# Lecture 1 - GWAS Statistics + Polygenic Risk Scores (PRS)

Wed, Mar 20, 9-10AM
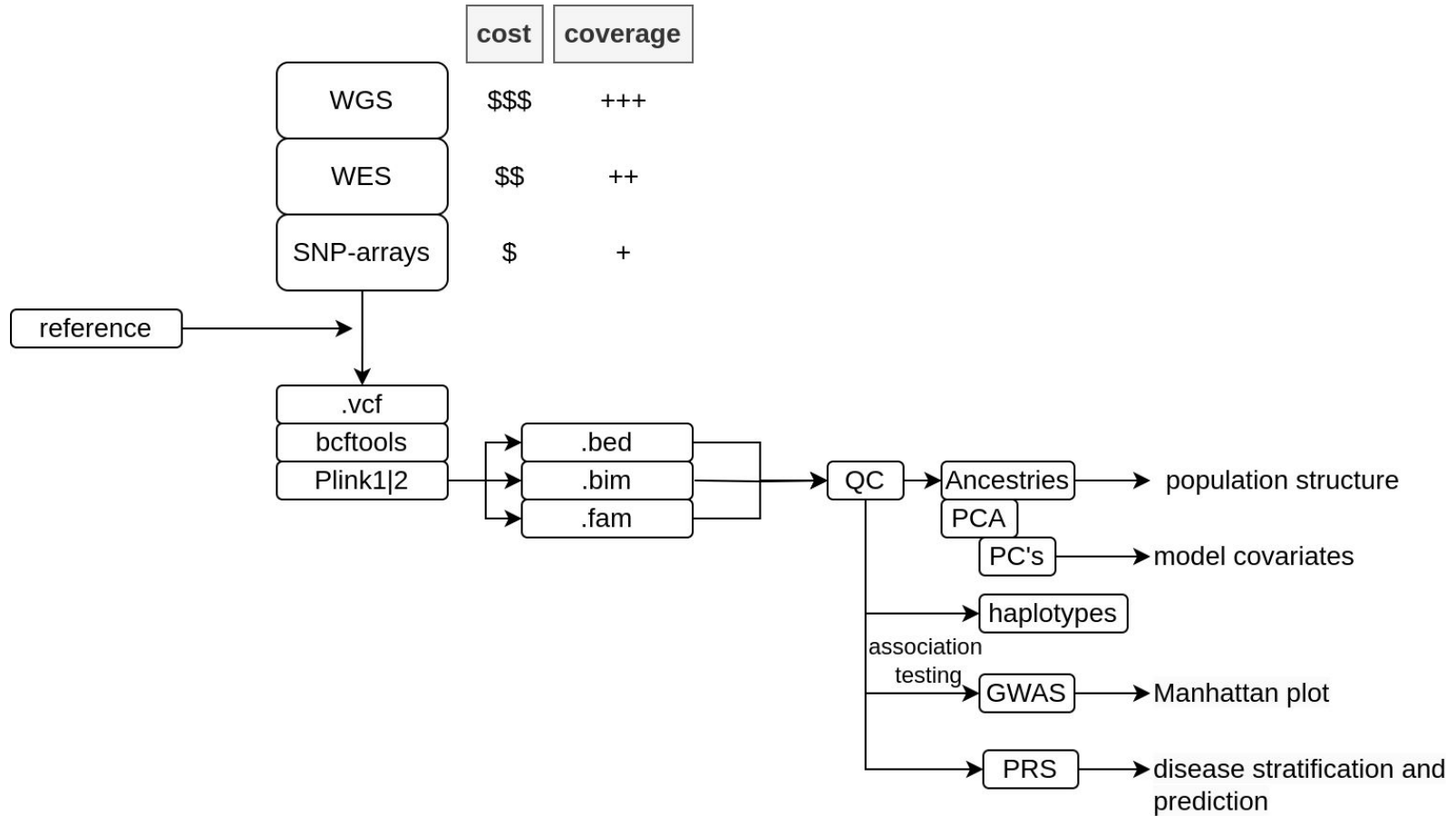
# Recap

HPC | GWAS

# Recap: HPC + SSH

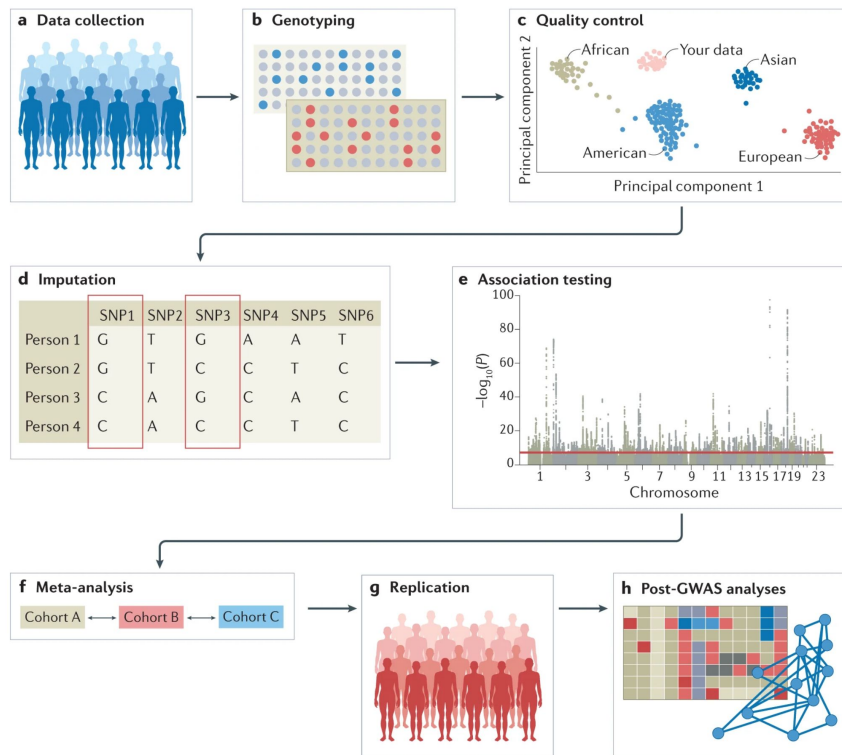| | long-term cost | scalability |
|---|---|---|
| laptop | $ | + |
| HPC | $$ | ++ |
| cloud | $$$ | +++ |

laptop —SSH→ HPC → cloud

on-site
VPN

@kennedy

| Shell | → | Bash |
| Job scheduler | → | SLURM |
| Scripts | → | Python |
| | → | R |
| | → | Plink1|2 |

# Recap: GWAS

| | cost | coverage |
|---|---|---|
| WGS | $$$ | +++ |
| WES | $$ | ++ |
| SNP-arrays | $ | + |

reference →

.vcf
bcftools
Plink1|2

→ .bed
→ .bim
→ .fam

→ QC →

Ancestries → population structure
PCA
PC's → model covariates

→ haplotypes

association testing

→ GWAS → Manhattan plot

→ PRS → disease stratification and prediction

# GWAS

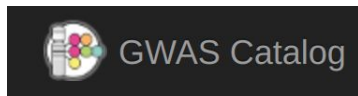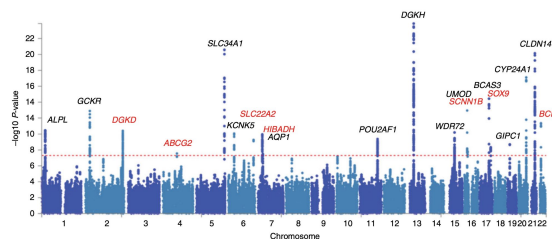Genome-wide Association Studies

# Genome-wide Association studies (GWAS)

**Single nucleotide polymorphism (SNP)**: This is a variation in a single nucleotide (i.e., **A**, **C**, **G**, or **T**) that occurs at a specific position in the genome. A SNP usually exists as two different forms (e.g., **A** vs. **T**). These different forms are called alleles. A SNP with two alleles has three different genotypes (e.g., **AA**, **AT**, and **TT**).
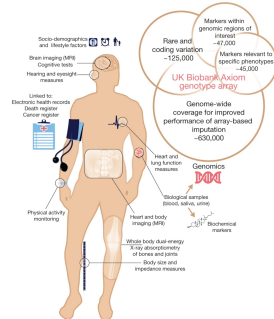
# Data sources & repositories



Summary statistics

general GWAS data repository

specific for GWAS Chronic Kidney Disease

Individual-level

# Testing for associations

## Genetic models

| | AA | AG | GG |
|---|---|---|---|
| **Additive model** | 0 | 1 | 2 |
| **Dominant model** | 0 | 1 | 1 |
| **Recessive model** | 0 | 0 | 1 |

- Additive model (ADD)
- Dominant model (DOM)
- Recessive model (REC)

biallelic SNP whose reference allele is **A** and the alternative allele is **G**.

# Testing for associations

## Contingency table

| genotype | AA | AG | GG | Total |
|----------|-----|-----|------|-------|
| **case** | 800 | 400 | 800 | 2000 |
| **control** | 1000 | 500 | 500 | 2000 |
| Total | 1800 | 900 | 1300 | 4000 |

### **Dominant** model

| genotype | **control:AA** | **case:AG/GG** | Total |
|----------|----------------|----------------|-------|
| **case** | 800 | 1200 | 2000 |
| **control** | 1000 | 1000 | 2000 |
| Total | 1800 | 2200 | 4000 |

### **Recessive** model

| genotype | **control:AA/AG** | **case:GG** | Total |
|----------|-------------------|-------------|-------|
| **case** | 1200 | 800 | 2000 |
| **control** | 1500 | 500 | 2000 |
| Total | 2700 | 1300 | 4000 |

### **Additive** model

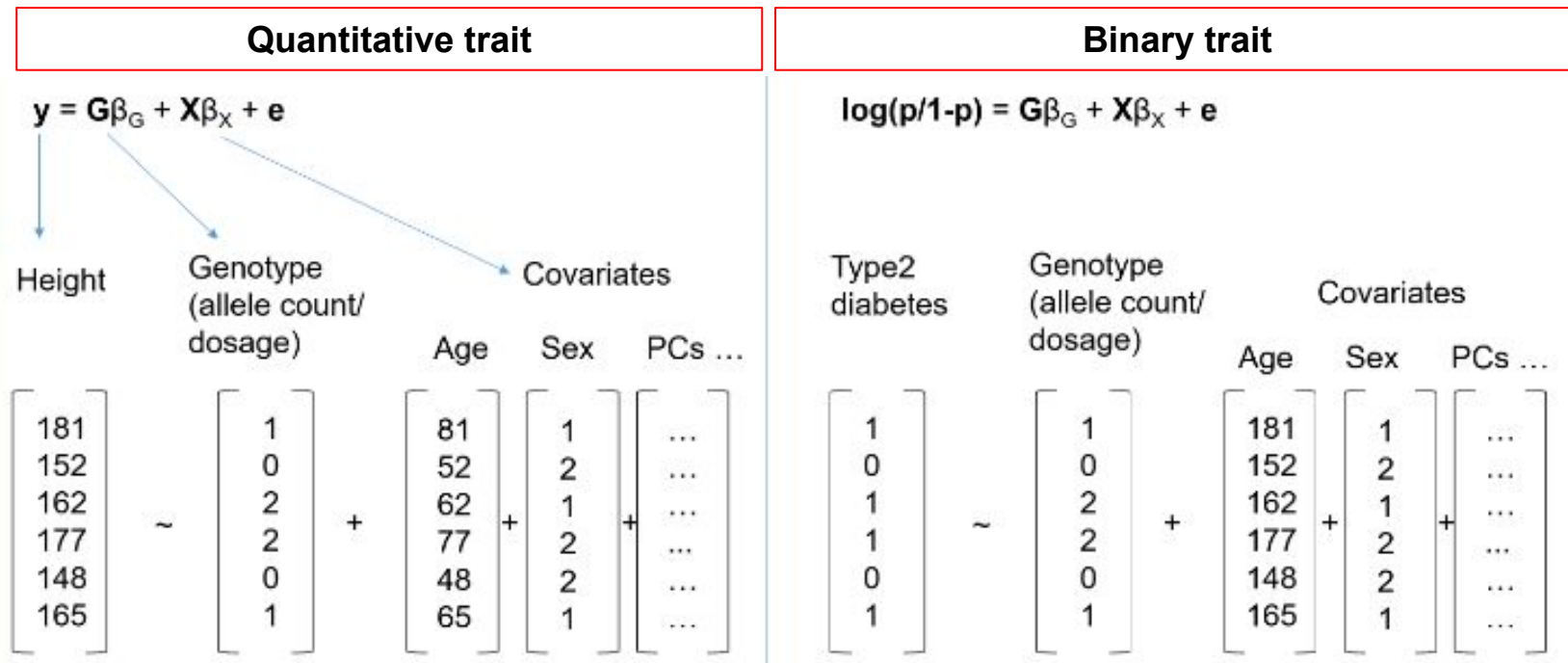| Allele (G) | **0** | **1** | **2** | Total |
|------------|-------|-------|-------|-------|
| **case** | 800 | 400 | 800 | 2000 |
| **control** | 1000 | 500 | 500 | 2000 |
| Total | 1800 | 900 | 1300 | 4000 |

# Testing for associations

**Quantitative** traits

$$y = G\beta_G + X\beta_X + e$$

- $G$ is the genotype matrix.
- $\beta_G$ is the effect size for variants.
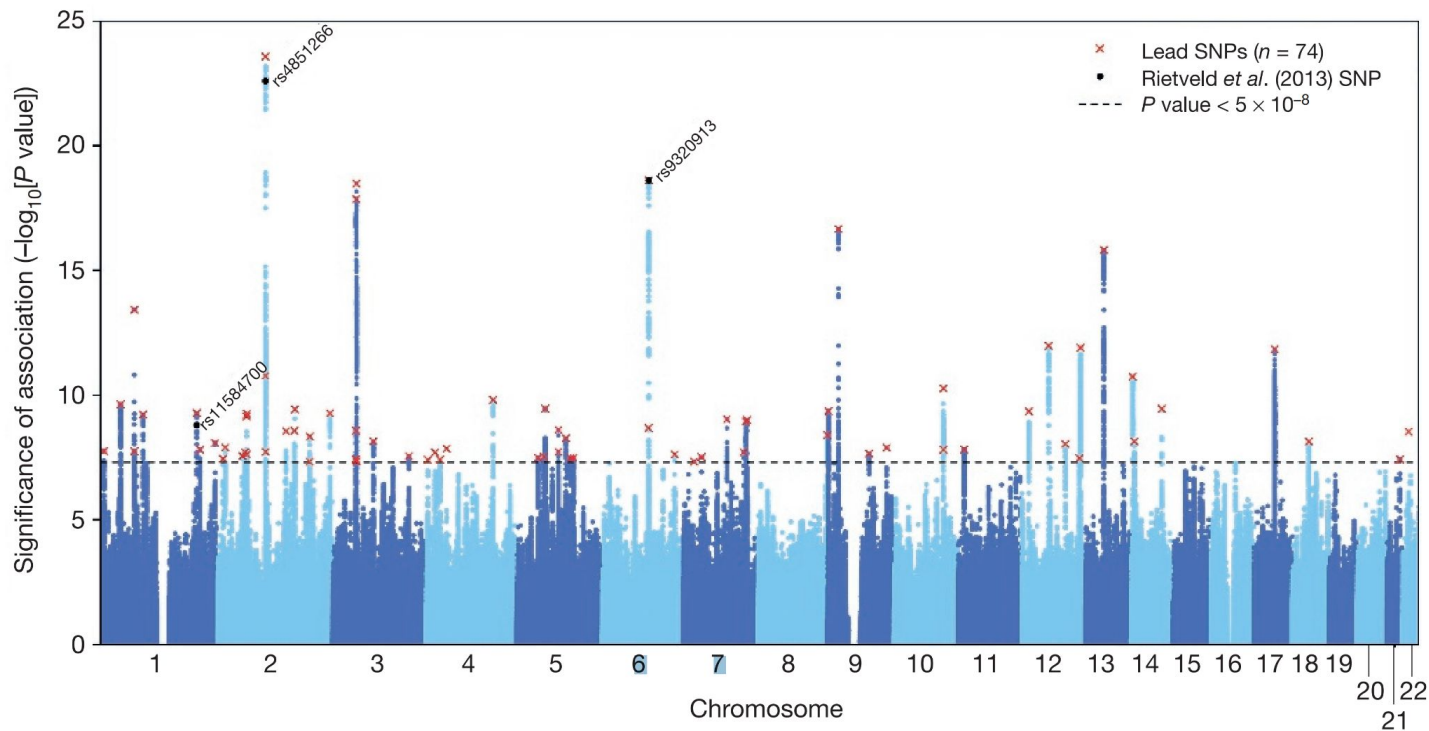- $X$ and $\beta_X$ are covariates and their effects.
- $e$ is the error term.
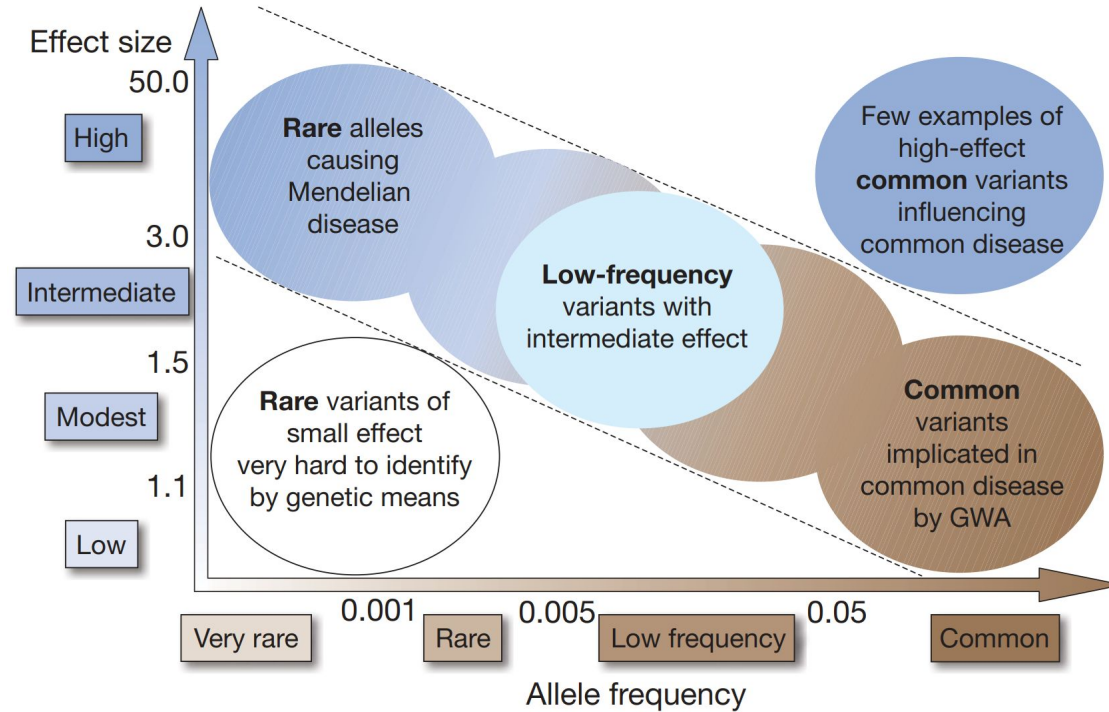
**Binary** traits

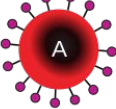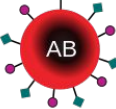$$logit(p) = G\beta_G + X\beta_X + e$$

# Testing for associations



| Quantitative trait | Binary trait |

$$y = G\beta_G + X\beta_X + e$$

Height ~ Genotype (allele count/dosage) + Covariates (Age, Sex, PCs ...)

$$\log(p/1-p) = G\beta_G + X\beta_X + e$$

Type2 diabetes ~ Genotype (allele count/dosage) + Covariates (Age, Sex, PCs ...)
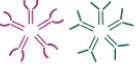
# Manhattan plot for EduYears associations (n = 293,723)

# Rare and common variants

# Haplotypes: ABO serological groups



| Blood group antigen | Tag SNP | Effect allele/non-effect allele |
|---|---|---|
| $A_1$ | rs507666 | A/G |
| $A_2$ | rs8176704 | A/G |
| B | rs8176746 | T/G |
| O | rs687289 | G/A |

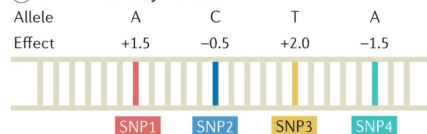# Polygenic Risk Scores (PRS)

stratification & disease trajectories

# Common workflow - single-trait PRS

**Polygenic risk scores (PRS)**



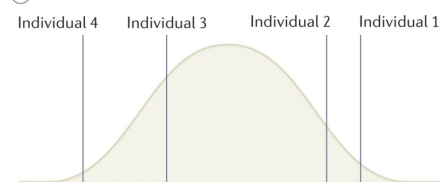score that summarizes the effect sizes of genetic variants on a certain disease or trait.

$$PRS_i = \sum_{j \in J} \beta_j \, G_{ij}$$

*i*-th individual
*j*-th variant
*G*: genotype
*β*: effect size

# Common workflow - single-trait PRS



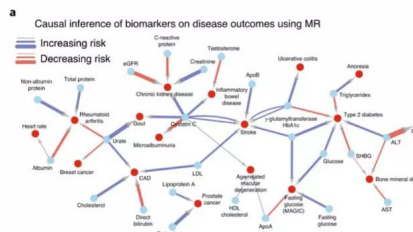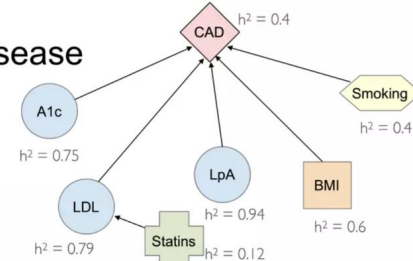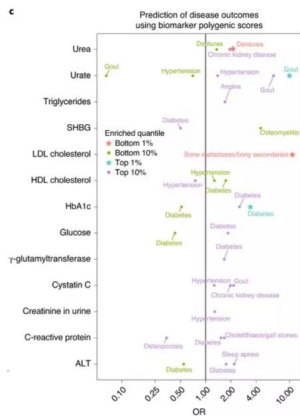| Category | Description | Methods/software |
|---|---|---|
| P value thresholding | P + T | C+T, PRSice, Plink |
| Beta shrinkage | genome-wide PRS model | LDpred |

**P+T** stands for Pruning + Thresholding, also known as Clumping and Thresholding (**C+T**)
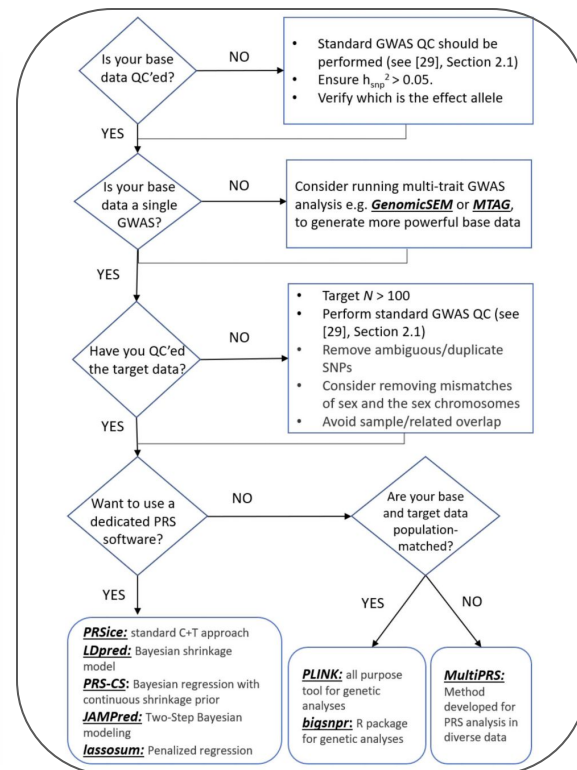
# Common workflow - multi-trait PRS

- Multiple observations suggest "biomarkers → disease" links
  - PRS-PheWAS analysis
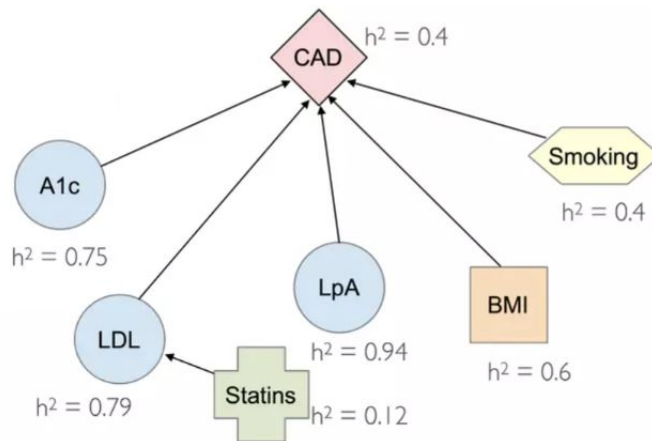  - Biomarkers are more heritable than disease
  - Mendelian Randomization

Prediction of disease outcomes using biomarker polygenic scores

Urea
Urate
Triglycerides
SHBG
LDL cholesterol
HDL cholesterol
HbA1c
Glucose
γ-glutamyltransferase
Cystatin C
Creatinine in urine
C-reactive protein
ALT

Enriched quantile
- Bottom 1%
- Bottom 10%
- Top 1%
- Top 10%

OR: 0.10 0.25 0.50 1.00 2.00 4.00 10.00

CAD $h^2 = 0.4$
A1c $h^2 = 0.75$
Smoking $h^2 = 0.4$
LpA $h^2 = 0.94$
BMI $h^2 = 0.6$
LDL $h^2 = 0.79$
Statins $h^2 = 0.12$

Causal inference of biomarkers on disease outcomes using MR
Increasing risk
Decreasing risk

Flowchart:

Is your base data QC'ed? — NO → 
- Standard GWAS QC should be performed (see [29], Section 2.1)
- Ensure $h_{snp}^2 > 0.05$.
- Verify which is the effect allele

YES ↓

Is your base data a single GWAS? — NO → Consider running multi-trait GWAS analysis e.g. *GenomicSEM* or *MTAG*, to generate more powerful base data

YES ↓

Have you QC'ed the target data? — NO → 
- Target $N > 100$
- Perform standard GWAS QC (see [29], Section 2.1)
- Remove ambiguous/duplicate SNPs
- Consider removing mismatches of sex and the sex chromosomes
- Avoid sample/related overlap

YES ↓

Want to use a dedicated PRS software? — NO → Are your base and target data population-matched?

YES ↓

**PRSice:** standard C+T approach
**LDpred:** Bayesian shrinkage model
**PRS-CS:** Bayesian regression with continuous shrinkage prior
**JAMPred:** Two-Step Bayesian modeling
**lassosum:** Penalized regression

YES → 
**PLINK:** all purpose tool for genetic analyses
**bigsnp:** R package for genetic analyses

NO → 
**MultiPRS:** Method developed for PRS analysis in diverse data

- Multi-PRS is a weighted sum of PRSs
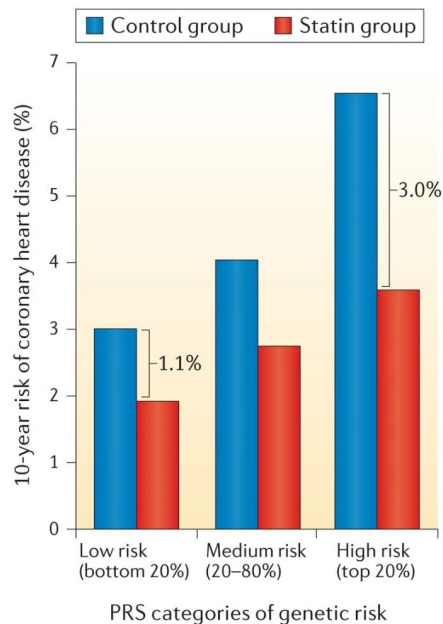  i.e. $w_1(PRS_1) + w_2(PRS_2) + w_3(PRS_3) + \ldots$

# Common workflow - multi-trait PRS
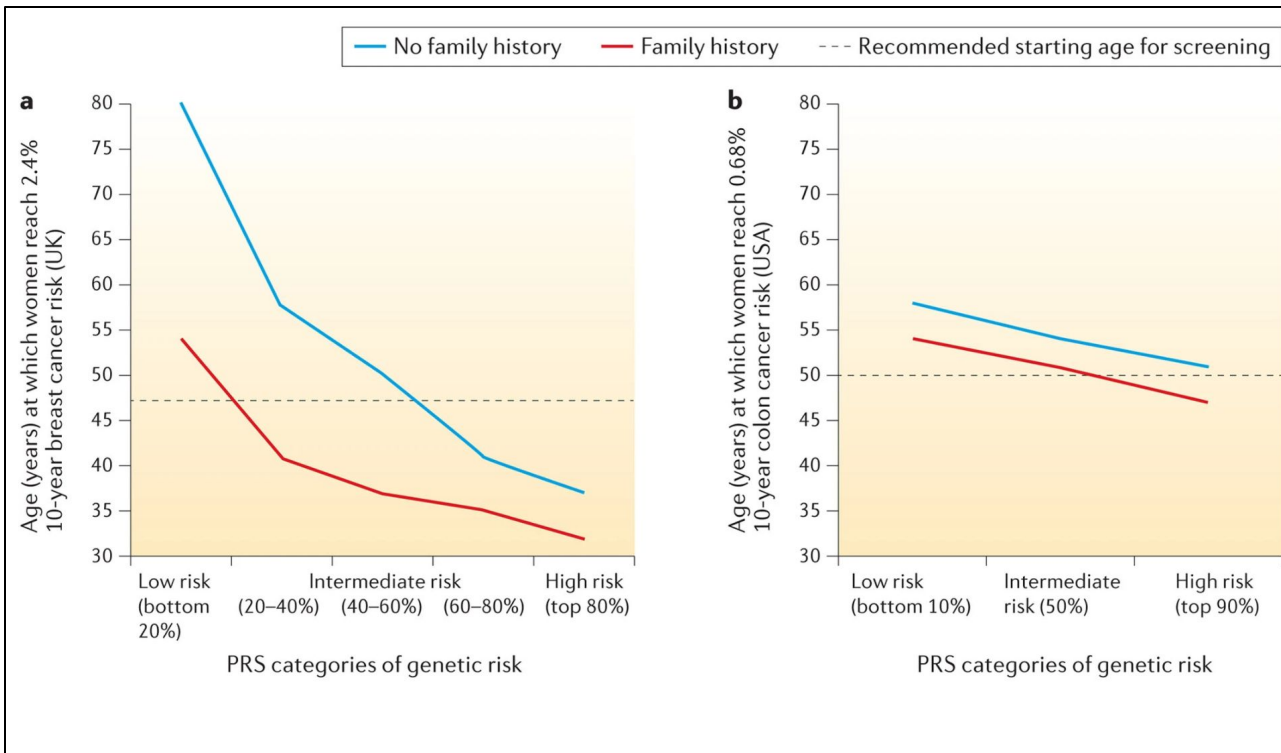
- Multi-PRS is a weighted sum of PRSs
  i.e. $w_1(PRS_1) + w_2(PRS_2) + w_3(PRS_3) + \ldots$
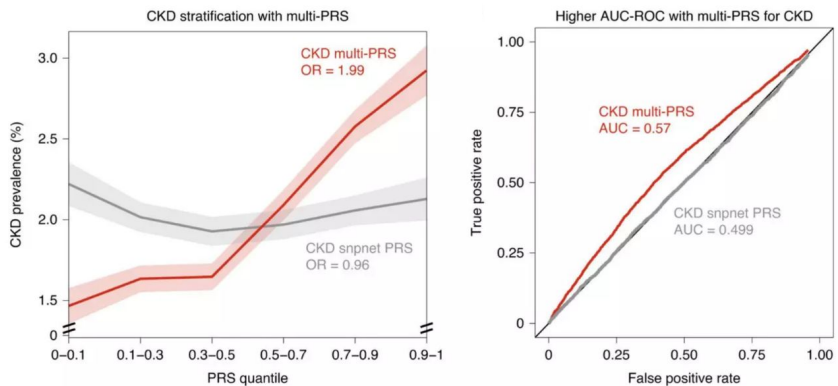
# Case study - disease stratification



| ARR (%) | 1.1 | 1.3 | 3.0 |
|---------|-----|-----|-----|
| RRR | 0.36 | 0.32 | 0.46 |

# Case study - multi-trait PRS improves disease prevalence prediction

# Use of PRS for trait / disease prediction



Prediction of disease outcomes using biomarker polygenic scores

Take extreme in PRS for biomarkers

Compare odds ratio for disease outcome relative to 40-60%ile bin

Applied PheWAS for ~160 diseases

# Use of PRS for trait / disease prediction



Prediction of disease outcomes using biomarker polygenic scores

Take extreme in PRS for biomarkers

Identify diseases with biomarker PRS associations

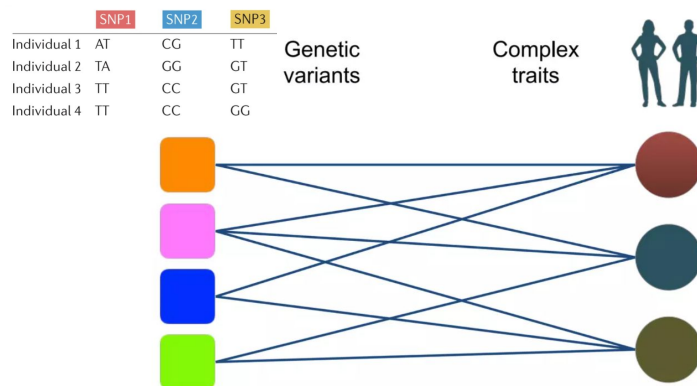Compare odds ratio for disease outcome relative to 40-60%ile bin

Applied PheWAS for ~160 diseases
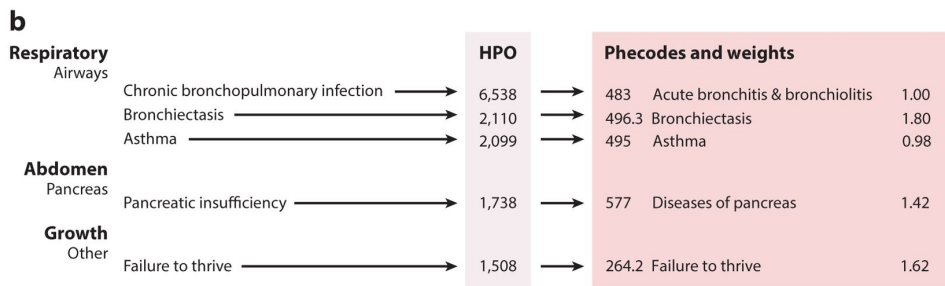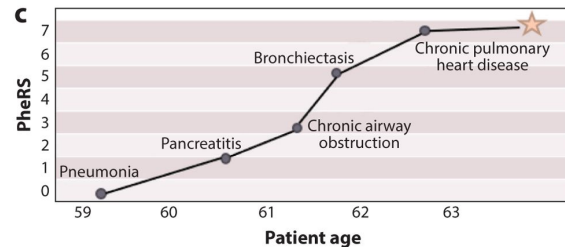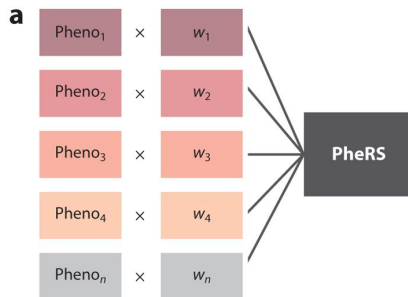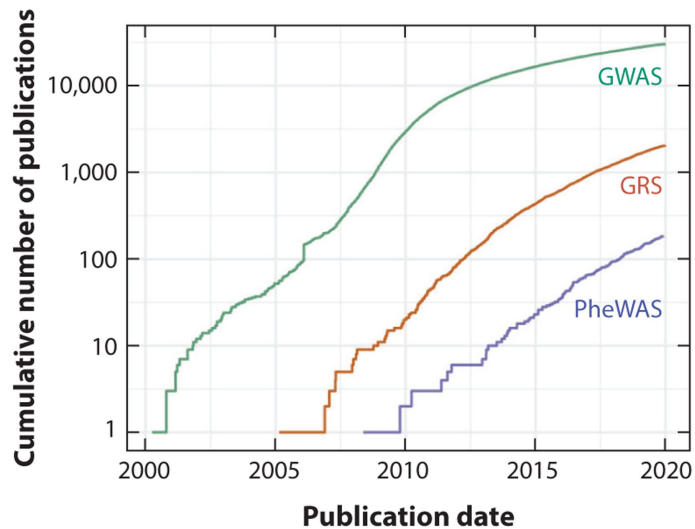
# Limitations - ethnicity / ancestry

# Limitations - polygenicity & pleiotropy

- Polygenicity: many variants - one trait
- Pleiotropy: one variant - many traits

| | SNP1 | SNP2 | SNP3 |
|---|---|---|---|
| Individual 1 | AT | CG | TT |
| Individual 2 | TA | GG | GT |
| Individual 3 | TT | CC | GT |
| Individual 4 | TT | CC | GG |

Genetic variants

Complex traits

- Large number of associations in population-based cohorts
- Can we group them together for enhanced interpretation?

# Going further - EHR's & PRS

# Summary

Two complementary approaches to improve predictive performance:

- Sample size → increase in **statistical power**
  - Multi-trait PRS analysis

Why does multi-PRS work?

- Quantitative traits have more power
- **Genetic correlation** between biomarkers and disease

The multi-trait PRS model:

Genetics →Biomarkers (molecular traits) →Disease