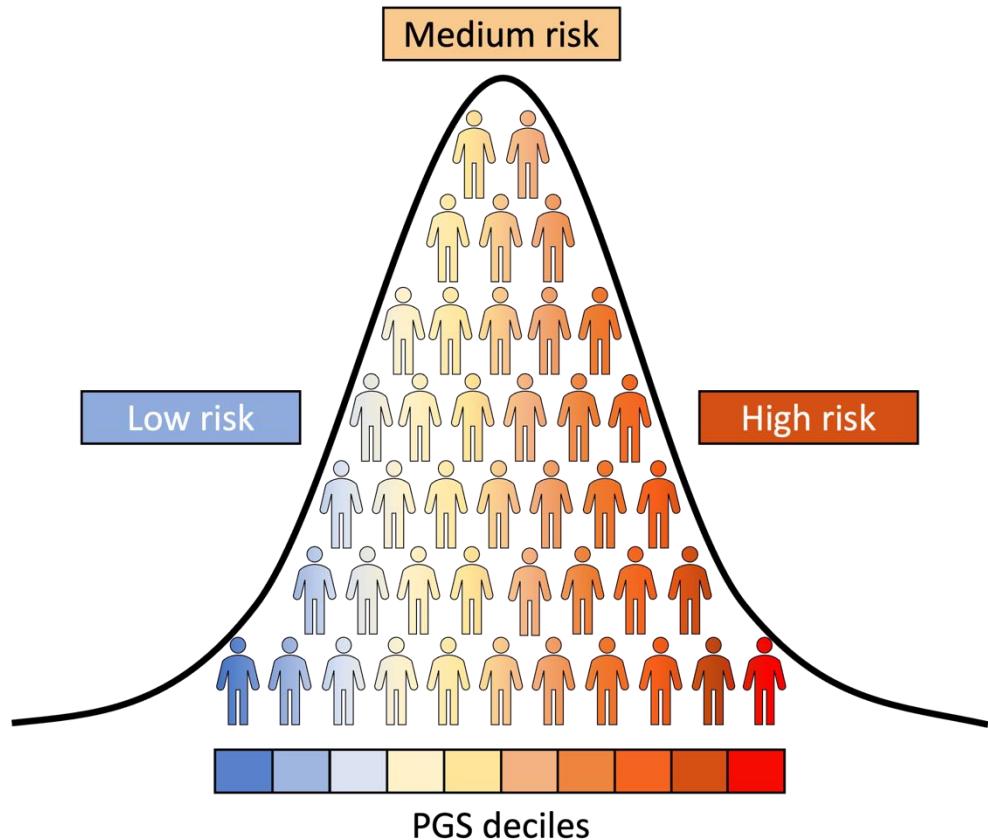
Odds ratio (OR)

Cut the distribution into deciles

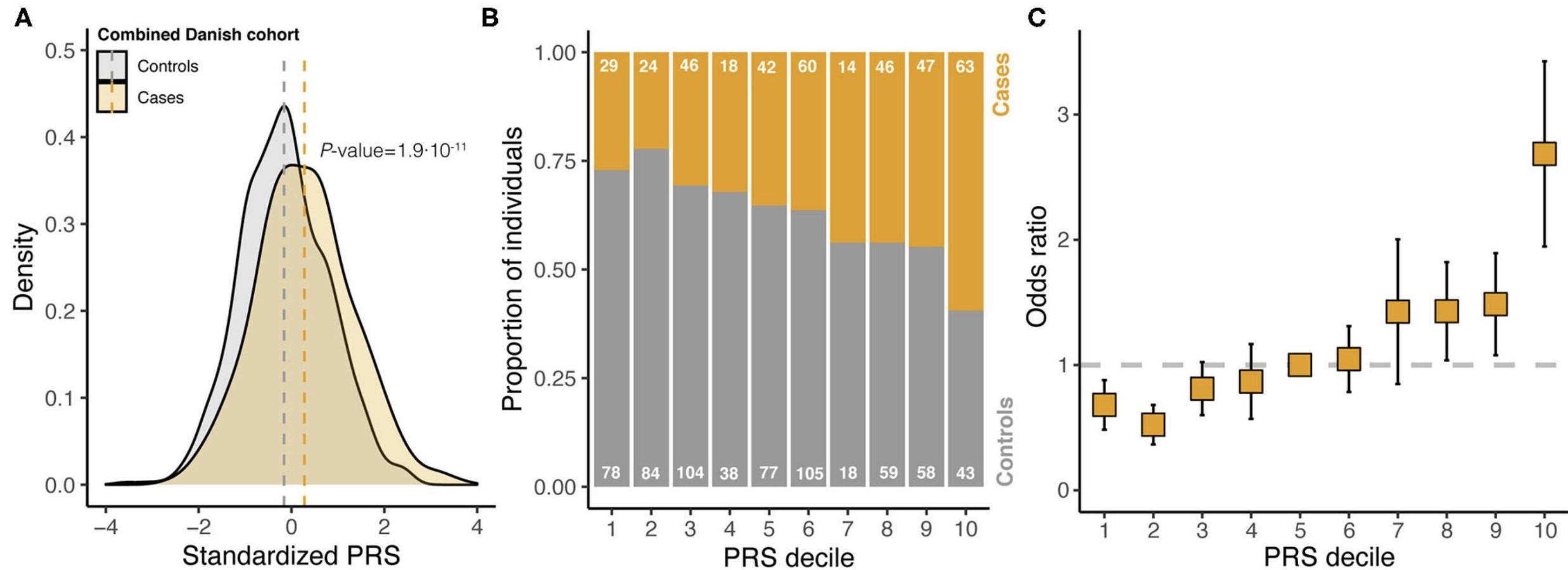
Each decile will include both cases and controls

Odds of being a case in each decile

Odds ratio for each decile compared to the first decile

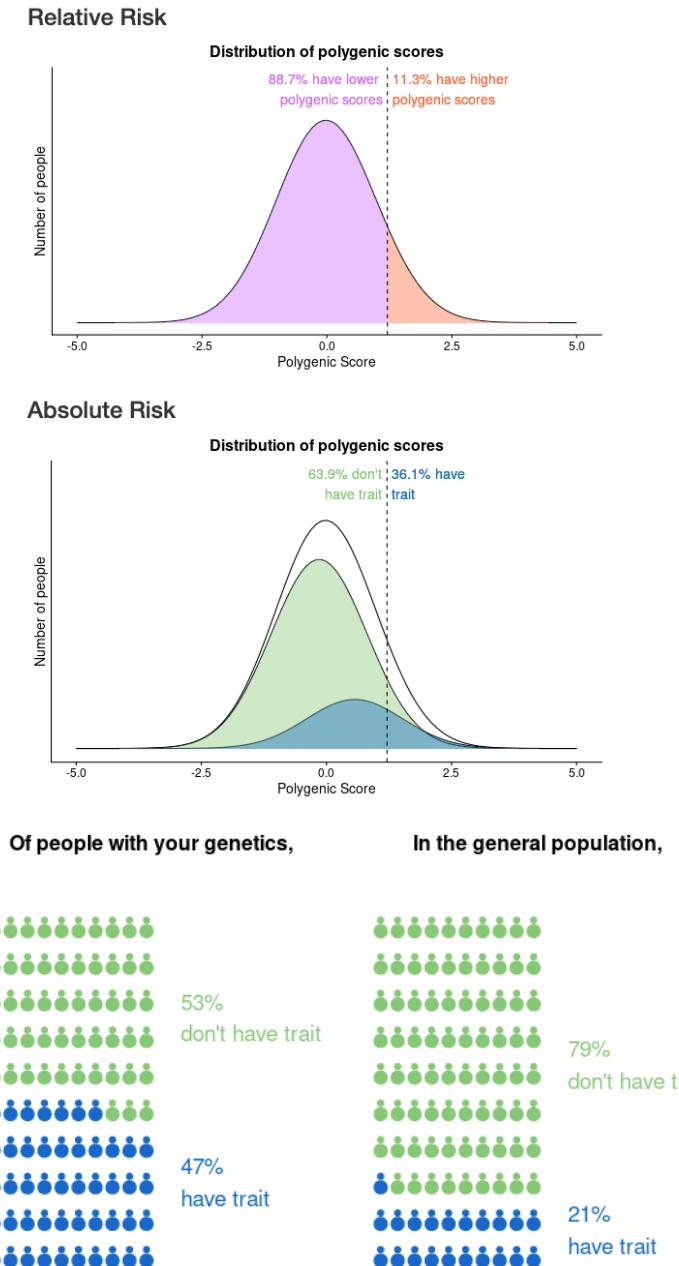


# A 14-SNP PGS

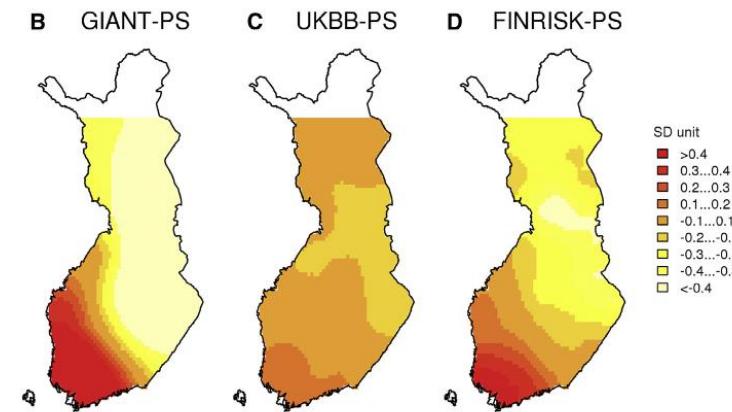
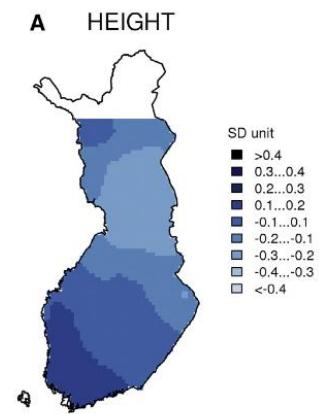
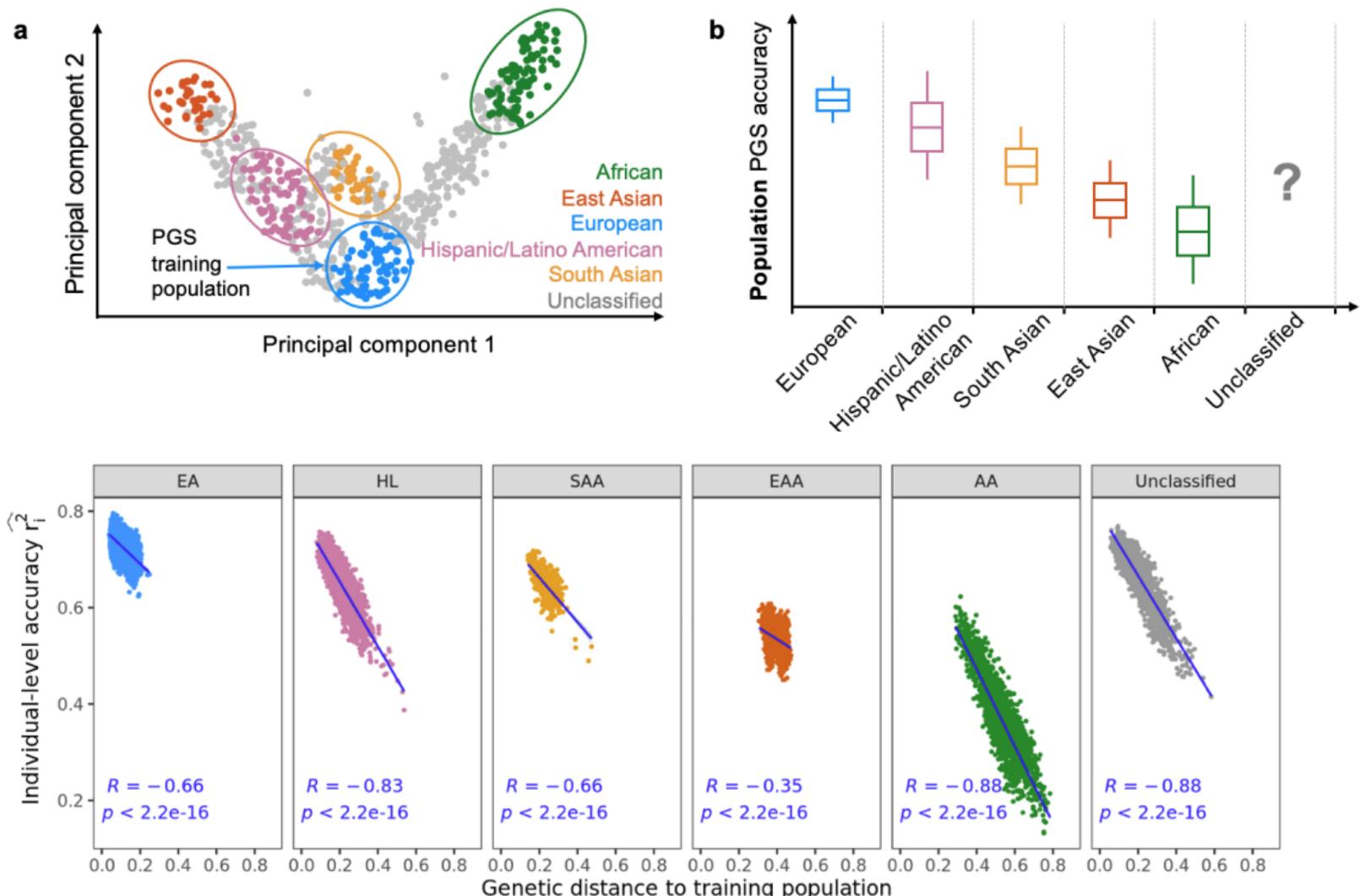


# INTERPRETABILITY & RISK COMMUNICATION

- The risk associated with the PGS is a relative risk
  - People have suggested methods to convert relative PGS risk to absolute risks ([https://opain.github.io/GenoPred/PRS\\_to\\_Abs\\_tool.html](https://opain.github.io/GenoPred/PRS_to_Abs_tool.html))
- The relative risk may sound high, but the absolute risk is low
- Hard to use meaningfully in clinical decisions without baseline risk
- Effective for population stratification, less so for individual prediction
- Lack of population transferability



# LACK OF TRANSFERABILITY



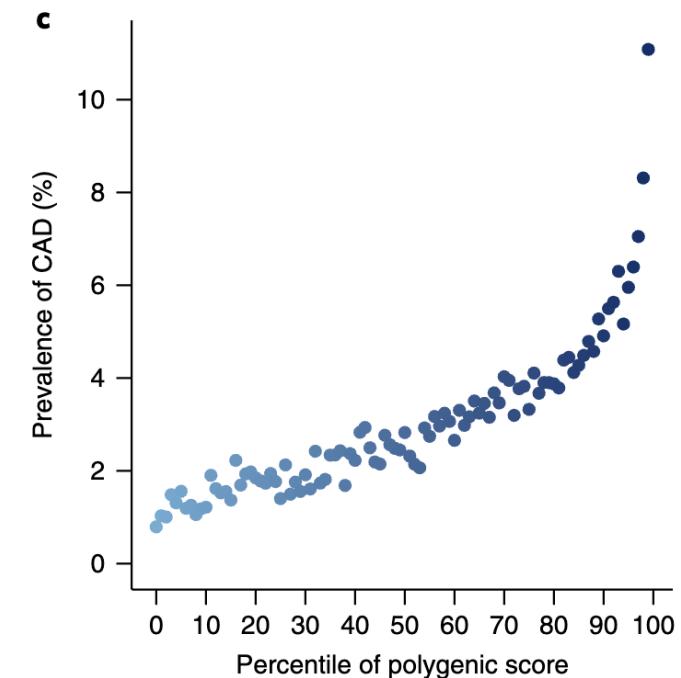
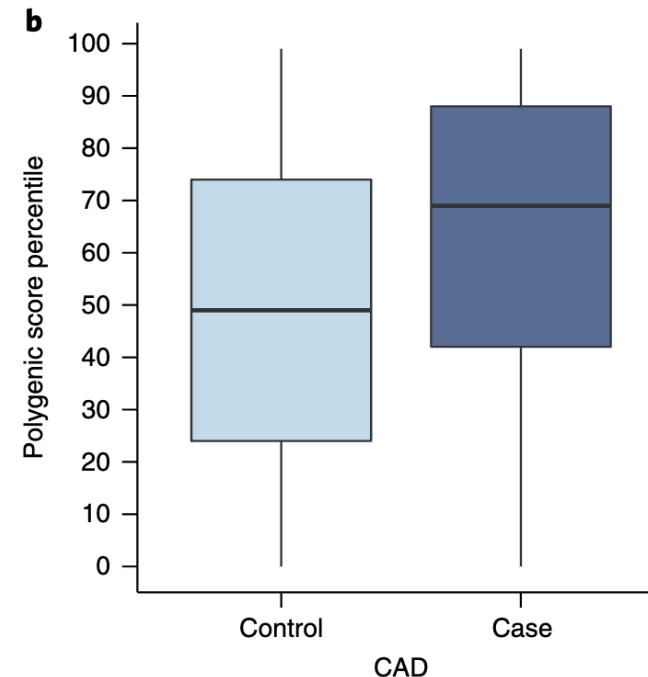
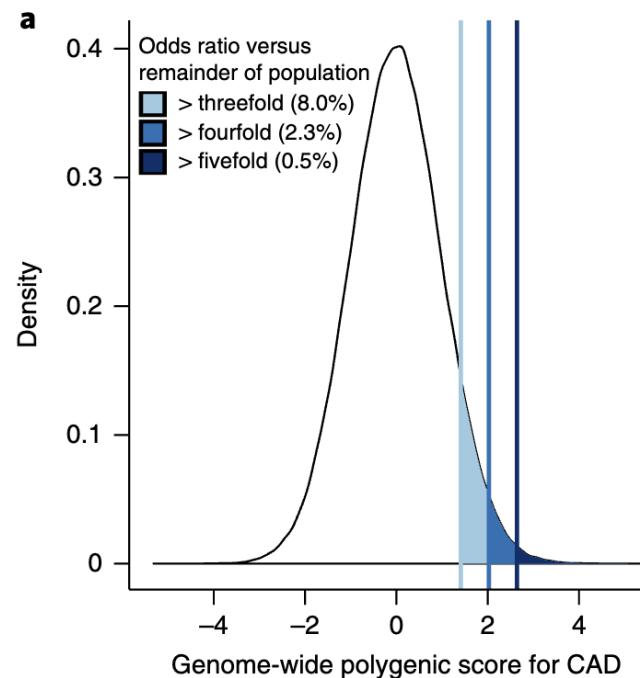
# APPLICATIONS?



# PGS FOR CAD

**Table 3 | Prevalence and clinical impact of a high GPS**

High GPS definition	Reference group	Odds ratio	95% CI	P value
<b>CAD</b>				
Top 20% of distribution	Remaining 80%	2.55	2.43-2.67	$<1 \times 10^{-300}$
Top 10% of distribution	Remaining 90%	2.89	2.74-3.05	$<1 \times 10^{-300}$
Top 5% of distribution	Remaining 95%	3.34	3.12-3.58	$6.5 \times 10^{-264}$
Top 1% of distribution	Remaining 99%	4.83	4.25-5.46	$1.0 \times 10^{-132}$
Top 0.5% of distribution	Remaining 99.5%	5.17	4.34-6.12	$7.9 \times 10^{-78}$

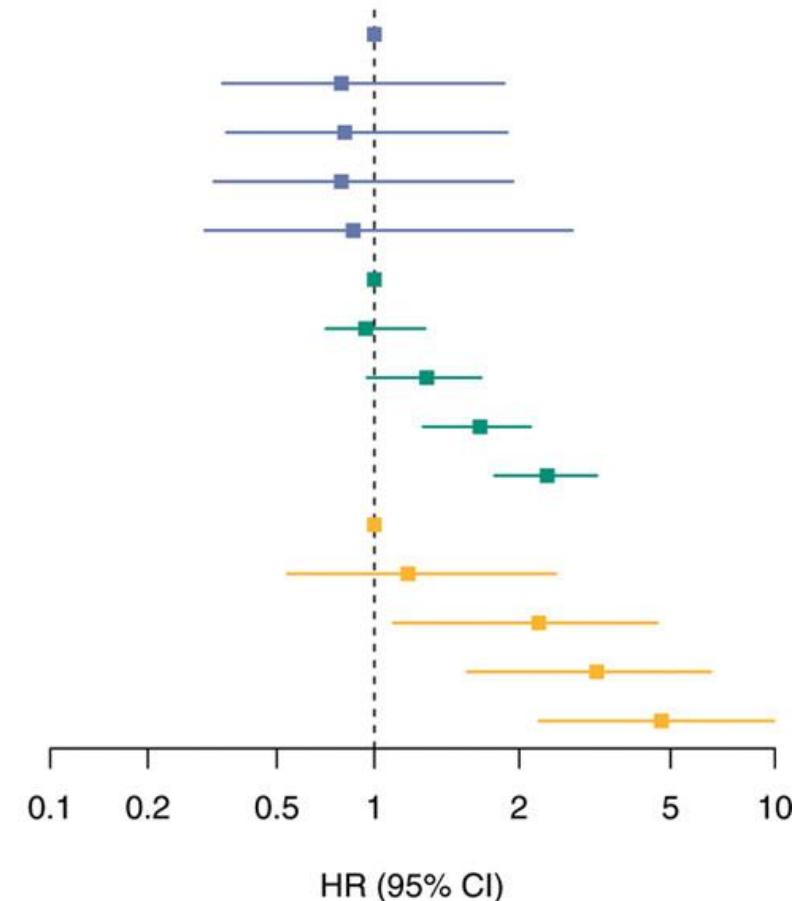


**Table 2 | Proportion of the population at three-, four- and fivefold increased risk for each of the five common diseases**

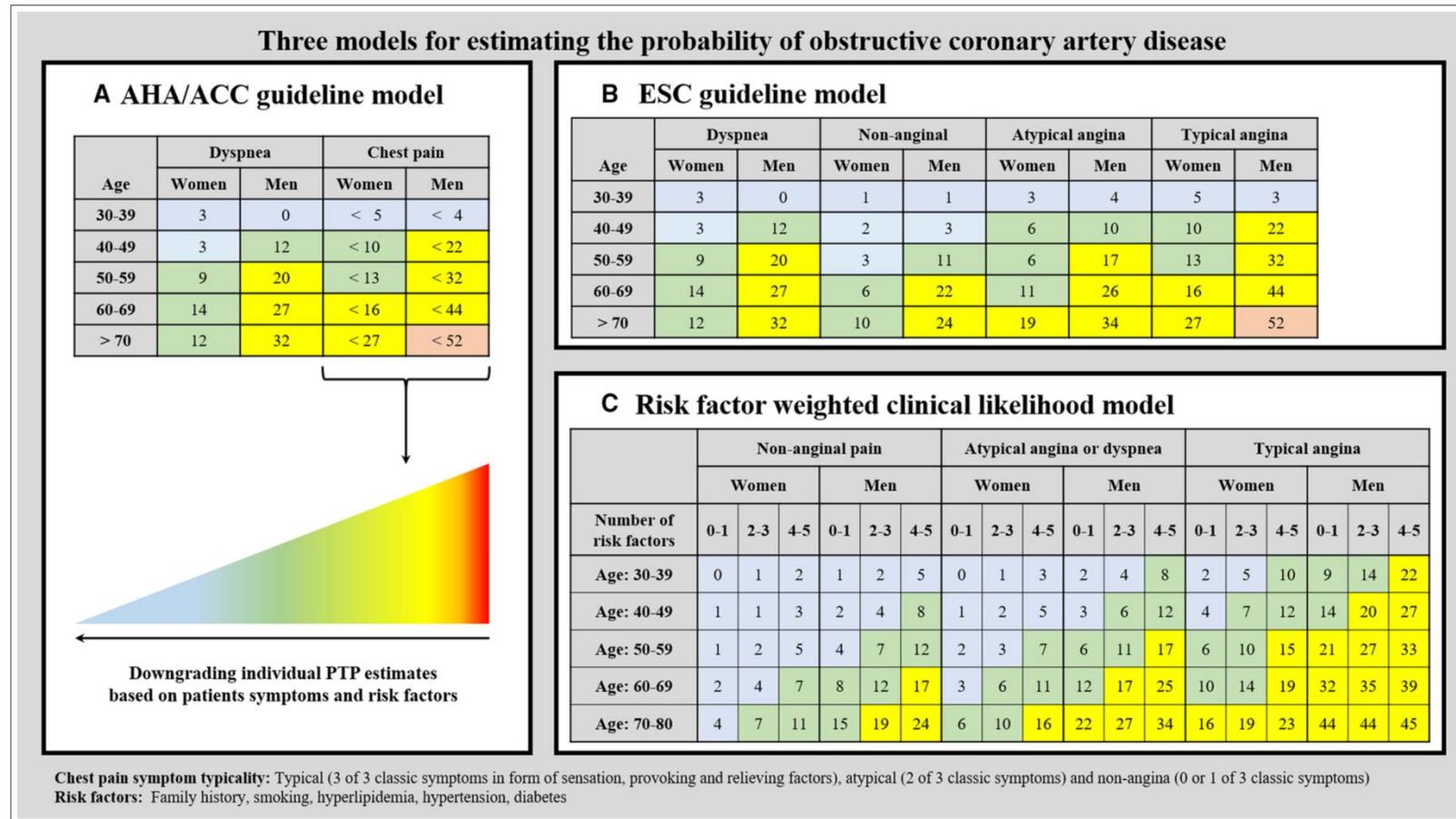
High GPS definition	Individuals in testing dataset (n)	% of individuals
<b>Odds ratio <math>\geq 3.0</math></b>		
CAD	23,119/288,978	8.0

# STRATIFIED USE OF PGS

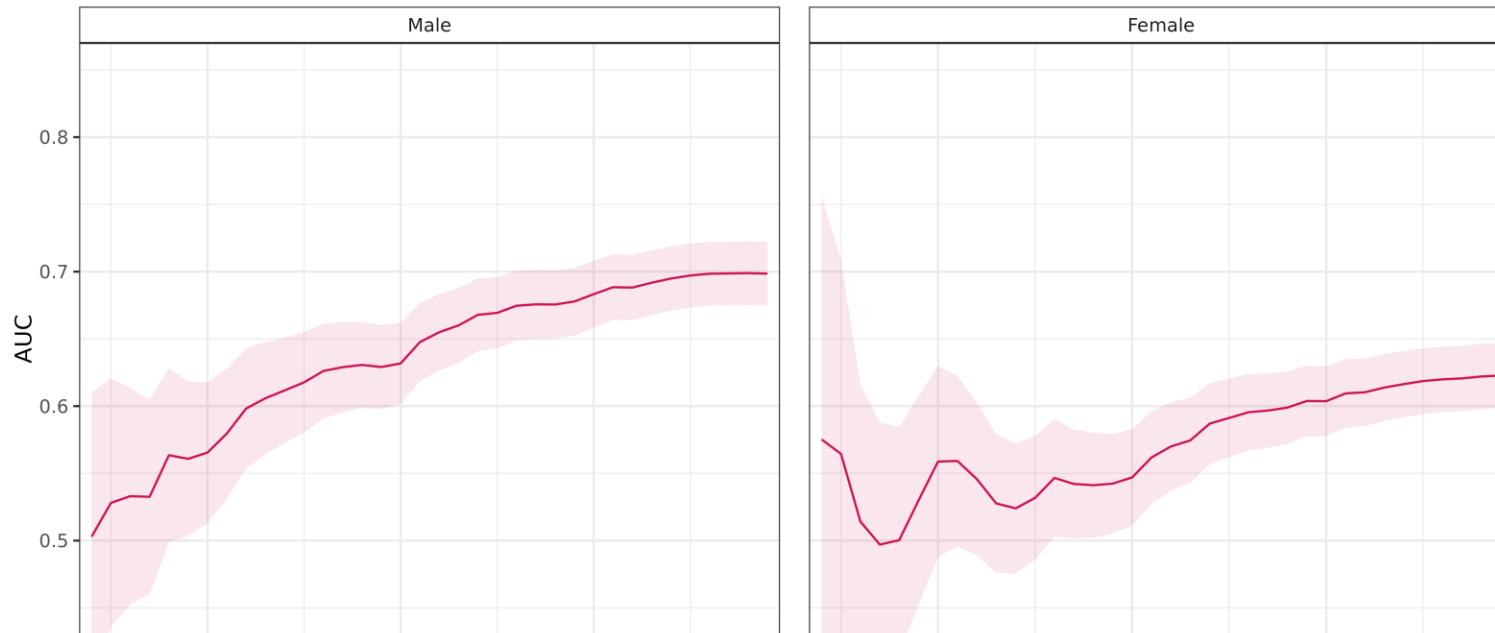
Group	LDL-C (mg/dL)	HR (95%CI)	P value	Total Inds / Cases
Low PRS	<100	1	[Reference]	1142 / 7
	100 – <130	0.79 (0.34–1.84)	0.58	4194 / 26
	130 – <160	0.81 (0.35–1.87)	0.62	4570 / 38
	160 – <190	0.79 (0.32–1.93)	0.6	2264 / 21
	≥ 190	0.86 (0.3–2.71)	0.86	701 / 7
Int PRS	<100	1	[Reference]	6824 / 61
	100 – <130	0.94 (0.71–1.24)	0.66	28352 / 283
	130 – <160	1.25 (0.95–1.63)	0.11	38352 / 615
	160 – <190	1.62 (1.23–2.13)	<0.001	21576 / 483
	≥ 190	2.34 (1.75–3.14)	<0.001	7867 / 246
High PRS	<100	1	[Reference]	624 / 8
	100 – <130	1.15 (0.54–2.46)	0.72	3082 / 44
	130 – <160	2.23 (1.08–4.59)	0.03	4808 / 151
	160 – <190	3.14 (1.52–6.5)	0.002	3091 / 144
	≥ 190	4.71 (2.23–9.94)	<0.001	1267 / 81



# RISK ASSESSMENT FOR OBSTRUCTIVE CAD

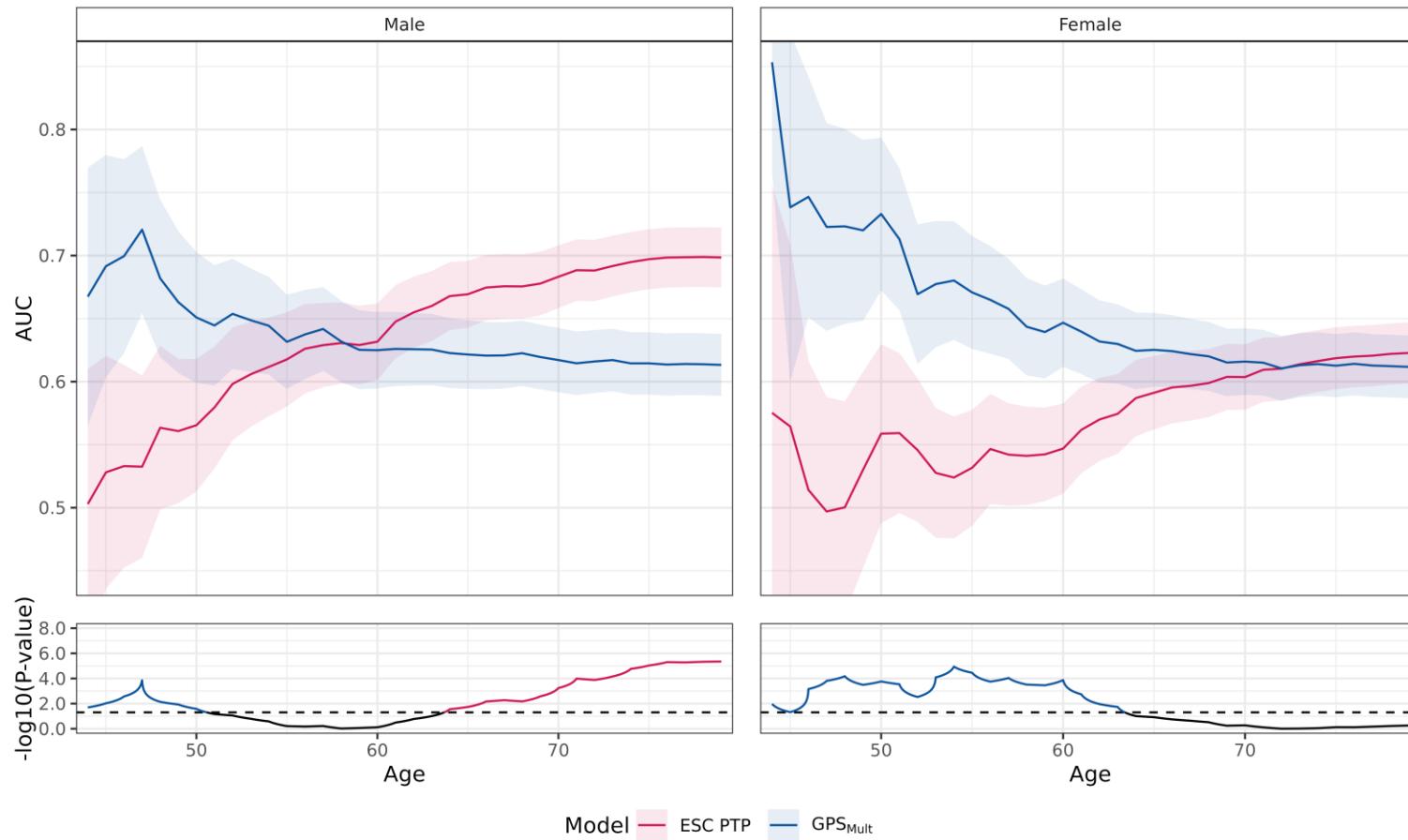


# STRATIFIED USE OF PGS

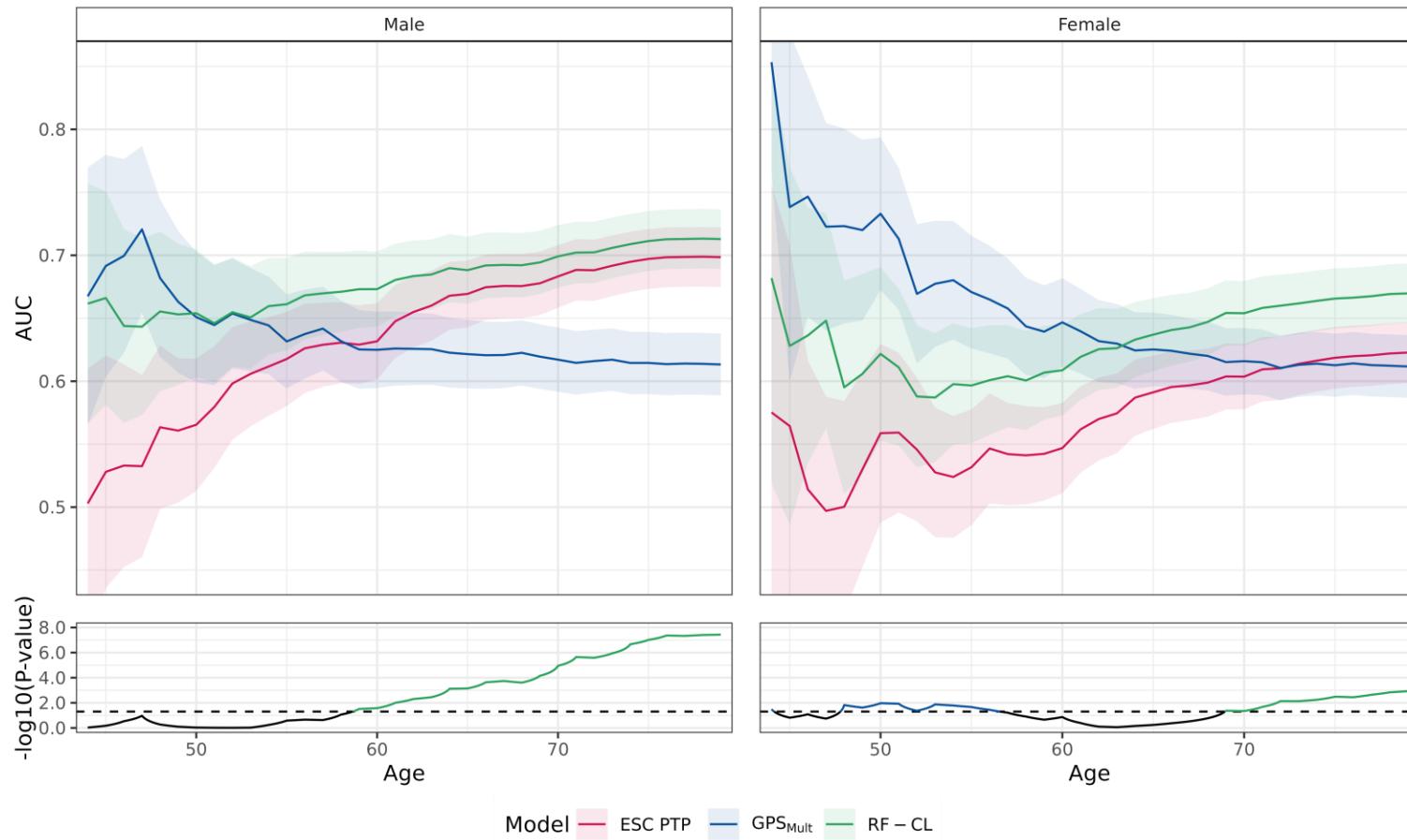


Model — ESC PTP

# STRATIFIED USE OF PGS



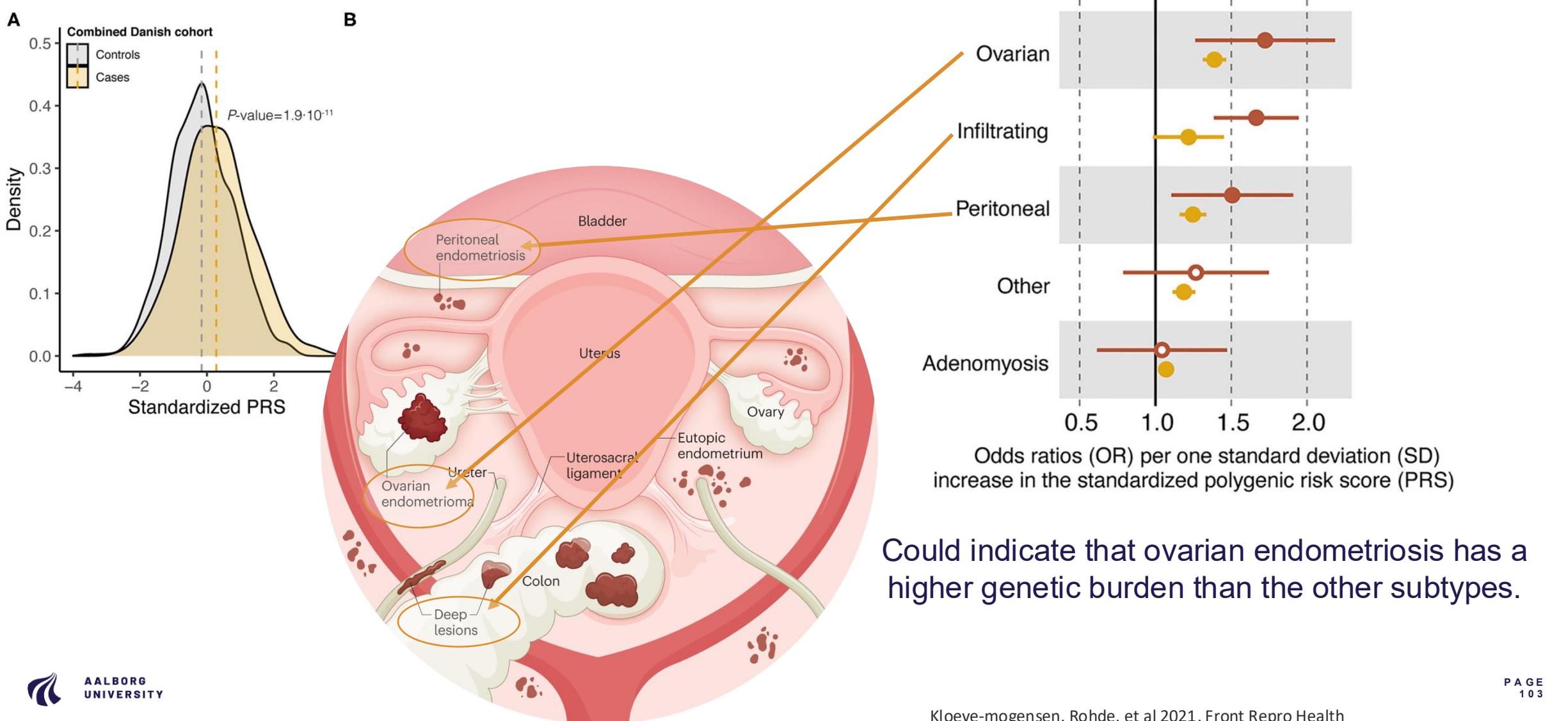
# STRATIFIED USE OF PGS



# STRATIFIED USE OF PGS

	<b>Group</b>	<b>Sample size</b>	<b>Cases</b>	<b>Referrals</b>	<b>Correct referrals</b>	<b>Cases referred</b>	<b>Cases missed</b>
ESC PTP	Women < 55	282 (13%)	27 (10%)	0 (0%)	0 (0%)	0%	27 (15%)
	Women ≥ 55	689 (33%)	146 (21%)	156 (23%)	40 (26%)	27%	106 (59%)
	Men < 55	395 (19%)	73 (18%)	189 (48%)	42 (22%)	58%	31 (17%)
	Men ≥ 55	733 (35%)	298 (41%)	661 (90%)	283 (43%)	95%	15 (8%)
	Total	2099 (100%)	544 (26%)	1006 (48%)	365 (36%)	67%	179 (100%)
ESC PTP + PGS	Women < 55	282 (13%)	27 (10%)	10 (4%)	2 (20%)	7%	25 (15%)
	Women ≥ 55	689 (33%)	146 (21%)	153 (22%)	48 (31%)	33%	98 (58%)
	Men < 55	395 (19%)	73 (18%)	179 (45%)	47 (26%)	64%	26 (15%)
	Men ≥ 55	733 (35%)	298 (41%)	620 (85%)	279 (45%)	94%	19 (11%)
	Total	2099 (100%)	544 (26%)	962 (46%)	376 (39%)	69%	168 (100%)

# LOOK INTO THE BIOLOG



# SESSION 3

- How to measure 'accuracy'?
- Interpretability and risk communication
- Lack of transferability
- Applications...?





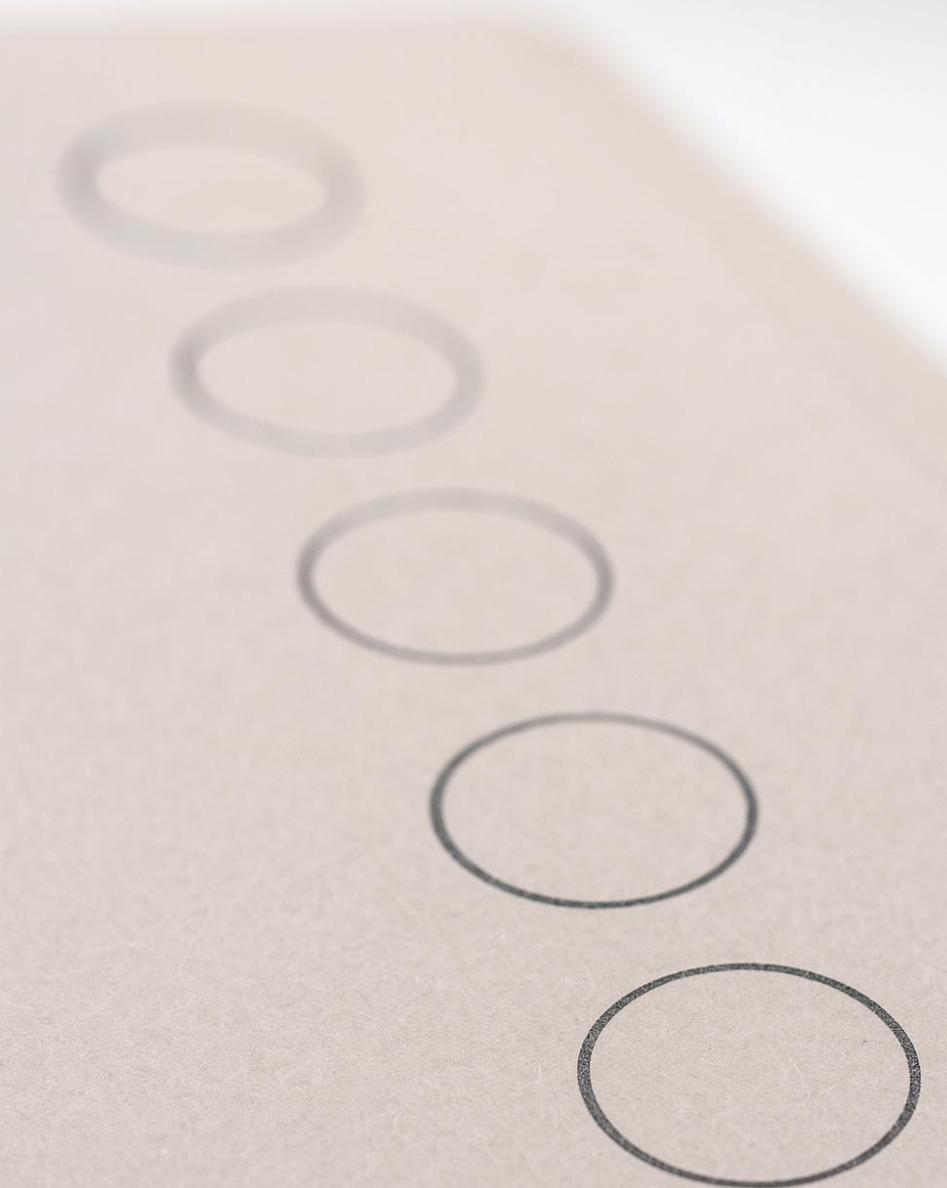
**BREAK**

# AGENDA

08:00 – 08:30	Welcome and common introductions
08:30 – 09:10	Session 1: Introduction to Polygenic Scores (PGS)
09:10 – 09:20	Break
09:20 – 10:00	Session 2: Data Sources and Computational Methods
10:00 – 10:10	Break
10:10 – 10:40	Session 3: Evaluating and Interpreting Polygenic Scores
10:40 – 11:00	Break
<b>11:00 – 11:45</b>	<b>Session 4: Advanced Applications and Future Directions</b>
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today

# SESSION 4

- Many challenges – how to circumvent these?
- Enhancing biological understanding?
- PGS in combination with proteomics



# CHALLENGES WITH PGS

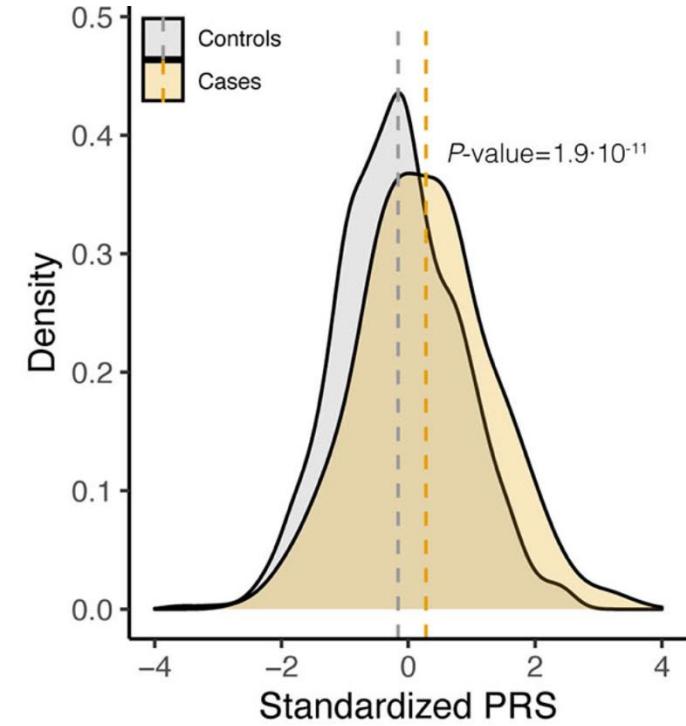
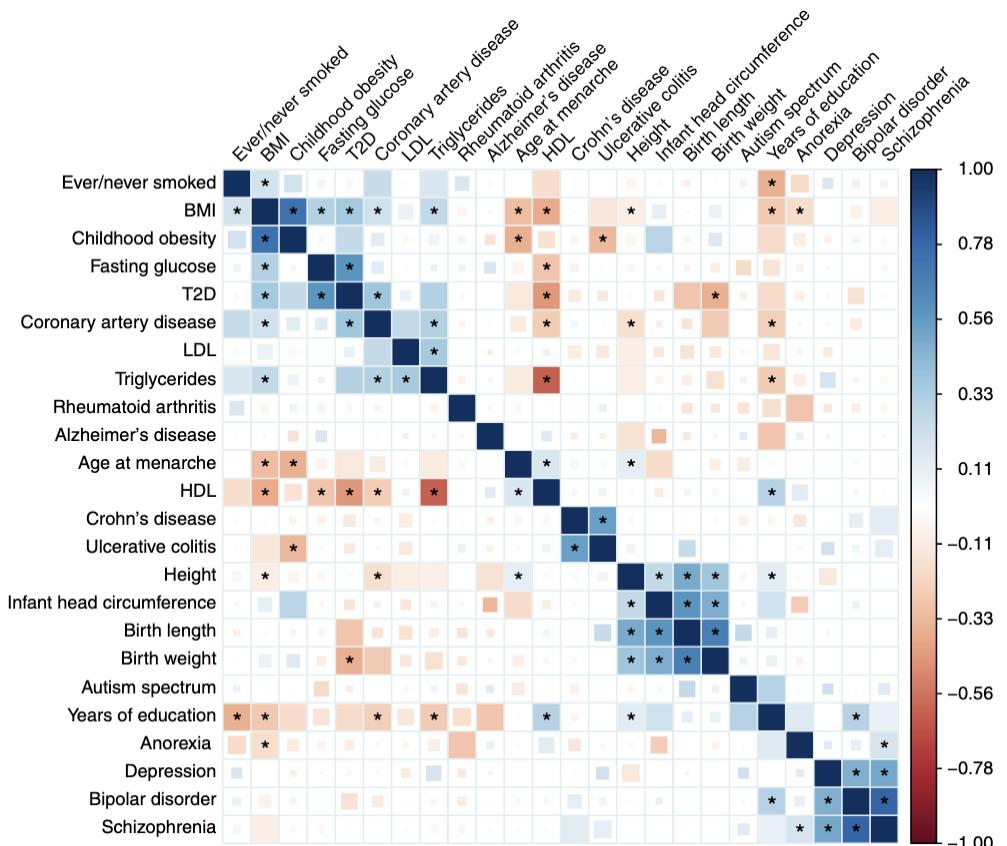
- 1) Too many individuals have 'average' PGS
- 2) Only capture additive effects
- 3) Only common variants
- 4) Transferability
- 5) Understanding biology

## Improving PGS and Future directions

- Improve discriminative ability using multi-trait models and functional annotation
- Ancestry-aware models
- Integration with single cell and spatial transcriptomics
- Integrative genomics (combining PGS with other 'omics)

# IMPROVING PGS

## BETTER DISCRIMINATIVE ABILITY

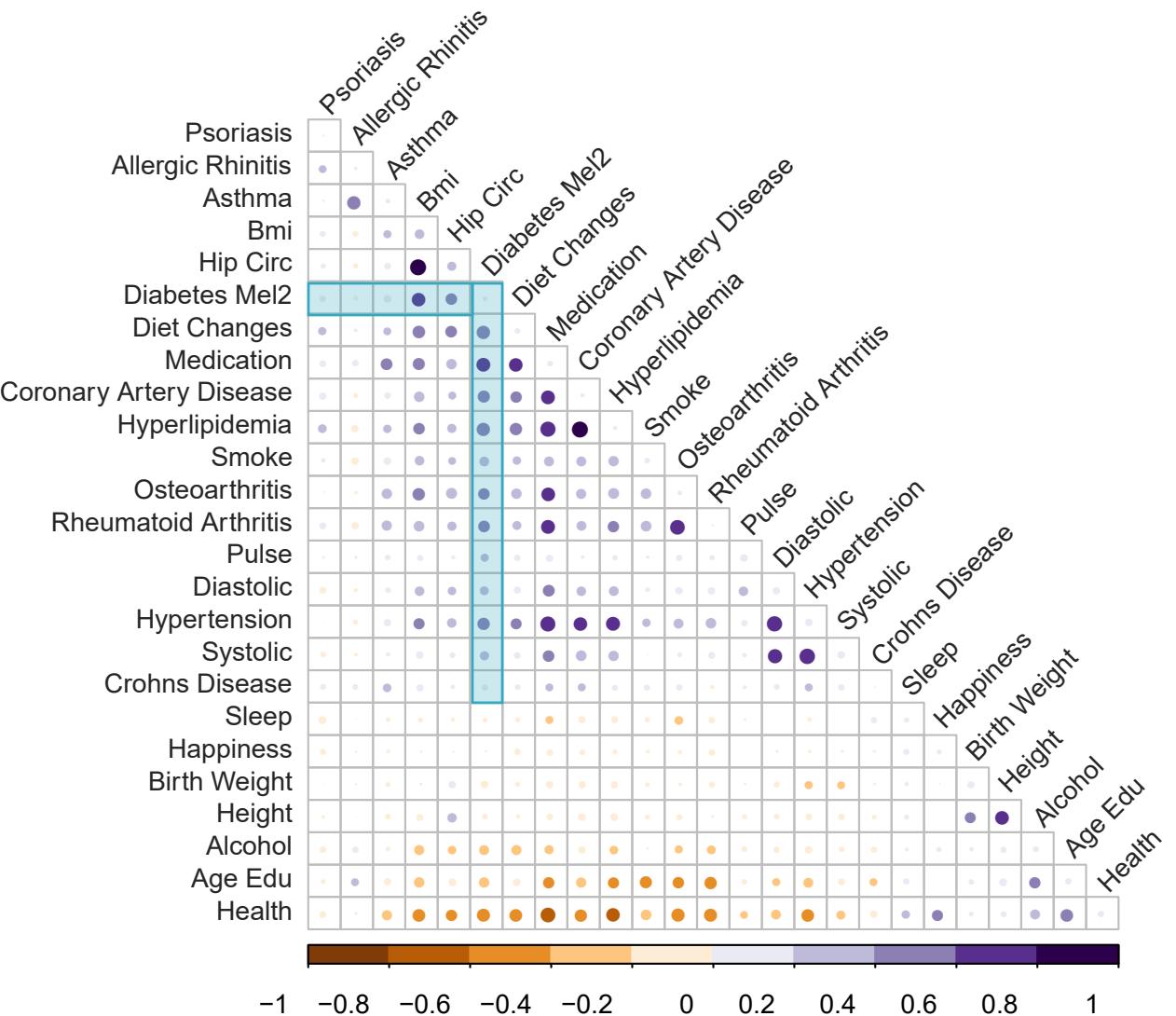


Identifying genetic correlations between complex traits and diseases can provide useful etiological insights and help prioritize likely causal relationships

# MULTI-TRAIT PGS FOR TYPE 2 DIABETES

T2D is strongly correlated with a range of complex diseases and traits, such as overweight, cardiovascular diseases, hypertension, chronic kidney disease etc.

Can we improve prediction accuracy by leveraging shared genetic signatures?



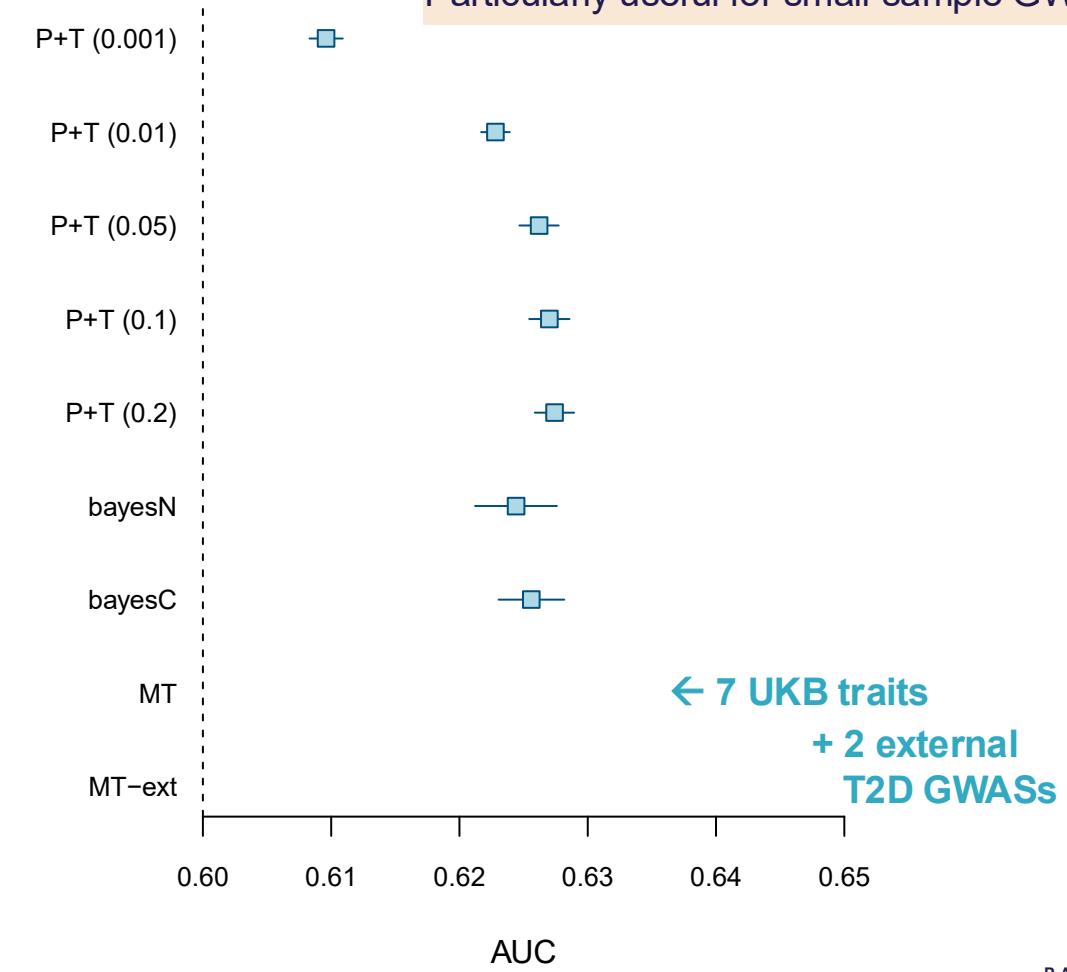
# MULTI-TRAIT PGS FOR TYPE 2 DIABETES

$$PGS = \sum X_i \hat{\beta}_i$$

Each marker effect ( $\hat{\beta}$ ) is weighted by  $w$ , which is found by using selection index theory

$$w = \begin{bmatrix} \frac{h_1^2}{M} + \frac{1}{N_1} & \dots & \frac{r_g h_1 h_k}{M} \\ \vdots & \ddots & \vdots \\ \frac{r_g h_k h_1}{M} & \dots & \frac{h_k^2}{M} + \frac{1}{N_k} \end{bmatrix}^{-1} \begin{bmatrix} \frac{h_1^2}{M} \\ \vdots \\ \frac{r_g h_1 h_k}{M} \end{bmatrix}$$

$$MT-PGS = \sum_{i=1}^m X_i \hat{\beta}_{wMTi}$$

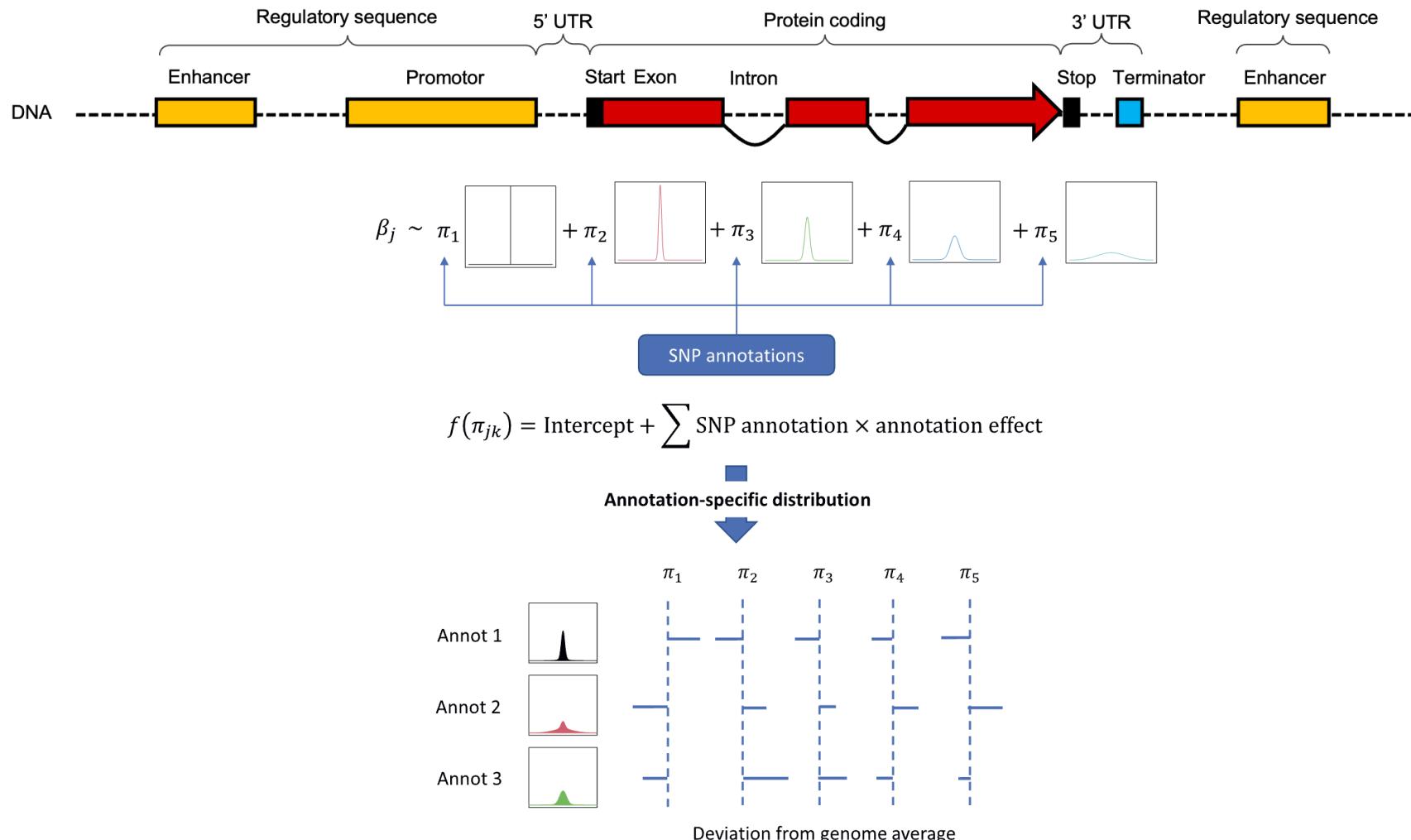


The PGS accuracy depends on the power of the GWAS (=sample size).

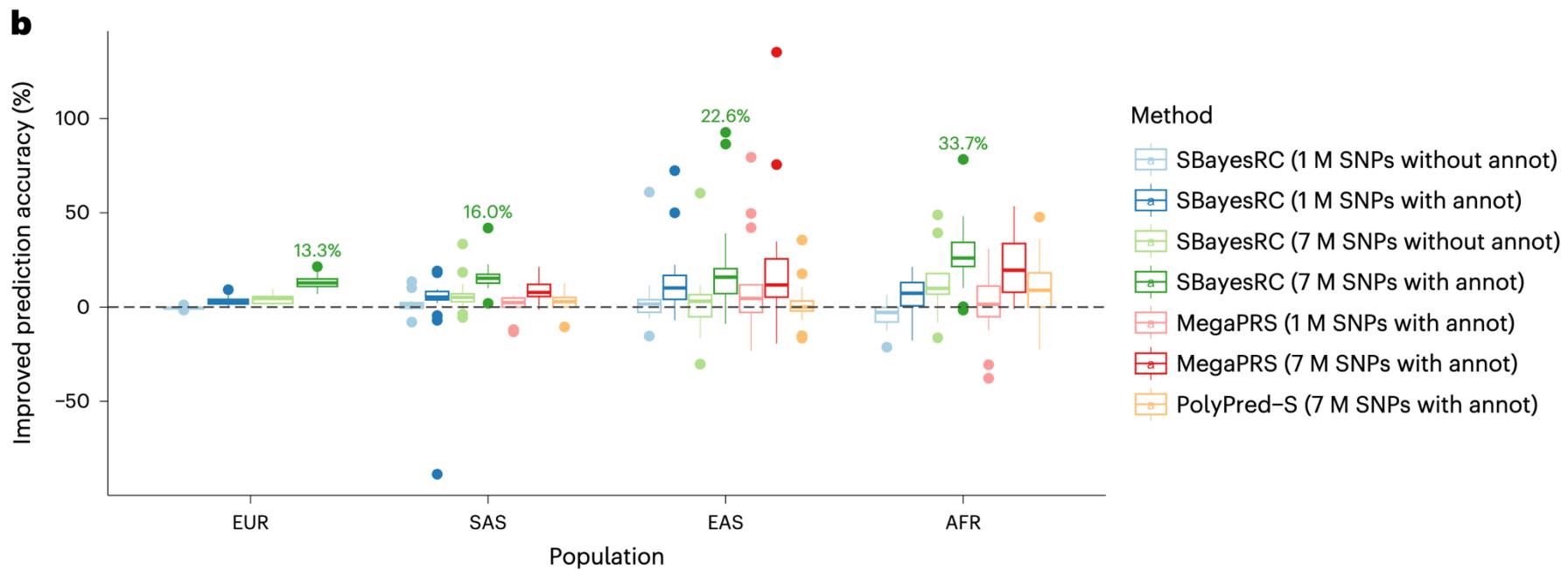
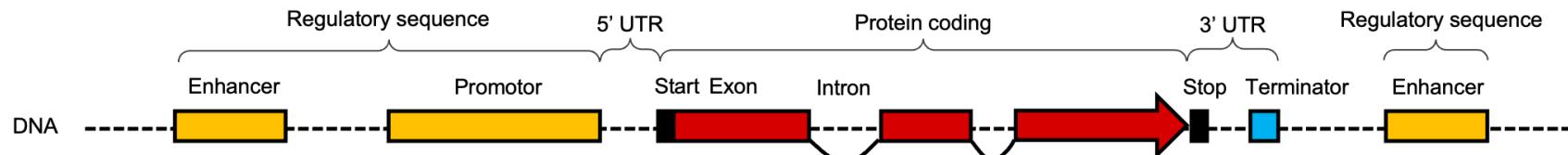
Utilising the genetic correlation among phenotypes is equivalent to increasing GWAS sample size.

Particularly useful for small-sample GWAS.

# ENHANCING PGS BY ANNOTATION



# ENHANCING PGS BY ANNOTATION

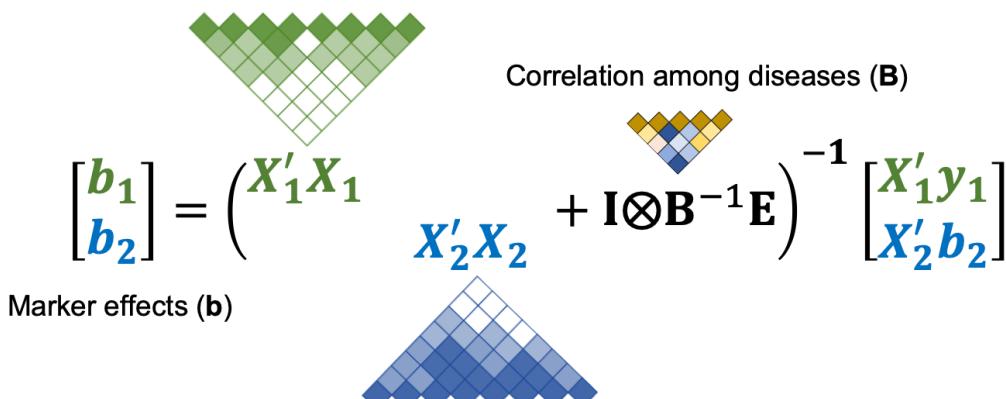
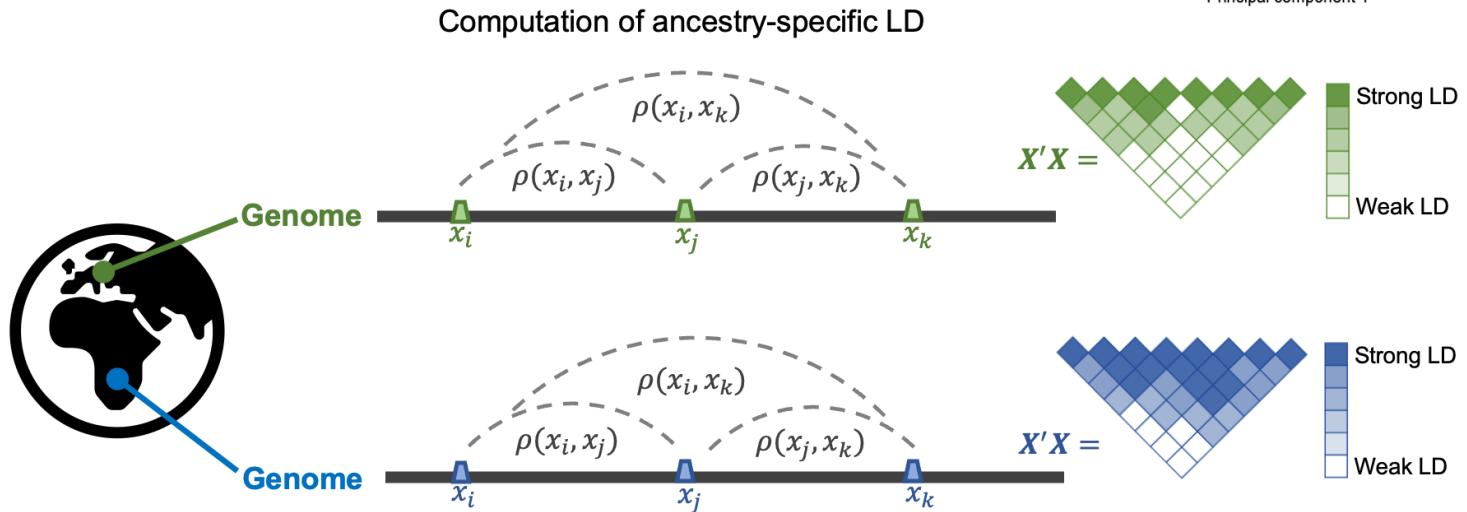


# CROSS-ANCESTRY PGS

$$b = \left( X'X + I \frac{\sigma_e^2}{\sigma_b^2} \right)^{-1} X'y$$

$$X'X = D^{0.5} BD^{0.5}$$

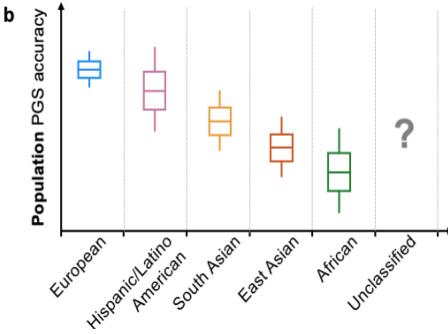
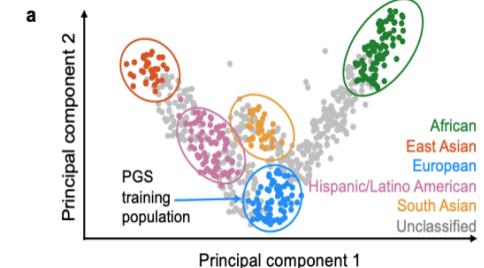
$$X'y = Db_m$$



MT-BLR

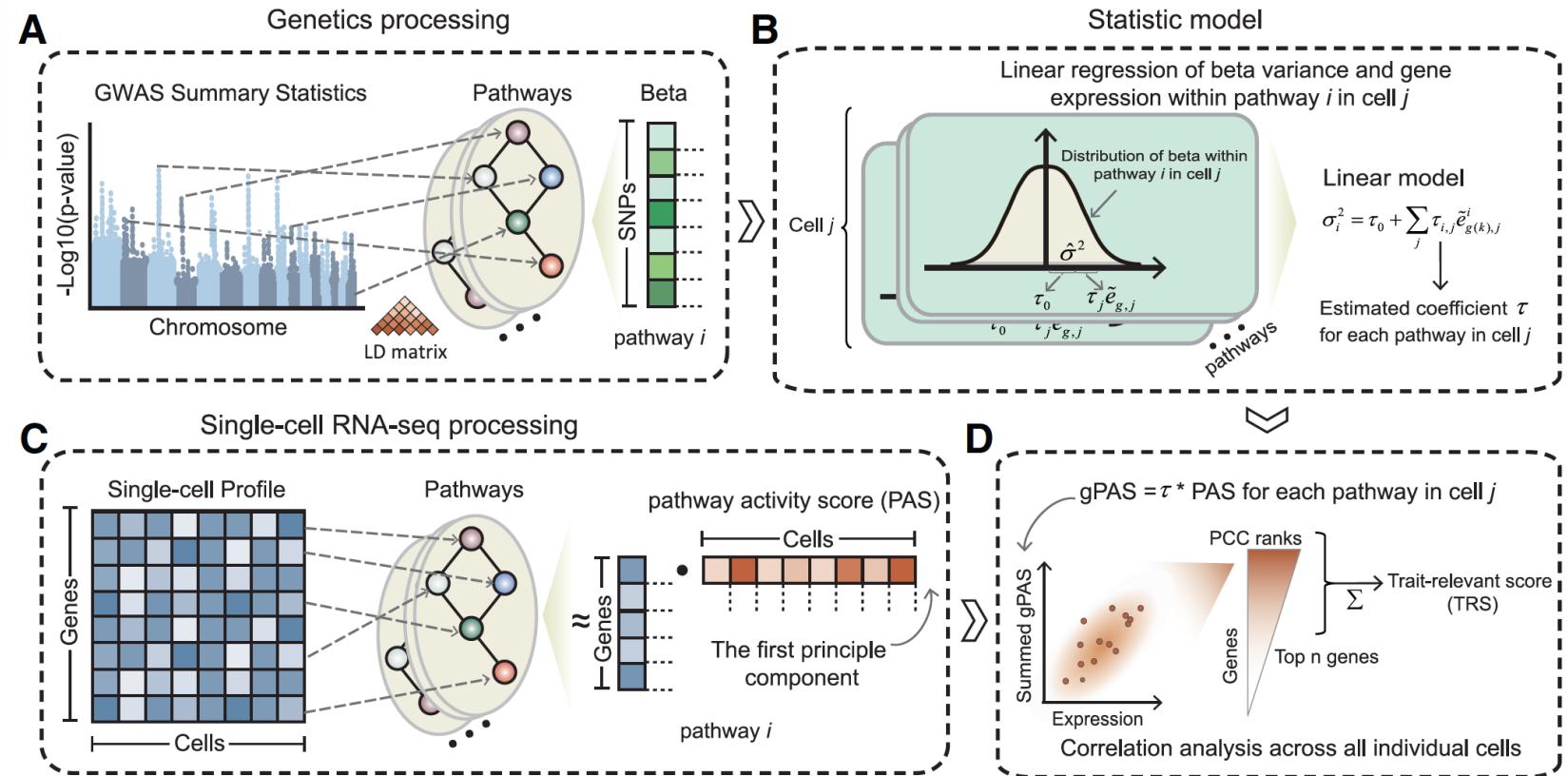
$$\begin{aligned} PGS_1 &= \sum X_1 b_1 \\ PGS_2 &= \sum X_2 b_2 \end{aligned}$$

Polygenic scores (PGS) combining ancestry-aware genetic risk factors.

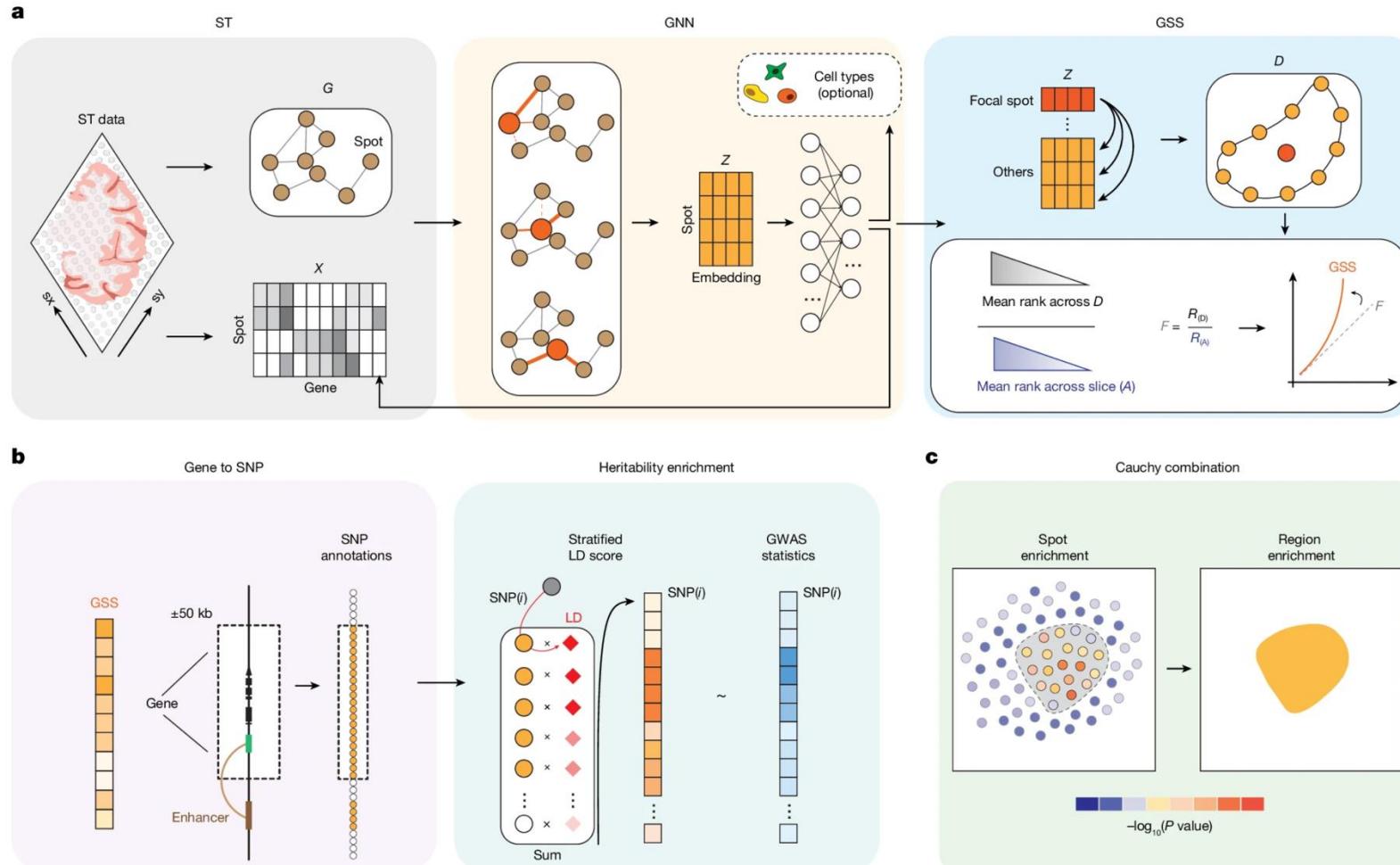


# UTILISING PGS WITH SINGLE CELL DATA

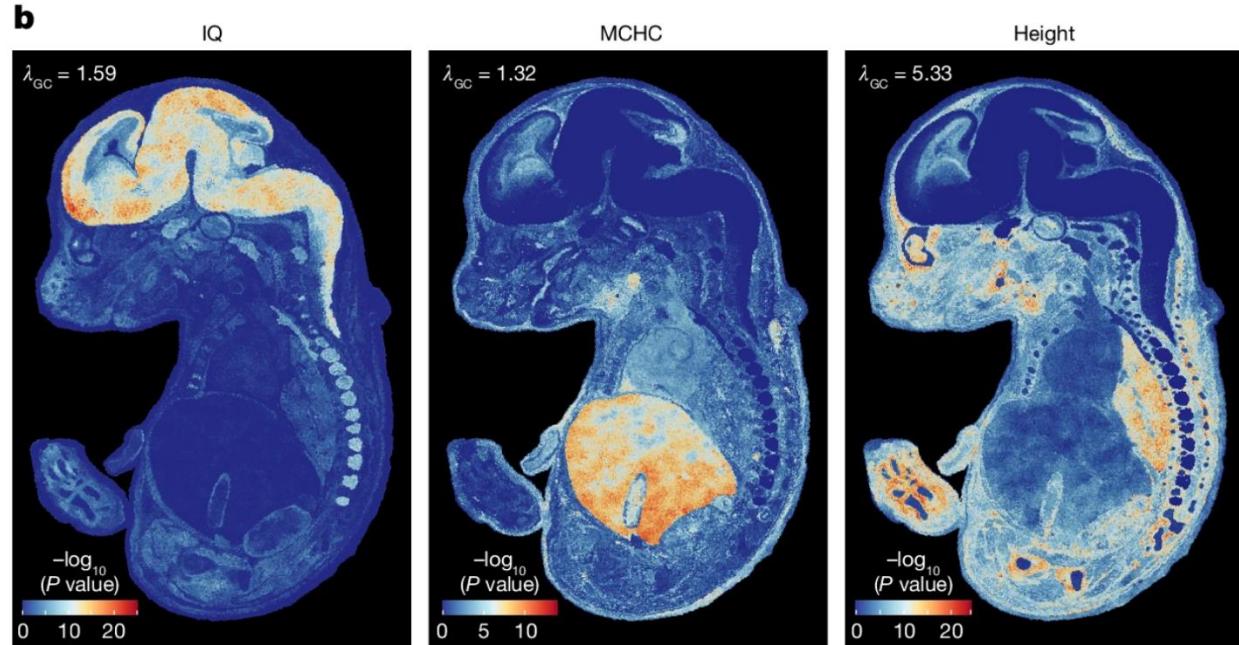
Integrate scRNA-seq data with GWAS data to discover cellular context for complex diseases



# UTILISING PGS WITH SPATIAL TRANSCRIPTOMICS



# UTILISING PGS WITH SPATIAL TRANSCRIPTOMICS

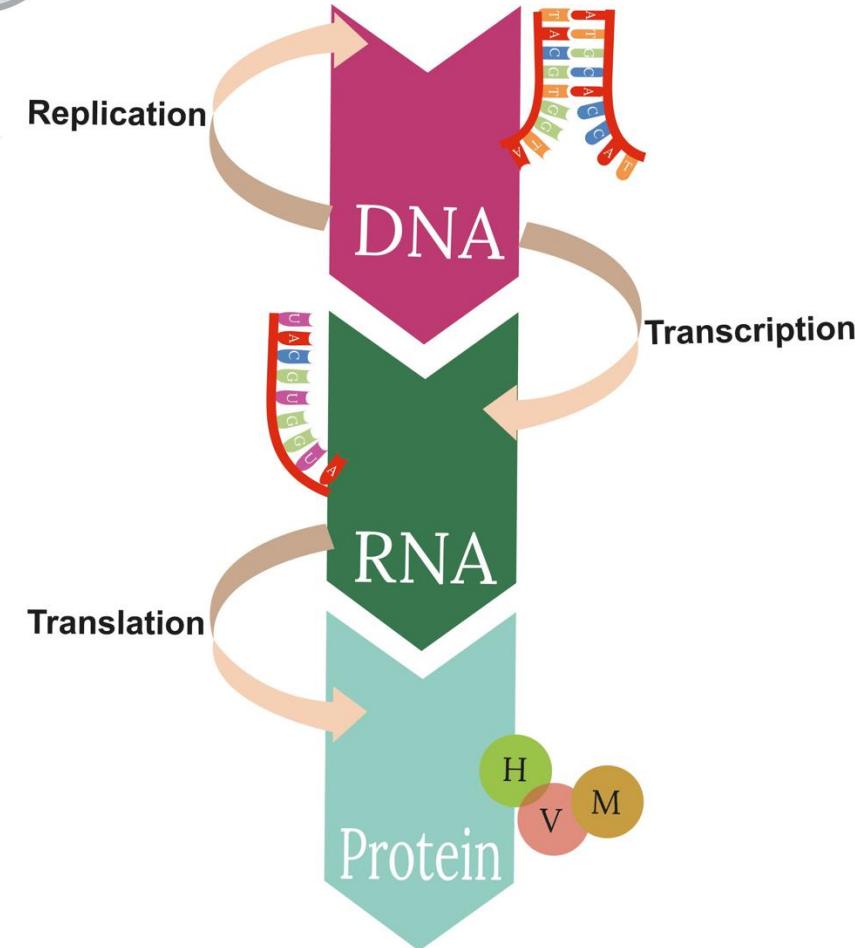
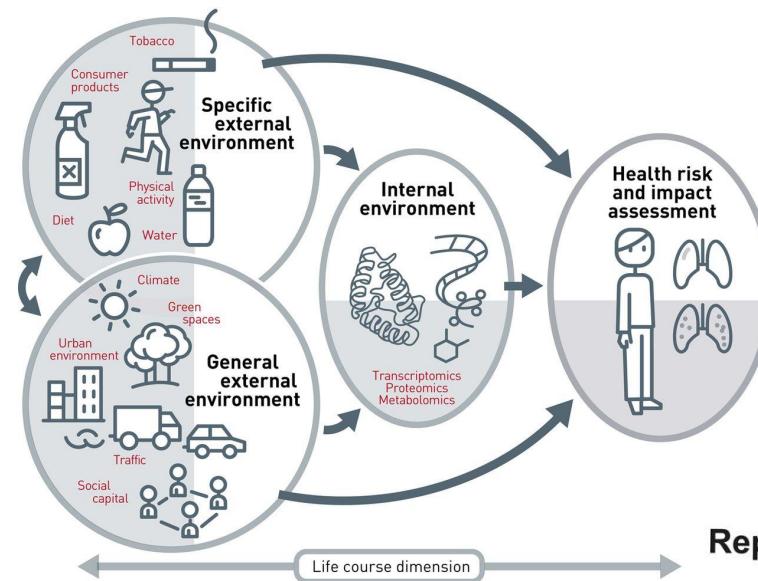


# Multifactorial traits have a **genetic component** and an **environmental component**

The **genetic profile is stable** throughout life – thus, a polygenic score may inform about our inherent disease risk.

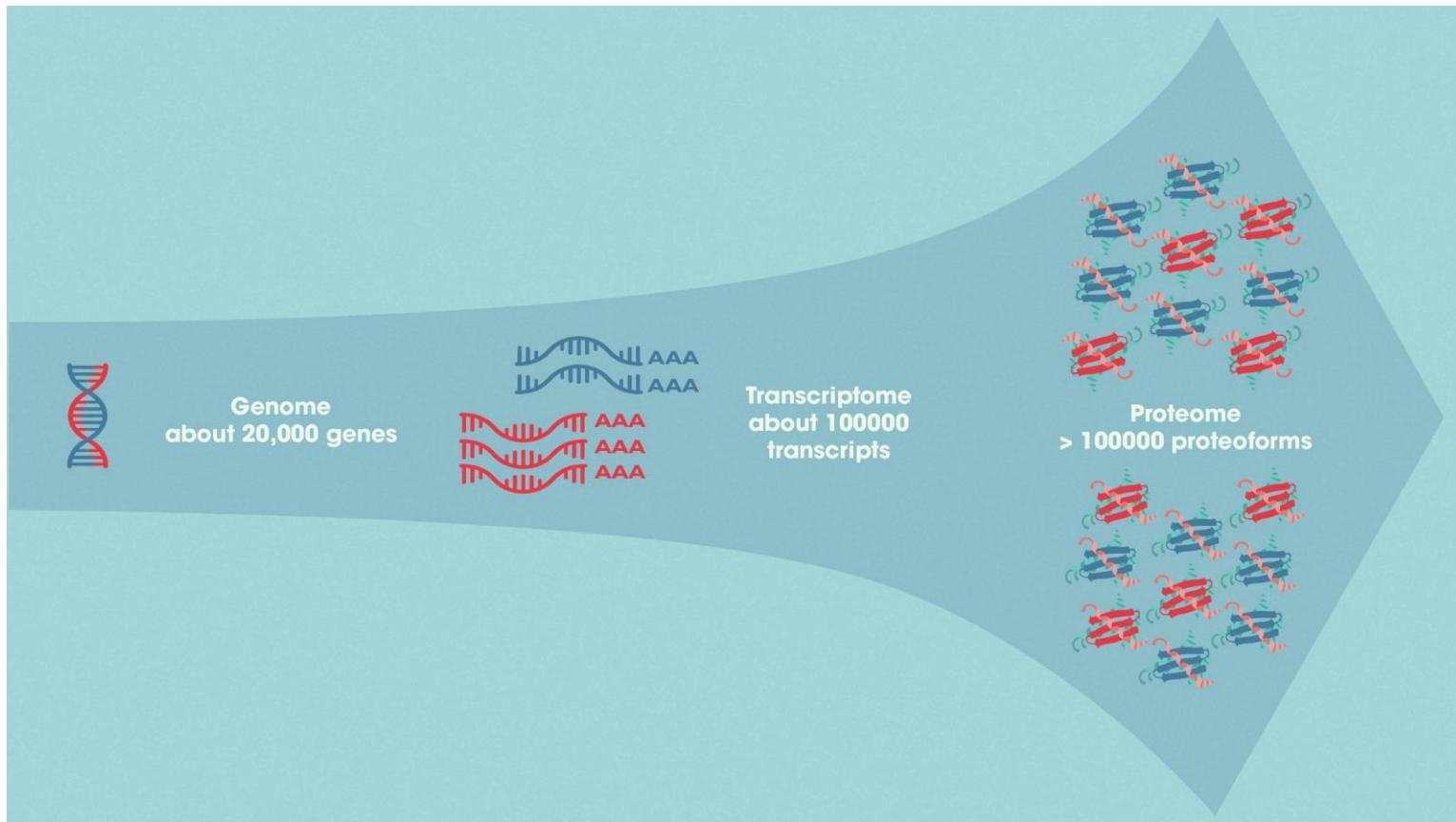
Environmental exposures affect **transcription and translation** **during life**, which hence might affect disease risk

**Therefore, other omic technologies** might inform about disease pathology and progression.



# GENOME VS PROTEOME

- Unlike the genome, the composition of the proteome is in a constant state of flux over time and throughout the organism
- Environmental exposures affect transcription of DNA, and translation of RNA
- Proteome may help us understand disease pathology



# Predicting presence of coronary plaques featuring high-risk characteristics using polygenic risk scores and targeted proteomics in patients with suspected coronary artery disease

Møller, P.L., Rohde, P. D., et al, Circulation: Genomics and Precision medicine, 2023  
Møller, P.L., Rohde, P. D., et al, Genome Medicine 2023



# PATIENT WITH CHEST PAIN



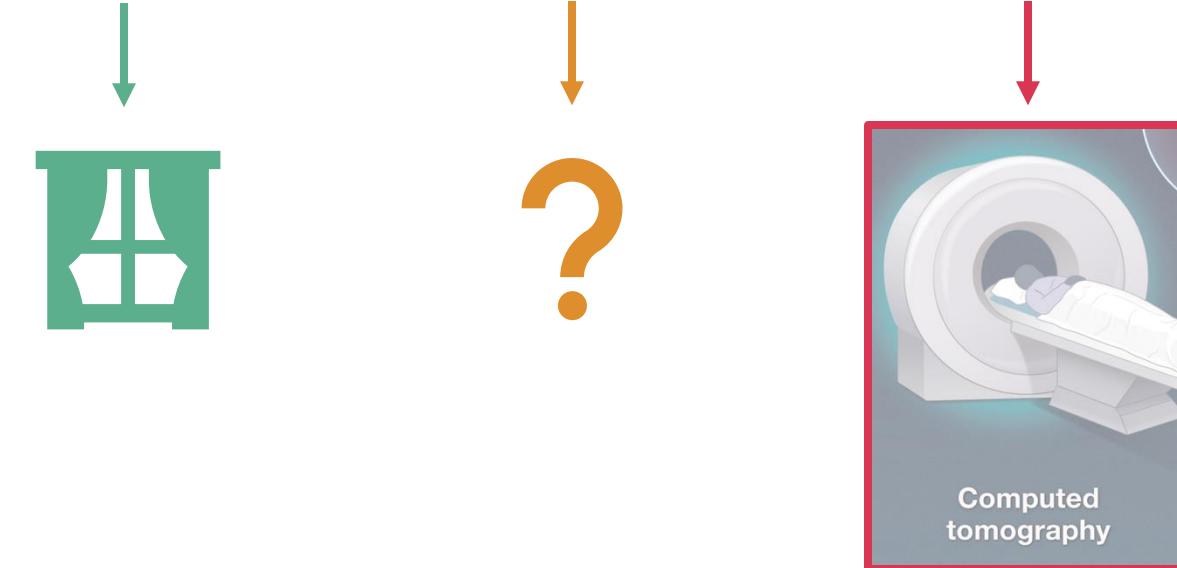
Patient with de-novo chest pain



Pretest probability (PTP) = gender, age and type of chest pain

Risk of obstructive **coronary artery disease** (CAD)

<5%                    5-15%                    >15%



# Dan-NICAD COHORT

Morten  
Böttcher



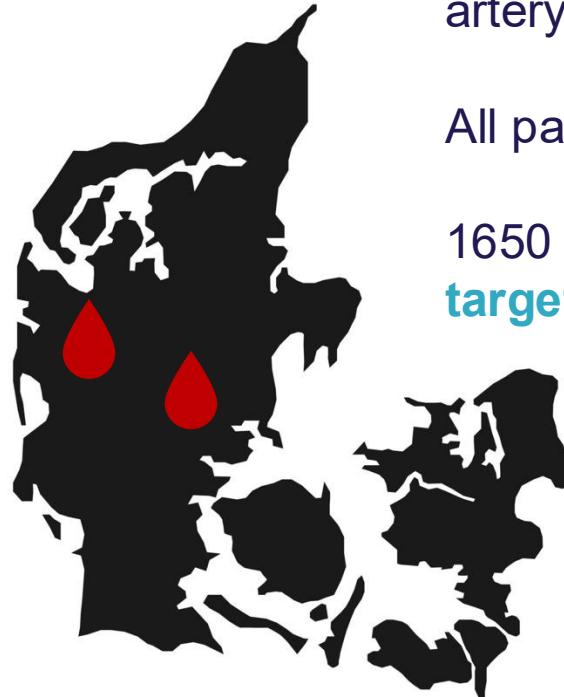
Simon  
Winther



~4000 patients with symptoms of obstructive coronary artery disease (CAD)

All patients undergo coronary CT

1650 patients have **plaque morphology**, **DNA** and **targeted proteomics** [368 proteins, Olink®]

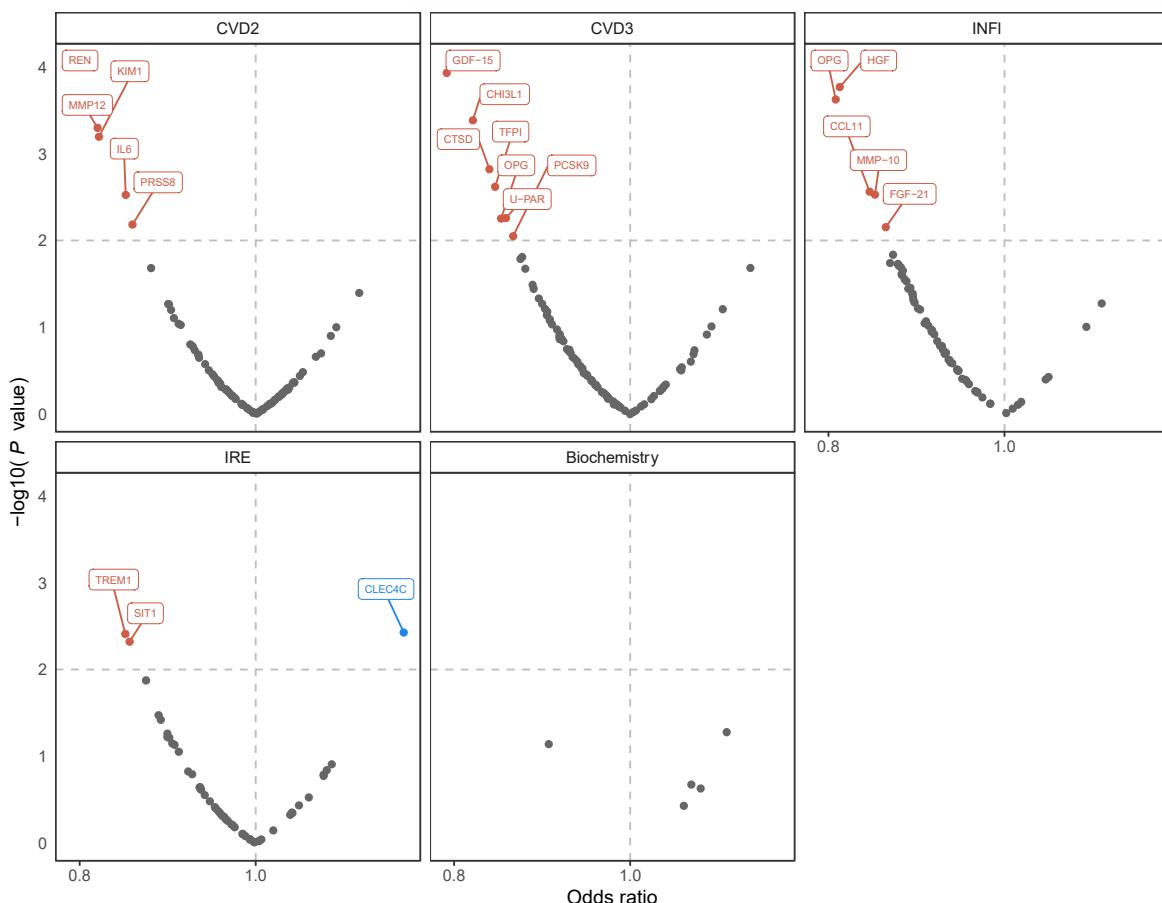
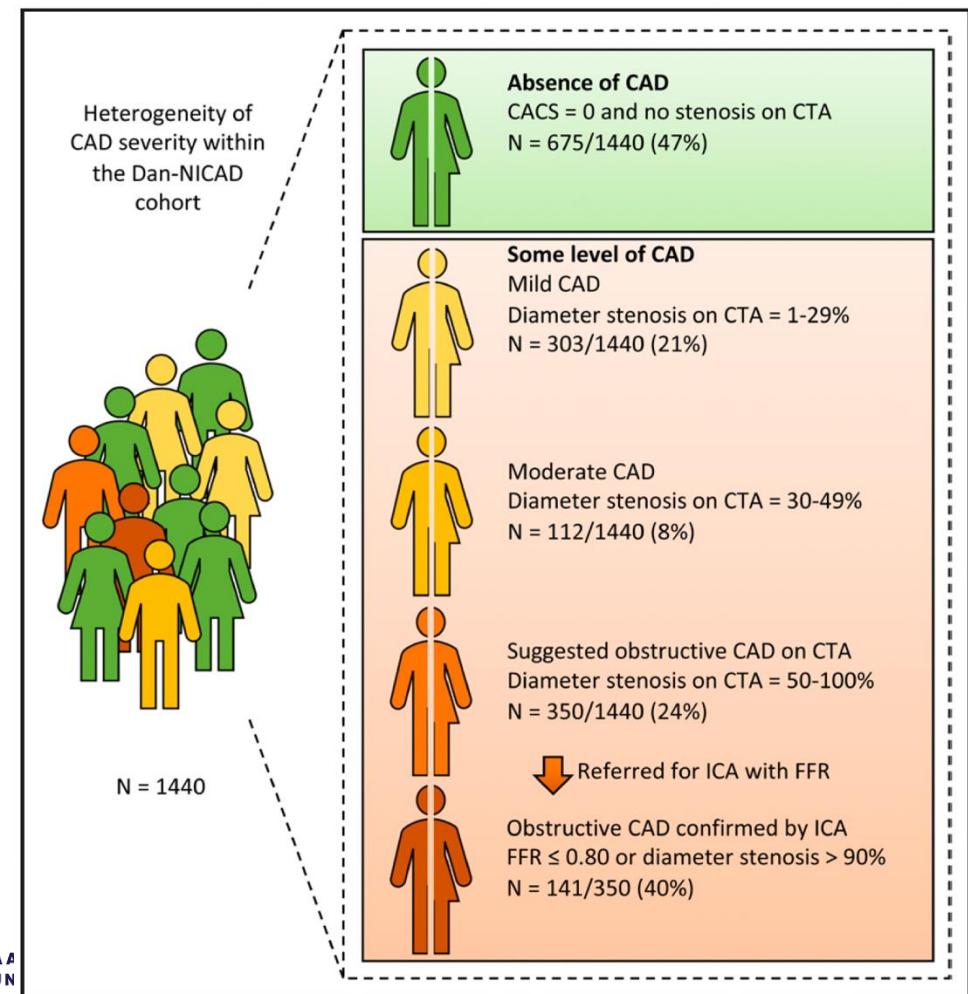


## ORIGINAL ARTICLE

# Combining Polygenic and Proteomic Risk Scores With Clinical Risk Factors to Improve Performance for Diagnosing Absence of Coronary Artery Disease in Patients With de novo Chest Pain

Peter Loof Møller<sup>1,2</sup>, MSc; Palle Duun Rohde<sup>3</sup>, MSc, PhD; Jonathan Nørtoft Dahl<sup>1,2</sup>, MD; Laust Dupont Rasmussen<sup>4</sup>, MD, PhD; Samuel Emil Schmidt<sup>2</sup>, MSc, PhD; Louise Nissen, MD, PhD; Victoria McGilligan<sup>2</sup>, PhD; Jacob F. Bentzon<sup>1,2</sup>, MD, PhD; Daniel F. Gudbjartsson, MSc, PhD; Kari Stefansson<sup>1,2</sup>, MD, PhD; Hilma Holm<sup>5</sup>, MD; Simon Winther, MD, PhD; Morten Böttcher<sup>2</sup>, MD, PhD; Mette Nyegaard<sup>1,2</sup>, MSc, PhD

# CAN WE IDENTIFY *HEALTHY* PATIENTS

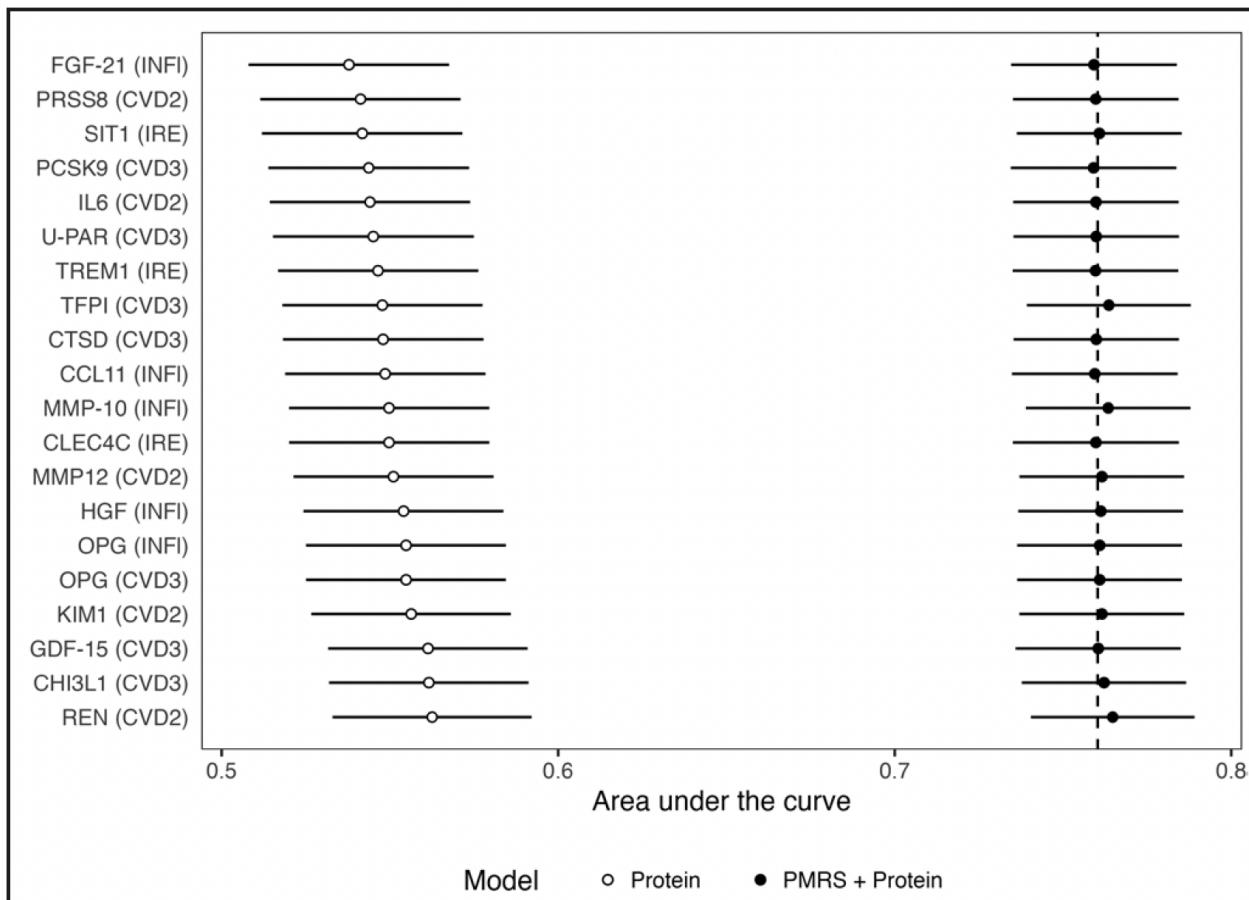
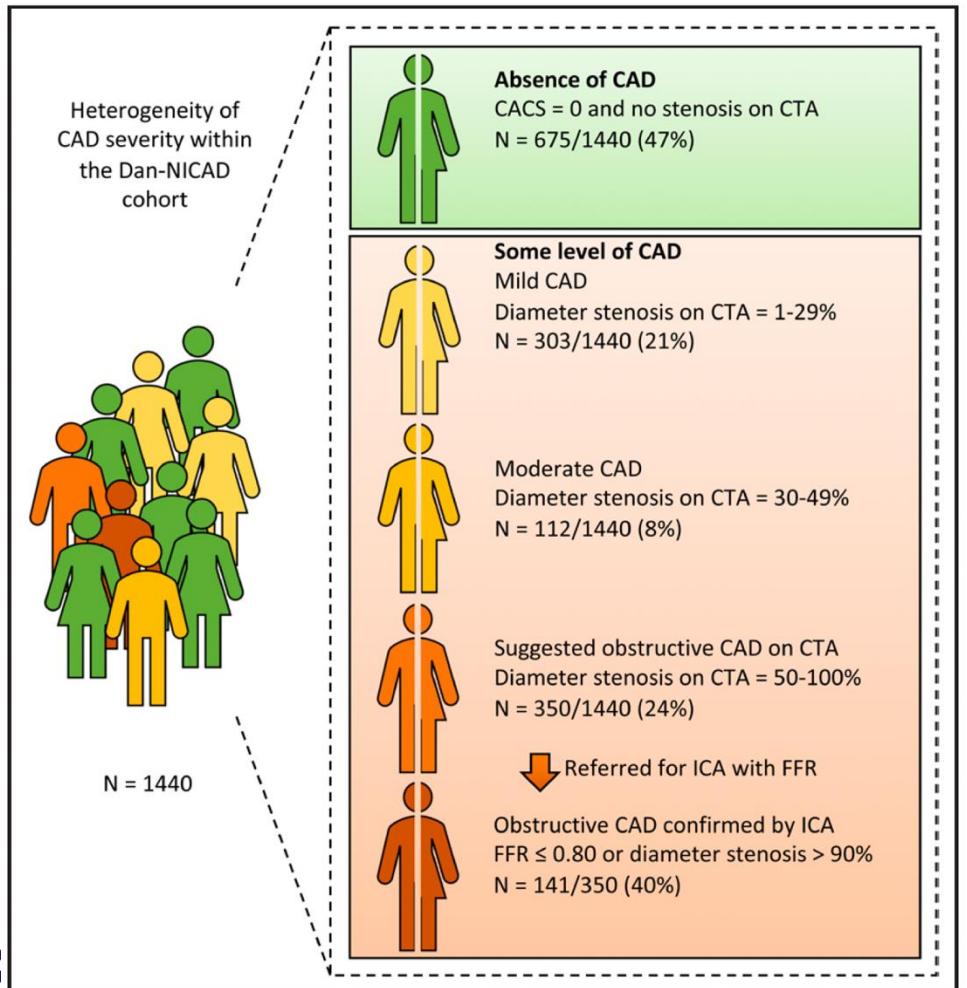


## ORIGINAL ARTICLE

# Combining Polygenic and Proteomic Risk Scores With Clinical Risk Factors to Improve Performance for Diagnosing Absence of Coronary Artery Disease in Patients With de novo Chest Pain

Peter Loof Møller<sup>1,2</sup>, MSc; Palle Duun Rohde<sup>3</sup>, MSc, PhD; Jonathan Nørtoft Dahl<sup>1,2</sup>, MD; Laust Dupont Rasmussen<sup>4</sup>, MD, PhD; Samuel Emil Schmidt<sup>1,2</sup>, MSc, PhD; Louise Nissen, MD, PhD; Victoria McGilligan<sup>5</sup>, PhD; Jacob F. Bentzon<sup>1,2</sup>, MD, PhD; Daniel F. Gudbjartsson, MSc, PhD; Kari Stefansson<sup>1,2</sup>, MD, PhD; Hilma Holm<sup>6</sup>, MD; Simon Winther, MD, PhD; Morten Böttcher<sup>2</sup>, MD, PhD; Mette Nyegaard<sup>1,2</sup>, MSc, PhD

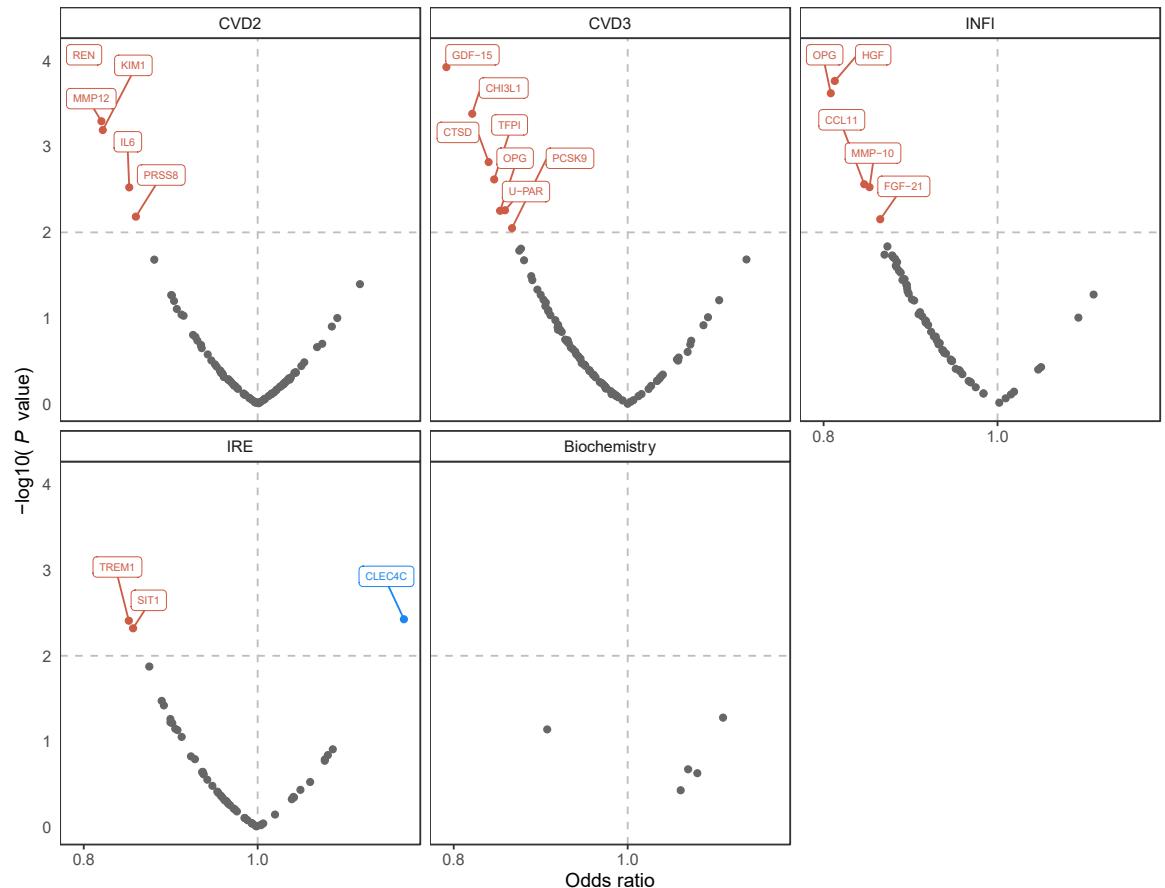
# CAN WE IDENTIFY *HEALTHY* PATIENTS



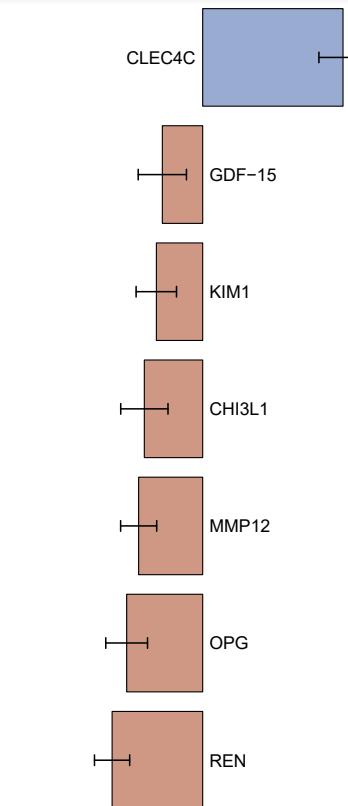
## ORIGINAL ARTICLE

# Combining Polygenic and Proteomic Risk Scores With Clinical Risk Factors to Improve Performance for Diagnosing Absence of Coronary Artery Disease in Patients With de novo Chest Pain

Peter Loof Møller<sup>1,2</sup>, MSc; Palle Duun Rohde<sup>3</sup>, MSc, PhD; Jonathan Nørtoft Dahl<sup>1,2</sup>, MD; Laust Dupont Rasmussen<sup>4</sup>, MD, PhD; Samuel Emil Schmidt<sup>1,2</sup>, MSc, PhD; Louise Nissen, MD, PhD; Victoria McGilligan<sup>5</sup>, PhD; Jacob F. Bentzon<sup>1,2</sup>, MD, PhD; Daniel F. Gudbjartsson, MSc, PhD; Kari Stefansson<sup>1,2</sup>, MD, PhD; Hilma Holm<sup>6</sup>, MD; Simon Winther, MD, PhD; Morten Böttcher<sup>7</sup>, MD, PhD; Mette Nyegaard<sup>1,2</sup>, MSc, PhD



Combining 300+  
proteins in ONE model

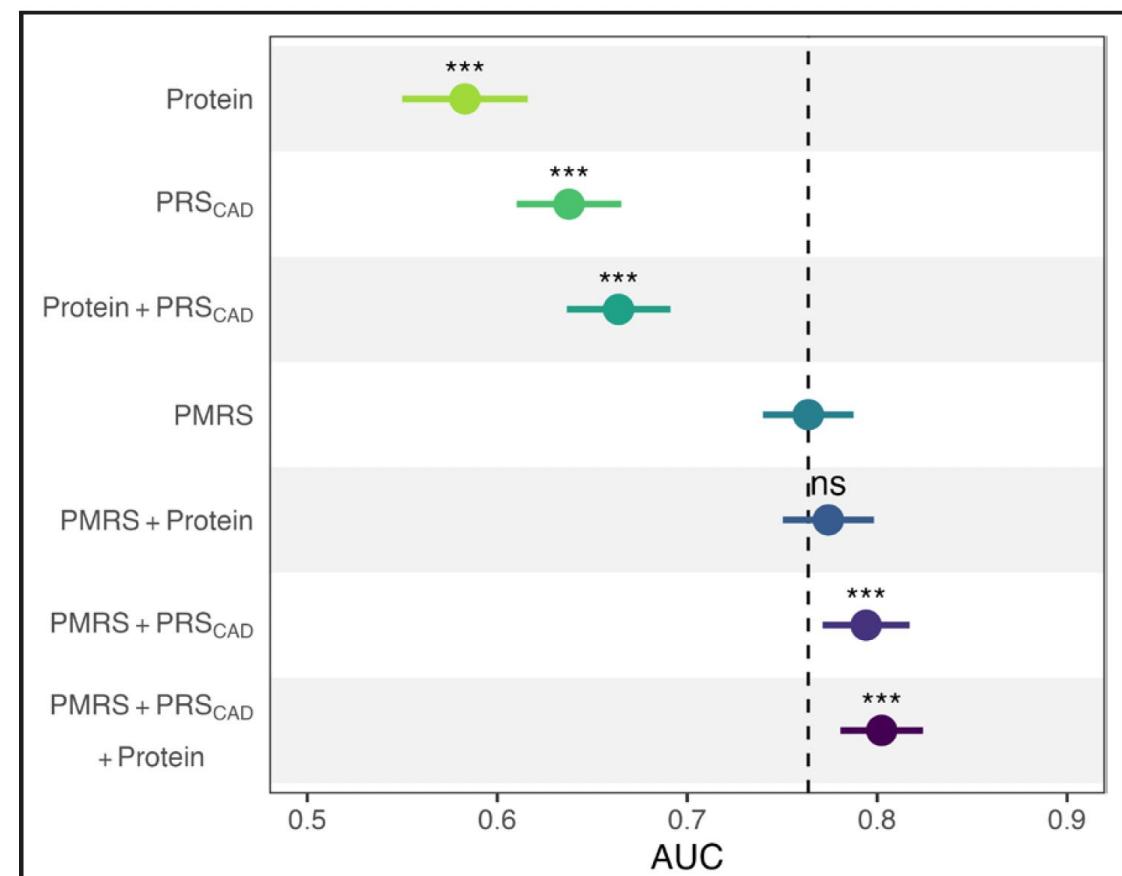
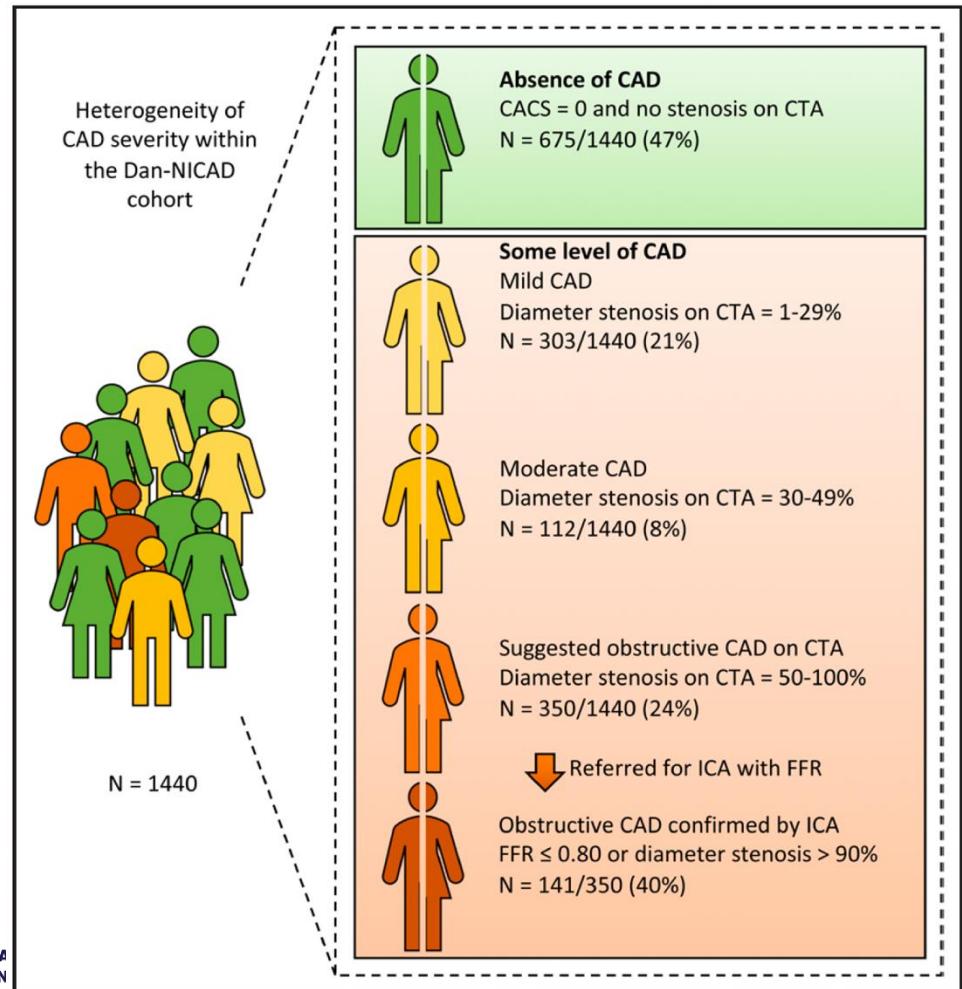


## ORIGINAL ARTICLE

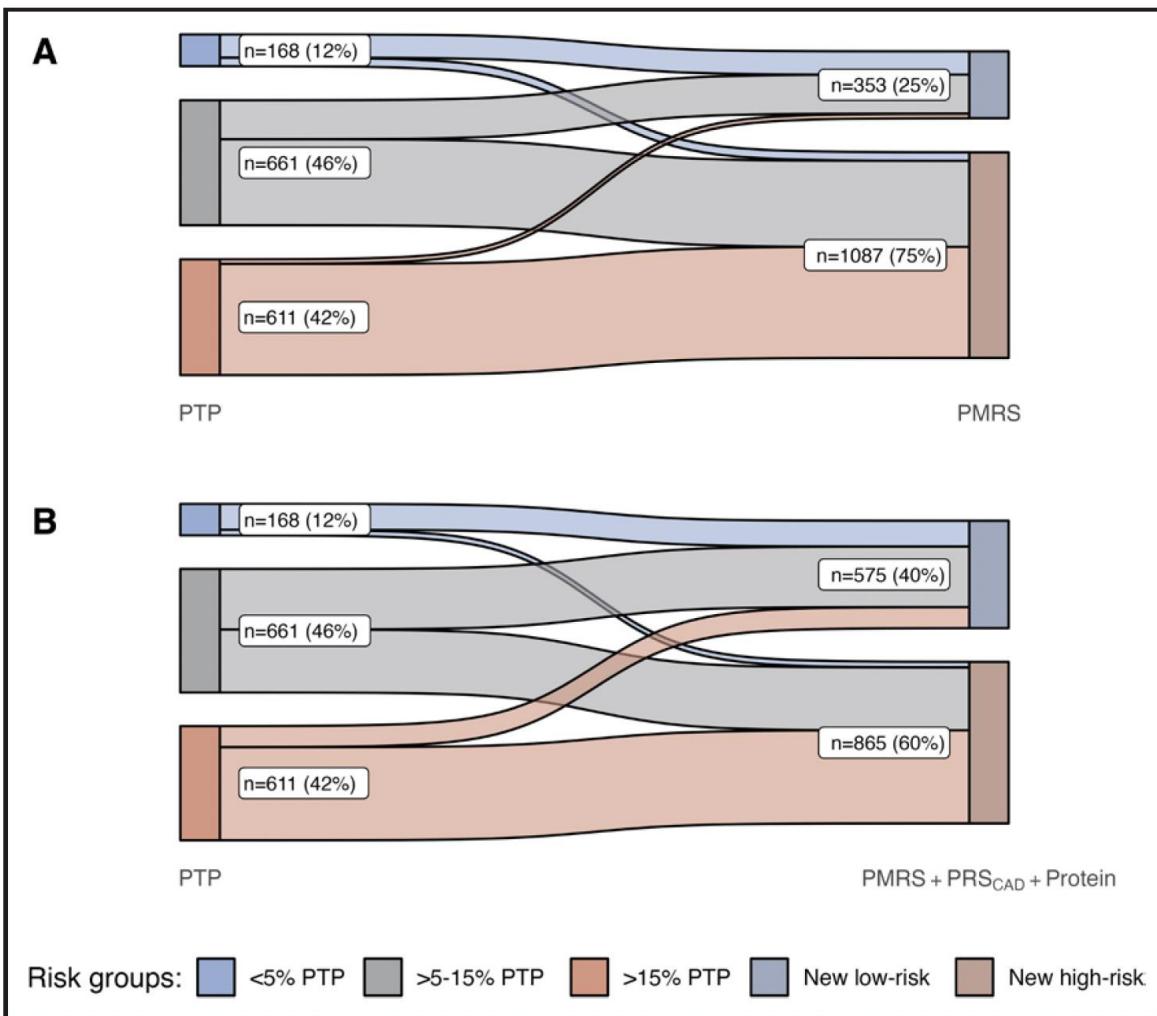
# Combining Polygenic and Proteomic Risk Scores With Clinical Risk Factors to Improve Performance for Diagnosing Absence of Coronary Artery Disease in Patients With de novo Chest Pain

Peter Loof Møller<sup>1</sup>, MSc; Palle Duun Rohde<sup>2</sup>, MSc, PhD; Jonathan Nørtoft Dahl<sup>1</sup>, MD; Laust Dupont Rasmussen<sup>3</sup>, MD, PhD; Samuel Emil Schmidt<sup>2</sup>, MSc, PhD; Louise Nissen, MD, PhD; Victoria McGilligan<sup>2</sup>, PhD; Jacob F. Bentzon<sup>1</sup>, MD, PhD; Daniel F. Gudbjartsson, MSc, PhD; Kari Stefansson<sup>1</sup>, MD, PhD; Hilma Holm<sup>2</sup>, MD; Simon Winther, MD, PhD; Morten Böttcher<sup>2</sup>, MD, PhD; Mette Nyegaard<sup>2</sup>, MSc, PhD

# CAN WE IDENTIFY *HEALTHY* PATIENTS



# CAN WE IDENTIFY *HEALTHY* PATIENTS



ORIGINAL ARTICLE

# Combining Polygenic and Proteomic Risk Scores With Clinical Risk Factors to Improve Performance for Diagnosing Absence of Coronary Artery Disease in Patients With de novo Chest Pain

Peter Loof Møller<sup>1</sup>, MSc; Palle Duun Rohde<sup>2</sup>, MSc, PhD; Jonathan Nørtoft Dahl<sup>1</sup>, MD; Laust Dupont Rasmussen<sup>3</sup>, MD, PhD; Samuel Emil Schmidt<sup>2</sup>, MSc, PhD; Louise Nissen, MD, PhD; Victoria McGilligan<sup>2</sup>, PhD; Jacob F. Bentzon<sup>2</sup>, MD, PhD; Daniel F. Gudbjartsson, MSc, PhD; Kari Stefansson<sup>2</sup>, MD, PhD; Hilma Holm<sup>2</sup>, MD; Simon Winther, MD, PhD; Morten Böttcher<sup>1</sup>, MD, PhD; Mette Nyegaard<sup>2</sup>, MSc, PhD

## Tradition method

## **Combining proteomics, genomics and clinical risk factors has better discriminative ability**

# PATIENT WITH CHEST PAIN



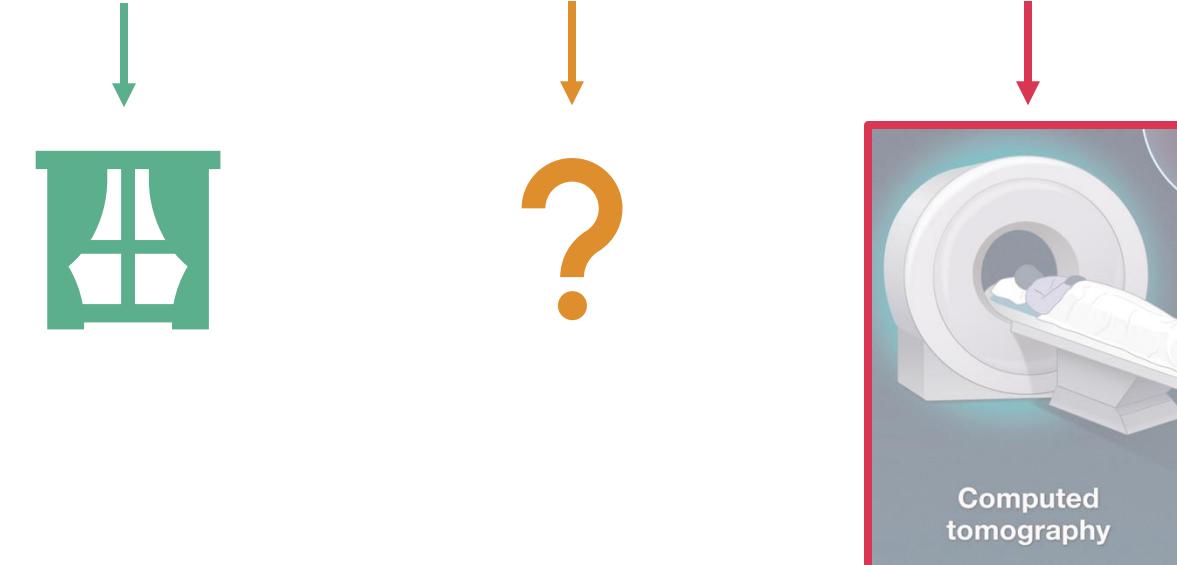
Patient with de-novo chest pain



Pretest probability (PTP) = gender, age and type of chest pain

Risk of obstructive **coronary artery disease** (CAD)

<5%                    5-15%                    >15%



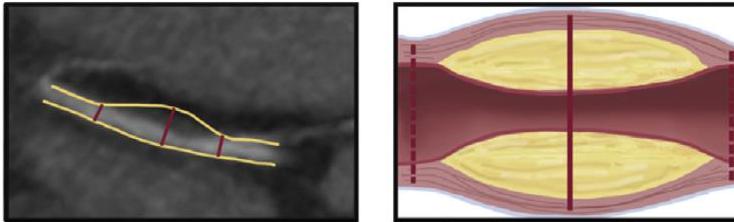
# LOW PTP

Some patients have a low PTP-score, but still have chest pain, and are at risk of a cardiac event.

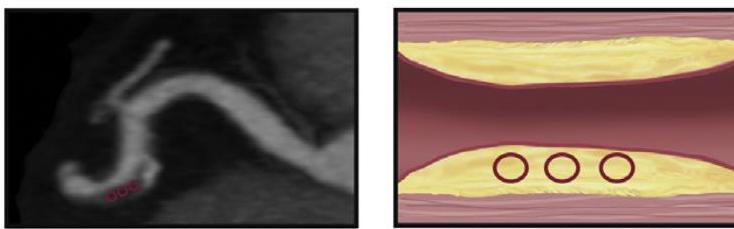


# HIGH RISK CORONARY PLAQUES [HRP]

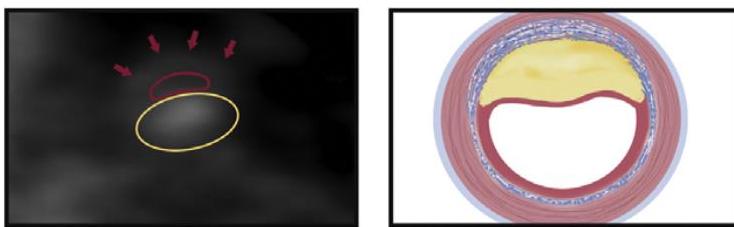
Positive Remodeling



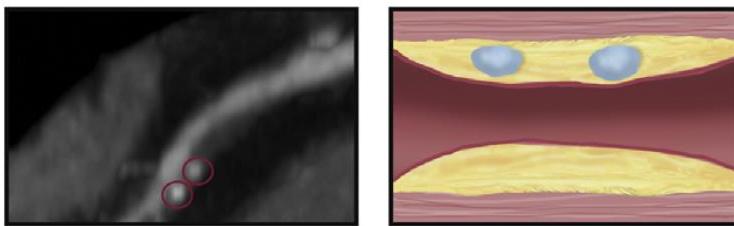
Low HU



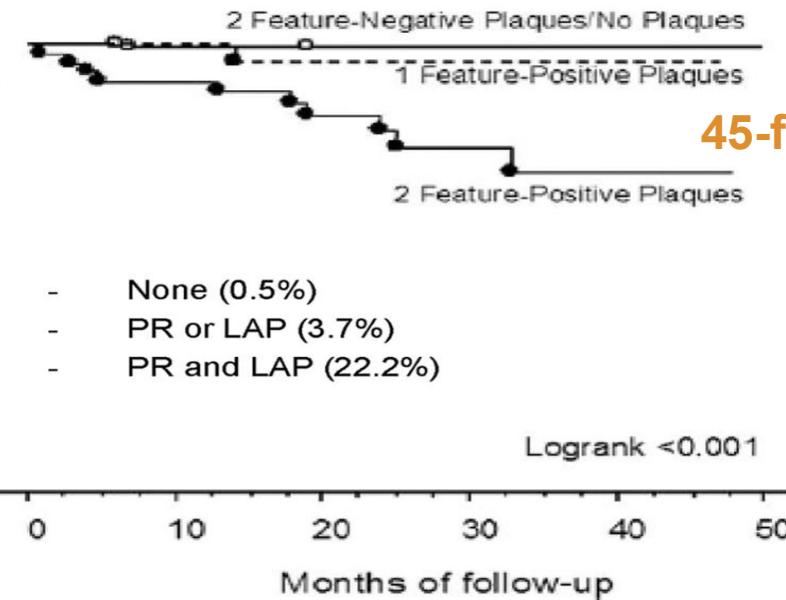
Napkin Ring Sign



Spotty Calcium



Cumulative Event Free Rate



45-fold higher likelihood of cardiac events

Can DNA and targeted proteomics predict the presence of HRP in low risk patients?



# PREDICTING HRP

Patients with chest pain undergoing CCTA  
n=1462

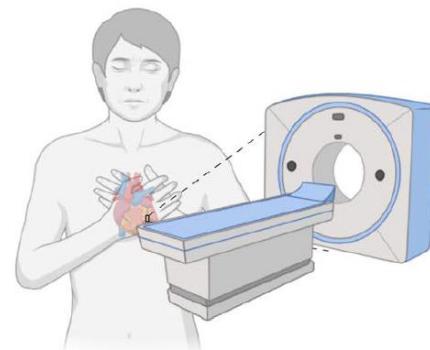
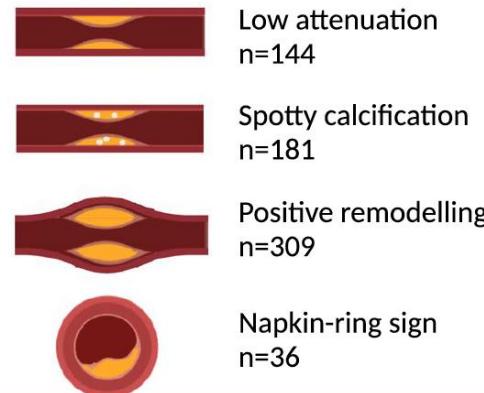


Image analysis to determine high-risk plaque characteristics



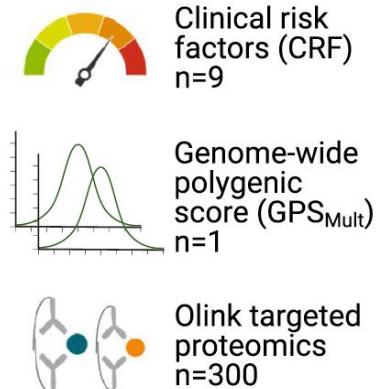
**High-risk plaque (HRP):** a plaque with two or more high-risk characteristics n=165

## RESEARCH

Predicting the presence of coronary plaques featuring high-risk characteristics using polygenic risk scores and targeted proteomics in patients with suspected coronary artery disease

Peter Loof Møller<sup>1,2</sup>, Palle Duun Rohde<sup>2</sup>, Jonathan Nørtoft Dahl<sup>3,4</sup>, Laust Dupont Rasmussen<sup>3,10</sup>, Louise Nissen<sup>3,4</sup>, Samuel Emil Schmidt<sup>2</sup>, Victoria McGilligan<sup>5</sup>, Daniel F. Gudbjartsson<sup>6,7</sup>, Kari Stefansson<sup>6,8</sup>, Hilma Holm<sup>6</sup>, Jacob Fog Bentzon<sup>4,9</sup>, Morten Böttcher<sup>3,4</sup>, Simon Winther<sup>3,4</sup> and Mette Nyegaard<sup>2\*</sup>

Features for prediction of high-risk plaque



**Fig. 1** Study design. 1462 patients underwent coronary computed tomography angiography (CCTA), followed by image analysis of high-risk plaque (HRP) characteristics. Finally, nine clinical risk factors, one multi-trait multi-ancestry genome-wide polygenic score ( $GPS_{Mult}$ ), and 300 proteins were used to predict HRP presence

# PREDICTING HRP

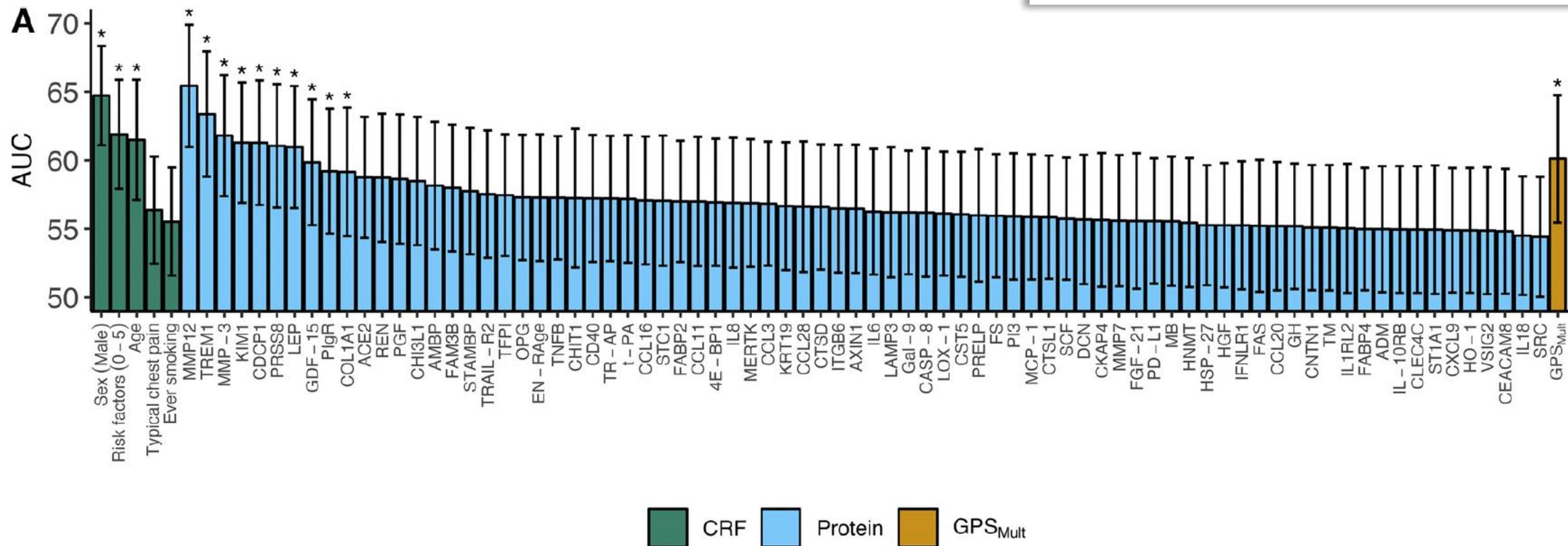
## RESEARCH

## Open Access



Predicting the presence of coronary plaques featuring high-risk characteristics using polygenic risk scores and targeted proteomics in patients with suspected coronary artery disease

Peter Loof Møller<sup>1,2</sup>, Palle Duun Rohde<sup>2</sup>, Jonathan Nørtoft Dahl<sup>3,4</sup>, Laust Dupont Rasmussen<sup>3,10</sup>, Louise Nissen<sup>3,4</sup>, Samuel Emil Schmidt<sup>2</sup>, Victoria McGilligan<sup>5</sup>, Daniel F. Gudbjartsson<sup>6,7</sup>, Kari Stefansson<sup>6,8</sup>, Hilma Holm<sup>6</sup>, Jacob Fog Bentzon<sup>4,9</sup>, Morten Böttcher<sup>3,4</sup>, Simon Winther<sup>3,4</sup> and Mette Nyegaard<sup>2\*</sup>



# PREDICTING HRP

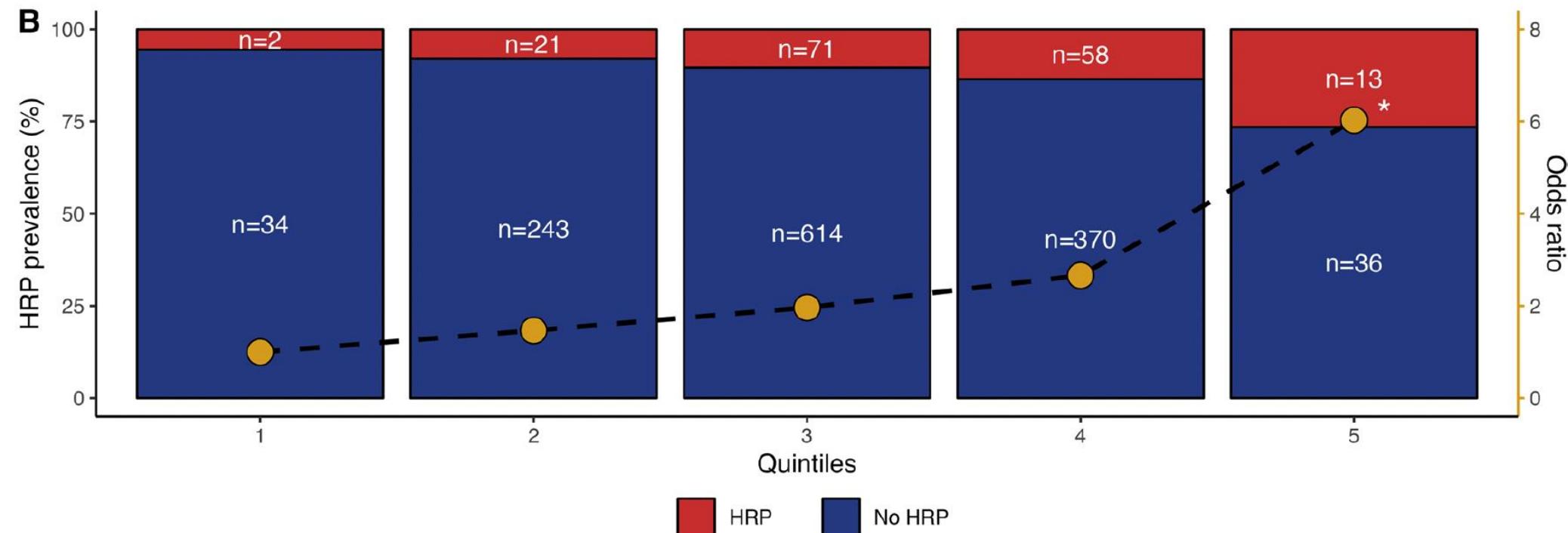
## RESEARCH

## Open Access



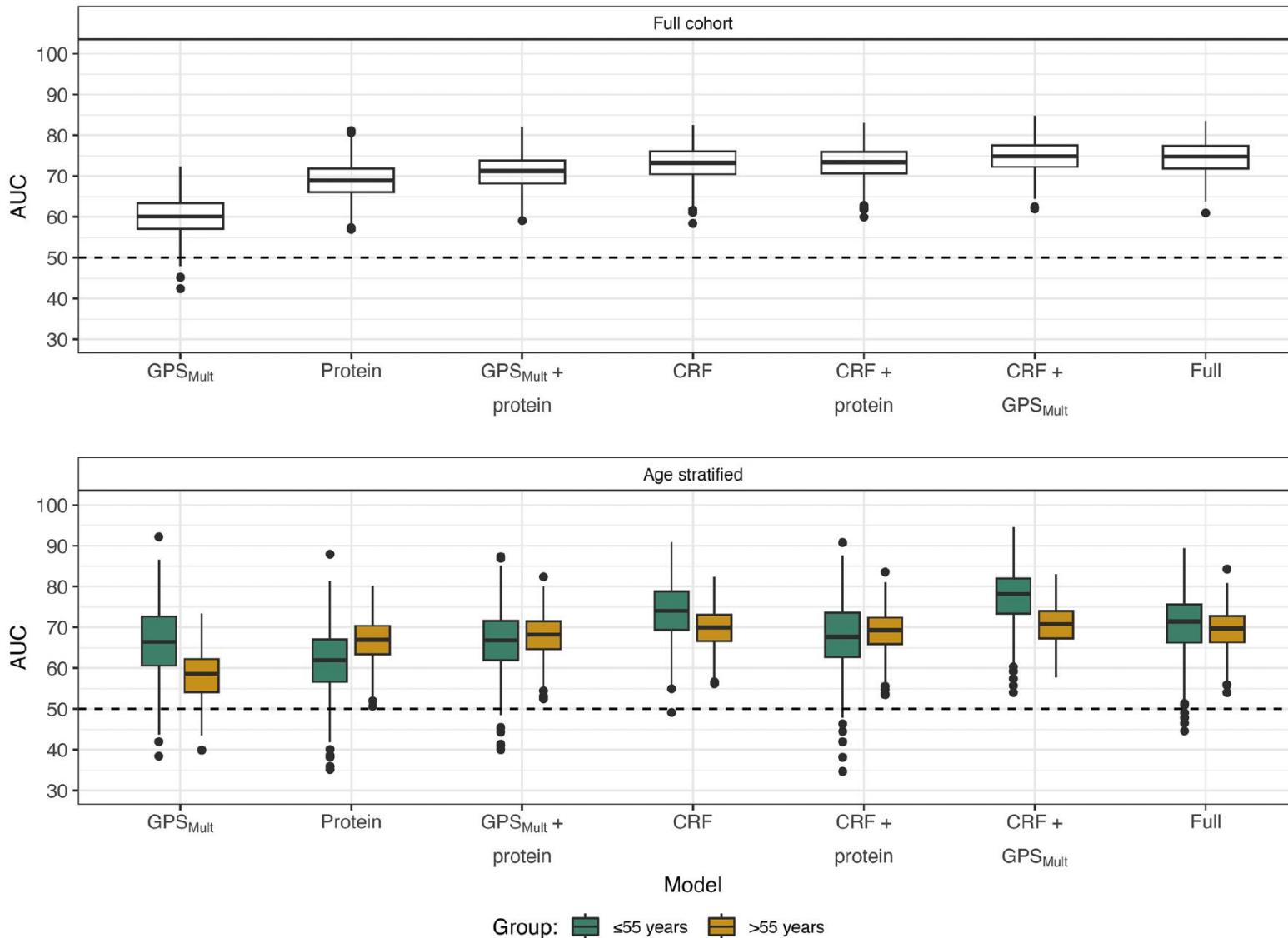
Predicting the presence of coronary plaques featuring high-risk characteristics using polygenic risk scores and targeted proteomics in patients with suspected coronary artery disease

Peter Loof Møller<sup>1,2</sup>, Palle Duun Rohde<sup>2</sup>, Jonathan Nørtoft Dahl<sup>3,4</sup>, Laust Dupont Rasmussen<sup>3,10</sup>, Louise Nissen<sup>3,4</sup>, Samuel Emil Schmidt<sup>2</sup>, Victoria McGilligan<sup>5</sup>, Daniel F. Gudbjartsson<sup>6,7</sup>, Kari Stefansson<sup>6,8</sup>, Hilma Holm<sup>6</sup>, Jacob Fog Bentzon<sup>4,9</sup>, Morten Böttcher<sup>3,4</sup>, Simon Winther<sup>3,4</sup> and Mette Nyegaard<sup>2\*</sup>



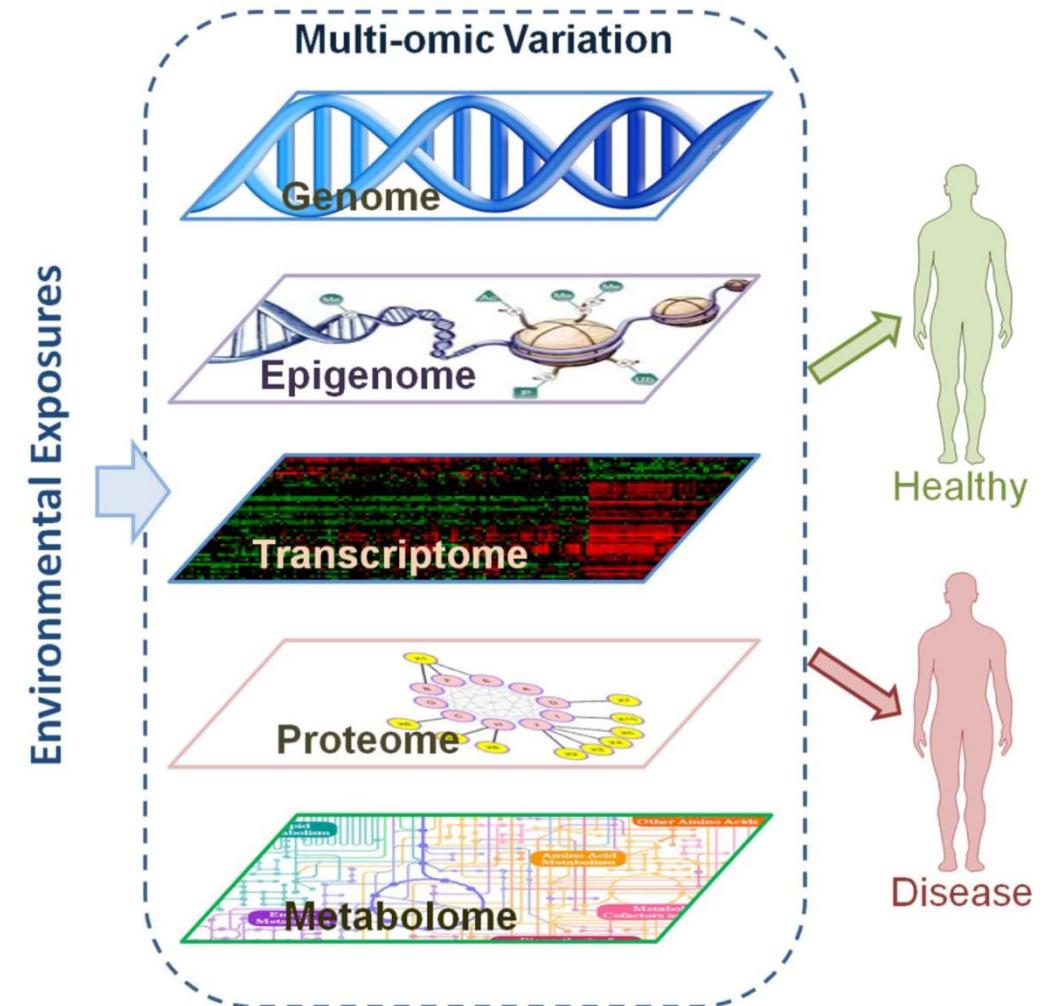


# PREDICTING HRP



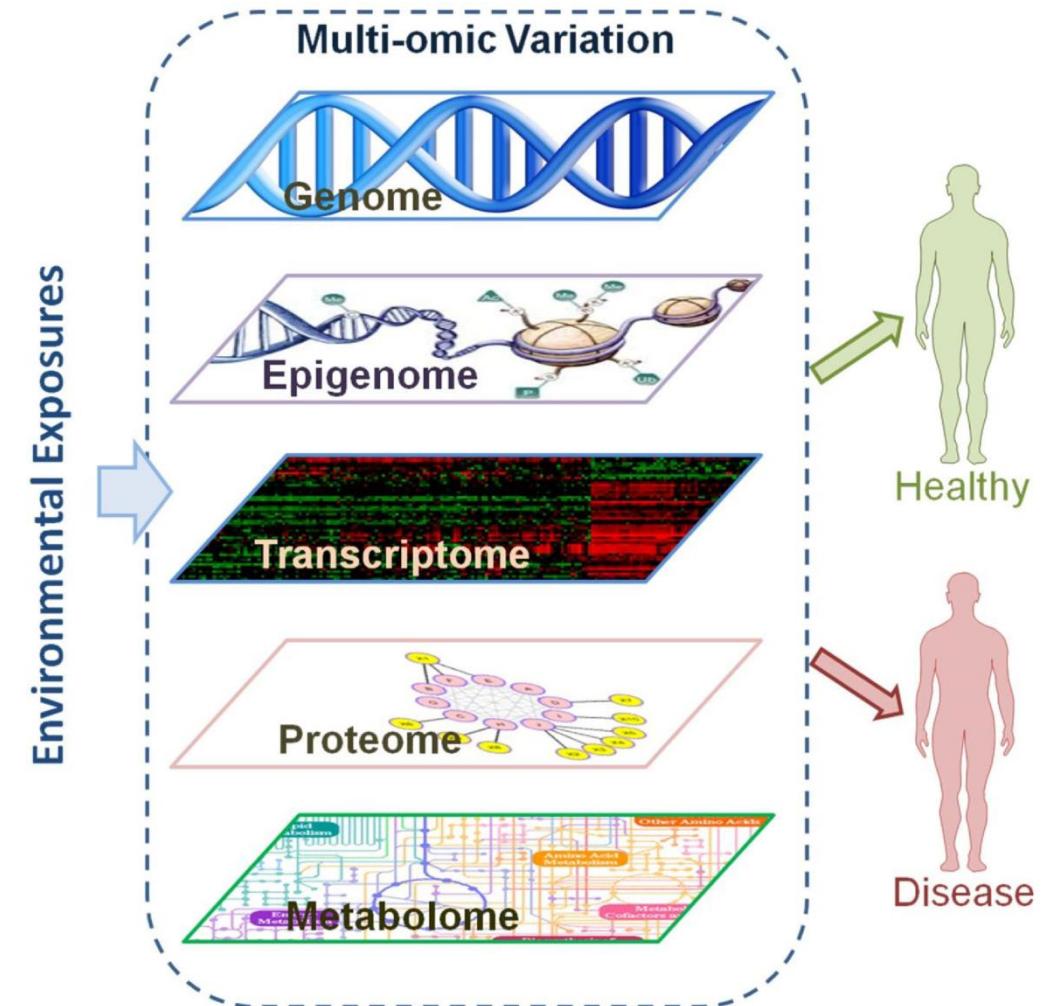
# VARIATION AT MULTIPLE LAYERS

- ❖ Biological data are now being collected in large scale across the different layers of molecular organisation.



# VARIATION AT MULTIPLE LAYERS

- ❖ Biological data are now being collected in large scale across the different layers of molecular organisation.
- ❖ Variation within each layer has generated knowledge about human complex diseases.



# SESSION 4

- Many challenges – how to circumvent these?
- Enhancing biological understanding?
- PGS in combination with proteomics





# A toolbox for statistical genetic analyses of complex traits

*Bioinformatics*, 36(8), 2020, 2614–2615  
doi: 10.1093/bioinformatics/btz955  
Advance Access Publication Date: 27 December 2019  
Applications Note

OXFORD

Genetics and population analysis  
**qgg: an R package for large-scale quantitative genetic analyses**

Palle Duun Rohde \*†, Izel Fourie Sørensen and Peter Sørensen\*

Department of Molecular Biology and Genetics, Aarhus University, Tjele, Denmark

\*To whom correspondence should be addressed.  
†Present address: Department of Chemistry, Aarhus University, Tjele, Denmark  
Associate Editor: Russell Schwartz

*Bioinformatics*, 2023, 39(11), btad656  
<https://doi.org/10.1093/bioinformatics/btad656>  
Advance Access Publication Date: 26 October 2023  
Applications Note

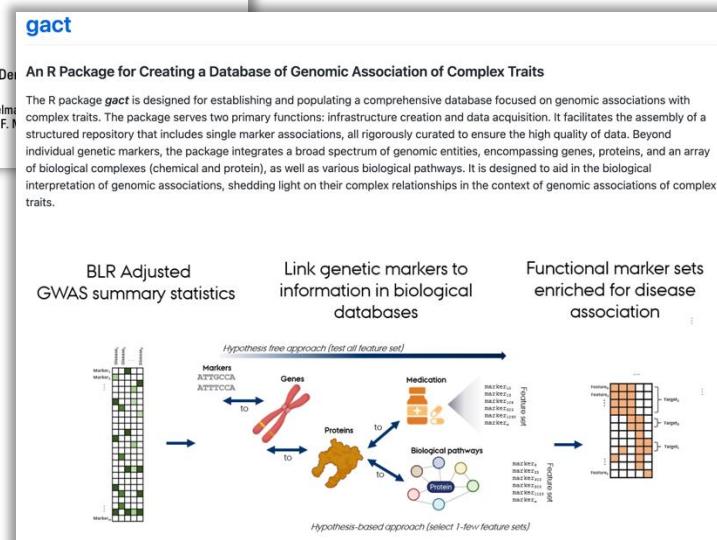
OXFORD

Genetics and population analysis  
**Expanded utility of the R package, qgg, with applications within genomic medicine**

Palle Duun Rohde 1,\*<sup>1</sup>, Izel Fourie Sørensen<sup>2</sup>, Peter Sørensen<sup>2,\*</sup>

<sup>1</sup>Genomic Medicine, Department of Health Science and Technology, Aalborg University, 9260 Gistrup, Denmark  
<sup>2</sup>Center for Quantitative Genetics and Genomics, Aarhus University, 8000 Aarhus, Denmark

\*Corresponding authors: Genomic Medicine, Department of Health Science and Technology, Aalborg University, Selmer Brøchner Building, 9260 Gistrup, Denmark. E-mail: palledr@hs.aau.dk (P.D.R.); Center for Quantitative Genetics and Genomics, Aarhus University, C. F. Møller Building, 8000 Aarhus, Denmark. E-mail: psq@qgg.au.dk (P.S.)  
Associate Editor: Christina Kendziorski



## Tutorials

Below is a set of tutorials used for the qgg package:

This tutorial provides a brief introduction to R package qgg using small simulated data examples.

### [Practicals\\_brief\\_introduction](#)

This tutorial provides an introduction to R package qgg using 1000G data.

### [Practicals\\_1000G\\_tutorials](#)

This tutorial provide a simple introduction to polygenic risk scoring (PRS) of complex traits and diseases using simulated data. The practical will be a mix of theoretical and practical exercises in R that are used for illustrating/applying the theory presented in the corresponding lecture notes on polygenic risk scoring.

### [Practicals\\_human\\_example](#)

In this tutorial we will be analysing quantitative traits observed in a mice population. The mouse data consist of phenotypes for traits related to growth and obesity (e.g. body weight, glucose levels in blood), pedigree information, and genetic marker data.

### [Practicals\\_mouse\\_example](#)

## Notes

Below is a set of notes for the quantitative genetic theory, statistical models and methods implemented in the qgg package:

### [Quantitative Genetics Theory](#)

### [Estimation of Genetic Predisposition](#)

### [Estimation of Genetic Parameters](#)

### [Linear Mixed Models](#)

### [Best Linear Unbiased Prediction Models](#)

### [REstricted Maximum Likelihood Methods](#)

### [Gene Set Enrichment Analysis](#)

### [Bayesian Linear Regression Models](#)

# THE PURPOSE OF TODAY

- ❖ Give an introduction to polygenic scores (PGS)
- ❖ Provide an introduction to complex trait genetics
  - Monogenic vs multifactorial aetiology
- ❖ How we can utilize genomic data to elucidate molecular genetic aetiology underlying complex traits
  - Genome-wide association studies (GWAS)
- ❖ Stratify a population/cohort by their inherent genetic load towards common complex diseases
  - Polygenic scores (PGS)
- ❖ Identify future projects of common interests



# AGENDA

08:00 – 08:30	Welcome and common introductions
08:30 – 09:10	Session 1: Introduction to Polygenic Scores (PGS)
09:10 – 09:20	Break
09:20 – 10:00	Session 2: Data Sources and Computational Methods
10:00 – 10:10	Break
10:10 – 10:40	Session 3: Evaluating and Interpreting Polygenic Scores
10:40 – 11:00	Break
11:00 – 11:45	Session 4: Advanced Applications and Future Directions
11:45 – 12:30	Lunch and short walk
12:30 – 15:30	Identification of 2-3 projects of common interest
15:30 – 16:00	Next steps and thank you for today