

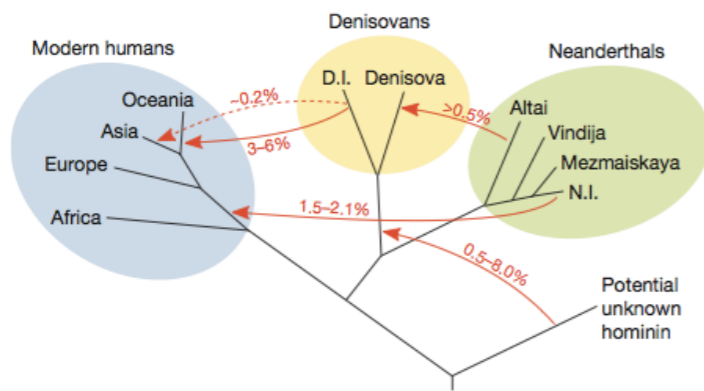
Probabilistic models for large-scale human genomic data

Sriram Sankararaman

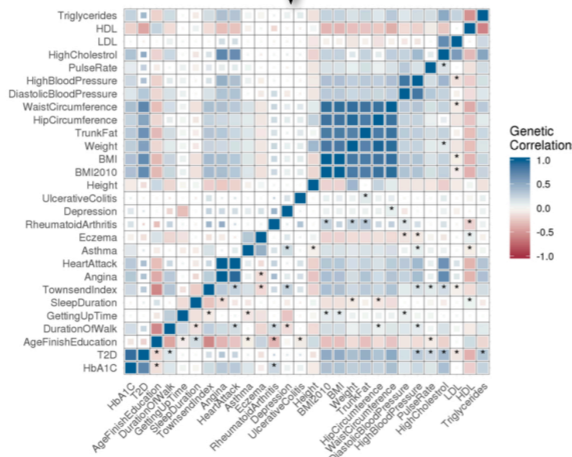
Computer Science, Human Genetics, Computational Medicine
UCLA

Machine learning for genomic data

Evolution



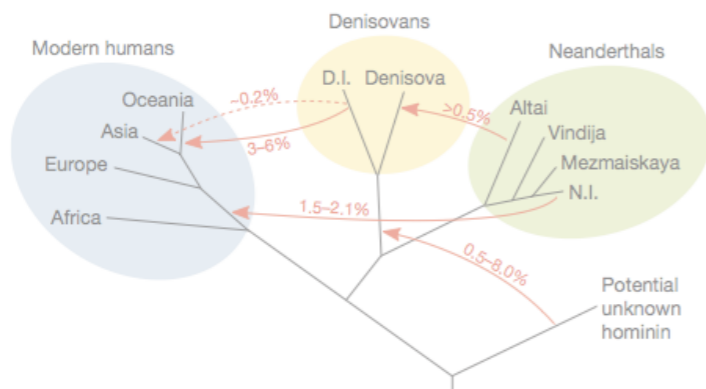
Clinical



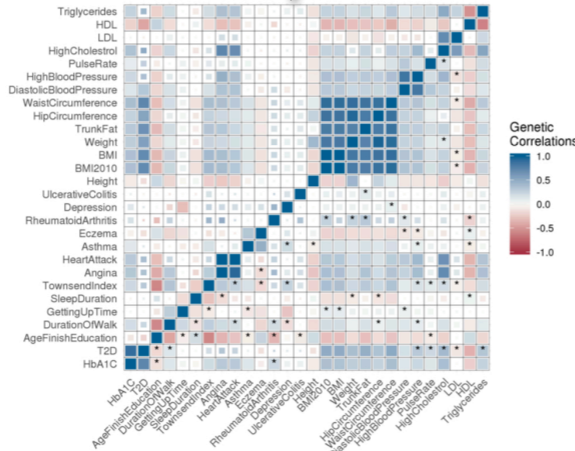
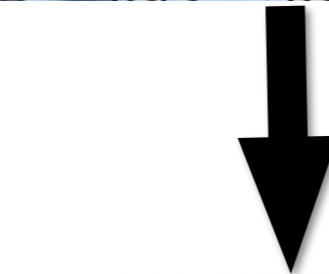
Complex traits

Machine learning for genomic data

Evolution



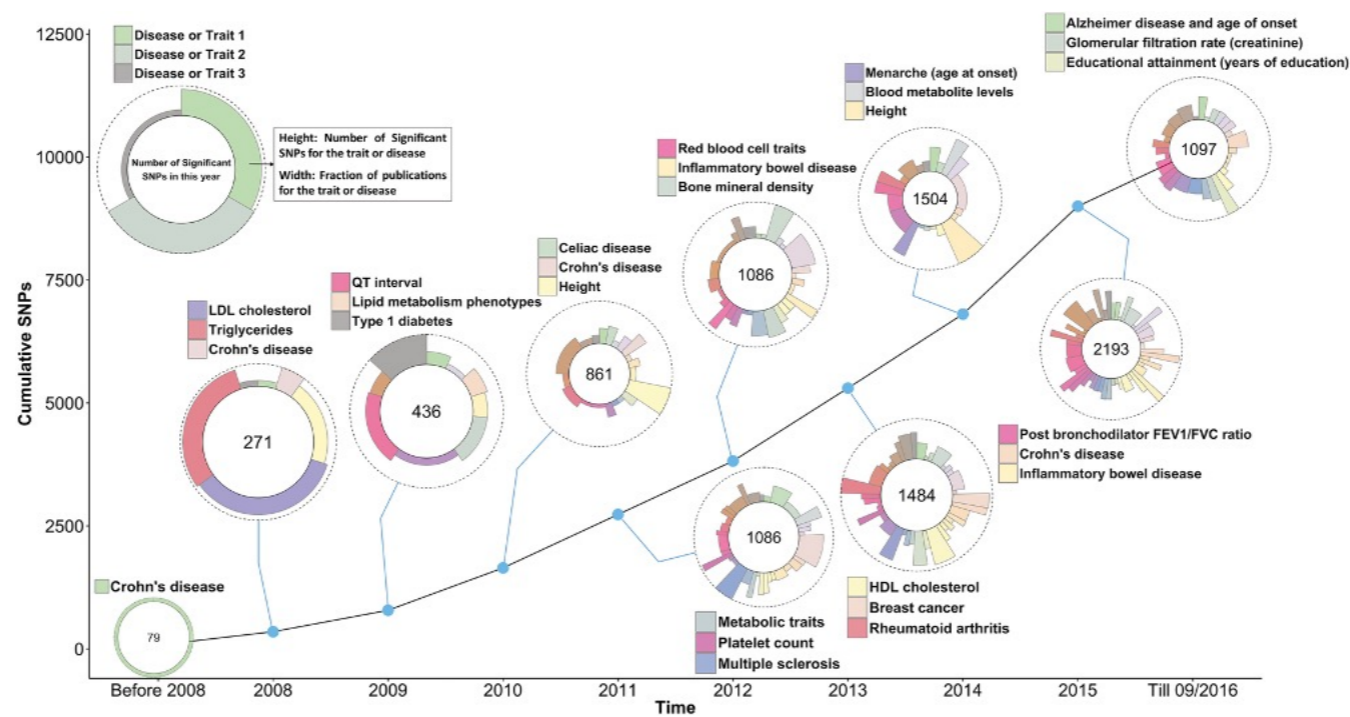
Clinical



Complex traits

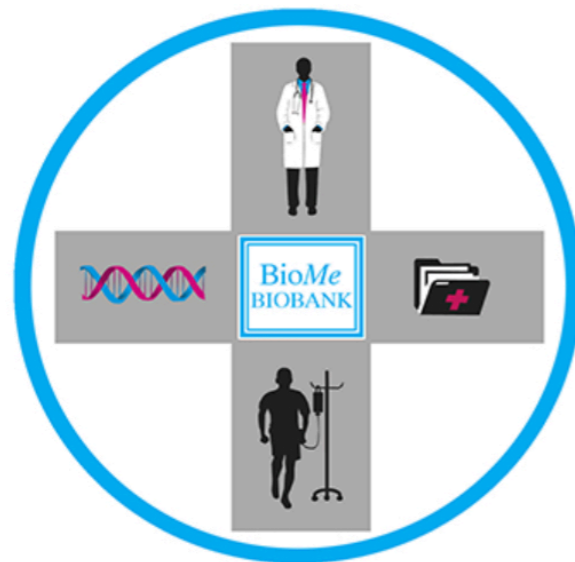
Genetic architecture of complex phenotypes

$$\text{Phenotype} = f(\text{Genotype}, \text{Environment})$$



Visscher et al. AJHG 2017

Growth of Biobanks



Machine Learning for Biobank-scale data

How can we learn about genetic architecture of complex traits and diseases from datasets that contain millions of genomes and thousands of traits ?

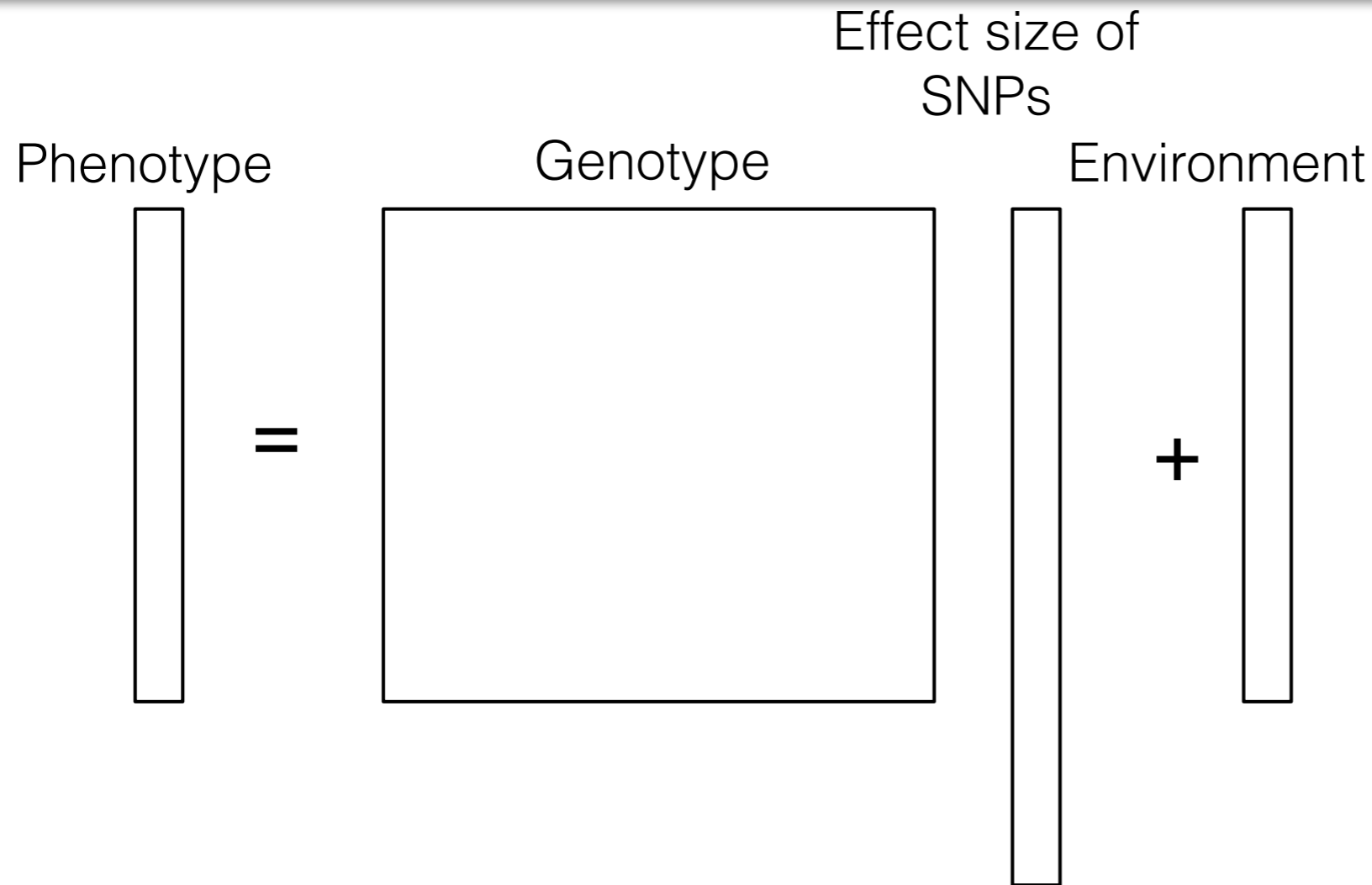
Statistical

Privacy

Computational

Interpretability

(Narrow-sense) Heritability



$$h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2}$$

SNP heritability

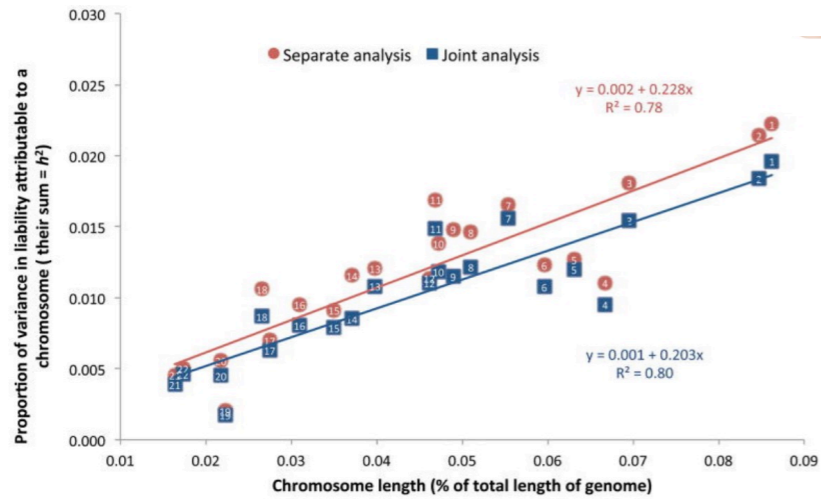
$$Y = X\beta + \epsilon$$

$$\beta_m \sim \mathcal{N}\left(0, \frac{\sigma_g^2}{M}\right) \quad \epsilon_n \sim \mathcal{N}(0, \sigma_e^2)$$

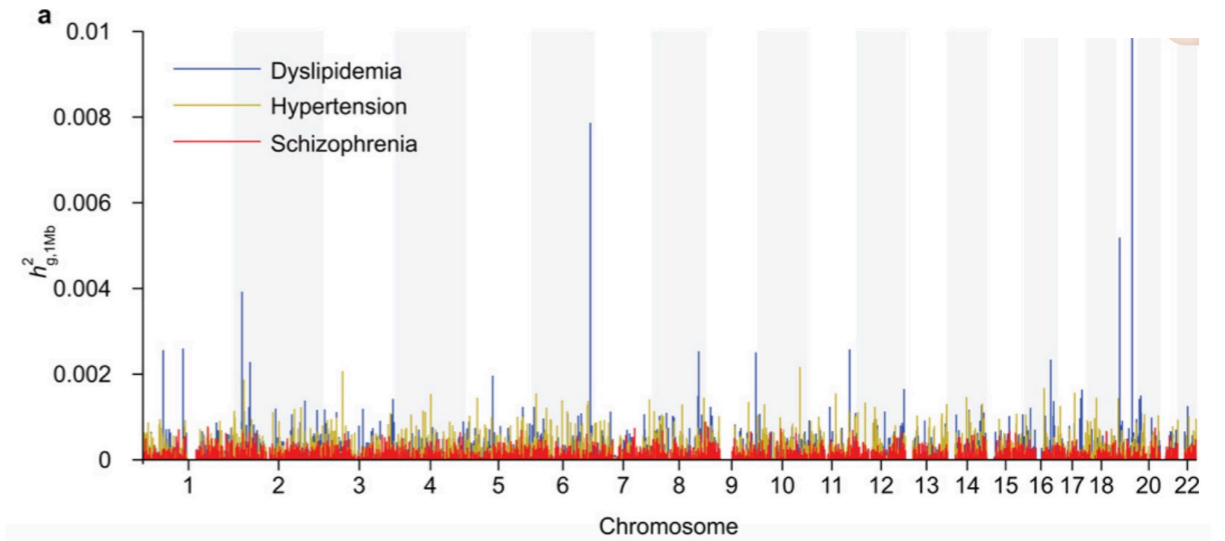
Genetic variance component Environmental variance component

Genetic variance component

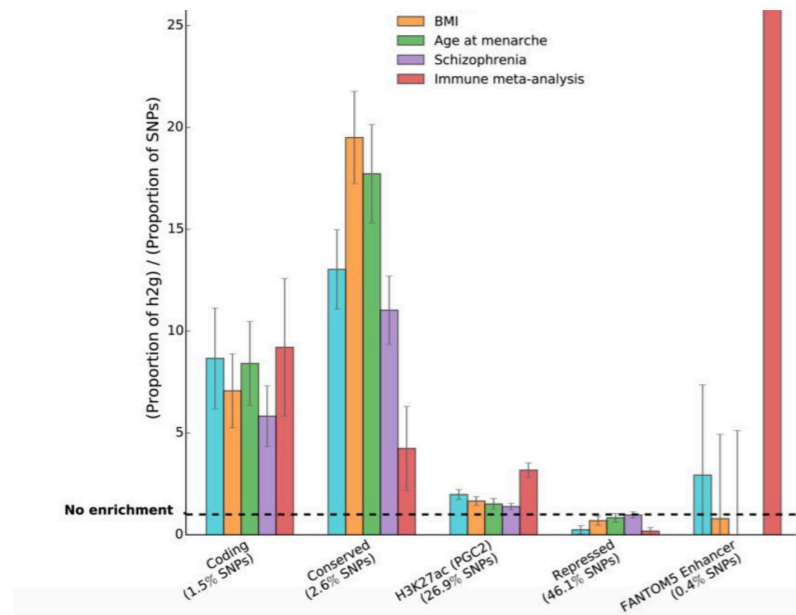
Beyond Heritability



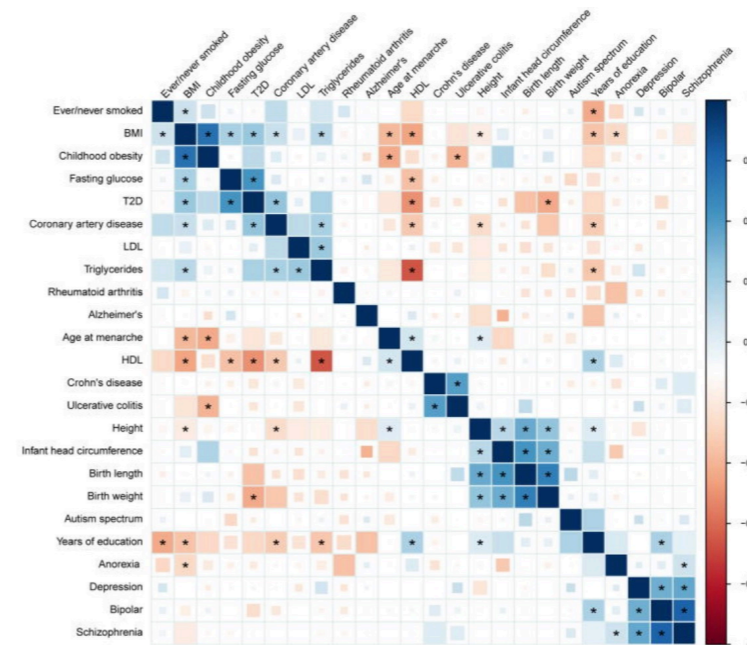
Lee et al. 2012



Loh et al. 2015

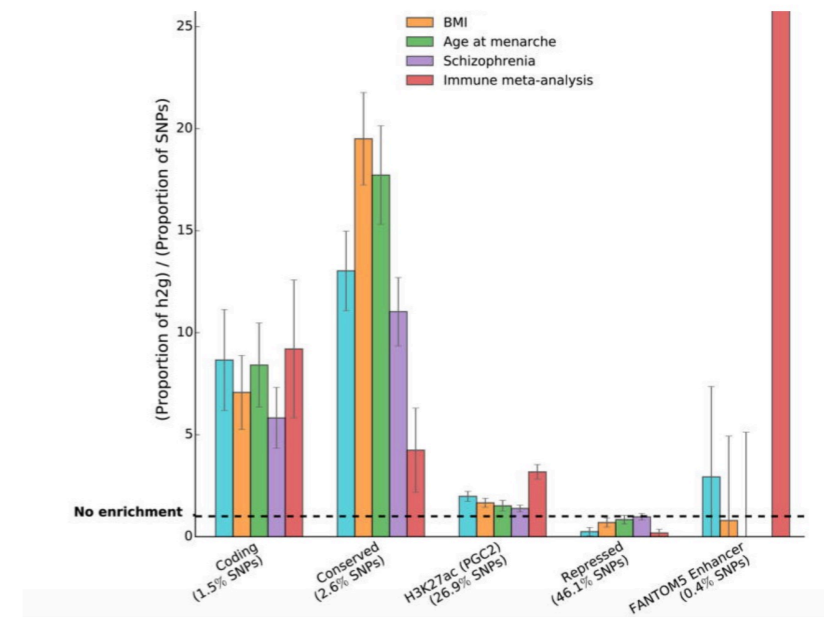
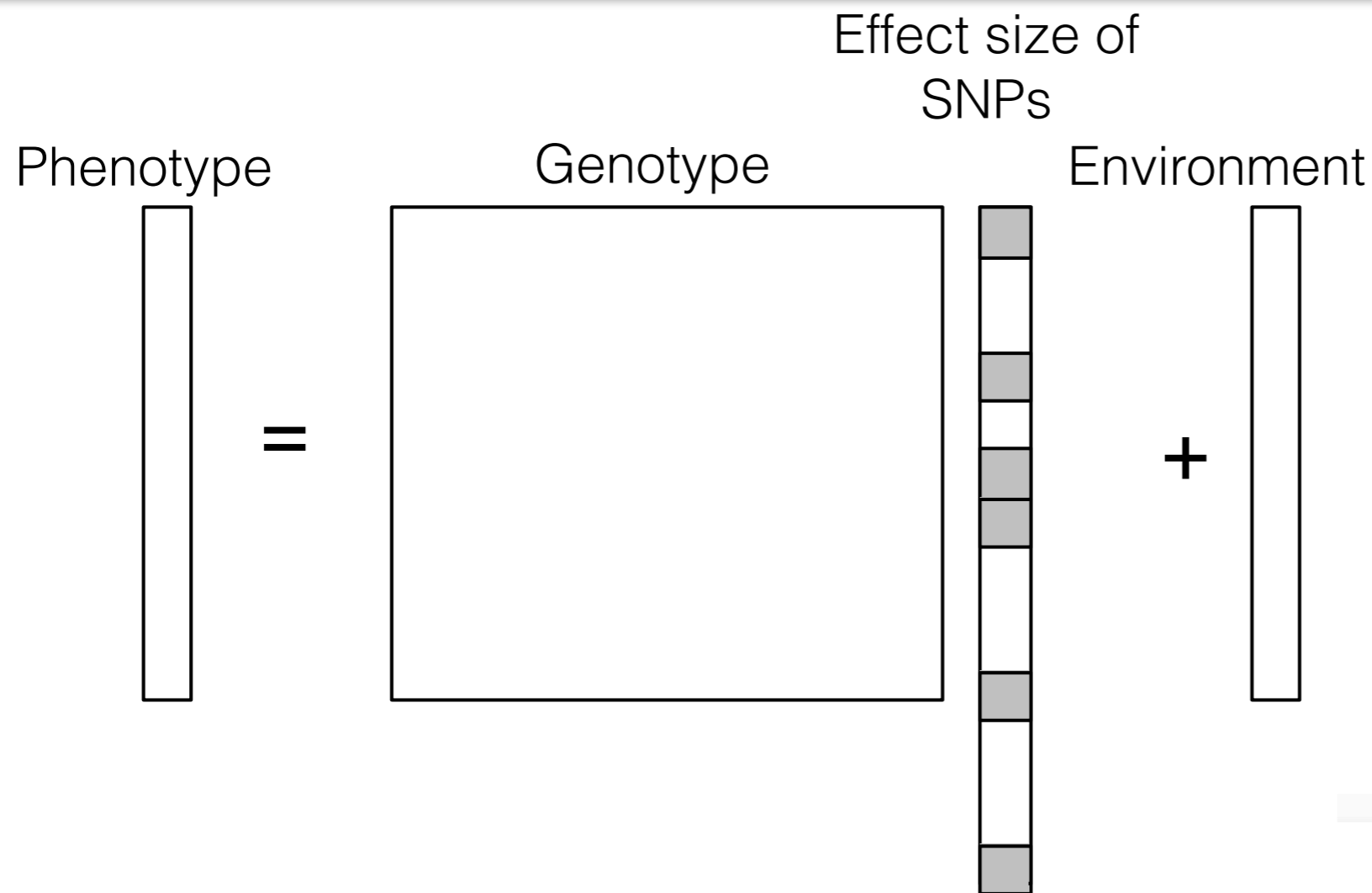


Finucane et al. 2015



Bulk-Sullivan et al. 2015

Variance components model



Finucane et al. Nature Genetics 2015

$$Y = X\beta + \epsilon$$

$$\beta_m \sim \mathcal{N}(0, \sigma_{Gene}^2) \text{ if } m \text{ is in Gene}$$

$$\beta_m \sim \mathcal{N}(0, \sigma_{NonGene}^2) \text{ otherwise}$$

Variance components model

$$\mathbf{y} = \sum_{k=1}^K \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\epsilon}$$

$$\boldsymbol{\beta}_k \sim \mathcal{N}(\mathbf{0}, \frac{\sigma_k^2}{M_k} \mathbf{I}_{M_k}), k \in \{1, \dots, K\}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}_N)$$

Goal

Estimate variance components $(\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2, \sigma_e^2)$

Estimating variance components

Maximum likelihood

$$\begin{aligned}(\hat{\sigma}_1^2, \hat{\sigma}_2^2, \dots, \hat{\sigma}_K^2, \hat{\sigma}_e^2) &= \operatorname{argmax}_{(\sigma_1^2, \dots, \sigma_K^2, \sigma_e^2)} \mathcal{LL}(\sigma_1^2, \dots, \sigma_K^2, \sigma_e^2) \\ &= \operatorname{argmax}_{(\sigma_1^2, \dots, \sigma_K^2, \sigma_e^2)} P(\mathbf{y} | \mathbf{X}_1, \dots, \mathbf{X}_K, \sigma_1^2, \dots, \sigma_K^2, \sigma_e^2)\end{aligned}$$

Computationally expensive

Scales as $\mathcal{O}(N^3)$

Challenging to apply to Biobank-scale data

Lippert et al. Nature Methods 2012
Zhou and Stephens, Nature Genetics 2012
Loh et al. Nature Genetics 2015

Alternate estimator

Method of Moments (HE-regression)

$$\begin{aligned} \text{cov}(\mathbf{y}) &= \sum_k \sigma_k^2 \frac{1}{M_k} \mathbf{X}_k \mathbf{X}_k^T + \sigma_e^2 \mathbf{I}_N \\ &= \sum_k \sigma_k^2 \mathbf{K}_k + \sigma_e^2 \mathbf{I}_N \\ &\approx \mathbf{y} \mathbf{y}^T \end{aligned}$$

Closed-form solution

Still needs computing GRMs $\mathbf{K}_k = \frac{1}{M_k} \mathbf{X}_k \mathbf{X}_k^T$

Scales as $\mathcal{O}(N^2 M)$

Haseman and Elston, AJHG 1972
Zhou Annals of Applied Statistics 2017

Randomized HE-regression (RHE-mc)

Avoid computing the GRM

$$\begin{aligned} \text{tr}(\mathbf{K}_l \mathbf{K}_k) &\approx \frac{1}{B} \sum_b z_b^T \mathbf{K}_l \mathbf{K}_k z_b \\ &= \frac{1}{M^2 B} \sum_b z_b^T \mathbf{X}_l \mathbf{X}_l^T \mathbf{X}_k \mathbf{X}_k^T z_b \end{aligned}$$

Multiply the genotype matrix with B random vectors

Hutchinson 1989

Wu et al. Bioinformatics 2018

Pazokitoroudi et al. Nature Communication 2020

Randomized HE-regression (RHE-mc)

Combines randomization with a method-of-moments estimator

Work with a “sketch” of the genotype

Multiply the genotype matrix with B random vectors

Efficiency depends on B : $\mathcal{O}\left(\frac{MNB}{\log_3(\max(N, M))}\right)$

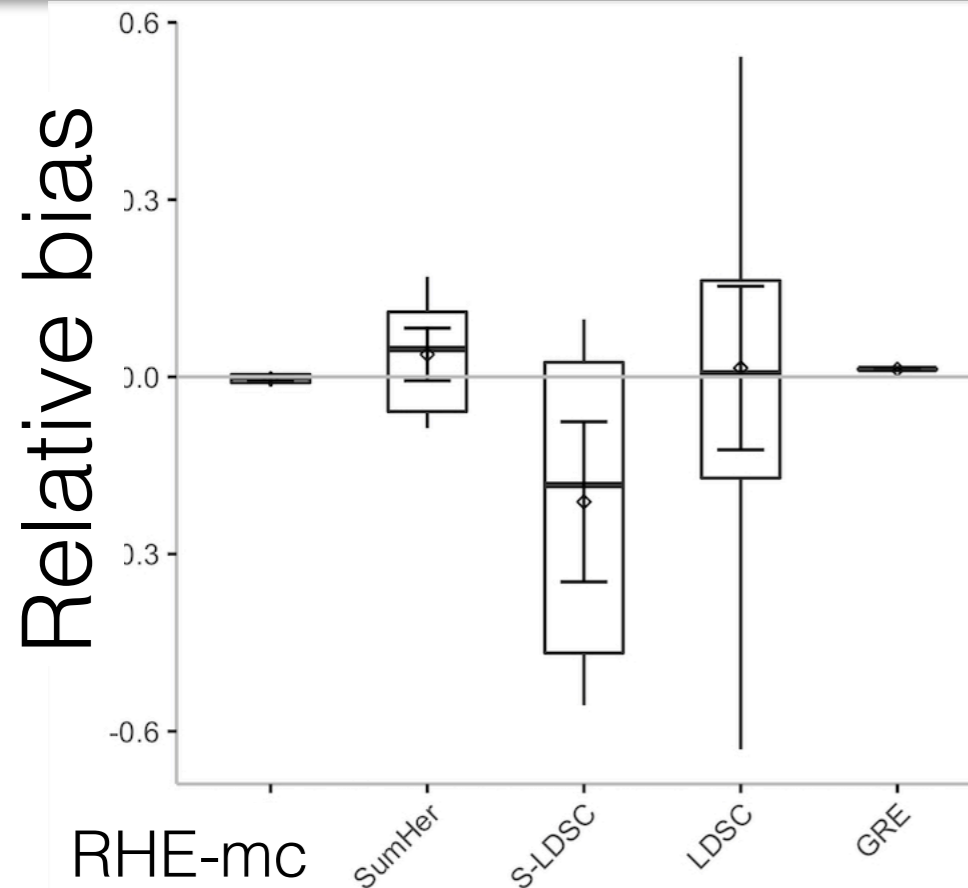
Accurate for B as small as 10

Hutchinson 1989

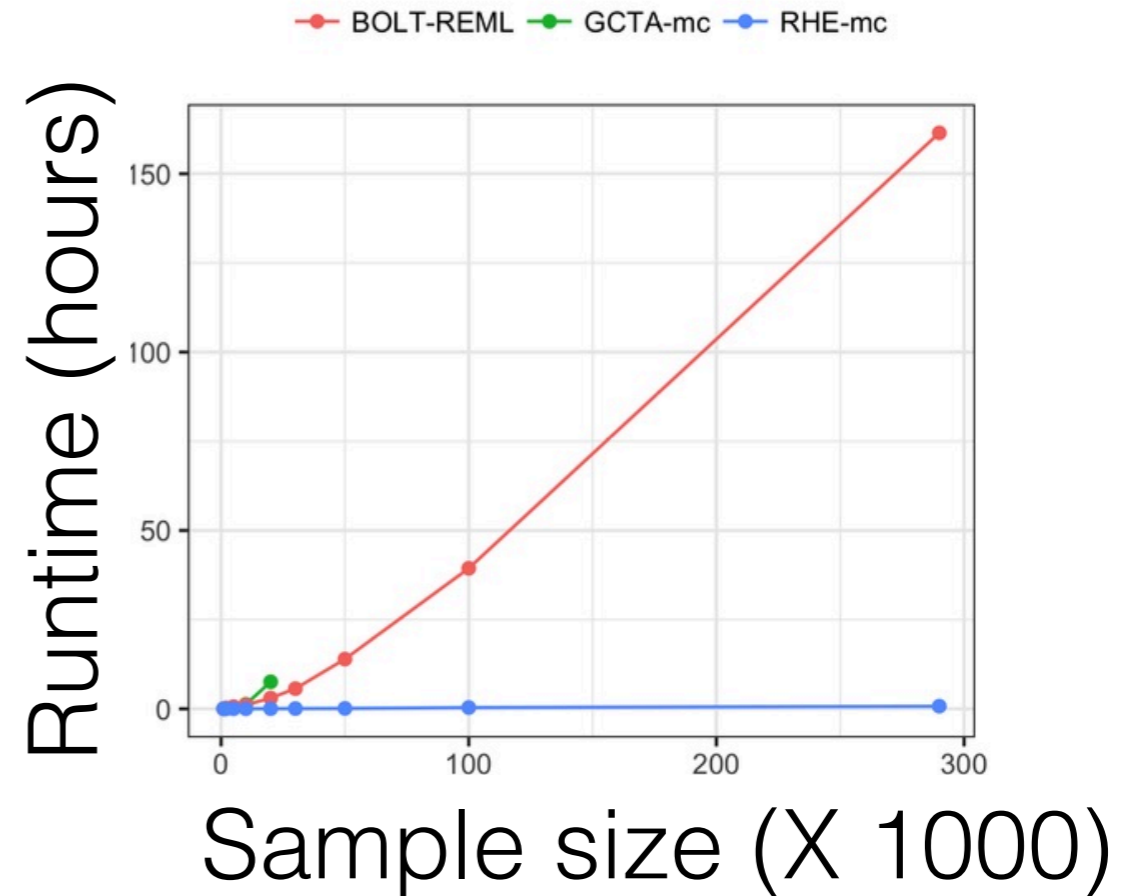
Wu et al. Bioinformatics 2018

Pazokitoroudi et al. RECOMB 2019, Nature Communication 2020

RHE-mc is accurate and scalable



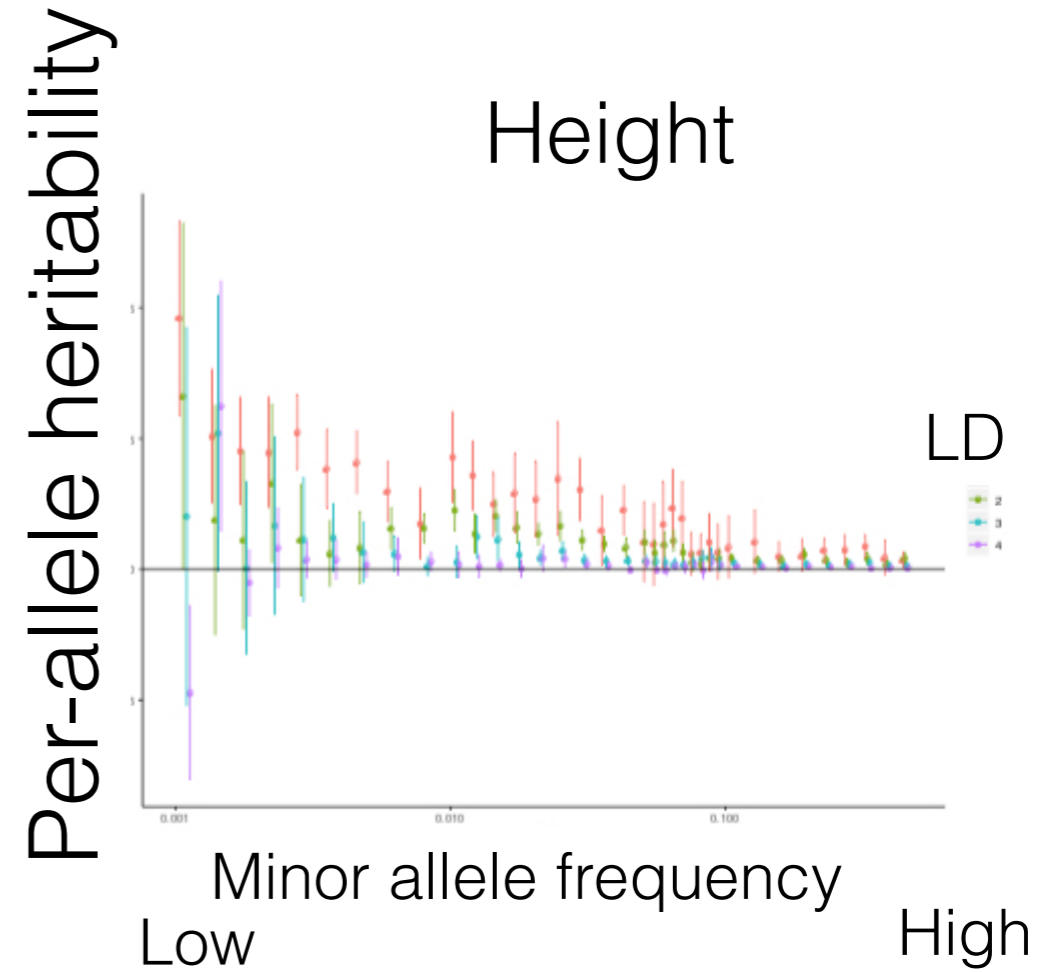
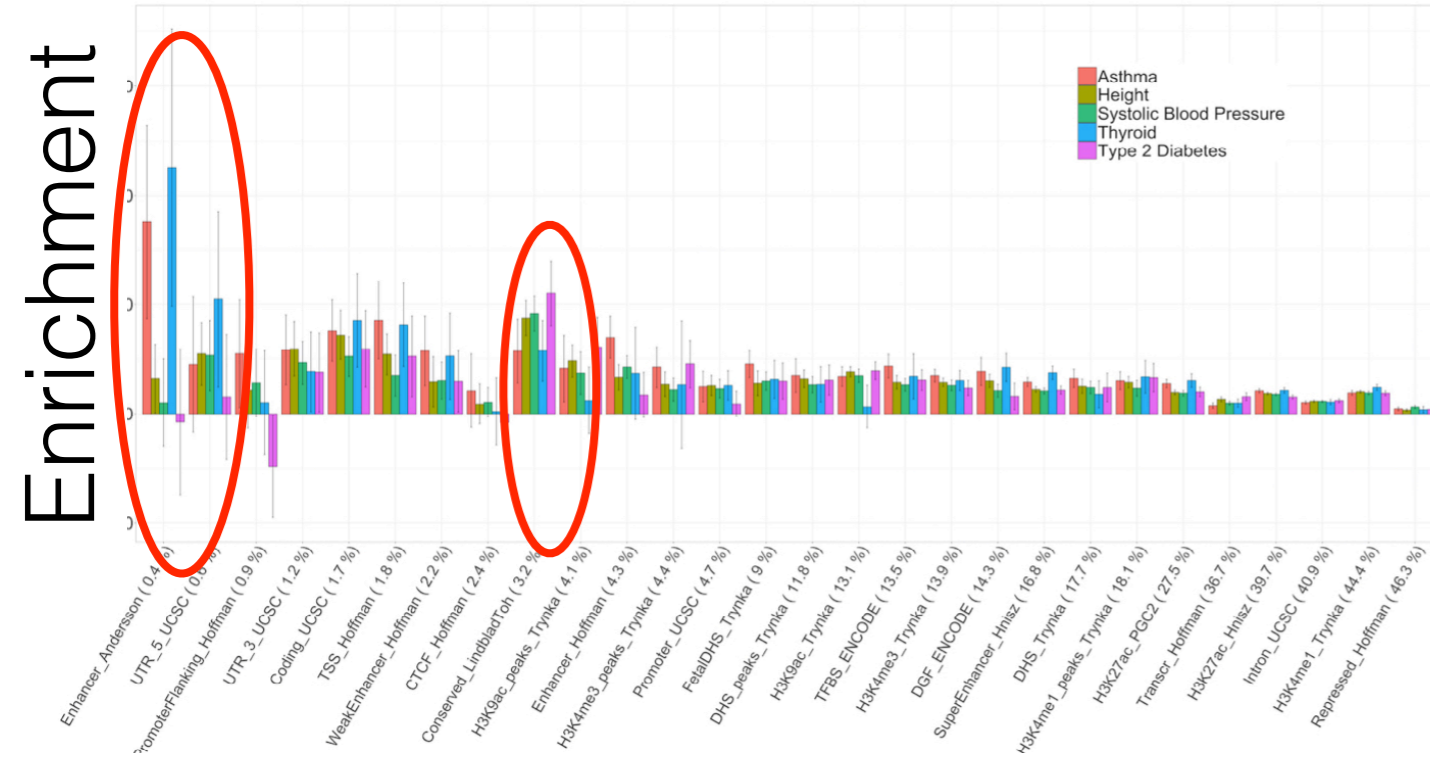
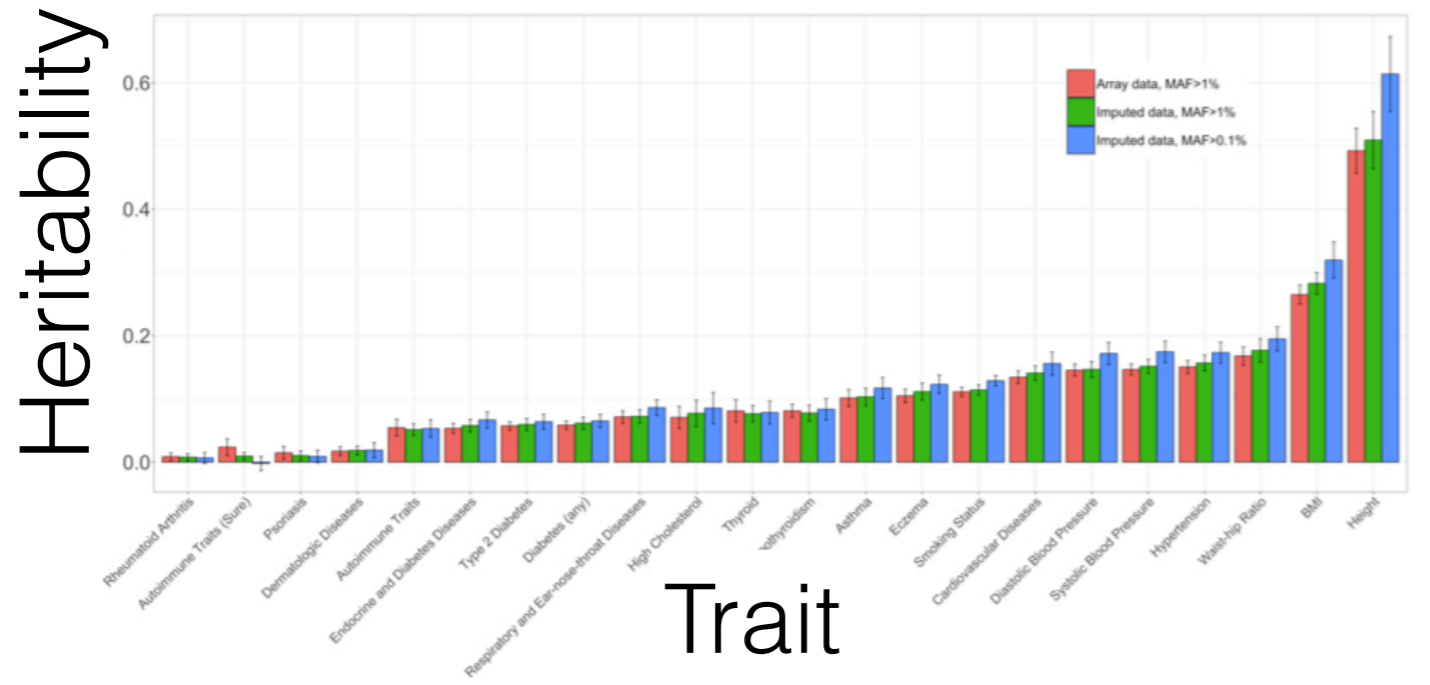
Summary-statistic methods



N	Parameters		Running time (hour)		
	M	K	RHE-mc	GCTA-mc	BOLT-REML
10,000	459,792	22	< 1	1.3	1
100,000	459,792	22	< 1	-	40
291,273	459,792	22	< 1	-	162
291,273	459,792	300	< 1	-	-
291,273	4,824,392	8	3.2	-	-
1,000,000	1,000,000	8	3	-	-
1,000,000	1,000,000	100	12.4	-	-

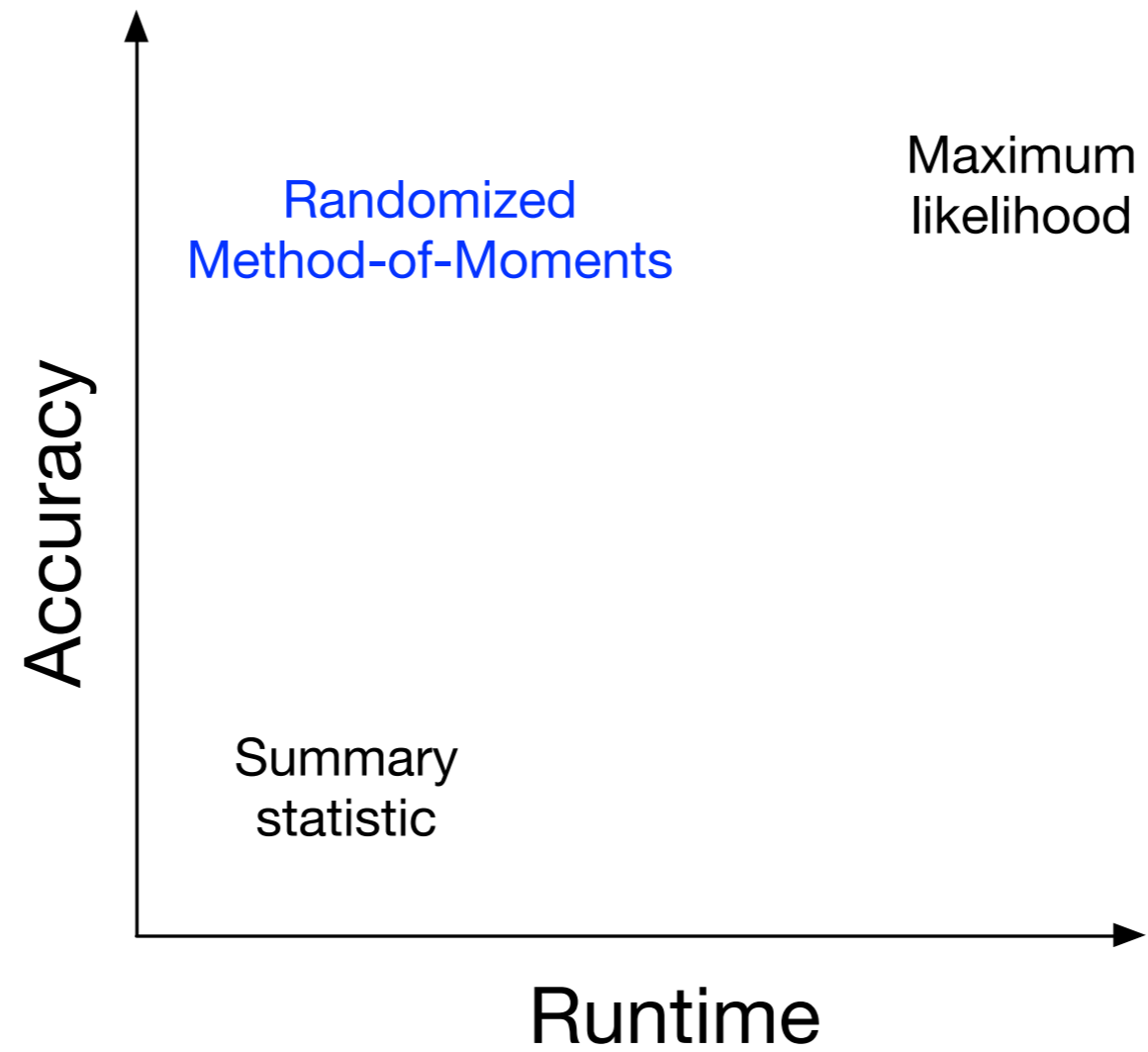


Insights from Biobank-scale analysis



Pazokitoroudi et al. Nature Communication 2020

RHE-mc



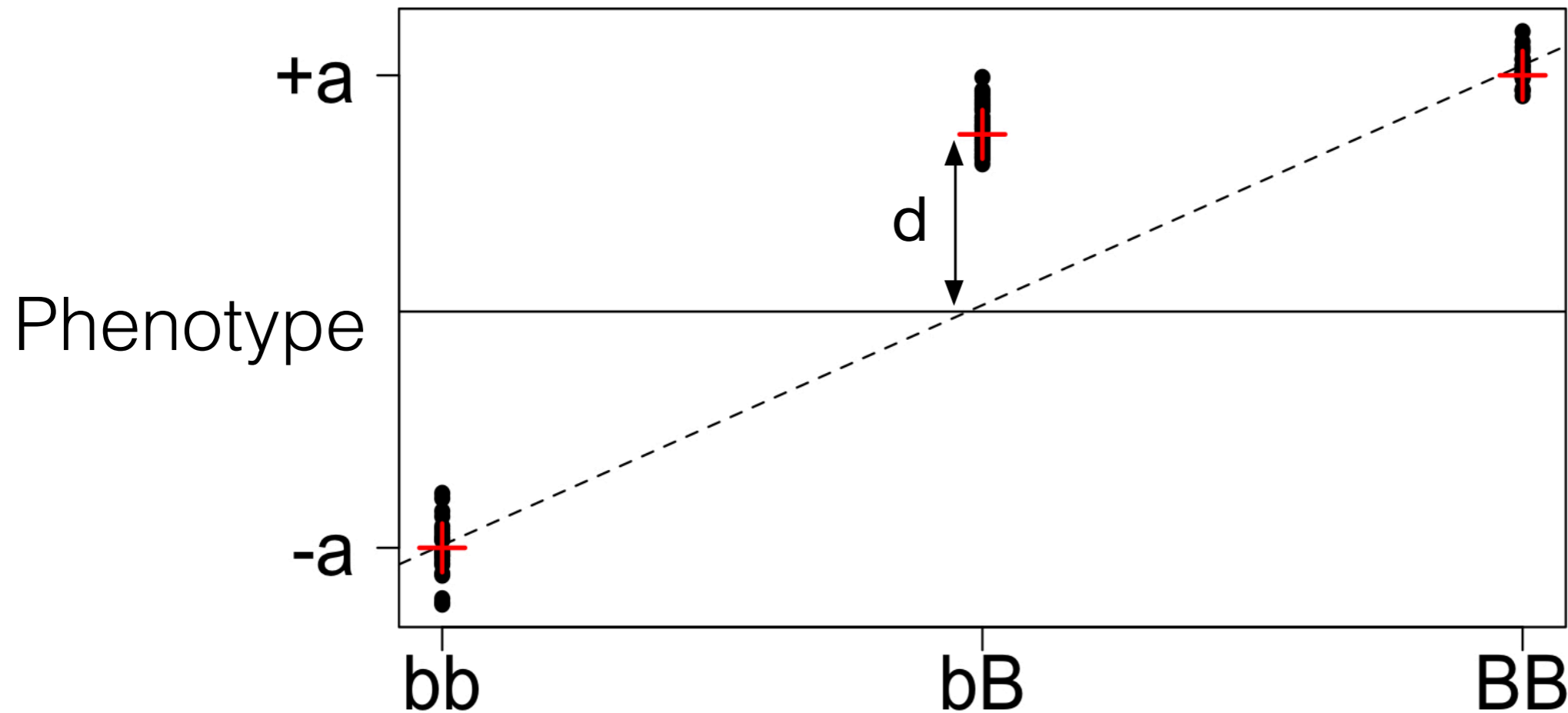
Beyond heritability

What is the contribution of non-linear effects ?

What is the contribution of environmental interactions ?

How are genetic effects shared across traits ?

Dominance deviation effects



Dominance deviation effects

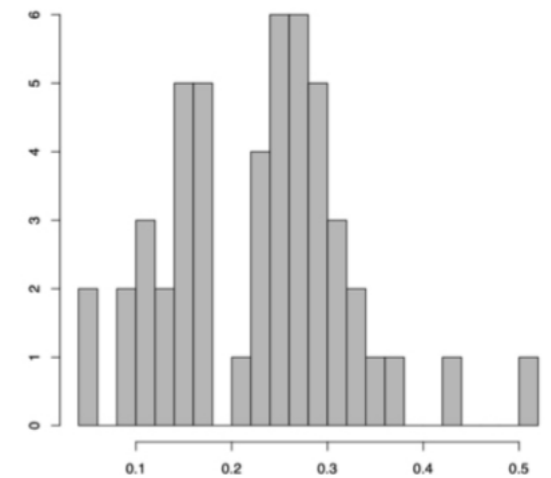
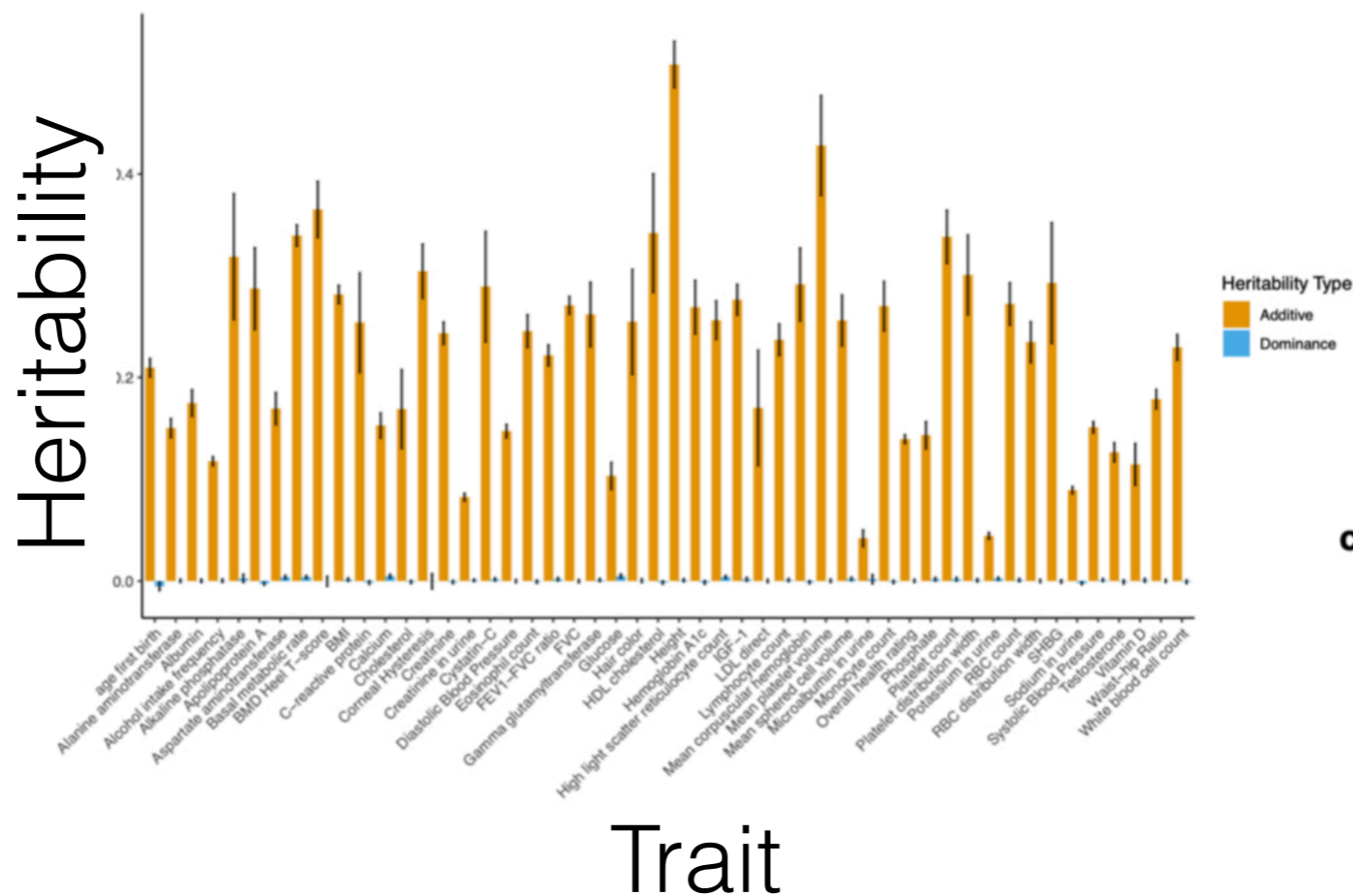
Additive variance component

$$y = X\beta + D\gamma + \epsilon$$

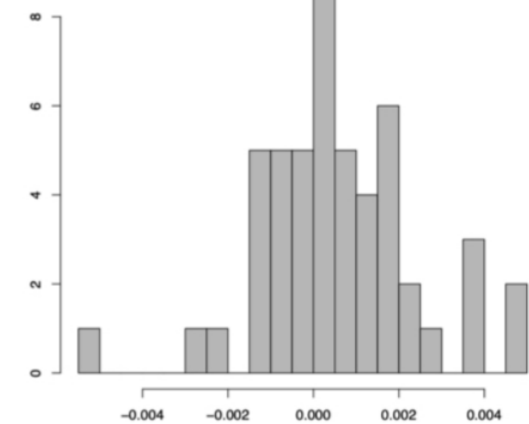
$$\beta \sim \mathcal{N}(\mathbf{0}, \frac{\sigma_a^2}{M} \mathbf{I}_M)$$

Dominance variance component

$$\gamma \sim \mathcal{N}(\mathbf{0}, \frac{\sigma_d^2}{M} \mathbf{I}_M)$$



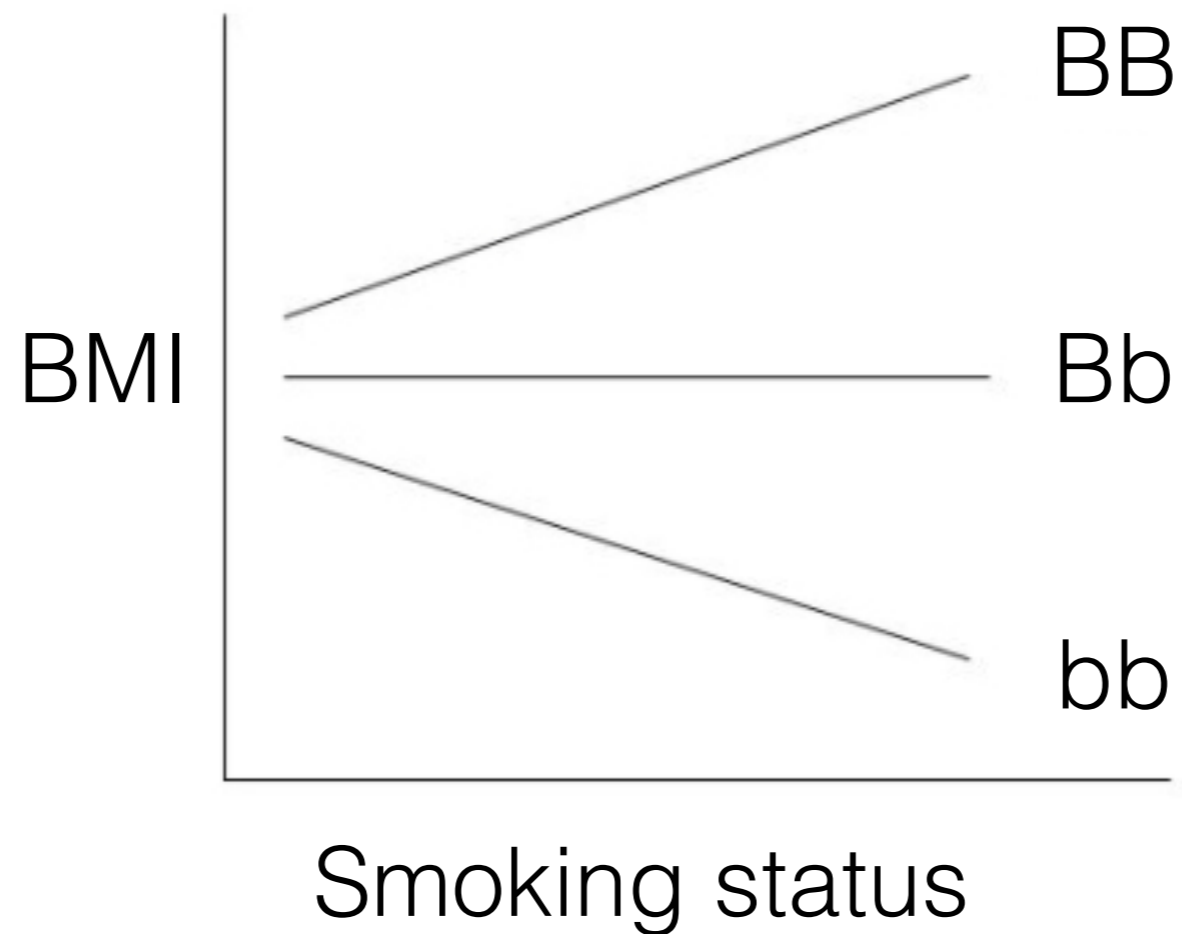
Additive heritability



Dominance heritability

Hivert et al. AJHG 2021
Pazokitoroudi et al. AJHG 2021

Gene-environment interactions (GxE)



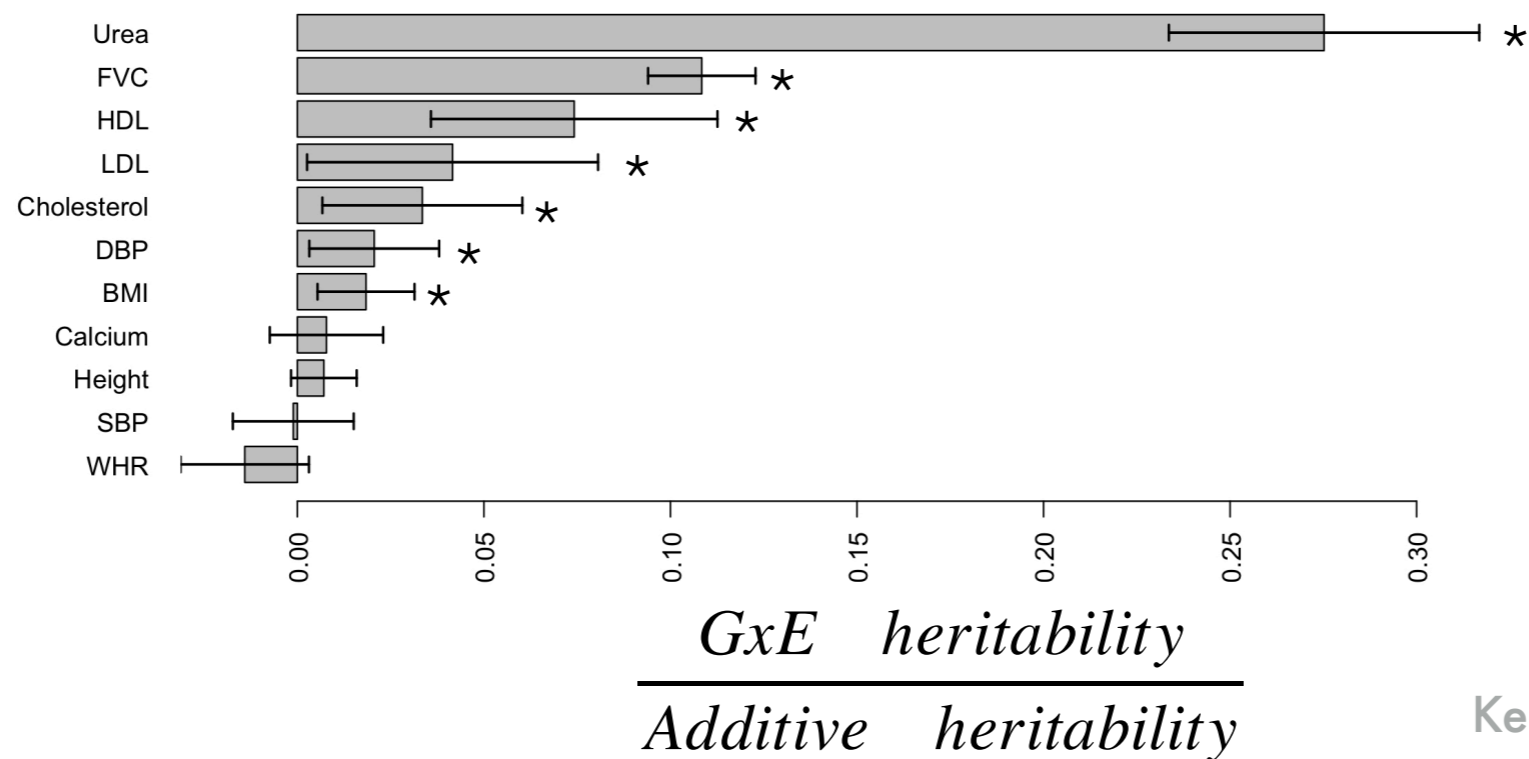
Gene-environment interactions (GxE)

$$y = X\beta + X \odot E\delta + \epsilon$$

$$\delta \sim \mathcal{N}\left(0, \frac{\sigma_{GE}^2}{ML} I\right)$$

GxE variance component

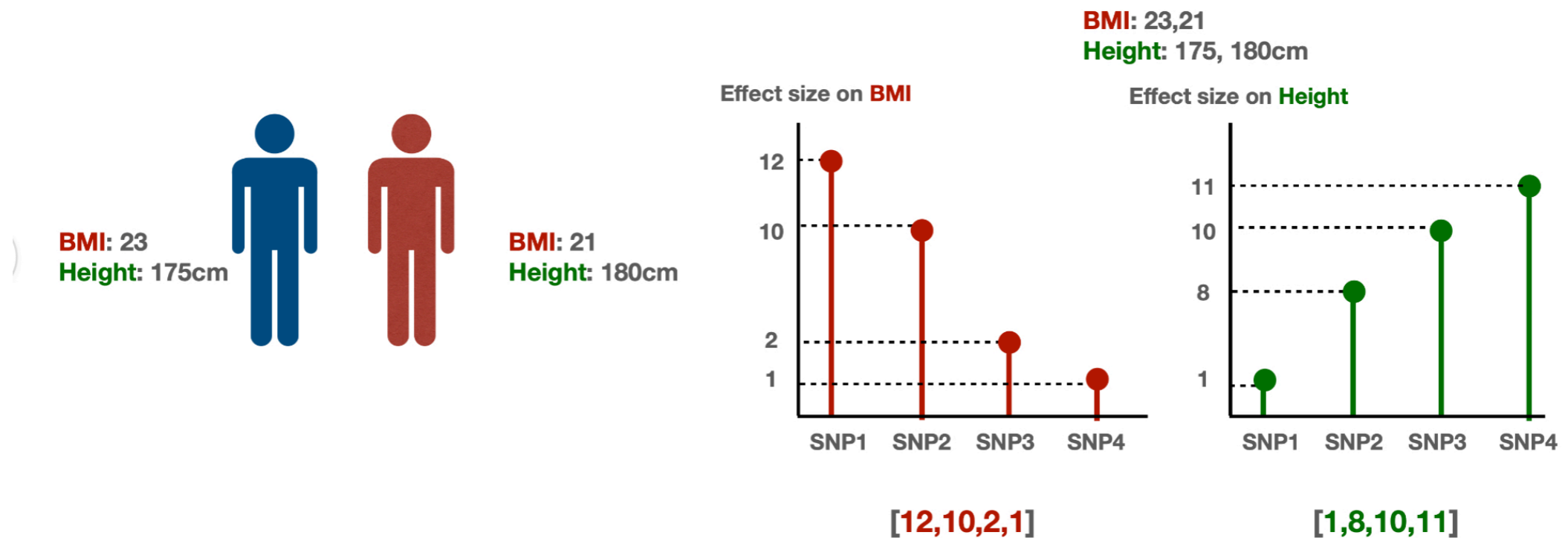
E = Smoking



Kerrin and Marchini AJHG 2020
Pazokitoroudi et al. RECOMB 2021

Genetic effects shared across traits

Genetic correlation



Estimating genetic correlation

$$y_1 = X_1\beta_1 + \epsilon_1$$

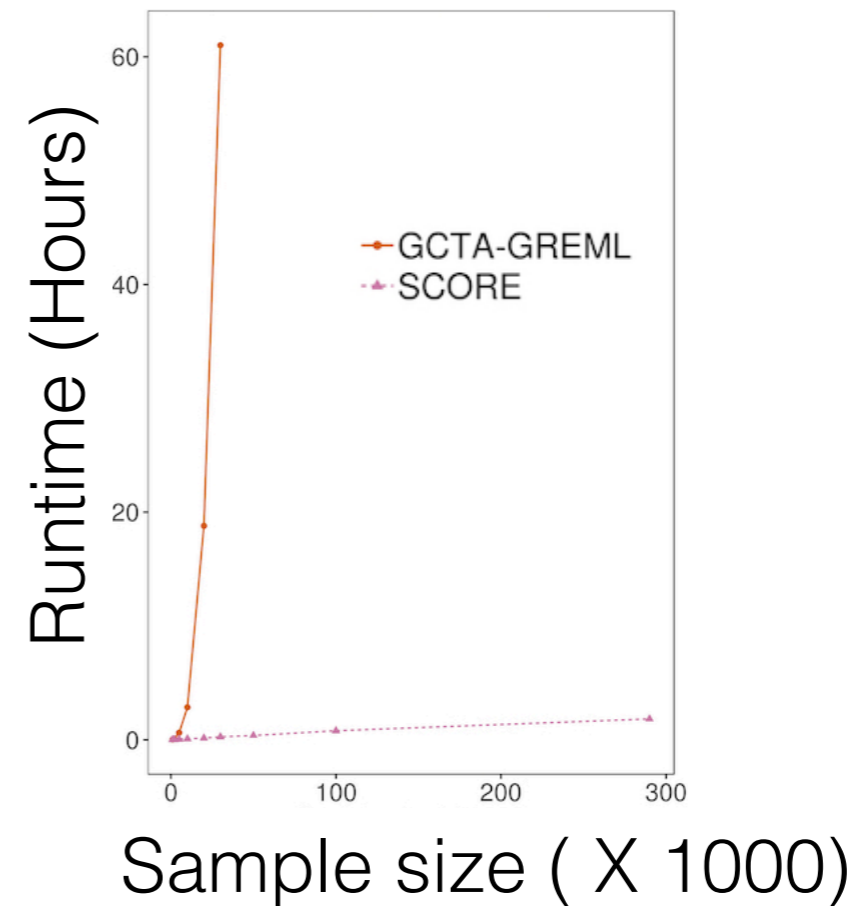
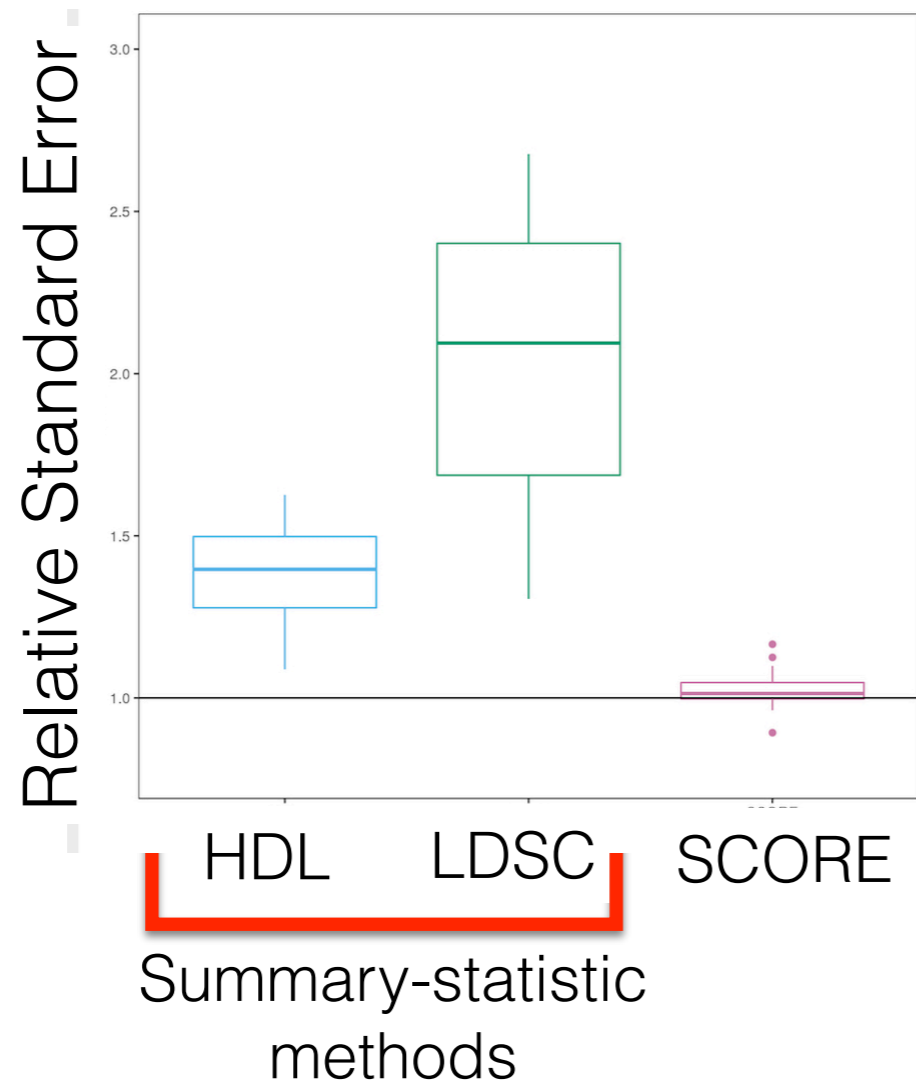
$$y_2 = X_2\beta_2 + \epsilon_2$$

$$\epsilon_{1n}, \epsilon_{2n} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \frac{\sigma_{e1}^2}{N} & \frac{\gamma_e}{N} \\ \frac{\gamma_e}{N} & \frac{\sigma_{e2}^2}{N} \end{bmatrix}\right)$$

$$\beta_{1m}, \beta_{2m} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \frac{\sigma_{g1}^2}{M} & \frac{\gamma_g}{M} \\ \frac{\gamma_g}{M} & \frac{\sigma_{g2}^2}{M} \end{bmatrix}\right)$$

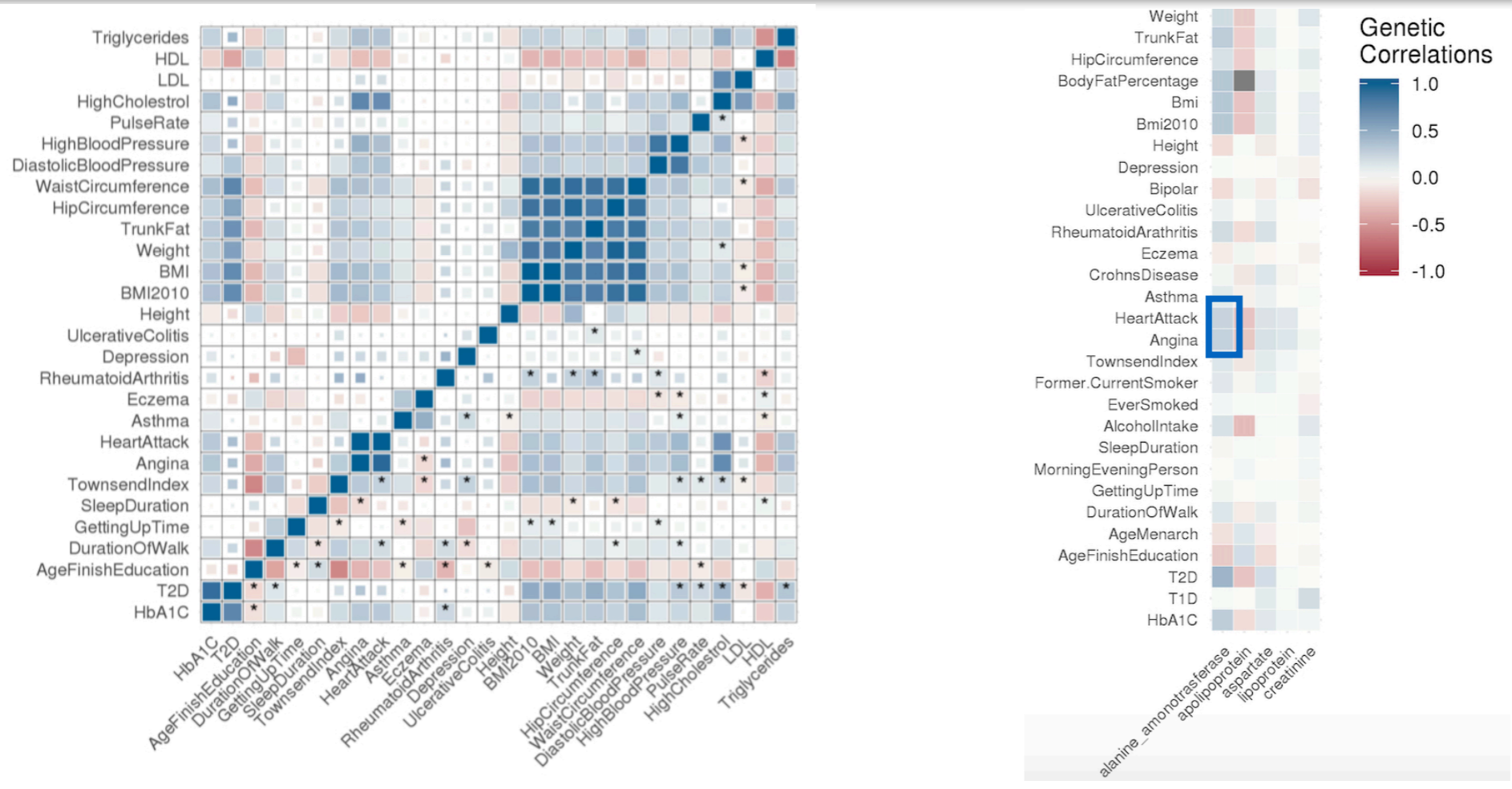
Coupling of genetic effects

Scalable genetic correlation estimator (SCORE)



Wu et al. RECOMB 2019, BioRxiv 2020

Genetic correlation



Wu et al. RECOMB 2019, BioRxiv 2020

Promises and challenges of Biobank-scale datasets

Challenging problems of statistical inference

Variance components analysis

Heritability estimation and partitioning

Non-linear contributions: Dominance and GxE

Genetic correlation

Promises and challenges of Biobank-scale datasets

Challenging problems of statistical inference

Variance components analysis

Action of negative selection on complex traits

Substantial additive heritability, negligible dominance heritability with substantial GxE for specific gene-environment combinations

Promises and challenges of Biobank-scale datasets

Challenging problems of statistical inference

Variance components analysis

Ongoing work

Biobank-scale association testing and prediction

Multi-ethnic analyses

Combining data modalities

Acknowledgments

UCLA

Bogdan Pasaniuc
Kathryn Burch
Paivi Pajukanta
Noah Zaitlen

University of Chicago

Andy Dahl

FIMM

Andrea Ganna

Tel Aviv University

Saharon Rosset



Funding

