

# Polygenic Prediction using GWAS Summary Statistics

SISG – Module 15

Dr Loic Yengo

[l.yengo@imb.uq.edu.au](mailto:l.yengo@imb.uq.edu.au)

[Slides credit: Dr Jian Zeng (UQ) & Prof. Naomi Wray (UQ)]

Institute for Molecular Bioscience  
The University of Queensland

The aim of this lecture is to introduce a few general classes of methods for genomic prediction applicable when individual-level data is unavailable.

$$\text{Linear Predictor} = \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \widehat{\beta}_3 x_{i3} + \cdots = \sum_{j=1}^{n_{SNP}} \widehat{\beta}_j x_{ij}$$

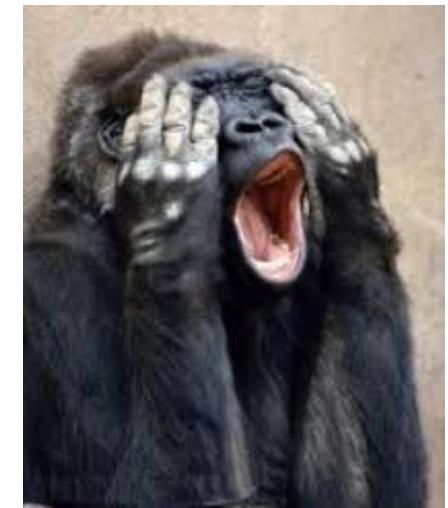
The diagram illustrates the components of the linear predictor. Blue arrows point from the text labels to the corresponding terms in the equation:

- An arrow points from "0, 1 or 2 alleles" to the term  $x_{i1}$ .
- An arrow points from "0, 1 or 2 alleles" to the term  $x_{i2}$ .
- An arrow points from "Which SNPs?" to the term  $\widehat{\beta}_j$ .
- An arrow points from "What weights?" to the term  $x_{ij}$ .

# How should we name those predictors?

Among Human geneticists

- **PRS** - Polygenic risk score
- **PGS** -Polygenic score
- **GRS** - Genetic risk score
- **Genetic score**
- **Genotypic score**
- **Allele score**
- **Profile score**
- ...



Among Animal Breeders...

- **GEBV**: Genomic Estimated Breeding Value

**In this lecture:** I'll use **PGS** or **Linear Predictor** as generic term

Credit: Prof. N. Wray

# Two families of methods...

- 1) Clumping + P-value Thresholding  
(GWAS-derived)
- 2) Summary Statistics-based  
Approximation of Individual-level data  
methods

# CP+T Method (nice and easy\*)

## Input

GWAS summary statistics

Validation sample (genotypes + phenotypes)

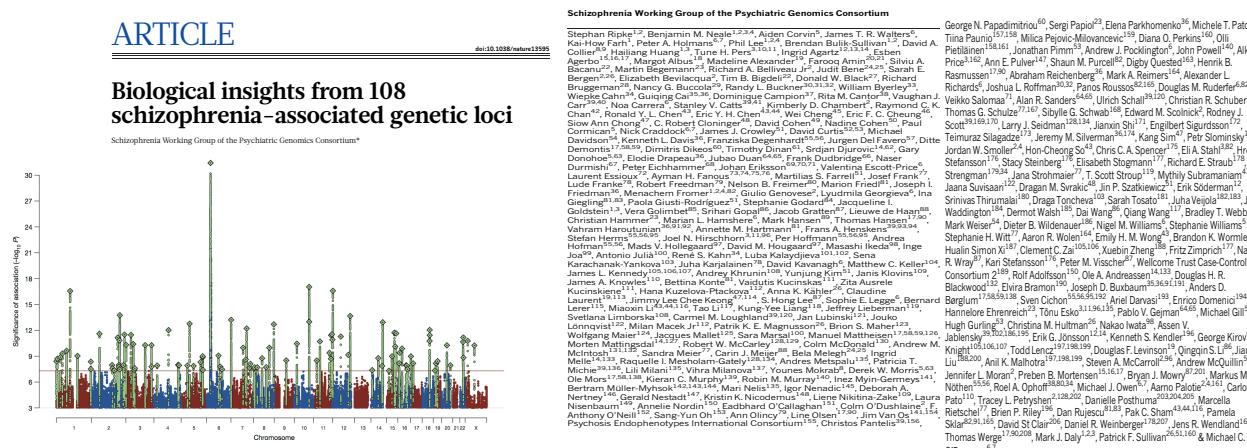
## Algorithm

- 1) Ascertain **independent** SNPs (**C**lumping) associated with the trait at a certain **P**-value threshold (**T**hresholding)
- 2) Calculate PGS for each set
- 3) Select the set with largest prediction accuracy in the validation sample

We'll use this method in the practical.

# Prediction using whole-genome variation

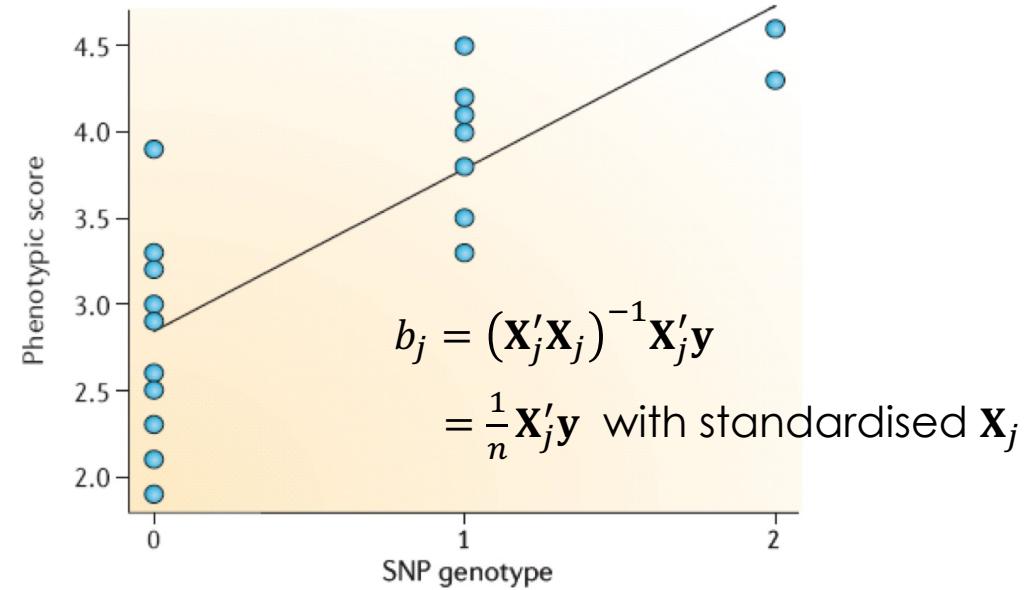
- Best linear unbiased prediction (**G**BLUP)
  - Bayesian alphabet methods (BayesA, B, C, R etc)
  - Require individual-level genotype and phenotype data.
  - Computationally demanding when # individuals and SNPs are large.
  - Data privacy and ethical considerations.



Results generated from a meta-analysis of 52 individual GWAS datasets.

# GWAS summary statistics (summary-level data)

- Marginal SNP effects
- Standard errors
- (Per-SNP) sample sizes
- Effect and other alleles (A1 and A2)
- Effect allele frequencies
- But we need more information...



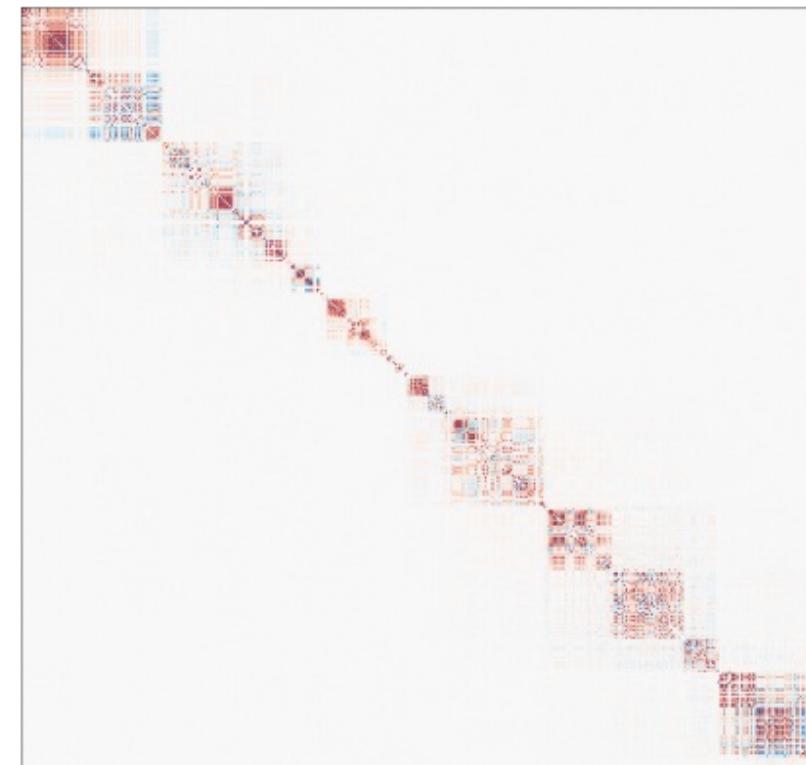
# Linkage disequilibrium (LD) data

- Usually obtained from a reference population
- LD correlation matrix

$$\mathbf{R} = \frac{1}{n} \mathbf{X}' \mathbf{X}$$

with standardised  $\mathbf{X}$

( $E[X_j] = 0$  and  $\text{var}[X_j] = 1$ )



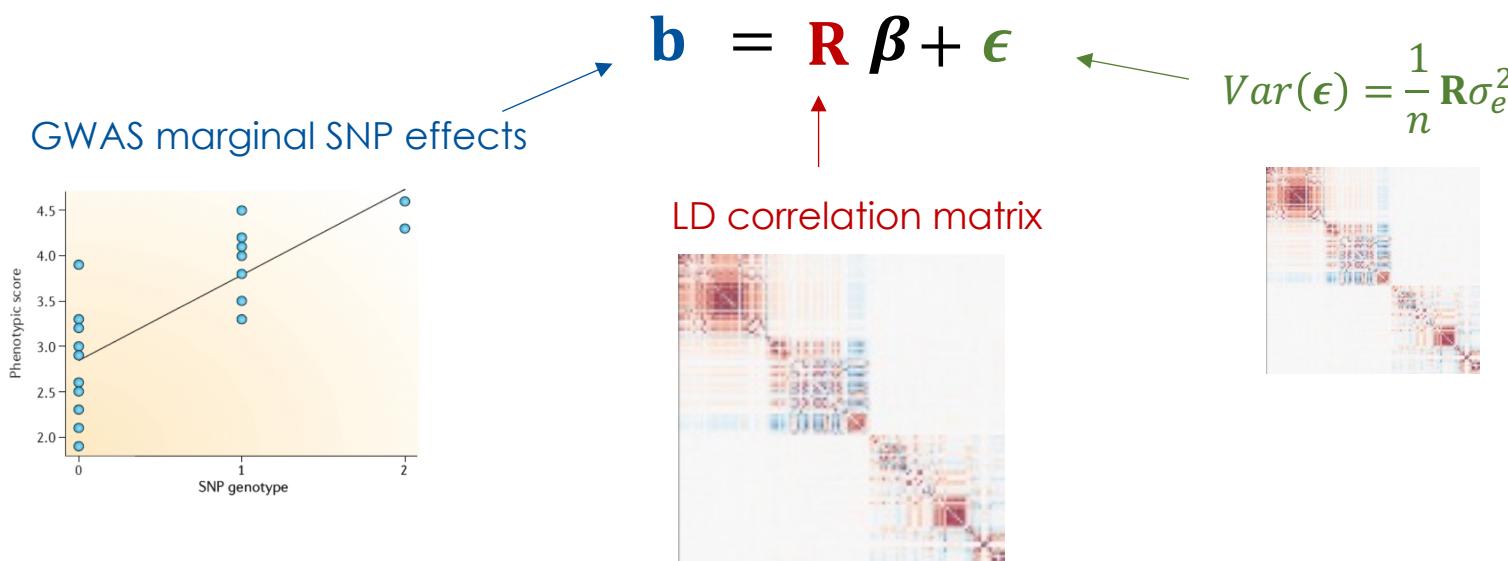
# From individual-level to summary-level model (*ideal conditions*)

Consider an individual-data model with a standardized genotype matrix  $\mathbf{X}$ :

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

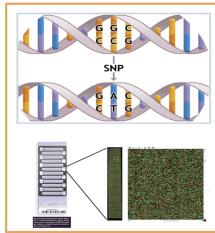
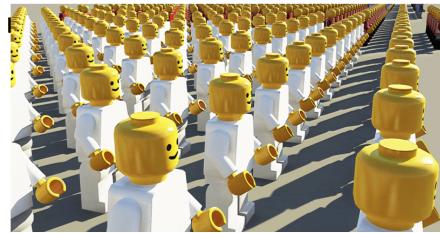
Multiply both sides by  $\frac{1}{n}\mathbf{X}'$  gives

$$\frac{1}{n}\mathbf{X}'\mathbf{y} = \frac{1}{n}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \frac{1}{n}\mathbf{X}'\mathbf{e}$$



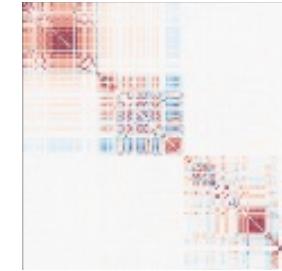
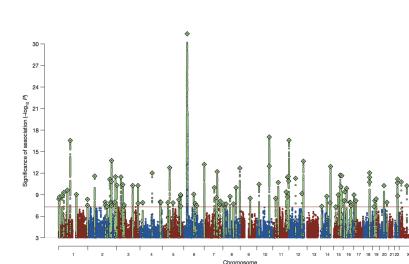
# Individual-level analysis

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$



# Summary-level analysis

$$\mathbf{b} = \mathbf{R}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$



BLUP\*



Bayes

SBLUP

SBayes

\*Best Linear Unbiased Predictor

# SNP-BLUP vs. SBLUP

## SNP-BLUP

- $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$
- $Var(\mathbf{e}) = \mathbf{I}\sigma_e^2$
- Mixed model equations (MME):  
$$[\mathbf{X}'\mathbf{X} + \mathbf{I}\lambda]\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$$
- $\lambda = \frac{\sigma_e^2}{\sigma_\beta^2}$
- $\hat{\boldsymbol{\beta}} = [\mathbf{X}'\mathbf{X} + \mathbf{I}\lambda]^{-1}\mathbf{X}'\mathbf{y}$

↑  
Genotype matrix

↑  
Phenotypes

## SBLUP

- $\mathbf{b} = \mathbf{R}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- $Var(\boldsymbol{\epsilon}) = \frac{1}{n}\mathbf{R}\sigma_e^2$
- Mixed model equations (MME):  
$$[\mathbf{R}'n\mathbf{R}^{-1}\mathbf{R} + \mathbf{I}\lambda]\hat{\boldsymbol{\beta}} = n\mathbf{R}'\mathbf{R}^{-1}\mathbf{b}$$
  
$$[n\mathbf{R} + \mathbf{I}\lambda]\hat{\boldsymbol{\beta}} = n\mathbf{b}$$

$$\hat{\boldsymbol{\beta}} = [n\mathbf{R} + \mathbf{I}\lambda]^{-1}n\mathbf{b}$$

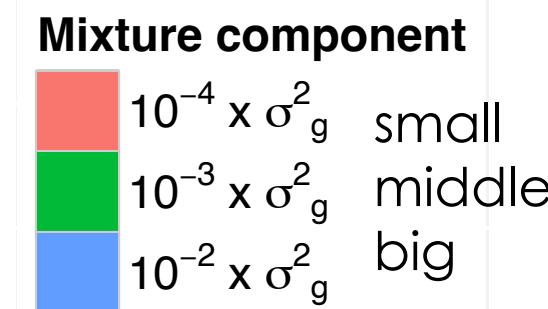
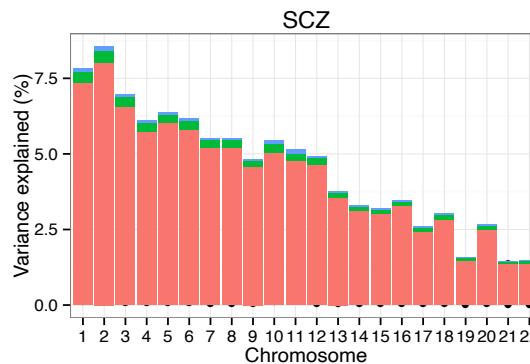
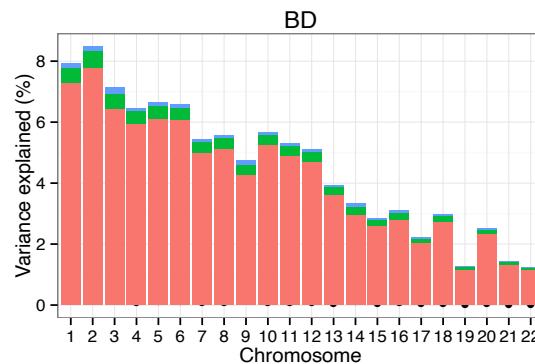
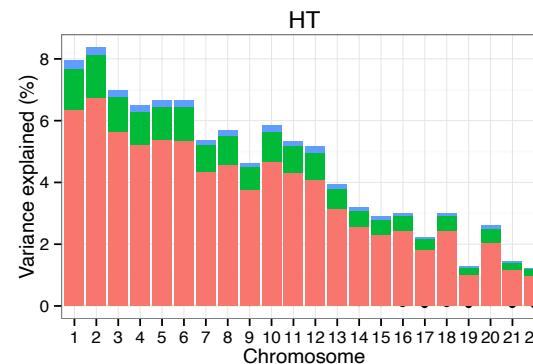
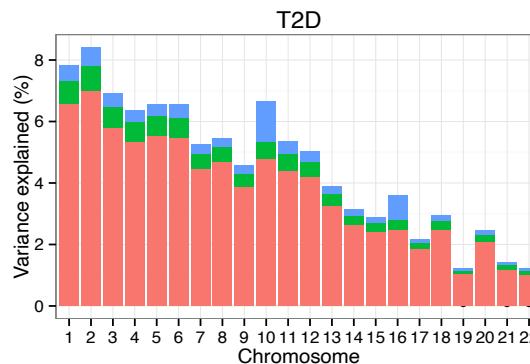
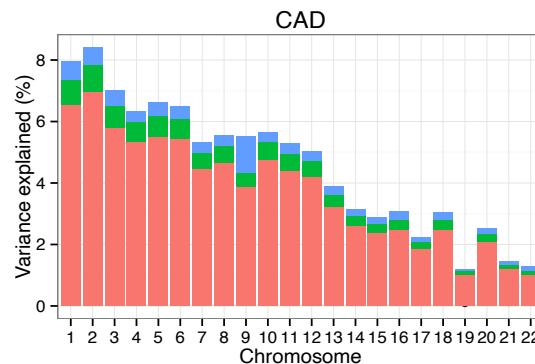
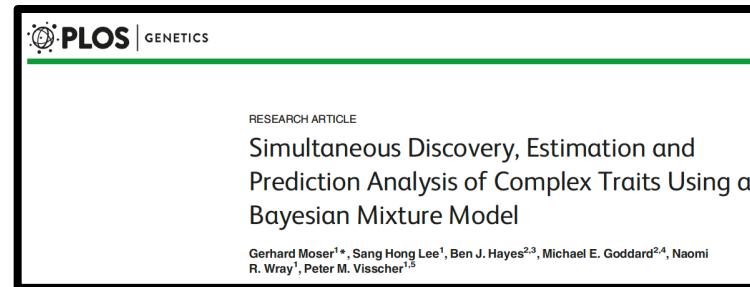
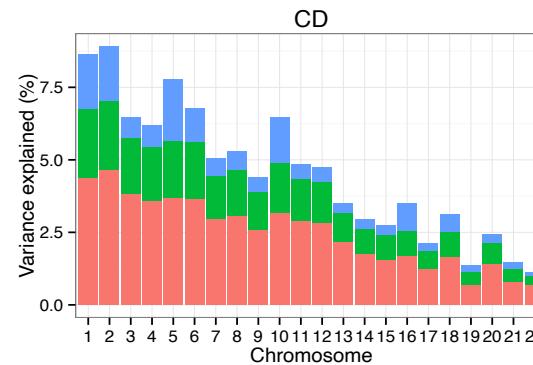
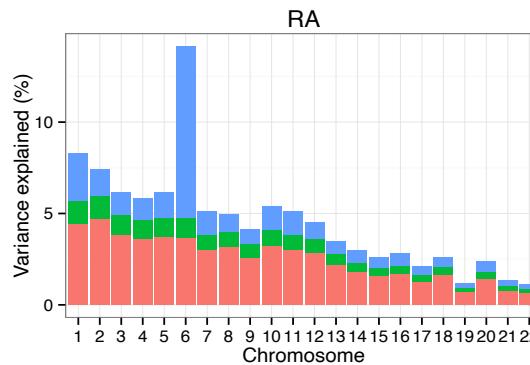
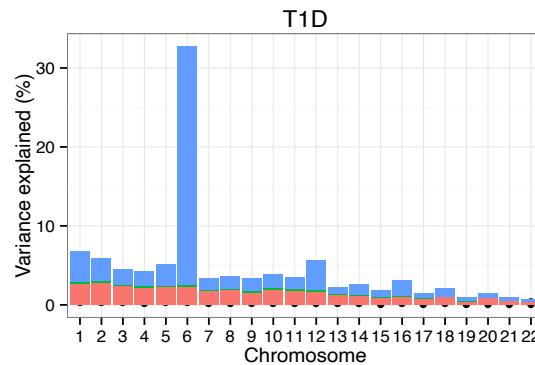
GWAS sample size

LD correlation matrix

Marginal SNP effects

Individual-level data  
Summary-level data

# Trait-specific genetic architecture



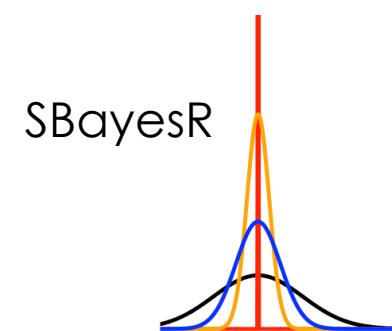
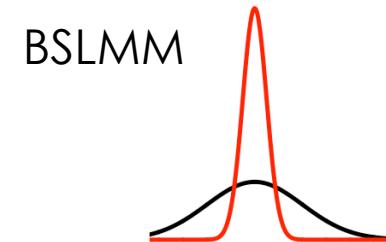
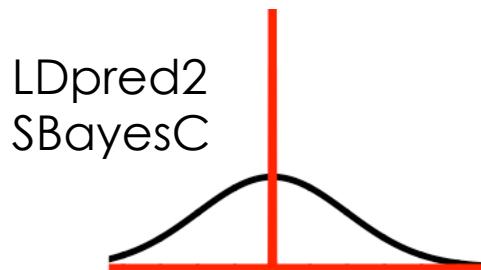
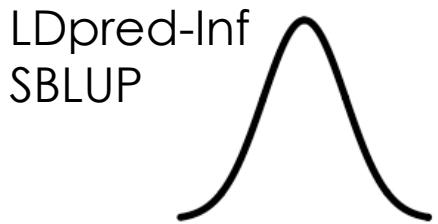
# SBayes

$$\mathbf{b} = \mathbf{R} \boldsymbol{\beta} + \boldsymbol{\epsilon}$$

GWAS marginal SNP effects      LD correlation matrix      SNP joint effects

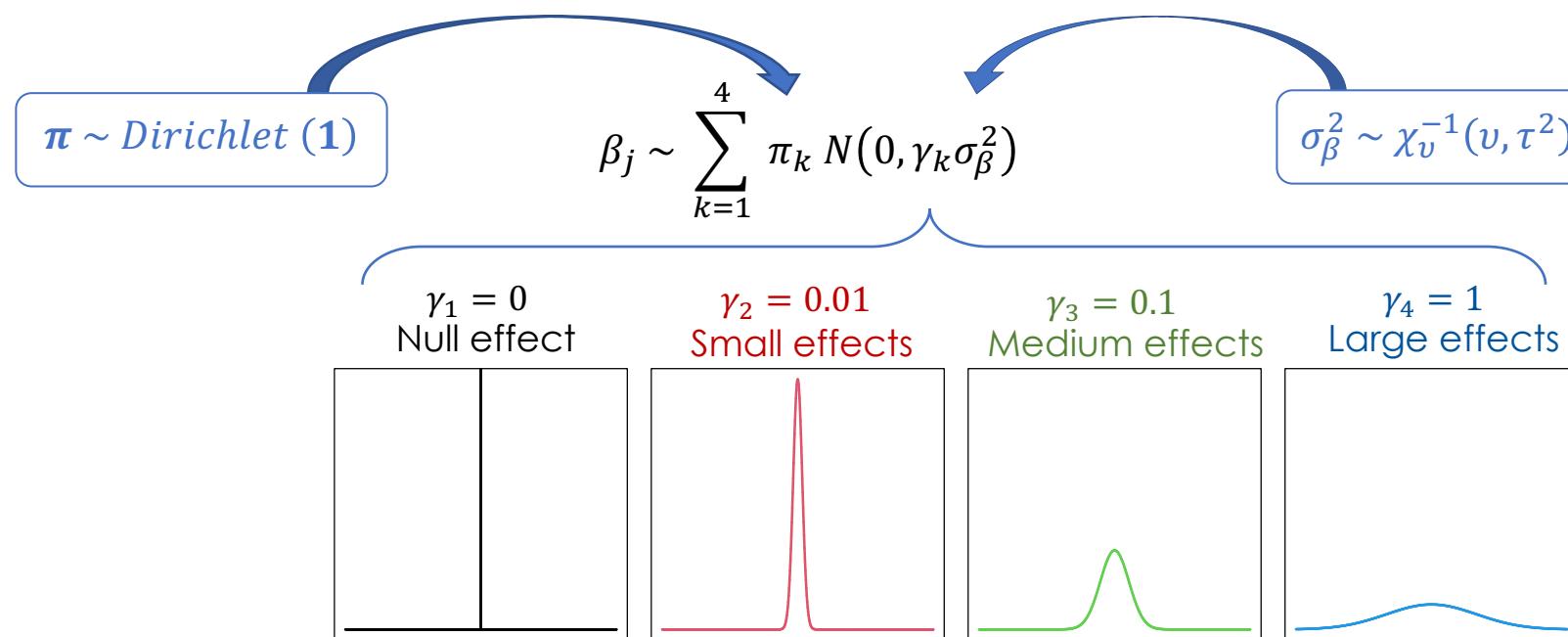
$$Var(\boldsymbol{\epsilon}) = \frac{1}{n} \mathbf{R} \sigma_e^2$$

- Prior distribution for each SNP effect



# SBayesR

- Each SNP effect follows a mixture distribution:



ARTICLE

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

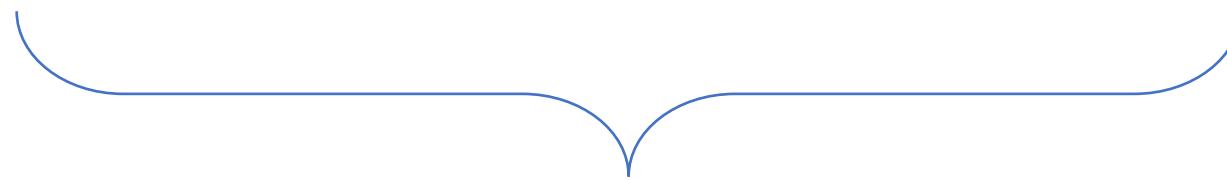
OPEN

Improved polygenic prediction by Bayesian multiple regression on summary statistics

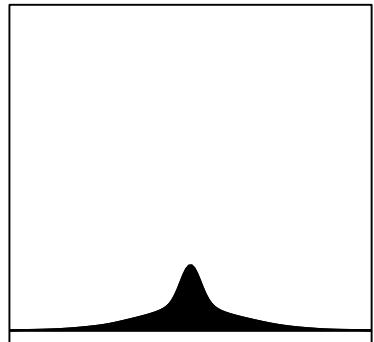
Luke R. Lloyd-Jones<sup>1,9\*</sup>, Jian Zeng<sup>10</sup>, Julia Sidorenko<sup>1,2</sup>, Loïc Yengo<sup>1</sup>, Gerhard Moser<sup>3,4</sup>, Kathryn E. Kemper<sup>1</sup>, Huanwei Wang<sup>10</sup>, Zhili Zheng<sup>1</sup>, Reedik Magi<sup>2</sup>, Tõnu Esko<sup>2</sup>, Andres Metspalu<sup>2,5</sup>, Naomi R. Wray<sup>1,6</sup>, Michael E. Goddard<sup>7</sup>, Jian Yang<sup>10</sup>,<sup>8\*</sup> & Peter M. Visscher<sup>10</sup>\*

# Account for various SNP effect distributions

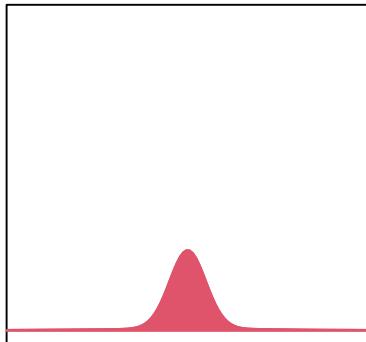
$$\beta_j \sim \pi_1 \quad \text{+} \quad \pi_2 \quad \text{+} \quad \pi_3 \quad \text{+} \quad \pi_4$$



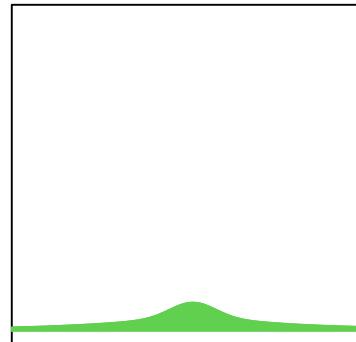
SNP 1



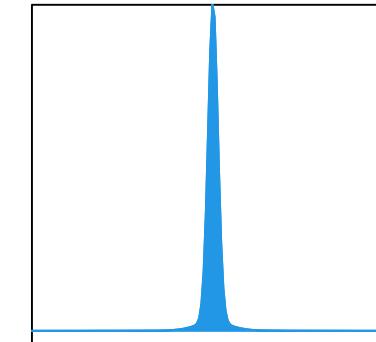
SNP 2



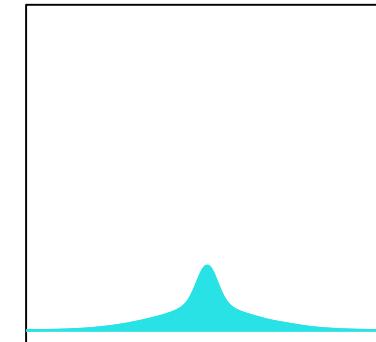
SNP 3



SNP 4



SNP 5



# The posterior distribution of SNP effects

Posterior  $\propto$  Likelihood  $\times$  Prior

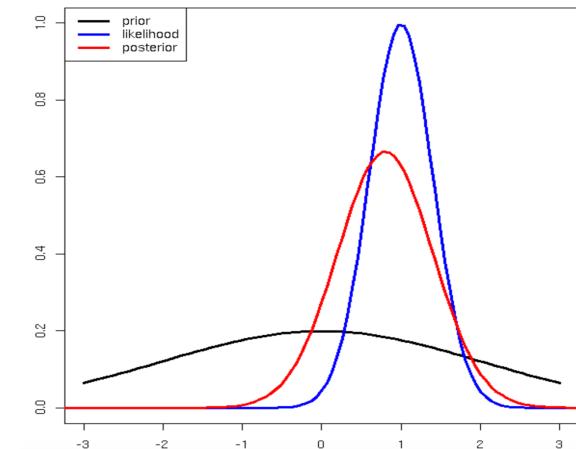
$$f(\beta | \text{Summary data}) \propto f(\text{Summary data} | \beta) \times f(\beta)$$

$$\beta | \mathbf{b} \sim N(\mathbf{C}^{-1} \mathbf{r}, \mathbf{C}^{-1} \sigma_e^2)$$

where

- $\mathbf{C} = n\mathbf{R} + \mathbf{G}^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} = \mathbf{X}'\mathbf{X} + \mathbf{G}^{-1} \frac{\sigma_e^2}{\sigma_\beta^2}$
- $\mathbf{r} = n\mathbf{b} = \mathbf{X}'\mathbf{y}$
- $\mathbf{G} = \text{diag}\{\gamma_j\}$

Summary-level data  
Individual-level data



# Single-site Gibbs sampling

Full conditional distribution for  $\beta_j$

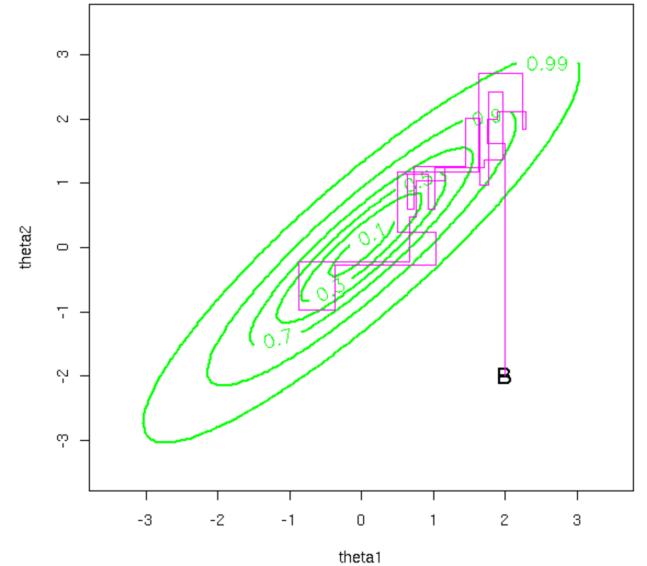
$$f(\beta_j \mid \mathbf{b}, \text{else}) = N\left(\frac{r_j}{C_j}, \frac{\sigma_e^2}{C_j}\right)$$

where

- $r_j = nb_j - \sum_{k \neq j} R_{jk} \beta_k = \mathbf{X}'_j (\mathbf{y} - \sum_{k \neq j} \mathbf{X}_k \beta_k)$

- $C_j = n + \frac{\sigma_e^2}{\gamma_j \sigma_\beta^2} = \mathbf{X}'_j \mathbf{X}_j + \frac{\sigma_e^2}{\gamma_j \sigma_\beta^2}$

Summary-level data  
Individual-level data



# Comparison between individual and summary level algorithms

## Algorithm 1 – Individual level data algorithm

---

Initialise parameters and read genotypes and phenotypes in PLINK binary format  
 Initialise  $\mathbf{y}^* = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$   
**for**  $i := 1$  to number of iterations **do**

- for**  $i := 1$  to  $p$  **do**
- Calculate  $r_j^* = \mathbf{x}_j' \mathbf{y}^*$
- Calculate  $r_j = r_j^* + \mathbf{x}_j' \mathbf{x}_j \beta_j^{(i-1)}$
- Calculate  $\sigma_\epsilon^2 = \sigma_\beta^2 \gamma_{\delta_j=c}$  for each of  $C$  classes (e.g., BayesR C=4 and  $\gamma = (0, 0.0001, 0.001, 0.01)$ )
- Calculate the left hand side  $l_{jc} = \mathbf{x}_j' \mathbf{x}_j + \frac{\sigma_\epsilon^2}{\sigma_\beta^2}$  for each of the  $C$  classes
- Calculate the log densities of given  $\delta_j = c$  using  $\log(\mathcal{L}_c) = -\frac{1}{2} \left[ \log \left( \frac{\sigma_\epsilon^2 l_{jc}}{\sigma_\beta^2} \right) - \frac{r_j^2}{\sigma_\beta^2 l_{jc}} \right] + \log(\pi_c)$ , where  $\pi_c$  is the current
- Calculate the full conditional posterior probability for  $\delta_j = c$  for  $C$  classes with  $\mathbb{P}(\delta_j = c | \theta, \mathbf{y}) = \frac{1}{\sum_{l=1}^C \exp[\log(\mathcal{L}_l) - \log(\mathcal{L}_c)]}$
- Using full conditional posterior probabilities sample class membership for  $\beta_j^{(i)}$  using categorical random variable sampler
- Given class sample SNP effect  $\beta_j^{(i)}$  from  $N \left( \frac{r_j}{l_{jc}}, \frac{\sigma_\epsilon^2}{l_{jc}} \right)$
- Given SNP effect adjust corrected phenotype side  $(\mathbf{y}^*)^{(i)} = (\mathbf{y}^*)^{(i-1)} - \mathbf{x}_j (\beta_j^{(i)} - \beta_j^{(i-1)})$

**od**

Sample update from full conditional for  $\sigma_\beta^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\beta = v_\beta + q$  and  $\tilde{S}^2_\beta = \frac{v_\beta S_\beta^2 + \sum_{j=1}^q \frac{\beta_j^2}{l_{jc}}}{v_\beta + q}$ , where  $q$  is the number of non-zero variants

Sample update from full conditional for  $\sigma_\epsilon^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\epsilon = n + v_e$  and scale parameter  $\tilde{S}_\epsilon^2 = \frac{SSE + v_e S_\epsilon^2}{n + v_e}$  and  $SSE = \mathbf{y}' \mathbf{y}$

Sample update from full conditional for  $\pi$ , which is Dirichlet( $C, \mathbf{c} + \boldsymbol{\alpha}$ ), where  $\mathbf{c}$  is a vector of length  $C$  and contains the counts of the number of variants in each variance class and  $\boldsymbol{\alpha} = (1, \dots, 1)$

Calculate genetic variance for  $h_{SNP}^2$  calculation using  $\sigma_g^2 = \text{Var}(\mathbf{X}\boldsymbol{\beta})$

Calculate  $h_{SNP}^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_\epsilon^2}$

---

## Algorithm 2 Summary data algorithm

---

Initialise parameters and read summary statistics  
 Reconstruct  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$  from summary statistics and LD reference panel  
 Calculate  $\mathbf{r}^* = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta}$   
**for**  $i := 1$  to number of iterations **do**

- for**  $i := 1$  to  $p$  **do**
- Calculate  $\mathbf{r}_j = \mathbf{r}_j^* + \mathbf{x}_j' \mathbf{x}_j \beta_j$
- Calculate  $\sigma_\alpha^2 = \sigma_\beta^2 \gamma_{\delta_j=c}$  for each of  $C$  classes (e.g., SBayesR C=4 and  $\gamma = (0, 0.01, 0.1, 1)'$ )
- Calculate the left hand side  $l_{jc} = \mathbf{x}_j' \mathbf{x}_j + \frac{\sigma_\epsilon^2}{\sigma_\beta^2}$  for each of the  $C$  classes
- Calculate the log densities of given  $\delta_j = c$  using  $\log(\mathcal{L}_c) = -\frac{1}{2} \left[ \log \left( \frac{\sigma_\epsilon^2 l_{jc}}{\sigma_\beta^2} \right) - \frac{r_j^2}{\sigma_\beta^2 l_{jc}} \right] + \log(\pi_c)$ , where  $\pi_c$  is the current
- Calculate the full conditional posterior probability for  $\delta_j = c$  for  $C$  classes with  $\mathbb{P}(\delta_j = c | \theta, \mathbf{y}) = \frac{1}{\sum_{l=1}^C \exp[\log(\mathcal{L}_l) - \log(\mathcal{L}_c)]}$
- Using full conditional posterior probabilities sample class membership for  $\beta_j^{(i)}$  using categorical random variable sampler
- Given class sample SNP effect  $\beta_j^{(i)}$  from  $N \left( \frac{r_j}{l_{jc}}, \frac{\sigma_\epsilon^2}{l_{jc}} \right)$
- Given SNP effect adjust corrected right hand side  $(\mathbf{r}^*)^{(i+1)} = (\mathbf{r}^*)^{(i)} - \mathbf{x}_j' (\beta_j^{(i+1)} - \beta_j^{(i)})$ .  $\mathbf{x}_j'$  is the  $j$ th column of  $\mathbf{X}'\mathbf{X}$ .

**od**

Sample update from full conditional for  $\sigma_\alpha^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\alpha = v_0 + q$  and  $\tilde{S}^2_\alpha = \frac{v_0 \tilde{v}_0 + \sum_{j=1}^q \frac{\beta_j^2}{l_{jc}}}{v_0 + q}$ , where  $q$  is the number of non-zero variants

Sample update from full conditional for  $\sigma_\epsilon^2$  from scaled inverse chi-squared distribution  $\tilde{v}_\epsilon = n + v_e$  and scale parameter  $\tilde{S}_\epsilon^2 = \frac{SSE + v_e S_\epsilon^2}{n + v_e}$  and  $SSE = \mathbf{y}' \mathbf{y} - \boldsymbol{\beta}' \mathbf{r}^* - \boldsymbol{\beta}' \mathbf{X}' \mathbf{y}$

Sample update from full conditional for  $\pi$ , which is Dirichlet( $C, \mathbf{c} + \boldsymbol{\alpha}$ ), where  $\mathbf{c}$  is a vector of length  $C$  and contains the counts of the number of variants in each variance class.

Calculate genetic variance for  $h_{SNP}^2$  calculation using  $\sigma_g^2 = MSS/n$ , where  $MSS = \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{y} - \hat{\boldsymbol{\beta}}' \mathbf{r}^*$

Calculate  $h_{SNP}^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_\epsilon^2}$

---

$\mathbf{X}'\mathbf{y}$  and  $\mathbf{X}'\mathbf{X}$  are sufficient statistics in the algorithm!

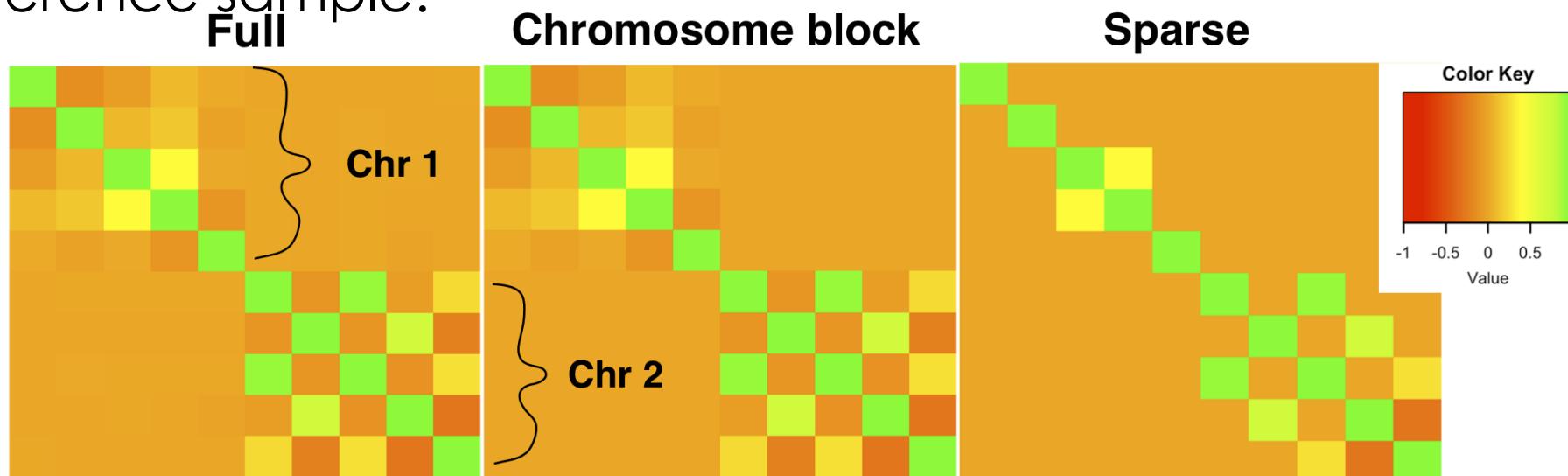
# Scaling the summary statistics

- Derivation is based on standardised genotypes.
- GWAS effects are usually in unit of allele count (per-allele effects).
- Need to rescale GWAS effects by  $b_j^* = s_j b_j$  where  $s_j = \sqrt{2p_j(1 - p_j)}$
- Allele frequency  $p_j$  could be reported with errors.
- Use of standard error is more robust:

$$\bullet SE_j = \sqrt{\frac{1}{\mathbf{x}'_j \mathbf{x}_j} \hat{\sigma}_e^2} = \sqrt{\frac{1}{nh_j} (\sigma_y^2 - s_j^2 b_j^2)} \quad \rightarrow \quad s_j = \sqrt{\frac{\sigma_y^2}{nSE_j^2 + b_j^2}} \approx \sqrt{\frac{\sigma_y^2}{nSE_j^2}}$$

# LD matrix from a reference

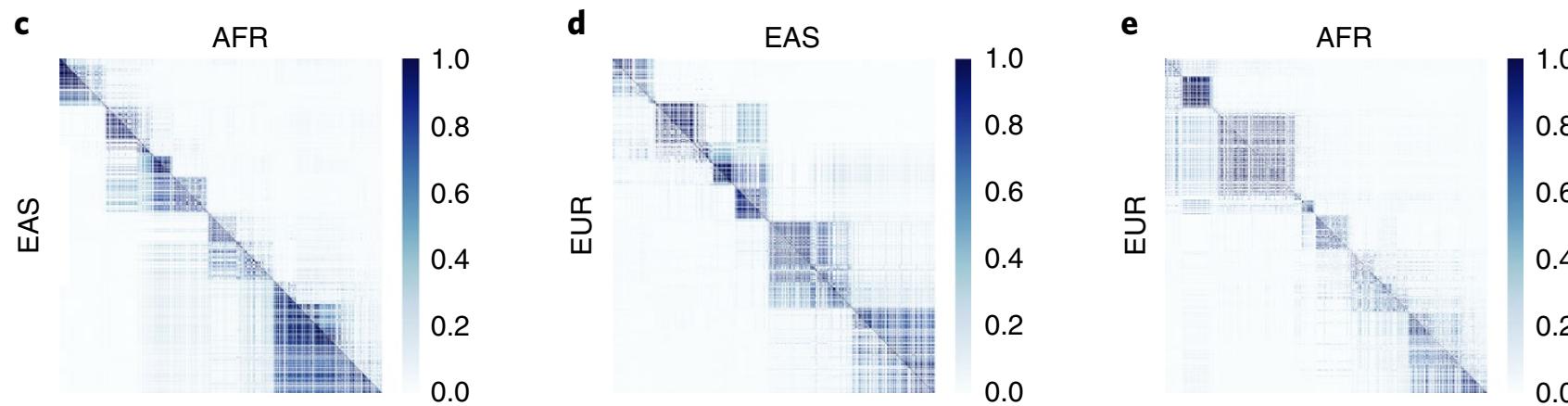
- Often cannot use genome-wide full LD matrix from the GWAS sample.
- Use a reduced (banded, sparse, or shrunk) LD matrix from a reference sample.



Challenging to extend to non-human cases,  
where cross-chromosomal LD is NOT negligible

# Implicit assumptions

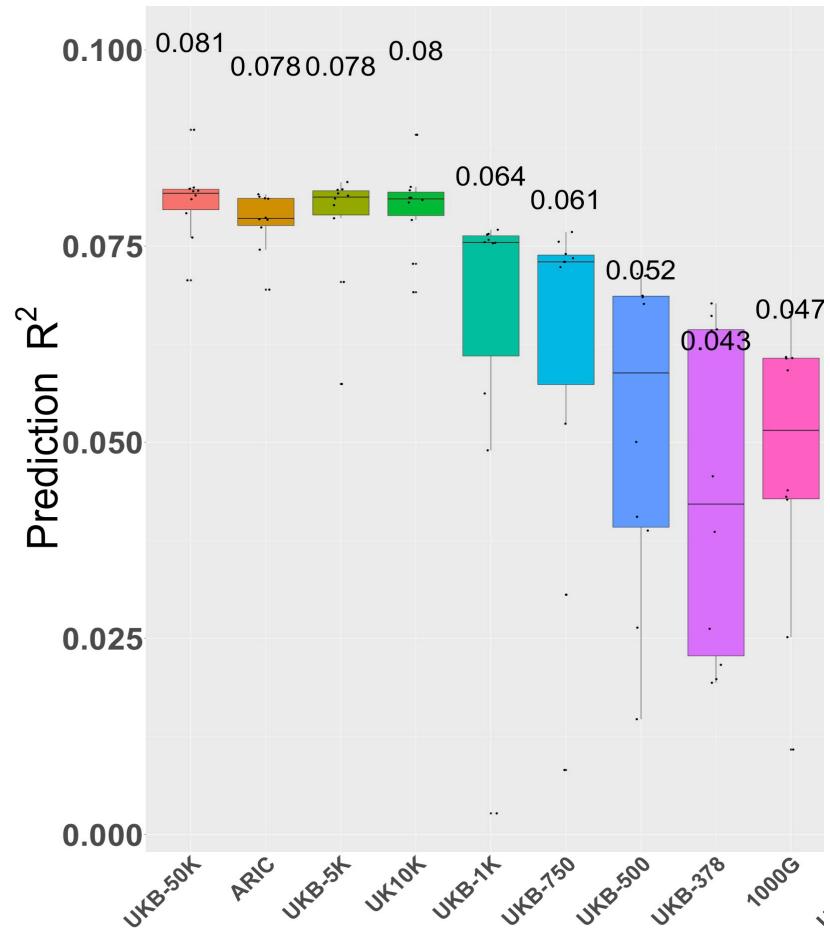
- LD reference population matches with GWAS population in genetics
  - No systematic differences in LD → same ancestry and population structure
  - Minimum sampling variance in LD → LD ref sample size cannot be too small



# Implicit assumptions

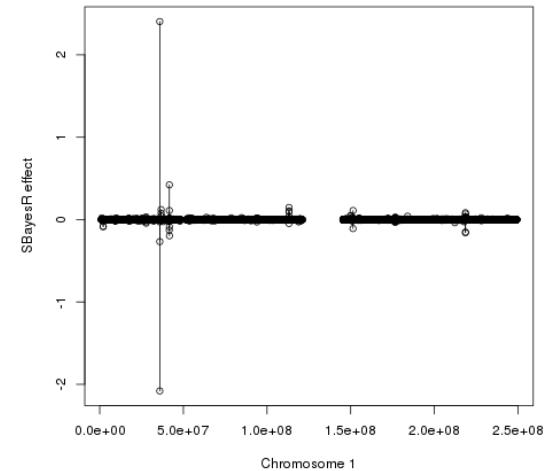
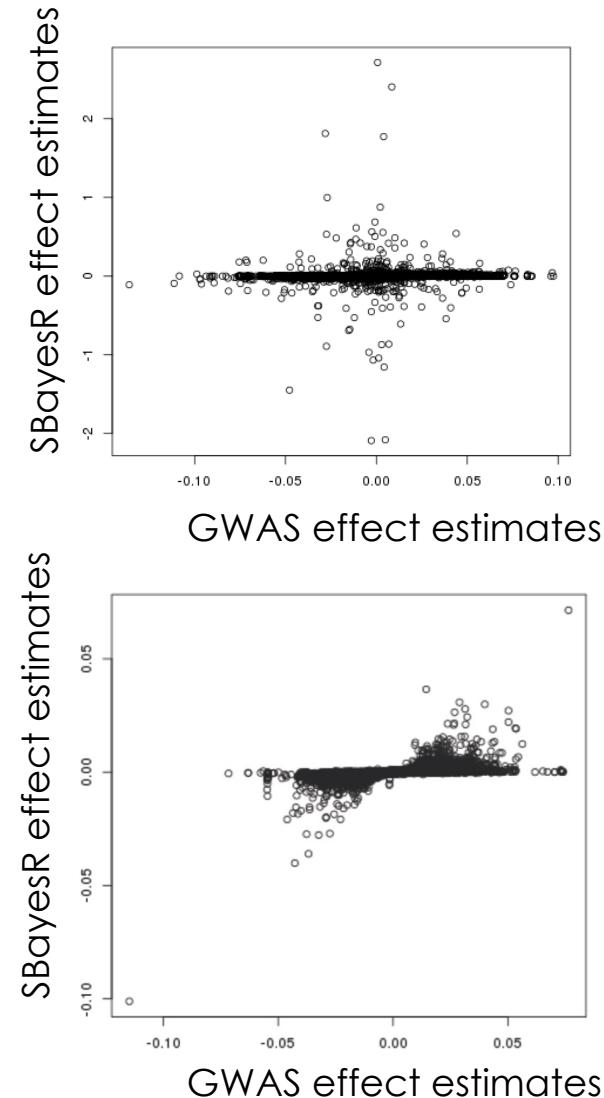
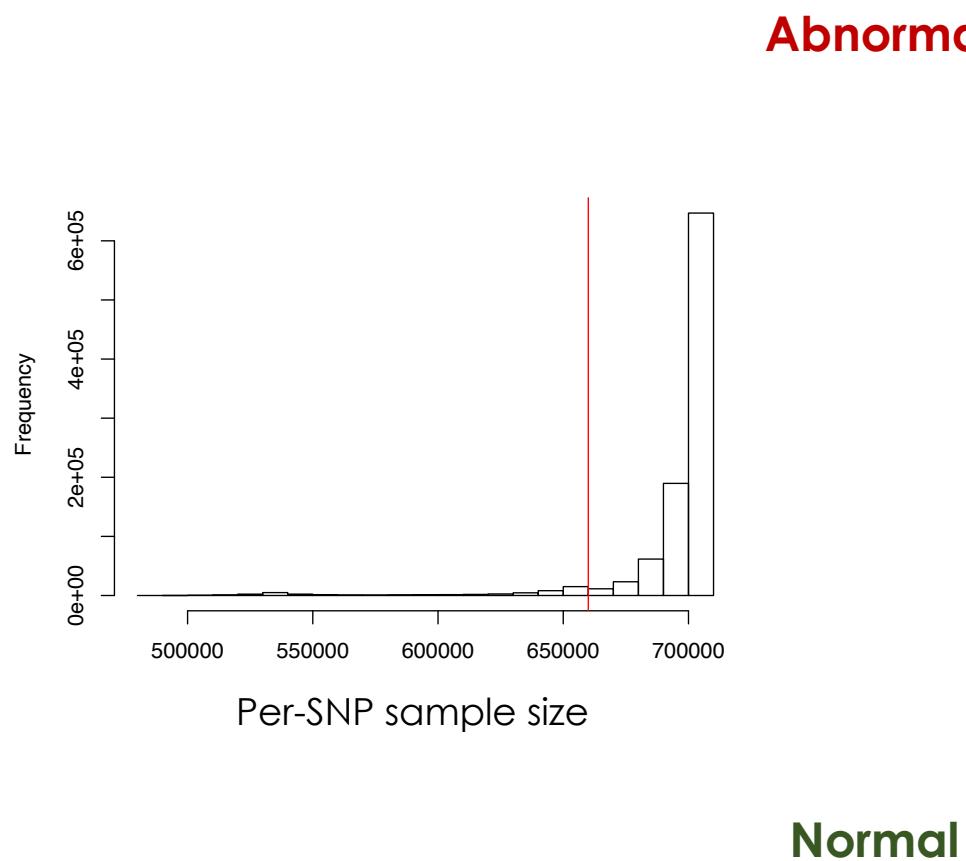
- LD reference population matches with GWAS population in genetics
  - No systematic differences in LD → same ancestry and population structure
  - Minimum sampling variance in LD → LD ref sample size cannot be too small
- GWAS data are collected on the same set of individuals
  - Often an issue in GWAS meta-analysis
  - Consistent genotyping platforms or imputation panels across cohorts
  - Remove SNP outliers in per-SNP sample size
- Violation these assumptions can cause model misspecification, resulting in attenuated prediction accuracy or even failure to reach convergence.

# Influence of choice of LD dataset



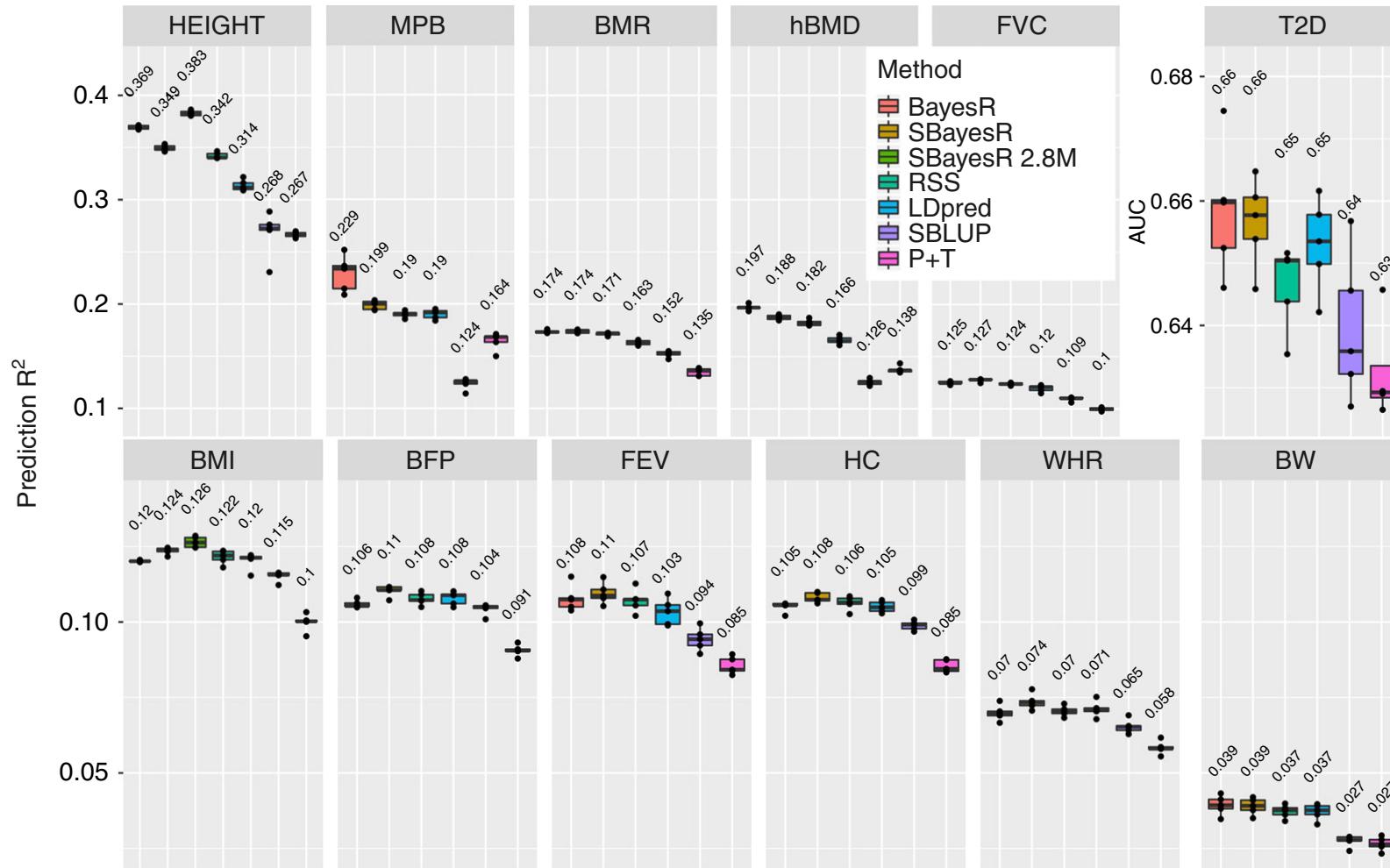
Lloyd-Jones & Zeng et al. 2019 NC

# Influence of heterogeneity in per-SNP sample size



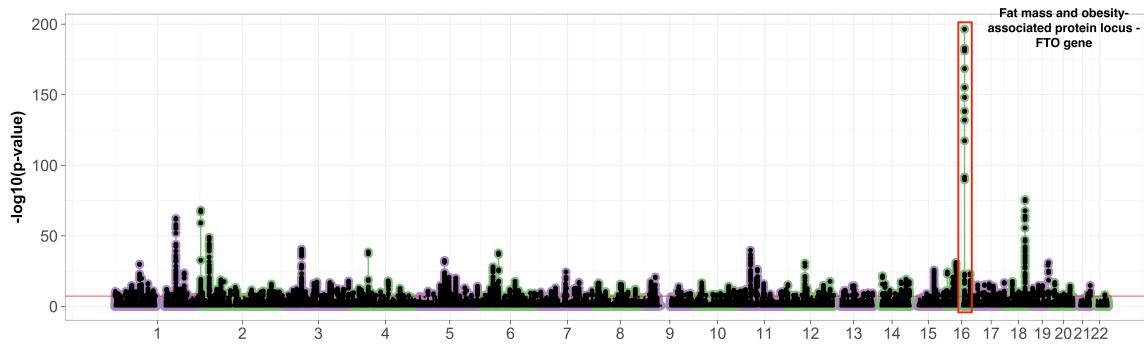
[https://cnsgenomics.com/software/gctb/  
#FAQ](https://cnsgenomics.com/software/gctb/#FAQ)

# Method comparison

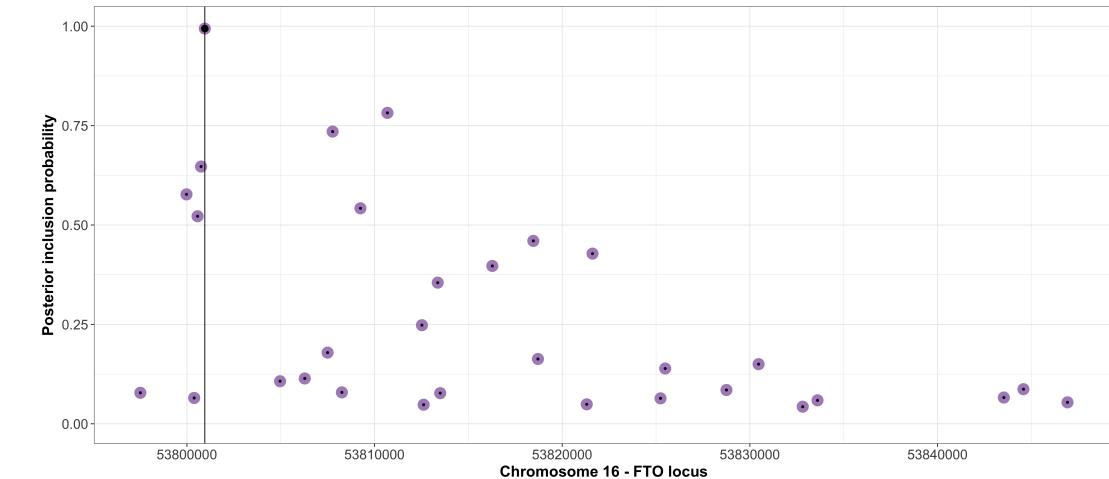
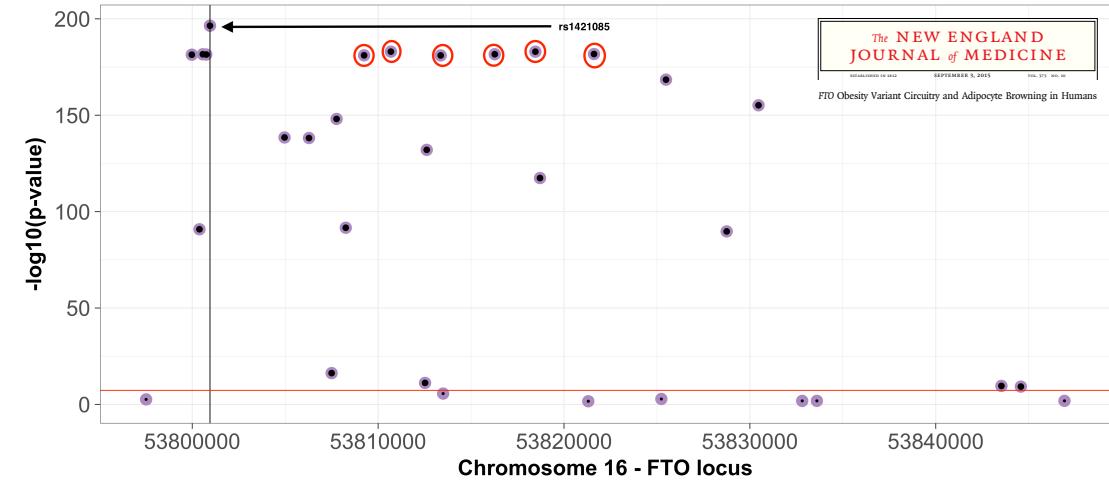


# Fine-mapping

Real data - body mass index

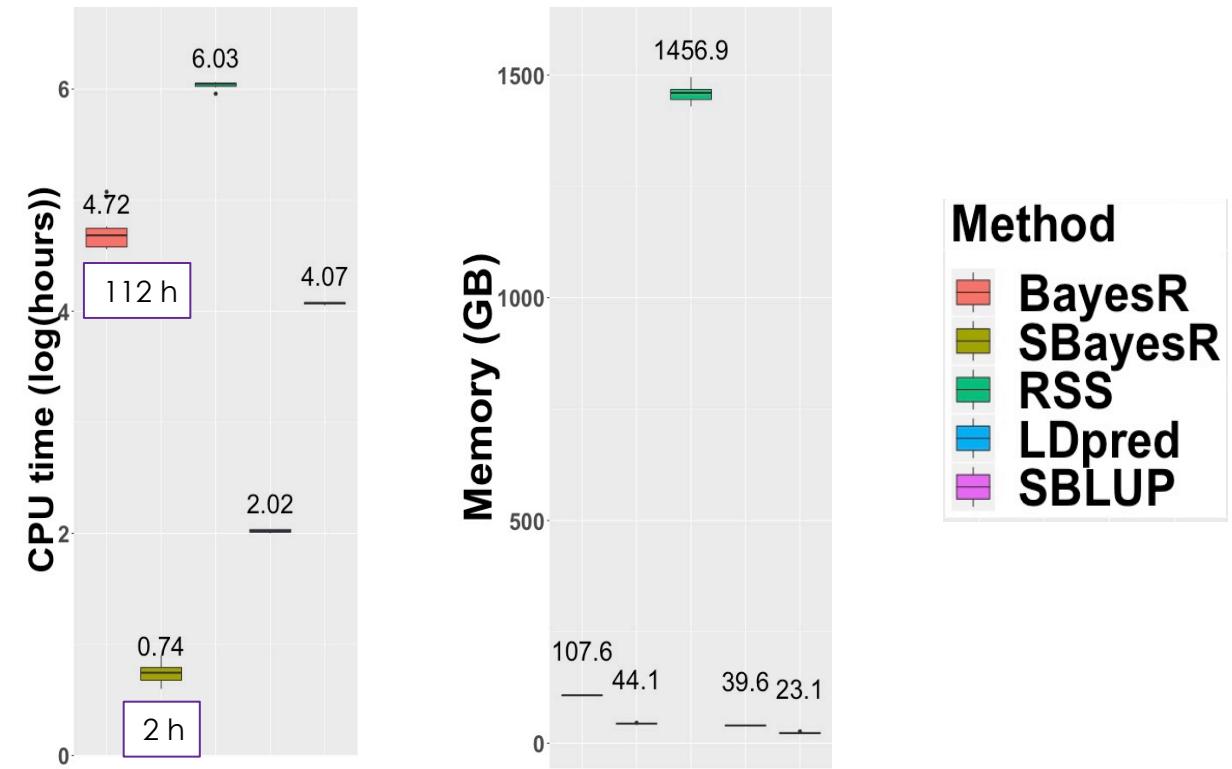


Real data - FTO - GWAS



# Computational efficiency

- 100k individuals
- 1M SNPs
- Resource consumption for summary-data-based methods is independent of sample size once GWAS summary statistics are obtained.



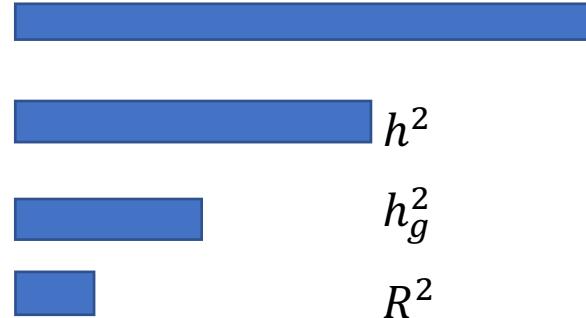
# Summary

- Summary-level methods exploit the full power of GWAS of large sample sizes for polygenic prediction.
- Approximations to the individual-level methods.
- Free from limitation of data accessibility.
- Computationally efficient.
- Flexible to incorporate other information, such as gene expression.

# Limits of PGS

# Understanding PGS limits

- ❖ PGS have a **theoretical** upper limit dependent on the **heritability of the trait** (how much of the variance of trait values between people is attributed to genetic factors).
- ❖ PGS have a **technical** upper limit associated with the proportion of **variance tagged** by the DNA variants measured.
- ❖ PGS have a **practical** upper limit dependent on the **sample size of the discovery sample** used to estimate effect sizes of risk alleles, and the **quality** of the discovery sample.
- ❖ PGS can be pushed closer to the technical upper limit by the **statistical methodology** used to generate the optimal weighting given to the risk alleles, and new methods integrate new biological data.



## Schizophrenia

**Max:**  
25% Liability  
AUC 0.84

**Current:**  
11% Liability  
AUC 0.74

Polygenic scores cannot be highly accurate predictors of phenotypes

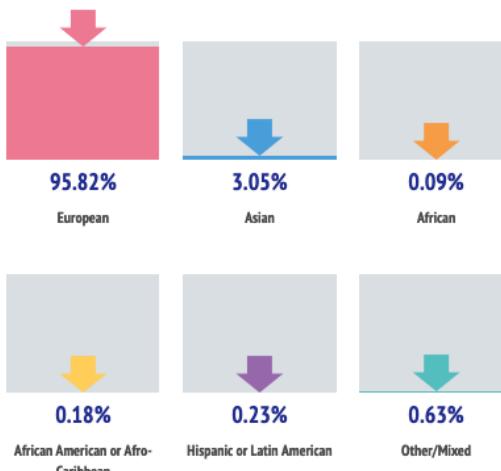
# **Loss of prediction accuracy of polygenic predictors across populations**

Additional Information

METRIC  By ParticipantsSTAGE  Discovery

## Total GWAS participants diversity

Version 1.0.0. Last check for data: 2022-01-18 09:34:08 .



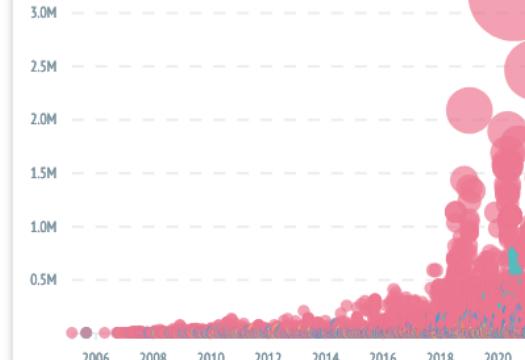
## Ancestry over time by parent term

Discovery Stage

All parent terms

OR

Search for one or more traits



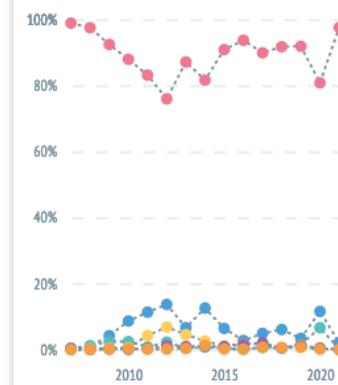
VIEW ALL

- European
- Asian
- African
- African American or Afro-Caribbean
- Hispanic or Latin American
- Other/Mixed

## Participants across all parent terms

Discovery Stage

All ancestries

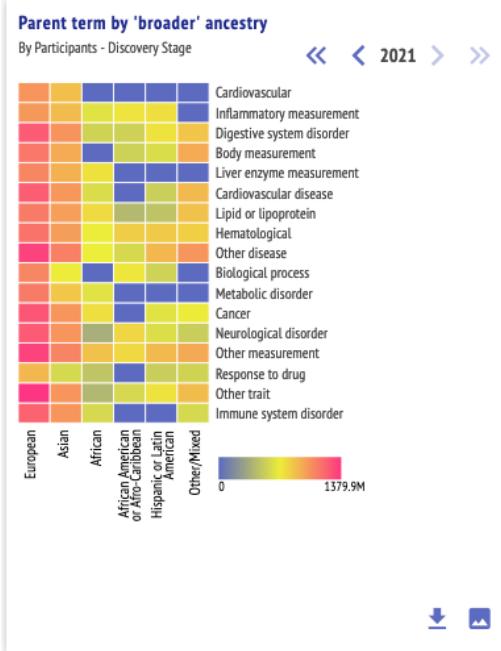
 Include not recorded

In 2022...

&gt;95% European

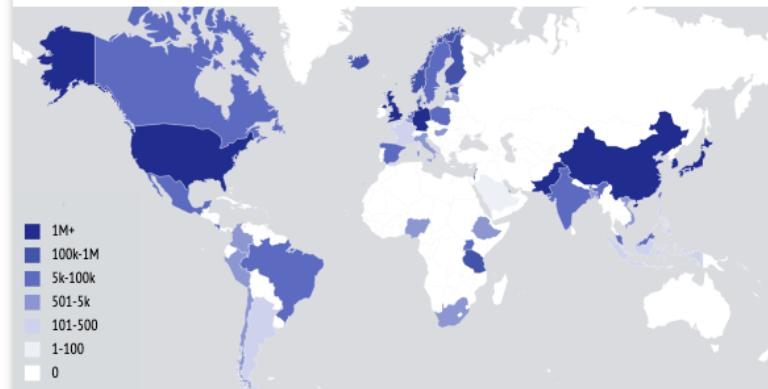
&lt;4% Asian

&lt;0.1% African



## Participants by country (all parent terms)

Both Stages



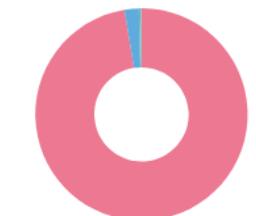
« « 2021 » »

## Participants by ancestry

Discovery Stage

Click to show associations discovered

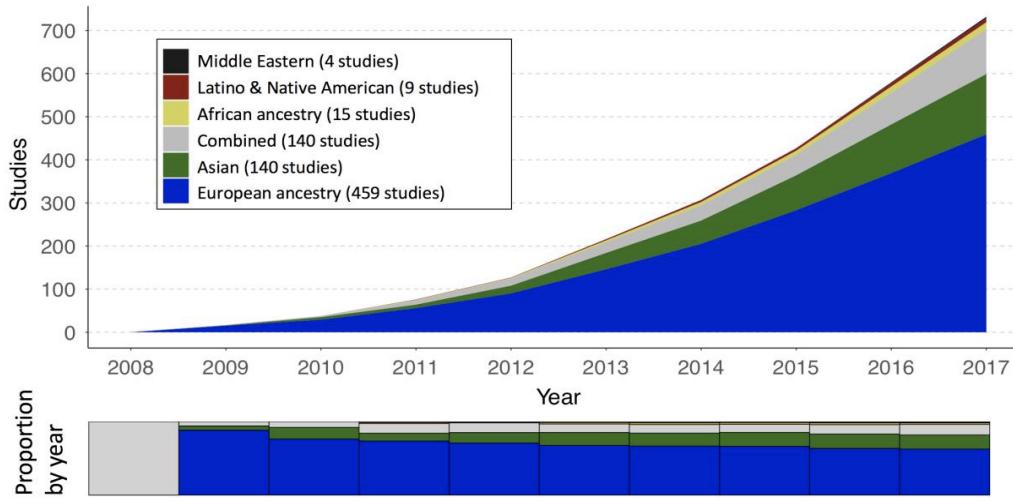
All parent terms



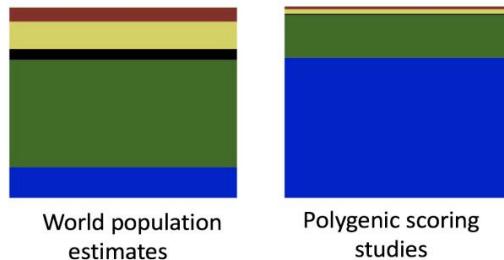
« « 2021 » »

# Euro-centric GWASs bias

A.

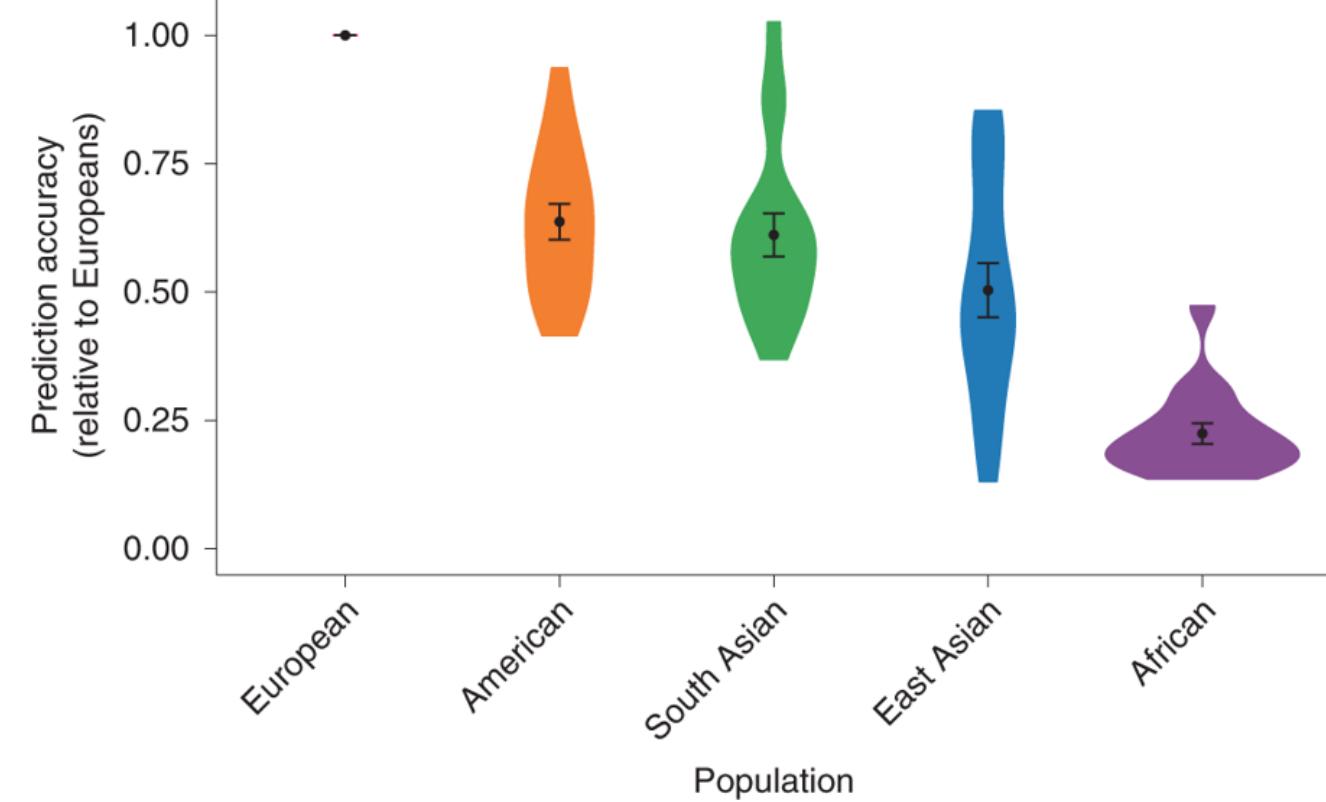


B.



Duncan et. al., *Nature Communications*, 2019

# Loss of prediction accuracy



Martin et. al., *Nature Genetics*, 2019

We are missing GWAS in non-European ancestries populations

Ancestry-specific causal variants (worse case scenario)

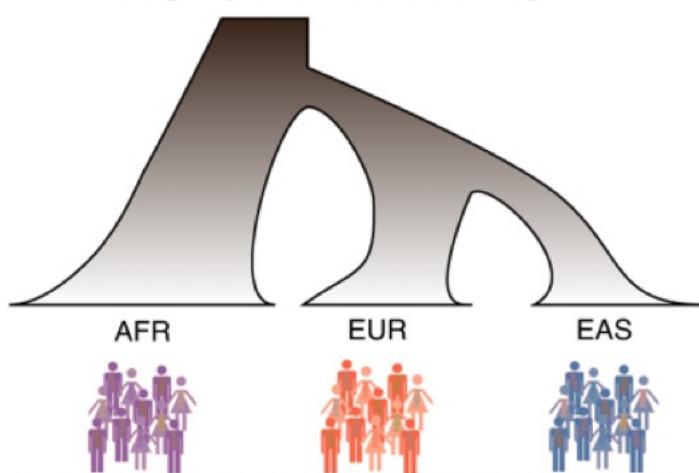
Same causal variants but different effect sizes

Same causal variants, same effect sizes, same heritability, but different haplotype frequency

...

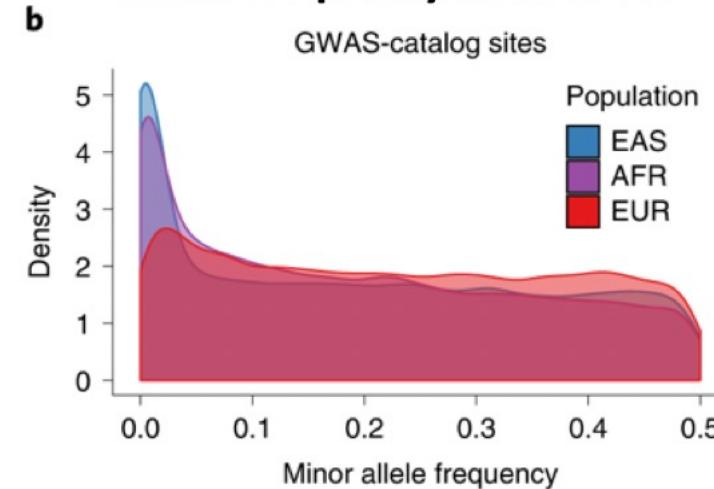
# Factors affecting PGS accuracy disparity

## a Demographic relationships

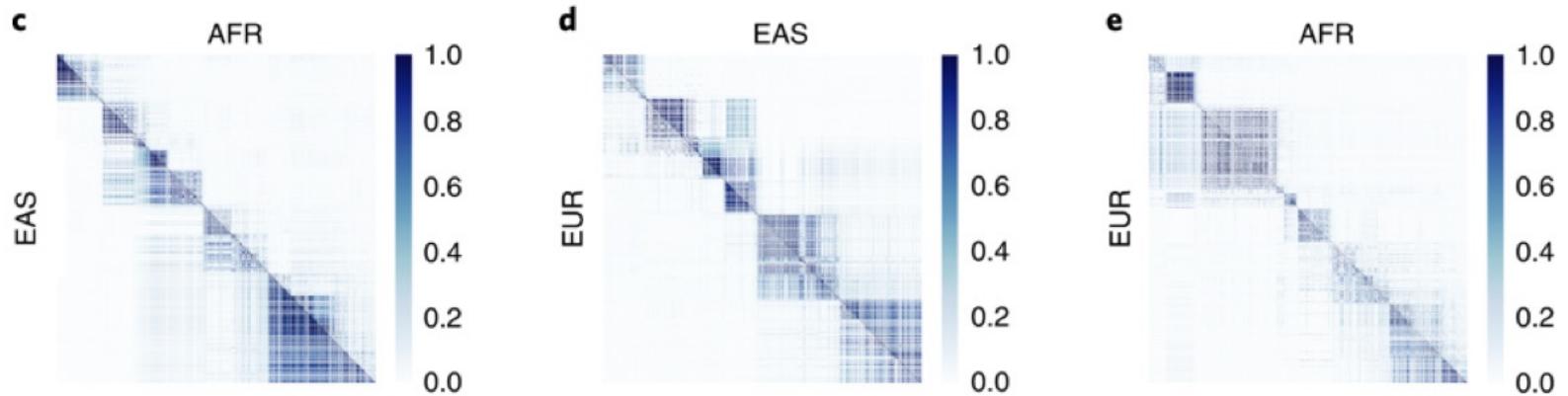


Martin et. al., *Nature Genetics*, 2019

## b Allele frequency differences



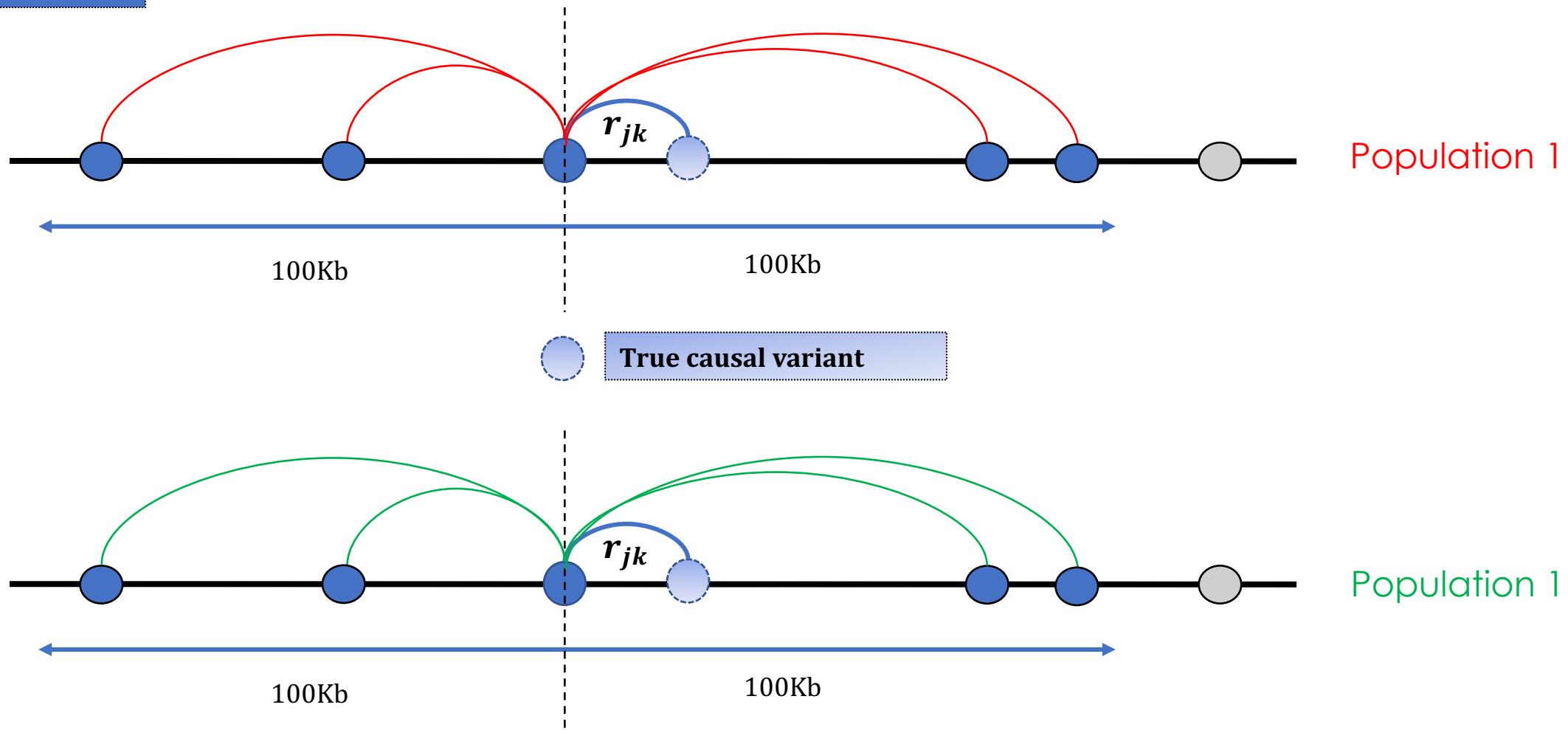
## c Local LD patterns between population pairs



# LD difference



SNPs in PGS



# Factors affecting PGS accuracy disparity



How to **quantify** the **loss of accuracy**  
attributable to MAF and LD?



ARTICLE

<https://doi.org/10.1038/s41467-020-17719-y>

OPEN

Theoretical and empirical quantification of the accuracy of polygenic scores in ancestry divergent populations

Ying Wang  <sup>1</sup>, Jing Guo<sup>1</sup>, Guiyan Ni<sup>1</sup>, Jian Yang  <sup>1,2</sup>, Peter M. Visscher  <sup>1</sup> & Loic Yengo  <sup>1</sup>✉

# 1.1 Theory

Relative Accuracy (RA)



◻  $\frac{R_2^2}{R_1^2}$  = Function of

$$r_b^2 h_2^2 / h_1^2,$$

$$p_{k1}^{(t)}, p_{k2}^{(t)}, \frac{\text{var}(\hat{y}_1)}{\text{var}(\hat{y}_2)}$$

and  $r_{jk,1}, r_{jk,2}$

$r_b^2$ : squared genetic correlation

$r_{jk,l}$ : LD correlation between  $j$ -th causal variant

$h_l^2$ : trait heritability in Population  $l$

and  $k$ -th tagSNP in Population  $l$

$p_{k,l}$ : MAF in Population  $l$

$\text{var}(\hat{y}_l)$ : variance of PGS in Population  $l$

# Conclusions

- (1) Relative accuracy is lower than expected if assuming GWAS SNPs are causal variants (=> wrong conclusion)
- (2) LD and MAF differences between populations account for the majority of the loss of prediction accuracy
- (3) Empirical data are consistent with causal variants being largely shared across populations.

Feature Review

# Genetic prediction of complex traits with polygenic scores: a statistical review

Ying Ma<sup>1</sup> and Xiang Zhou  <sup>1,2,\*</sup>

# More methods...

# Methods = What SNPs + What Weights (Linear Models)

- Clumping + Thresholding (Historical – less popular now)
- COJO + Thresholding
- Random Regression (Genome-wide) Methods
  - Individual-level data: **BLUP**, Bayesian Regression (e.g., **BayesR**, BOLT-LMM)
  - Summary-statistics methods



# SS Parametric Methods\*

- SBayes- family: no tuning but need to specify the model (e.g., S, C, R, etc.)
- SBLUP (tuning parameter  $\lambda = M(1 - h^2)/(Nh^2)$  is fixed)
- LDpred infinitesimal  $\Leftrightarrow$  SBLUP  
(tuning parameter  $\lambda$  is selected on a grid – **need tuning sample**)
- LDpred  $\Leftrightarrow$  SBayesC parameters  $\pi$  and  $h^2$  either estimated with MCMC or selected on a grid – parametrization slightly differs from SBayesC
- LASSOSUM (penalizes the  $L_1$  norm): tuning parameter  $\lambda$  and  $s$  (shrinkage on the LD matrix) selected on a grid – **need tuning sample** or GWAS summary statistics.
- Deterministic Bayesian Sparse Linear Mixed Model (DBSLMM) – combines C+T, COJO, LDSC and SBLUP

Brisbane Teams  
Melbourne Teams  
Arhus Teams  
(Speed + Vihhjalmsson)

Arhus Teams  
(Speed + Vihhjalmsson + Prive)

Pak Sham Group (2017)

Zhou Group (2020)

# SS Nonparametric Methods

- Continuous Shrinkage Prior (Dirichlet Prior)  $\Leftrightarrow$  infinite number of components in mixture
  - DPR method (2017)  
Zhou Group [Michigan]
  - PRS CS (2019, 2022)  
Ge Group [Boston]
- Non-parametric Shrinkage [Partitioned GWAS summary stats]
  - Sunyaev Group (2020)  
[Boston]
- Semi-Parametric: MEGA-PRS (combines multiple methods and learn weights from validation sample)

# Parametric vs Nonparametric

## Bayesian Parametric

- Quest for the “Best” prior -- flexibility
- Flexibility has cost
- Modelling flexibility to incorporate additional information
- MCMC can be expensive

## Nonparametric

- No need to specify THE best prior
- Not straightforward to include external information
- Bayesian => MCMC expensive
- Non-Bayesian => need tuning data

## Empirical Bayes Parametric (e.g., SBLUP)

- Strong assumptions on effect sizes distribution
- Inference is fast if no MCMC
- May require tuning sample

...The eternal tension between Bias (explanation) and Variance (Prediction)

# Notable extensions

- Multi-trait
  - Functional annotations
  - Tuning using summary statistics (LASSO-SUM)
  - Combine individual-level and summary data
  - A few Deep Learning Applications...
- Improving genetic prediction by leveraging genetic correlations among human diseases and traits**
- [Robert M. Maier](#)✉, [Zhihong Zhu](#), [Sang Hong Lee](#), [Maciej Trzaskowski](#), [Douglas M. Ruderfer](#), [Eli A. Stahl](#), [Stephan Ripke](#), [Naomi R. Wray](#), [Jian Yang](#), [Peter M. Visscher](#)✉ & [Matthew R. Robinson](#)✉
- [Nature Communications](#) 9, Article number: 989 (2018) | [Cite this article](#)
- Incorporating functional priors improves polygenic prediction accuracy in UK Biobank and 23andMe data sets**
- [Carla Márquez-Luna](#)✉, [Steven Gazal](#), [Po-Ru Loh](#), [Samuel S. Kim](#), [Nicholas Furlotte](#), [Adam Auton](#), [23andMe Research Team](#) & [Alkes L. Price](#)✉
- [Nature Communications](#) 12, Article number: 6052 (2021) | [Cite this article](#)

Article

Leveraging both individual-level genetic data and GWAS summary statistics increases polygenic prediction

Clara Albiñana<sup>1, 2</sup>✉, Jakob Grove<sup>1, 3, 7, 8</sup>, John J. McGrath<sup>2, 4, 5</sup>, Esben Agerbo<sup>1, 2</sup>, Naomi R. Wray<sup>6, 5</sup>, Cynthia M. Bulik<sup>9, 10, 11</sup>, Merete Nordentoft<sup>1, 12, 13</sup>, David M. Hougaard<sup>1, 14</sup>, Thomas Werge<sup>1, 15, 13, 16</sup>, Anders D. Børglum<sup>1, 3, 7</sup>, Preben Bo Mortensen<sup>1, 2</sup>, Florian Privé<sup>1, 2, 17</sup>, Bjarni J. Vilhjálmsson<sup>1, 2, 8, 17</sup>✉

stacked clumping and thresholding (SCT) and meta-PRS. We find that, when large individual-level data are available, the linear combination of PRSs (meta-PRS) is both a simple alternative to meta-GWAS and often more accurate.

# Incorporating functional annotation data

Credit: slides from Dr Jian Zeng and Shouye Liu

## nature communications

Explore content ▾ About the journal ▾ Publish with us ▾

nature > nature communications > articles > article

Article | Open Access | Published: 18 October 2021

### Incorporating functional priors improves polygenic prediction accuracy in UK Biobank and 23andMe data sets

Carla Márquez-Luna , Steven Gazal, Po-Ru Loh, Samuel S. Kim, Nicholas Furlotte, Adam Auton, 23andMe Research Team & Alkes L. Price 

#### LDpredFunct method

### Exploiting biological priors and sequence variants enhances QTL discovery and genomic prediction of complex traits

I. M. MacLeod , P. J. Bowman, C. J. Vander Jagt, M. Haile-Mariam, K. E. Kemper, A. J. Chamberlain, C. Schrooten, B. J. Hayes & M. E. Goddard 

BMC Genomics 17, Article number: 144 (2016) | Cite this article

6209 Accesses | 146 Citations | 9 Altmetric | Metrics

#### BayesRC method

## PLOS COMPUTATIONAL BIOLOGY

 OPEN ACCESS  PEER-REVIEWED

RESEARCH ARTICLE

### Leveraging functional annotations in genetic risk prediction for human complex diseases

Yiming Hu , Qiongshi Lu , Ryan Powles, Xinwei Yao, Can Yang, Fang Fang, Xinran Xu, Hongyu Zhao 

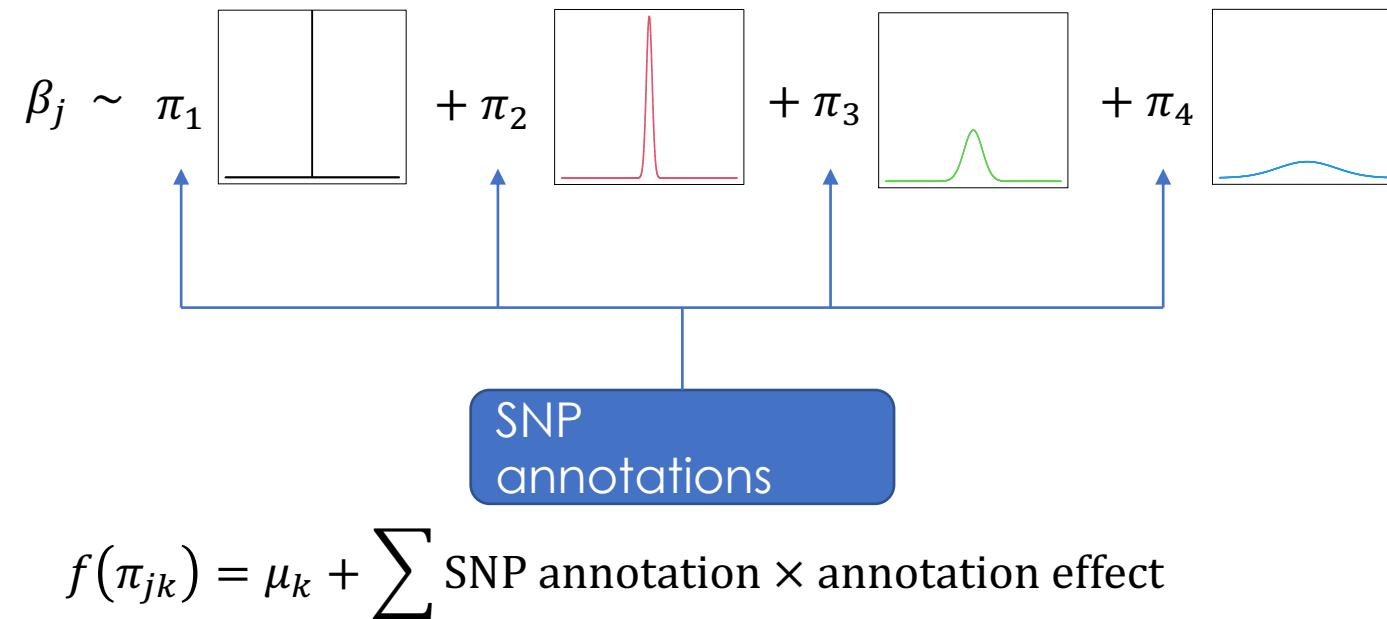
#### AnnoPred method

### Winner's Curse Correction and Variable Thresholding Improve Performance of Polygenic Risk Modeling Based on Genome-Wide Association Study Summary-Level Data

Jianxin Shi , Ju-Hyun Park, Jubao Duan, Sonja T. Berndt, Winton Moy, Kai Yu, Lei Song, William Wheeler, Xing Hua, Debra Silverman, Montserrat Garcia-Closas, Chao Agnes Hsiung, Jonine D. Figueroa, [...] Nilanjan Chatterjee  [ view all ]

#### P+T-funct-LASSO method

# Model annotation effects



## Assumption

- Annotation effects are additive at the GLM scale.

## Pros

- Estimation of conditional effects.
- Allow annotation overlap.
- Interpretation.

## Cons

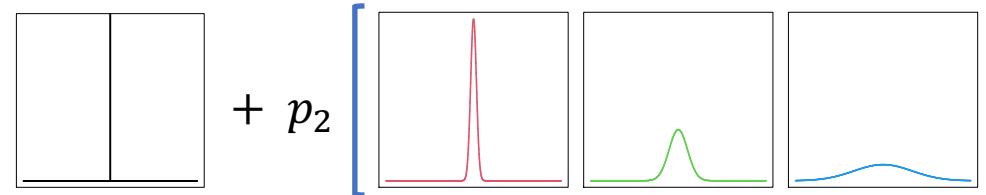
- # annotation effect parameters  $\times 4$ .
- $\pi_{j1} + \pi_{j2} + \pi_{j3} + \pi_{j4} = 1$ .

# Model annotation effects

- Three independent 2-component models:

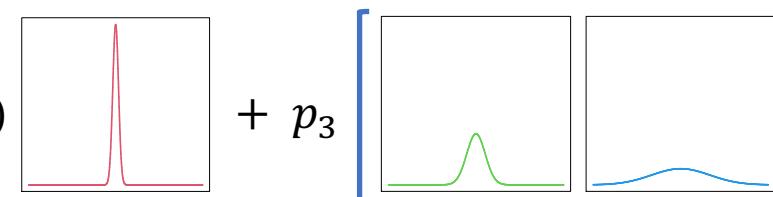
- For all SNPs

$$\beta_j \sim (1 - p_2)$$



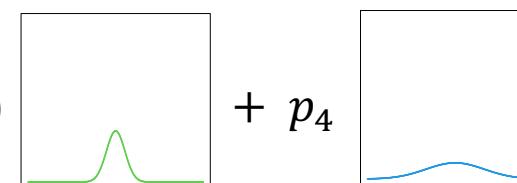
- For SNPs with nonzero effects

$$\beta_j \sim (1 - p_3)$$



- For SNPs with at least medium effects

$$\beta_j \sim (1 - p_4)$$



$$f(p_{jk}) = \mu_k + \sum \text{SNP annotation} \times \text{annotation effect}$$

# GCTB software

<https://cnsgenomics.com/software/gctb>

**GCTB**  
A tool for Genome-wide Complex Trait Bayesian analysis

GCTA SMR GSMR OSCA GCTB Program in CTG CTG forum

**Overview**

**Download**

Executable files  
Tutorial data  
LD matrices  
Source code  
Update log

Basic options  
Bayesian alphabet  
Summary Bayesian Alphabet

**Download**

**Executable files**

[gctb\\_2.03beta\\_Linux.zip](#) (Lastest version) In this version, we introduced a more robust parameterisation for SBayesR and SBayesS to address the convergence issue which sometimes occurs when the GWAS summary statistics are from a meta-analysis. With the same command line code, the program will start with the original model but if a convergence issue is detected (indicated by a negative sampled value of residual variance) it will restart the MCMC process with a more robust parameterisation. The robust model can also be directed invoked by `--robust`. The robust parameterisation for SBayesR and SBayesS is inspired by the parameterisation in LDpred ([Vilhjálmsson et al. 2015](#)). Details will soon be available on our website.

[gctb\\_2.02\\_Linux.zip](#)  
[gctb\\_2.0\\_Linux.zip](#)  
[gctb\\_1.0\\_Linux.zip](#)  
[gctb\\_1.0\\_Mac.zip](#)

**Tutorial data**

[GCTB tutorial data](#)

**LD matrices**

The following LD matrices were computed based on 1.1 million common SNPs in a random sample of 50K unrelated individuals of European ancestry in UK Biobank dataset unless otherwise noted.

- [Shrunk sparse matrix](#)
- [Shrunk sparse LD matrix \(2.8 million common SNPs\)](#)

In the shrunk sparse matrices, described in [Lloyd-Jones et al. \(2019\)](#), the observed LD correlations computed from a reference sample were shrunk toward the expected values defined by a [genetic map](#), following the algorithm in [Wen and Stephens \(2010\)](#). After shrinkage, LD correlations smaller than a threshold (default 1e-5) were set to be zero to give a sparse format, which is more efficient in storage and computation.

- [Sparse matrix \(including MHC regions\)](#)

The sparse matrices described in [Zeng et al. \(2021\)](#) were computed by setting the likely chance LD to zero based on a chi-squared test (default threshold at chi-squared test statistic of 10).

- [Banded matrix \(including MHC regions\)](#)

While the shrunk sparse matrices were used in our original SBayesR paper, [Prive et al. \(2021\)](#) found that using a banded matrix with a window size of 3 cM per SNP can improve prediction accuracy. Therefore, we have created such a LD matrix in GCTB format for SBayesR analysis.

**Source code**

[GCTB 2.0 standard version](#)  
[GCTB 2.0 MPI version](#)  
[GCTB 1.0 standard version](#)  
[GCTB 1.0 MPI version](#)

The MPI version implements a distributed computing strategy that scales the analysis to very large sample sizes. A significant improvement in computing time is expected for a sample size > 10,000. The MPI version needs to be compiled on user's machine. See [README.html](#) in the tarball for instructions of compilation and usage. A testing dataset is also included in each tarball.

# Tutorial

<http://cnsgenomics.com/software/gctb/#Tutorial>

## GCTB

A tool for Genome-wide Complex Trait Bayesian analysis

GCTA SMR GSMR OSCA GCTB Program in CTG CTG forum

### Overview

### Download

### Basic options

### Bayesian alphabet

### Summary Bayesian Alphabet

Options

### Tutorial

### FAQ

## Tutorial

This updated version of the GCTB software (version 2.0) includes summary-data based versions of the individual-data Bayesian linear mixed models previously implemented. These methods require summary statistics from genome-wide association studies, which typically include the estimated univariate effect for each single nucleotide polymorphism (SNP), the standard error of the effect, the sample size used to estimate the effect for each SNP, the allele frequency of an allele, and an estimate of LD among SNPs, all of which are easily accessible from public databases.

This tutorial will outline how to use the summary data based methods and accompanies the manuscript [Improved polygenic prediction by Bayesian multiple regression on summary statistics](#). To recreate the tutorial you will need the [PLINK 2](#) software and the updated version of the [GCTB](#) software [compiled](#) from source or the binary executable in your path. The data for the tutorial are available at [Download](#). The tutorial is designed so that each part can be reconstructed or picked up at any point with the following directory structure.

```
tutorial
gctb # Static binary executable for Linux 64-bit systems
data # 1000 Genomes data in PLINK format
pheno_generation # R scripts for generating simulated phenotypes
pheno # Simulated phenotypes
gwas # PLINK GWAS summary statistics
ldm # LD correlation matrices and scripts to calculate them
ma # Summary statistics in GCTB compatible format
sbayesr # Scripts for running and SBayesR analysis using GCTB and subsequent results
```

### Data

The tutorial will run through all aspects of the software using genotype data from chromosome 22 of Phase 3 of the [1000 Genomes Project \(1000G\)](#). The genotype data have been filtered to exclude variants with minor allele frequency (MAF) < 0.01, which left 15,938 SNPs available for analysis. These data include a subset of 378 individuals of European ancestry from the CEU, TSI, GBR and FIN populations.

N.B. The results generated from this example tutorial using data from the 1000G are designed to be lightweight and should be only used to explore how to run summary data based analyses using GCTB. The results do not make the best use of GCTB capabilities given the small sample size (N = 378) in the example.

### Phenotype simulation and GWAS

Using these genotypes, phenotypes were generated under the multiple regression model  $y_i = \sum_{j=1}^p w_{ij}\beta_j + \epsilon_i$ , where  $w_{ij} = (x_{ij} - 2q_j)/\sqrt{2q_j(1-q_j)}$  with  $x_{ij}$  being the reference allele count for the  $i$ th individual at the  $j$ th SNP,  $q_j$  the allele frequency of the  $j$ th SNP and  $\epsilon_i$  was sampled from a normal distribution with mean 0 and variance  $\text{Var}(W\beta)(1/h_{\text{SNP}}^2 - 1)$  such that  $h_{\text{SNP}}^2 = 0.1$  for each of 20 simulation replicates, which is much larger than the contribution to the genome-wide SNP-based heritability ( $h_{\text{SNP}}^2$ ) estimate for chromosome 22 for most quantitative traits. All phenotypes were generated using the R programming language with scripts used available in the [pheno\\_generation](#) subdirectory.

For each scenario replicate, we randomly sampled a new set of 1,500 causal variants. The genetic architecture simulated contained two causal variants of large effect explaining 3% and 2% of the phenotypic variance respectively and a polygenic tail of 1,498 causal variants sampled from a  $N(0, 0.05^2 / 1,498)$  distribution such that the expected total genetic variance explained by all variants was 0.1.

For each of the 20 simulation replicates, simple linear regression for each variant was run using the PLINK 2 software to generate summary statistics. Code for running the GWAS and the output is available in the [gwas](#) subdirectory.

### GCTB summary statistics input format

The GCTB summary-based methods have inherited the [GCTA-COJO .ma](#) format.

```
SNP A1 A2 freq b se p N
rs1001 A G 0.8493 0.0024 0.0055 0.6653 129853
rs1002 C G 0.0306 0.0034 0.0115 0.7659 129799
rs1003 A C 0.5128 0.0045 0.0038 0.2319 129830
```

Columns are SNP identifier, the effect allele, the other allele, frequency of the effect allele, effect size, standard error, p-value and sample size. The headers are not keywords and will be omitted by the program. Important: "A1" needs to be the effect allele with "A2" being the other allele and "freq" should be the frequency of "A1".

[Presentation Title] | [Date]

The transformation of your summary statistics file to the [.ma](#) format will depend on the software used to analyse your data. The [ma](#) subdirectory contains a simple R script for constructing [.ma](#) files from PLINK 2 [--linear](#) output. There is no need to filter your summary statistics to match the LD reference as GCTB will perform this data management for you.