

# Efficient variance components analysis across millions of genomes

Ali Pazokitoroudi<sup>1</sup>, Yue Wu<sup>1</sup>, Kathryn S. Burch<sup>2</sup>, Kangcheng Hou<sup>3</sup>, Bogdan Pasaniuc<sup>4,5,6</sup>, Sriram Sankararaman<sup>1,5,6,\*</sup>

<sup>1</sup>Department of Computer Science, UCLA, Los Angeles, California <sup>2</sup>Bioinformatics Interdepartmental Program, UCLA, Los Angeles, California <sup>3</sup>Department of Pathology and Laboratory Medicine, David Geffen School of Medicine, UCLA, Los Angeles, California <sup>4</sup>College of Computer Science and Technology, Zhejiang University, Hangzhou, Zhejiang, China <sup>5</sup>Department of Human Genetics, David Geffen School of Medicine, UCLA, Los Angeles, California <sup>6</sup>Department of Computational Medicine, David Geffen School of Medicine, UCLA, Los Angeles, California

\* Contact: [sriram@cs.ucla.edu](mailto:sriram@cs.ucla.edu)

## Abstract

While variance components analysis has emerged as a powerful tool in complex trait genetics, existing methods for fitting variance components do not scale well to large-scale datasets of genetic variation. Here we present a method for variance components analysis that is accurate and efficient: capable of estimating one hundred variance components on a million individuals genotyped at a million SNPs in a few hours. We illustrate the utility of our method in estimating and partitioning variation in a trait explained by genotyped SNPs (SNP-heritability). Analyzing 22 traits with genotypes from 300,000 individuals across about 8 million common and low-frequency SNPs, we observe that pre-allele effect sizes increase with decreasing minor allele frequency (MAF) and linkage disequilibrium (LD) consistent with the action of negative selection. Partitioning heritability across 28 functional annotations, we observe enrichment of heritability in PANTOM5 enhancers in asthma, eczema, thyroid and autoimmune disorders.

## Multi-component Linear Mixed Models (LMMs)

$$\begin{aligned} \mathbf{y} &= \sum_{k=1}^K \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\epsilon} \\ \boldsymbol{\epsilon} &\sim N(0, \sigma_e^2 I_N) \\ \boldsymbol{\beta}_k &\sim N(0, \frac{\sigma_k^2}{M_k} I_{M_k}), k \in \{1, \dots, K\} \end{aligned} \quad (1)$$

• Each of the  $M$  SNPs is assigned to one of  $K$  non-overlapping categories.

•  $\boldsymbol{\beta}_k$ :  $M_k$ -vector of SNP effect sizes for the  $k$ -th category.

•  $\sigma_k^2$ : the residual variance and  $\sigma_k^2$ : the variance component of the  $k$ -th category.

The total SNP heritability :

$$h_{SNP}^2 = \frac{\sum_{k=1}^K \sigma_k^2}{(\sum_{k=1}^K \sigma_k^2) + \sigma_e^2} \quad (2)$$

The SNP heritability of category  $k$ :

$$h_k^2 = \frac{\sigma_k^2}{(\sum_{k=1}^K \sigma_k^2) + \sigma_e^2}, k \in \{1, \dots, K\} \quad (3)$$

## Method of Moments

$$\text{cov}(\mathbf{y}) = E[\mathbf{y}\mathbf{y}^T] - E[\mathbf{y}]E[\mathbf{y}^T] = \sum_k \sigma_k^2 \mathbf{K}_k + \sigma_e^2 I_N \quad (4)$$

•  $\mathbf{K}_k = \frac{\mathbf{X}_k \mathbf{X}_k^T}{M_k}$ : the genetic relatedness matrix (GRM) computed from all SNPs of  $k$ -th category.

## Method of Moments objective function

Using  $\mathbf{yy}^T$  as our estimate of the empirical covariance, solve the following problem to find the variance components:

$$(\hat{\sigma}_1^2, \dots, \hat{\sigma}_K^2, \hat{\sigma}_e^2) = \underset{k}{\operatorname{argmin}} (\sigma_1^2, \dots, \sigma_K^2, \sigma_e^2) \| \mathbf{yy}^T - \sum_k \sigma_k^2 \mathbf{K}_k + \sigma_e^2 I_N \|_F^2 \quad (5)$$

The method of moments estimate satisfies the normal equations:

$$\begin{bmatrix} T & b \\ b^T & N \end{bmatrix} \begin{bmatrix} \sigma_1^2 \\ \vdots \\ \sigma_K^2 \\ \sigma_e^2 \end{bmatrix} = \begin{bmatrix} c \\ \mathbf{y}^T \mathbf{y} \end{bmatrix} \quad (6)$$

•  $T$ :  $K \times K$  matrix with entries  $T_{k,l} = \text{tr}(\mathbf{K}_k \mathbf{K}_l)$ ,  $k, l \in \{1, \dots, K\}$

•  $b$ :  $K$ -vector with entries  $b_k = \text{tr}(\mathbf{K}_k) = N$

•  $c$ :  $K$ -vector with entries  $c_k = \mathbf{y}^T \mathbf{K}_k \mathbf{y}$ .

• Computational complexity: GRM  $\mathcal{O}(N^2 M_k)$ . Computing  $T_{k,l}$ ,  $c_k$ ,  $k, l \in \{1, \dots, K\}$ :  $\mathcal{O}(N^2)$ . Solving the normal equation:  $\mathcal{O}(K^3)$

## Scalable Method-of-Moments estimator (RHE-mc)

The key bottleneck in solving the normal equation: the computation of  $K_{k,l}$ ,  $T_{k,l}$ ,  $c_k$ . Estimate  $T_{k,l}$ ,  $k, l \in \{1, \dots, K\}$  using an unbiased estimator:

$$T_{k,l} = \text{tr}(\mathbf{K}_k \mathbf{K}_l) \approx \widehat{T}_{k,l} = \frac{1}{B} \frac{1}{M_k M_l} \sum_b \mathbf{z}_b^T \mathbf{X}_k \mathbf{X}_k^T \mathbf{X}_l \mathbf{X}_l^T \mathbf{z}_b \quad (7)$$

•  $\mathbf{z}_1, \dots, \mathbf{z}_B$ :  $B$  independent random vectors with zero mean and covariance  $I_N$ .

• No explicit GRM computation.

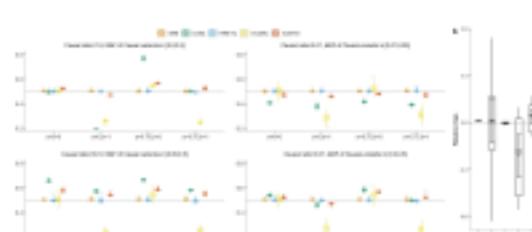
• Estimating the values  $T_{k,l}$  is  $\mathcal{O}(NM^2)$ .

• Can improve the time complexity to  $\mathcal{O}(\frac{NM^2}{\max(\log_2 N, \log_2 M)})$ .

Overall time complexity:  $\mathcal{O}(\frac{NM^2}{\max(\log_2 N, \log_2 M)}) + K^3$ .

## Results

### RHE-mc yields relatively unbiased estimates of total SNP heritability



## RHE-mc is scalable

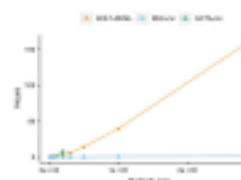
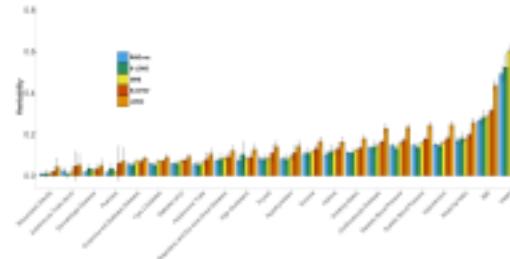
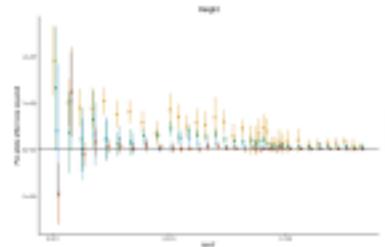


Table 1 Comparison of running time of RHE-mc, GCTA-mc, and REML-mc		
Processors	N	Time (min)
1	10,000	0.02
1	20,000	0.04
1	30,000	0.06
1	40,000	0.08
1	50,000	0.10
1	60,000	0.12
1	70,000	0.14
1	80,000	0.16
1	90,000	0.18
1	100,000	0.20
2	10,000	0.02
2	20,000	0.04
2	30,000	0.06
2	40,000	0.08
2	50,000	0.10
2	60,000	0.12
2	70,000	0.14
2	80,000	0.16
2	90,000	0.18
2	100,000	0.20
4	10,000	0.02
4	20,000	0.04
4	30,000	0.06
4	40,000	0.08
4	50,000	0.10
4	60,000	0.12
4	70,000	0.14
4	80,000	0.16
4	90,000	0.18
4	100,000	0.20

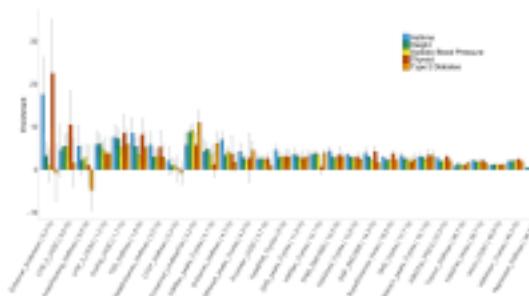
## Estimates of genome-wide SNP heritability in the UK Biobank



## Per-allele squared effect size of height as a function of MAF and LD (see [1] for other traits)



## RHE-mc estimates of heritability enrichment across functional annotations in the UK Biobank



## References

- [1] Ali Pazokitoroudi, Yue Wu, Kathryn S. Burch, Kangcheng Hou, Bogdan Pasaniuc, and Sriram Sankararaman. Efficient variance components analysis across millions of genomes. *Nature Communications*, 11(4020), 2020.

## Acknowledgements

This research was conducted using the UK Biobank Resource under applications 33127 and 33297. We thank the participants of UK Biobank for making this work possible. We thank Rob Brown for feedback on this manuscript. This work was funded by NIH grants R01HG0309120 (B.P. and K.S.B.), R35GM125055 (S.S.) an Alfred P. Sloan Research Fellowship (S.S.), and a NSF grant III-1705121 (Y.W. and S.S.).