

1 Title Page

2 Title

3 A versatile, fast and unbiased method for estimation of gene-by-environment interaction effects  
4 on biobank-scale datasets

5

6 Author names

7 Mohammad Khan, MSc<sup>1,2</sup>, Matteo Di Scipio, BSc<sup>1,2</sup>, Conor Judge, BMBS, BEng<sup>1</sup>, Nicolas  
8 Perrot, PhD<sup>1</sup>, Michael Chong, MSc<sup>1</sup>, Shihong Mao, PhD<sup>1</sup>, Shuang Di, BSc, LLB, MEd, MSc<sup>3,4</sup>,  
9 Walter Nelson, BSc<sup>3</sup>, Jeremy Petch, PhD<sup>1,2,3,5</sup>, \*Guillaume Paré, MD, MSc<sup>1,6,7,8</sup>

10

11 Affiliations

12 (1) Population Health Research Institute, David Braley Cardiac, Vascular and Stroke

13 Research Institute, Hamilton Health Sciences and McMaster University, Hamilton,

14 Canada.

15 (2) Department of Medicine, Faculty of Health Sciences, McMaster University, Hamilton,

16 ON, Canada.

17 (3) Centre for Data Science and Digital Health, Hamilton Health Sciences, Hamilton, ON,

18 Canada.

19 (4) Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada.

20 (5) Institute of Health Policy, Management and Evaluation, University of Toronto, Toronto,

21 ON, Canada

22 (6) Thrombosis and Atherosclerosis Research Institute, David Braley Cardiac, Vascular and  
23 Stroke Research Institute, Hamilton, Canada.

24 (7) Department of Pathology and Molecular Medicine, McMaster University, Michael G.  
25 DeGroot School of Medicine, Hamilton, Canada.

26 (8) Department of Health Research Methods, Evidence, and Impact, McMaster University,  
27 Hamilton, Canada.

28

## 29 Contact information

30 Dr Guillaume Paré, Population  
31 Health Research Institute,  
32 Hamilton Health Sciences and  
33 McMaster University, Hamilton,  
34 ON L8L 2X2, Canada  
35 pareg@mcmaster.ca

36

37

## 38 Figures/Tables

39 7 (6 figures; 1 table);  
40 7; supplemental

41

## 42 Keywords

43 Gene-environment interactions, genome-wide linear regression, cardiometabolic biomarkers,  
44 interaction polygenic score, waist-hip-ratio

## 45 Abstract

46 Current methods to evaluate gene-by-environment (GxE) interactions on biobank-scale datasets  
47 are limited. MonsterLM enables multiple linear regression on genome-wide datasets, does not  
48 rely on parameters specification and provides unbiased estimates of variance explained by GxE  
49 interaction effects. We applied MonsterLM to the UK Biobank for eight blood biomarkers  
50 (N=325,991), identifying significant genome-wide interaction variance with waist-to-hip ratio  
51 for five biomarkers, with variance explained by interactions ranging from 0.11 to 0.58. 48% to  
52 94% of GxE interaction variance can be attributed to variants without significant marginal  
53 association with the phenotype of interest. Conversely, for most traits, >40% of interaction  
54 variance was explained by less than 5% of genetic variants. We observed significant  
55 improvements in polygenic score prediction with incorporation of GxE interactions in four  
56 biomarkers. Our results imply an important contribution of GxE interaction effects, driven  
57 largely by a restricted set of variants distinct from loci with strong marginal effects.

58

## 59 Introduction

60 Identifying gene-by-environment (GxE) interactions is difficult because individual interaction  
61 effects are expected to be small<sup>1</sup>, the multiple hypothesis burden is considerable<sup>2,3</sup>, and the  
62 sample sizes needed are correspondingly large<sup>4</sup>. Most previous analyses have focused on  
63 identifying interactions with variants marginally associated with a phenotype of interest<sup>5,6</sup>.  
64 Hitherto, methods developed to estimate the overall effect of these interactions rely on variance  
65 component methods, due to the predictor ( $m$ ) > observation ( $n$ ) problem, where single nucleotide  
66 polymorphisms (SNPs) ( $m$ ) vastly outnumber the participants ( $n$ )<sup>7,8</sup>. These methods are  
67 advantageous for smaller datasets; however, they can be limited when applied to larger datasets  
68 due to computational burden<sup>7</sup>. Furthermore, variance component methods depend on strong  
69 assumptions about the underlying genetic model and often require *a priori* specification of  
70 parameters and/or hyper-parameters, such as polygenicity, minor allele frequency (MAF), and  
71 linkage disequilibrium (LD) dependence<sup>9-13</sup>. While never formally tested in the context of GxE  
72 interactions, it has previously been shown these assumptions can lead to important biases in  
73 heritability estimates<sup>9-11,14-18</sup>. Novel methods are thus needed to enable fast and unbiased  
74 calculations of the variance explained ( $R^2$ ) by GxE interactions in large samples, on multiple  
75 traits and without the need for genetic model assumptions.

76  
77 Our method is similar to the generalized random effects (GRE) model<sup>19</sup>, building on the  
78 observation that the multiple regression coefficient of determination can be used to accurately  
79 estimate heritability<sup>19</sup>. Extending this observation to include an environmental exposure variable  
80 and computing the interactions between genotypes and the environmental exposure allows us to  
81 examine the variance explained by genetic interactions with an environmental exposure. Using

82 linear regression across the genome presents a key problem: there are far more SNPs ( $m$ ) than  
83 participants ( $n$ ) in genome-wide studies, and thus it becomes difficult to estimate heritability and  
84 interaction variance<sup>20,21</sup>. By partitioning the genome into non-overlapping regions, it becomes  
85 possible to estimate genome-wide interactions with environmental exposures by reducing  $m$   
86 within each region to a size where  $m < n$ . However, partitioning the genome into large blocks  
87 still presents challenges. First, LD spillage at the junction of blocks can theoretically inflate  
88 heritability estimates if many such junctions exist<sup>9</sup>. Second, any residual population stratification  
89 effects would be amplified if heritability at each region is overestimated and this effect is  
90 expected to be proportional to the number of blocks<sup>22</sup>. Third, computing prediction  $R^2$  on large  
91 blocks with high dimensionality can be slow. By using the conjugate gradient method<sup>23</sup> with  
92 graphics processing unit (GPU) acceleration<sup>24</sup>, it is possible to perform multiple linear regression  
93 modelling efficiently on large (25,000 SNPs) blocks (Supplementary Table 1). The potential for  
94 residual population stratification effects and LD spills are minimized as only approximately 60  
95 blocks are typically needed for genome-wide analyses and variants are LD-pruned. A block size  
96 of 25,000 SNPs also ensures that  $n > 10m$  for accurate estimations.

97

98 We propose a novel method, MonsterLM, to estimate the proportion of variance explained by  
99 GxE interactions for continuous traits, in a fast, accurate, efficient and unbiased manner on  
100 biobank-scale datasets ( $N > 300,000$ ). We hypothesized that GxE interactions contribute  
101 significantly to complex trait variance. Our objective was to quantify and characterize these  
102 contributions for continuous traits. We illustrate an overview of our computational analyses in  
103 Figure 1.

104

## 105 Methods

### 106 UK Biobank

107 The UK Biobank is a large population-based study which includes over 500,000 participants  
108 living in the United Kingdom<sup>15,32</sup>. Men and women aged 40–69 years were recruited between  
109 2006 and 2010, and extensive phenotypic and genotypic data was collected. We selected 325,991  
110 unrelated British individuals from the UK Biobank with both genotype and biomarker data for  
111 inclusion in the analysis. This study used genetic variants from the ‘V3’ release of the UK  
112 Biobank data including those present in the Haplotype Reference Consortium and 1000 Genomes  
113 panels with imputation quality greater than 0.7, no deviation from Hardy-Weinberg equilibrium  
114 ( $P > 1 \times 10^{-10}$ ) and minor allele frequency greater than 1%<sup>15</sup>. Genotype data were filtered by  
115 removing highly correlated SNPs with a LD  $r^2$  value of more than 0.9 and removing SNPs with a  
116 MAF of less than 0.01, as the focus of this report is on common variants. After quality control  
117 filtering, there remained 1,031,135 SNPs and 325,991 individuals. Raw genotypes were  
118 normalized to have a mean of zero and variance of one. For the current analysis we examined  
119 eight biomarkers including Apolipoprotein B, Bilirubin, Total Cholesterol, C-reactive protein  
120 (CRP), HbA1c, HDL-Cholesterol, LDL-Cholesterol, Triglycerides and the environmental  
121 exposure, waist-to-hip ratio (WHR).

122

123 For secondary analyses, we randomly partitioned the UK Biobank participants into two sets: a  
124 discovery set containing 80% of the participants used for model building and a validation set  
125 containing the remaining 20% of the participants. This was done to remove the potential for  
126 overfitting that can occur when using derived models on the same datasets for prediction  
127 purposes<sup>33</sup>.

128

## 129 MonsterLM Estimations of Variance Explained by GxE Effects

130 The standard linear model for a phenotypic trait ( $Y$ ) when an interaction term is included can be  
131 expressed as:

$$Y = \beta_G G + \beta_E E + \beta_{GE} GE + \epsilon \quad (1)$$

132 Where  $G$  is the genotype matrix,  $E$  is the environmental exposure,  $GE$  is the Hadamard product  
133 between each genotype and environmental exposure, resulting in a matrix with the same  
134 dimensionality as  $G$ . The betas ( $\beta$ ) represent the true marginal effects associated with their  
135 respective term. To account for covariate effects such as age, sex,  $E$ , and population stratification  
136 we first regress  $Y$  onto the covariates and the first twenty genetic principal components and  
137 extract the residuals of the model ( $y_{residuals}$ ). The residuals ( $y_{residuals}$ ) become our phenotype  
138 used for analyses in MonsterLM:

$$y_{residuals} = \beta_G G + \beta_{GE} GE + \epsilon \quad (2)$$

139 Both phenotype and environmental exposures are quantile normalized after residualization, such  
140 that mean is zero and variance one. Through residualization of the environmental exposure, we  
141 can leave  $E$  out of the model. For simplicity, we denote the augmented matrix of  $G$  and  $GE$  as  $U$   
142 with dimension  $n \times 2m$ , where  $n$  is number of participants included and  $m$  is number of SNPs:

$$U = G \mid GE \quad (3)$$

143 And  $y_{residuals}$  becomes:

$$y_{residuals} = \beta_U U + \epsilon \quad (4)$$

144

145 The MonsterLM method enables multiple linear regression on biobank-scale datasets by  
146 parallelizing the calculation of least squares regression, including the interaction terms, between

147 the genotypes and environmental factors. The calculation is done such that the only practical  
 148 limitation is the inversion of the  $U$  matrix, but without any restriction on  $n$ . This limitation is  
 149 circumvented using the conjugate gradient method and GPU acceleration<sup>24</sup>. Importantly,  
 150 MonsterLM requires neither parametrization nor assumptions regarding the genetic architecture  
 151 of traits analyzed (such as polygenicity of effects, MAF and LD dependence). Genotypic data  
 152 was partitioned into blocks with a maximal size of 25,000 SNPs ( $m$ ) to minimize LD spillage  
 153 between blocks and to optimize speed of the matrix calculation.

154

155 Given a quantitative trait  $Y$ , the least squares estimate for  $\hat{\beta}_U$ , the estimated effects vector,  
 156 corresponding to the genotype and GxE interaction is:

$$\hat{\beta}_U = (U^T U)^{-1} U^T Y \quad (5)$$

157

158 After computing  $\hat{\beta}_U$  using conjugate gradient, the predicted values of  $Y$  denoted as  $\hat{y}$ , can be  
 159 computed as:

$$\hat{y} = \hat{\beta}_U U \quad (6)$$

160 This same method can be applied if we use the genotype matrix only ( $G$ ) instead of  $U$  to compute  
 161  $\hat{\beta}_G$  and  $\hat{y}$ . Once  $\hat{y}$  is calculated for each block (with and without interactions), we calculate the  
 162 variance explained for the full model ( $U$ ) and the model without interactions. Since  $R^2$  is a biased  
 163 estimator, the adjusted  $R^2$  ( $\overline{R^2}$ ) is used as our estimate for variance explained. Then, to calculate  
 164 the interaction variance explained ( $\overline{R^2}_{GE_i}$ ) we compute the difference in  $\overline{R^2}_{GE_i}$  as:

$$\overline{R^2}_{GE_i} = \overline{R^2}_{U_i} - \overline{R^2}_{G_i} \quad (7)$$

165



166 Since we remove SNPs in very high LD ( $r^2 > 0.9$ ), the remaining variant set will not be highly  
 167 correlated. We can then estimate the total contribution of variance by genome-wide environment  
 168 interaction ( $\overline{R^2}_{GWE}$ ) by taking the sum over all blocks:

$$\overline{R^2}_{GWE} = \sum_{i=1}^j \overline{R^2}_{GE_i} \quad (8)$$

169 Where  $j$  is the number of partitioned blocks used for analysis (i.e. 60 blocks for current analyses)  
 170 and  $i$  is the index of the current block.

171

172 The 95% confidence (CI) of the  $R^2_{GE_i}$  term can be estimated for each block using asymptotic  
 173 properties described by Graf and Alf<sup>34</sup>. The asymptotic variance for the difference between  
 174 the  $R^2_{GE_i}$  of two models is given by:

$$\widehat{Var}_\infty(R^2_{GE_i}) = \widehat{Var}_\infty(R^2_{U_i}) + \widehat{Var}_\infty(R^2_{G_i}) - 2\widehat{Cov}_\infty(R^2_{U_i, G_i}) \quad (9)$$

175

176 Where:

$$\widehat{Var}_\infty(R^2_{U_i}) = 4R^2_{U_i}(1 - R^2_{U_i})^2/n \quad (10)$$

$$\widehat{Var}_\infty(R^2_{G_i}) = 4R^2_{G_i}(1 - R^2_{G_i})^2/n \quad (11)$$

$$2\widehat{Cov}_\infty(R^2_{U_i, G_i}) = 8r_{U_i}r_{G_i}(0.5(2r_{U_i, G_i} - r_{U_i}r_{G_i})(1 - r_{U_i, G_i}^2 - r_{G_i}^2 - r_{U_i}^2) + r_{U_i, G_i}^3)/n \quad (12)$$

177

178

179 And:

$$r_{U_i, G_i} = \frac{r_{G_i}}{r_{U_i}} \quad (13)$$

180 The 95% CI for a single block can then be derived using the Wald estimate:

$$95\% \text{ CI} = R^2_{GE_i} \pm 1.96 \sqrt{\widehat{\text{Var}}_{\infty}(R^2_{GE_i})} \quad (14)$$

181 To estimate the 95% CI for our  $\overline{R^2}_{GWE}$  estimate, we calculate the total asymptotic variance as the  
 182 sum of the individual variances ( $\overline{R^2}_{GE_i}$ ) for each block, but since our estimates use  $\overline{R^2}_{GE}$  to  
 183 estimate GxE interactions, we also adjust the variance of  $R^2_{GE_i}$  similarly to obtain the adjusted  
 184 asymptotic variance of  $\overline{R^2}_{GWE}$ , which is then used to calculate the 95% CIs.

$$\widehat{\text{Var}}_{\infty}(\overline{R^2}_{GWE}) = \sum_{i=1}^j \left( \frac{n-1}{n-m-1} \right)^2 \widehat{\text{Var}}_{\infty}(R^2_{GE_i}) \quad (15)$$

185 Where  $n$  is the number of samples and  $m$  is the number of SNPs tested per block  $i$ . With the total  
 186 asymptotic variance estimated, we calculate the 95% CI for the  $\overline{R^2}_{GWE}$  as:

$$95\% \text{ CI} = \overline{R^2}_{GWE} \pm 1.96 \sqrt{\widehat{\text{Var}}_{\infty}(\overline{R^2}_{GWE})} \quad (16)$$

## 187 Simulations to Validate the MonsterLM Method

188 We tested MonsterLM with simulations using UK Biobank genotypes filtered as described  
 189 above. We used chromosome 18 to generate a single block of 8,913 SNPs (smallest block  
 190 allowing for efficient simulations). We then simulated the true, unobserved effects ( $\beta_G, \beta_E, \beta_{GE}$ )  
 191 from a normal distribution, assuming 20% of SNPs have a marginal effect associated with the  
 192 simulated trait of interest,  $Y_{sim}$  (i.e.  $\beta_G \neq 0$ ). We further assumed that 10% of the causal SNPs  
 193 (i.e. 2% of total SNPs) have an interaction effect (i.e.  $\beta_{GE} \neq 0$ ). The values were chosen based on  
 194 similar estimates with heritability of WHR through MonsterLM. The error was sampled from an

195 independent and identically distributed normal distribution. The simulated trait ( $Y_{sim}$ ) was then  
196 computed as:

$$Y_{sim} = \beta_G G + \beta_E E + \beta_{GE} GE + \epsilon \quad (17)$$

197 We divided the above case into three scenarios. The first scenario considered that  $E$  was not  
198 dependent on  $G$  and the genetic and interaction effects for all SNPs were randomly generated  
199 from a standard normal distribution. The next two scenarios considered that  $E$  was dependent on  
200  $G$ . In these scenarios,  $E$  was simulated to have 20% of its variance explained by  $G$  (i.e.  
201 heritability), as WHR was observed to have similar heritability empirically. Scenario 2 further  
202 assumed that the genetic effects could be zero when the interaction effect was non-zero for a  
203 specific SNP  $i$  ( $\beta_{G,i} = 0, \beta_{GE,i} \neq 0$ ) and the SNPs explaining  $E$  were the same as the SNPs with an  
204 interaction effect. Scenario 3 assumed that both the genetic and interaction effects were non-zero  
205 for a specific SNP  $i$  ( $\beta_{G,i} \neq 0, \beta_{GE,i} \neq 0$ ) and that the SNPs explaining  $E$  were not the same as the  
206 SNPs with an interaction effect. To ensure realistic scenarios were simulated, we varied the  
207 variance of the normal distributions to achieve pre-specified genetic, environment and interaction  
208 effects. The heritability ( $\overline{R^2}_G$ ) was set to 0.025, variance explained by the environmental  
209 exposure ( $\overline{R^2}_E$ ) was set to 0.2 and variance explained by the interactions ( $\overline{R^2}_{GE}$ ) was set to 0.005.  
210 We also considered 3 multi-block scenarios identical to the above scenarios; whereby  
211 chromosome 11 was split into 3 blocks of roughly 15,000 SNPs each. Each block had  $\overline{R^2}_G$  set to  
212 0.025,  $\overline{R^2}_E$  to 0.2, and  $\overline{R^2}_{GE}$  set to 0.005, such that the variance explained by interactions across  
213 the whole chromosome ( $\overline{R^2}_{GWE}$ ) was 0.015.  
214

## 215 Directionality of Effects Analysis

216 After computing  $\overline{R^2}_{GWE}$  for our eight biomarkers, we tested whether direction of effect was  
217 concordant between marginal and interaction regression coefficients for each SNP. Concordant  
218 direction of effects is defined as when  $\hat{\beta}_G$  has the same sign (+/+, -/-) as  $\hat{\beta}_{GE}$  for a single SNP and  
219 its associated interaction. Discordant direction of effects is defined as when the  $\hat{\beta}_G$  and  $\hat{\beta}_{GE}$  have  
220 a different sign (+/-, -/+) for a single SNP and its associated interaction. We used a subset of  
221  $\hat{\beta}_G$  and  $\hat{\beta}_{GE}$  coefficients that were in low LD ( $r^2 < 0.1$ ) and computed the direction of effect  
222 concordance for this subset. We then plotted the sign concordance twice: first as a function of  $\hat{\beta}_G$   
223  $P$  – values ( $P_G$ ), then as a function of  $\hat{\beta}_{GE}$   $P$  – values ( $P_{GE}$ ), which were computed from  
224 association of single SNPs and their respective interaction on the biomarker traits. Two-  
225 proportion Z-tests were used to compare the proportion of directionally concordant marginal and  
226 interaction effects for each biomarker in each threshold compared to a null count at a proportion  
227 of 0.50.

## 228 Stratification of Estimates by MAF and LD

229 SNPs were stratified by MAF and LD score into a total of 20 bins: 5 MAF bins ( $0.01 \leq 0.1$ ,  
230  $0.1 < \text{MAF} \leq 0.2$ ,  $0.2 < \text{MAF} \leq 0.3$ ,  $0.3 < \text{MAF} \leq 0.4$ , and  $0.4 < \text{MAF} \leq 0.5$ ) and 4 LD score quantiles  
231 ( $0 < \text{LD} \leq 0.25$ ,  $0.25 < \text{LD} \leq 0.50$ ,  $0.50 < \text{LD} \leq 0.75$ , and  $0.75 < \text{LD} \leq 0.9$ ). MAF and LD score were  
232 calculated using a subset of 5000 participants from the UKBiobank. We then computed the  
233 variance explained ( $\overline{R^2}_{Gbins}$ ,  $\overline{R^2}_{GWEIbins}$ ) and divided each estimate by the total number of SNPs  
234 in each bin to get an  $\overline{R^2}$  per SNP value that was compared between bins and to the total genetic  
235 and interaction variance estimates ( $\overline{R^2}_{Gbins}$ ,  $\overline{R^2}_{GWEIbins}$ ).

236

## 237 Polygenic Scores Analysis

238 To calculate polygenic scores ( $PS$ ) without interactions ( $PS_G$ ), we first selected SNPs based on  
239 univariate  $P_G$  derived from regression of each variant with biomarker concentration from the  
240 discovery set. We then combined the selected SNPs into a single block from the discovery set,  
241 and applied MonsterLM regression to obtain the multiple linear regression coefficients ( $\hat{\beta}_G$ ).  
242 Using these coefficients, we calculated the  $PS_G$  in the validation set as:

$$PS_{G,i} = \sum_j^O G_{i,j} \hat{\beta}_{G,j} \quad (18)$$

243 Where  $PS_{G,i}$  is the individual polygenic score of participant  $i$ ,  $j$  is the SNP number and  $O$   
244 represents the total number of SNPs included in this analysis. We then evaluated the  
245 predictiveness of each  $PS_G$  using  $\overline{R^2}$  in the validation set. We repeated the same process for four  
246 univariate  $P_G$  thresholds ( $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ ) for each biomarker.

247  
248 We define  $PS_{GE}$  as the  $PS$  with GxE interactions included. To include GxE interactions, we  
249 selected significant interactions based on  $P_{GE}$  obtained from regressing each variant and its  
250 associated GxE interaction with biomarker concentration in the discovery set. These interactions  
251 are selected from the subset of SNPs included in polygenic scores without interactions. The  
252 interactions passing the univariate  $P_{GE}$  thresholds ( $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ ) were then included with  
253 the SNPs to create a single block. We applied MonsterLM regression to obtain the multiple  
254 linear regression coefficients ( $\hat{\beta}_G$ ,  $\hat{\beta}_{GE}$ ). Using these coefficients, we calculate the  $PS_{GE}$  as:

$$PS_{GE,i} = \sum_j^O G_{i,j} \hat{\beta}_{G,j} + \sum_k^P (G \times E)_{i,k} \hat{\beta}_{GE,k} \quad (19)$$

255 Where  $PS_{GE,i}$  is the polygenic score with interactions incorporated for participant  $i$ , summed over  
256 each SNP ( $j$ ) and, if included, its associated interaction ( $k$ ).  $O$  represents the SNPs included in the  
257  $PS_{GE}$ , while  $P$  represents the interactions included, a subset of  $O$ . As with the  $PS_G$ , we evaluated  
258 the predictiveness of each polygenic score using  $\overline{R^2}$  in the validation set. We repeated for all  
259 pairwise combinations of the four  $P_G$  thresholds and the four  $P_{GE}$  thresholds, resulting in 16  $PS_{GE}$   
260 for each biomarker.

## 261 Results

### 262 Simulation results

263 We conducted 100 simulations for each of the three scenarios (Figure 2A). On average, the  
264  $\overline{R^2}_{GE}$  was observed to be close to 0.005, the true underlying  $\overline{R^2}_{GE}$  that was predefined for  
265 interactions. We compared the estimated  $\overline{R^2}_{GE}$  to the true  $\overline{R^2}_{GE}$  and found that the difference in  
266  $\overline{R^2}_{GE}$  was not significant ( $P>0.05$ ). After verifying MonsterLM for a single block, we conducted  
267 100 simulations using three contiguous blocks from chromosome 11 under the same three  
268 scenarios. The  $\overline{R^2}_{GWE}$  was observed to be close to 0.015, the true  $\overline{R^2}_{GWE}$  set for interactions  
269 (Figure 2B). Our calculated 95% CIs were also well-calibrated for our simulated data.

270

### 271 Estimation of Genome-Wide Environmental Interaction Effects

272 Next, we applied MonsterLM to estimate the variance explained by interactions between waist-  
273 hip-ratio (WHR) and genetic variants for eight blood biomarkers (Apolipoprotein B, Bilirubin,  
274 Total Cholesterol, CRP, HbA1c, HDL-Cholesterol, LDL-Cholesterol, Triglycerides) linked to  
275 cardio-metabolic diseases. WHR was selected as the environmental exposure because it is a  
276 measure of central obesity linked to a wide range of adverse metabolic consequences, including

277 diabetes and cardiovascular disease (CVD)<sup>25</sup>. As such, it represents an excellent marker of the  
278 effect of the modern obesogenic environment on metabolism. We observed significant variance  
279 explained by interaction effects for five of the eight biomarkers, with interaction  $\overline{R^2}$  ranging from  
280 0.11 to 0.58 (Figure 3A). As expected, all heritability estimates were significant and consistent  
281 with previous work<sup>16</sup>. Furthermore, we observed the presence of significant directionality for  
282 interaction effects at both  $P_G$  and  $P_{GE} < 10^{-3}$  significance threshold (Figure 3B; Supplementary  
283 Figure 1). When stratifying variants according to MAF and LDscore, there was a general  
284 tendency for SNPs with low MAF (i.e.  $0.01 < \text{MAF} < 0.1$ ) and higher LDscore to  
285 disproportionally contribute to interaction variance explained per SNP (Supplementary Figure 2)  
286  
287 The presence of significant gene-by-WHR (GxWHR) interactions prompted additional questions.  
288 First, do GxE interactions arise from SNPs strongly associated with the trait of interest, as has  
289 been commonly assumed, or are the variants contributing to GxE interactions independent from  
290 those with marginal effects? To address this question, we randomly split participants into a  
291 discovery set comprising 80% of participants (260,792 individuals) with the remaining 20%  
292 comprising the validation set. Using the five biomarkers with significant GxE interaction  
293 variance, we conducted linear regression on the discovery set using biomarker concentration as  
294 the outcome variable and a single SNP as the predictor variable, repeating this process for all  
295 SNPs and extracting  $P_G$ . We then selected SNPs according to six association  $P_G$  thresholds:  $<1$   
296 (i.e. all SNPs),  $< 10^{-1}$ ,  $<10^{-2}$ ,  $<10^{-3}$ ,  $<10^{-4}$ ,  $<10^{-5}$ . Each SNP set was then tested for association  
297 with the corresponding biomarker in the validation set, using the least number of blocks possible.  
298 We evaluated the total  $\overline{R^2}_G$  and  $\overline{R^2}_{GE}$  for each of the five SNP sets. The  $\overline{R^2}_G$  and  $\overline{R^2}_{GE}$  was then  
299 compared to the variance explained when including all SNPs (i.e.  $P_G < 1$ ) for the validation set

300 ( $\overline{R^2}_{G-val}$  and  $\overline{R^2}_{GWE-val}$ ). We estimated the proportion of  $\overline{R^2}_G$  recovered when including an  
301 increasing proportion of SNPs in the analysis (Figure 4; Supplementary Figure 3). We observed  
302 that between 51-86% of the original  $\overline{R^2}_{G-val}$  calculated in the validation set could be recovered  
303 only using SNPs with  $P_G < 10^{-3}$  from the discovery set (Supplementary Table 2). We then  
304 similarly estimated the proportion of variance explained by GxE interactions recovered when  
305 including an increasing proportion of SNPs, based on  $P_G$ . At a  $P_G$  threshold of  $< 10^{-3}$ , only 1-8%  
306 of total  $\overline{R^2}_{GWE-val}$  was recovered in the validation set (Figure 4), suggesting that a majority of  
307  $\overline{R^2}_{GWE-val}$  involves SNPs with  $P_G > 10^{-3}$ . At the  $P_G < 10^{-2}$  threshold, the interaction variance  
308 recovered ranged from 2-13% whereas the corresponding range was 0-58% at the  $P_G < 10^{-1}$   
309 threshold.

310

311 As our results showed that GxE interactions are largely derived from SNPs without strong  
312 marginal associations, we next sought to address whether a few strong GxE interactions are  
313 responsible for the large variance explained by interactions, or whether it is the result of many  
314 small interactions. We conducted regression on each SNP and its associated interaction from the  
315 discovery set. We selected interactions based on five discovery  $P_{GE}$  thresholds:  $< 1$  (i.e. all  
316 SNPs),  $< 10^{-1}$ ,  $< 10^{-2}$ ,  $< 10^{-3}$ ,  $< 10^{-4}$ ,  $< 10^{-5}$ . In other words, an interaction term was included in the  
317 validation sample analysis if it passed the  $P_{GE}$  threshold in the discovery set. Importantly, all  
318 SNPs were included in the analysis, irrespective of whether their corresponding interaction terms  
319 were included or not. The interaction  $\overline{R^2}_{GWE}$  were computed in the validation set and compared  
320 to the  $\overline{R^2}_{GWE-val}$  estimates (Figure 5; Supplementary Figure 4). We observed that up to 45% of  
321 the total  $\overline{R^2}_{GWE-val}$  was recovered at a discovery  $P_{GE}$  threshold  $< 10^{-3}$ , corresponding to 0.2-3.3%  
322 of the SNPs tested in our initial analyses (Supplementary Figure 4). Indeed, high recovery of



323 variance explained by interaction was also observed at the  $P_{GE} < 10^{-2}$  (range: 14-78%) and  
324  $P_{GE} < 10^{-1}$  (range: 48-94%) thresholds. To confirm the specificity of interaction effects, we  
325 conducted a sensitivity analysis using Apolipoprotein B (Supplementary Table 3). We randomly  
326 selected a set of interaction terms equal to the number of interactions included at the  $P_{GE} < 10^{-2}$   
327 threshold (62,904 SNPs out of 1.2 million SNPs tested). We then calculated the  $\overline{R^2}_{GWE-val}$   
328 using this set of randomly chosen interaction effects. The randomly selected SNPs had an  
329  $\overline{R^2}_{GWE-val}$  of 0.02, compared to an  $\overline{R^2}_{GWE-val}$  of 0.25 for the interaction terms with  $P_{GE} < 10^{-2}$  in  
330 the validation set.

331

### 332 Polygenic Scores Analysis

333 Finally, we examined if the predictiveness of polygenic score ( $PS$ ) could be improved by  
334 incorporating interactions. To select SNPs and interaction effects to be included in each  $PS$ , we  
335 used both  $P_G$  and  $P_{GE}$  thresholds of  $10^{-2}$ ,  $10^{-3}$ ,  $10^{-4}$ , and  $10^{-5}$  in the discovery set when testing  
336 either each SNP individually or both a single SNP and corresponding interaction, respectively.  
337 Each  $PS$  was then tested in the validation sample for association with its corresponding  
338 biomarker.  $PS$  prediction  $\overline{R^2}$  was modestly improved for the four biomarkers with the highest  
339 interaction variance by incorporating interaction effects (Figure 6), with the relative increase in  
340 prediction  $\overline{R^2}$  ranging from 0% to 8% across the biomarkers analyzed. Significant improvements  
341 in prediction of Apolipoprotein B, Bilirubin and HDL-Cholesterol levels were observed at the  
342 95% confidence level (for interaction significance thresholds of  $10^{-3}$ ,  $10^{-4}$ ,  $10^{-5}$ ; Supplementary  
343 Table 4). Notably, there was no improvement in the Total Cholesterol  $PS$  with  
344 interactions ( $PS_{GE}$ ) compared to their respective  $PS$  without interactions ( $PS_G$ ) values (Figure  
345 6), consistent with the  $P_{GE}$  results for Total Cholesterol (Figure 5).

## 346 Discussion

347 In this report, we developed a novel method, MonsterLM, to estimate variance explained by  
348 genome-wide interactions with environmental exposures. Using simulations, we verified that  
349 MonsterLM estimates the variance explained by interaction effects accurately and precisely.

350 Analysis of UK Biobank biomarker data demonstrated the presence of significant GxE  
351 interactions effects with WHR, a marker of metabolically deleterious adiposity. The interaction  
352 estimates for five of the eight biomarkers analysed were significant with estimates ranging  
353 between 0.11 to 0.58 of overall variance, prompting further analyses into these results.

354  
355 MonsterLM provides distinct advantages over current methods for GxE analysis (Table 1)<sup>26-30</sup>.

356 In most settings, inference methods for genome-wide SNP-heritability and GxE interactions  
357 make assumptions on genetic architecture. These assumptions are parametrized by polygenicity  
358 (the number of variants with effects) and MAF/LD-dependence (the coupling of effects with  
359 MAF, LD or other functional annotations). Since the true genetic architecture of any given trait  
360 is unknown, existing methods are susceptible to bias and often yield vastly different estimates  
361 even when applied to the same data<sup>10-12</sup>. This is also the case for the estimation of Genome-Wide  
362 Environment interactions, where different assumptions about the structure of interactions result  
363 in a variety of different estimates<sup>26-30</sup>. Although multi-component methods that stratify SNPs by  
364 LD/MAF can address these robustness issues, fitting multiple variance components to biobank  
365 scale data is highly resource intensive<sup>16</sup>, and this problem is compounded when considering  
366 interactions where the number of variables analyzed increases by two-fold. Alternate methods  
367 that explicitly model these dependencies are also sensitive to model misspecification<sup>9-13</sup>.

368 MonsterLM makes no assumption with respect to the genetic model and does not rely on

369 parametrization for underlying assumptions. Our partitioning approach combined with methods  
370 to accelerate computations allows for fast, unbiased genome-wide computations of heritability  
371 and GxE interactions for both small datasets and biobank-scale data. Our method also enables  
372 testing for interactions with specific environmental exposures instead of overall effects from  
373 multiple environmental outcomes. By reducing the assumptions required for computing  
374 heritability and GxE interactions, MonsterLM has the potential to uncover greater insights into  
375 the genetic architecture of GxE interactions.

376

377 Our analyses revealed the presence of significant GxE interactions for five of eight blood  
378 biomarkers with WHR. Interaction effects ranged from null to very strong, and in the cases of  
379 Apolipoprotein B, Bilirubin, and HDL-Cholesterol, explained a higher proportion of overall  
380 variance than heritability. These results have important implications for future research. First,  
381 our observations suggest that there are real interactions between genetics and exposures that  
382 contribute greatly to complex trait variance. Second, genetic associations are likely to be  
383 heterogenous when comparing populations with dramatically different obesogenic environmental  
384 exposures. The observation that GxE effects do not come from SNPs with strong marginal  
385 effects suggests this may not impact top GWAS hits excessively. We also observed the presence  
386 of significant directionality effects for strongly significant SNPs and their associated interaction  
387 effects, which suggest an overall greater impact of genetic variation under certain environmental  
388 conditions. There are also clinical implications for these observations. For instance,  
389 Apolipoprotein B is a *bona fide* risk factor for coronary artery disease (CAD)<sup>31</sup>. A strong  
390 interaction effect with WHR is observed, suggesting WHR is also an important modulator of  
391 genetic risk of CAD mediated through Apolipoprotein B.

392

393 Our results also provide some insights into why identification of GxE interactions has been  
394 challenging<sup>1</sup>. Many prior studies have reasonably focused the search for significant GxE  
395 interactions on variants with genome-wide significant marginal effects. However, our results  
396 show that only a small proportion of GxE interaction effects can be explained by such variants.  
397 Rather, the majority of GxE interaction effects are due to variants with unremarkable marginal  
398 effects. On the other hand, we also show that a relatively small minority of variants is  
399 responsible for a disproportionate contribution to GxE interactions. Altogether, these findings  
400 offer hope that the identification of specific interactions is possible. Indeed, we also show in a  
401 proof-of-concept experiment that incorporation of GxE interactions can significantly improve PS  
402 prediction, albeit modestly.

403

404 Some limitations are worth mentioning. First, we quantile normalized all traits before analysis,  
405 and while this protects against potential scaling effects, it could also bias results towards the null.  
406 Second, MonsterLM is not meant to identify specific GxE interactions but rather to quantify the  
407 overall, genome-wide contributions of GxE interactions to continuous traits. Another limitation  
408 includes the potential loss of information from LD pruning to account for high correlation in the  
409 genotype data and from filtering rare variants (MAF<1%).

410

411 In this report, we have established the presence of GxE interactions in cardiometabolic  
412 biomarkers. We observed that SNPs with strong marginal effects contribute weakly to the  
413 variance of GxE interaction effects, and that there is a disproportionate contribution from a  
414 relatively small minority of variants. Our results also highlight the potential for pathway

415 analysis, examining specific genes involved in GxE interactions. MonsterLM provides flexibility  
416 for any form of genetic architecture, environmental exposures and interaction models, and serves  
417 as the basis for more advanced future analyses into the specifics of genome-wide environmental  
418 interactions and importantly, the contribution of GxE interactions to dichotomous traits such as  
419 disease status.

420

421

## 422 Acknowledgements

423 The authors are thankful for all the UK Biobank participants.

## 424 Author Contributions

425 **MK**: data curation, software, formal analysis, investigation, visualization, writing (original  
426 draft); **MD**: formal analysis, visualization, writing (review and editing); **CJ**: formal analysis,  
427 visualization, writing (review and editing); **NP**: formal analysis; **MC**: data curation, analysis  
428 interpretation, writing (review and editing); **SM**: data curation, software; **SD**: software; **WN**:  
429 software **JP**: software; **GP**: conceptualization, supervision, funding acquisition, methodology,  
430 project administration, writing (review and edit).

## 431 Competing Interests statement

432 None of the authors report competing interests.

## 433 Author ORCIDs

434 Mohammad Khan: <https://orcid.org/0000-0001-5076-279X>

435 Matteo Di Scipio: <https://orcid.org/0000-0001-6280-3739>

436 Conor Judge: <https://orcid.org/0000-0001-9473-2920>

437 Nicolas Perrot: <https://orcid.org/0000-0002-2395-2333>

438 Michael Chong: <https://orcid.org/0000-0002-0555-4622>

439 Shihong Mao: <https://orcid.org/0000-0002-0881-2412>

440 Shuang Di: <https://orcid.org/0000-0003-1707-7613>

441 Guillaume Paré: <https://orcid.org/0000-0002-6795-4760>

442

443 **Code Availability Statement**

444 All custom code is available upon request.

## References

1. Aschard, H. A perspective on interaction effects in genetic association studies. *Genet. Epidemiol.* **40**, 678–688 (2016).
2. Dempfle, A. *et al.* Gene-environment interactions for complex traits: definitions, methodological requirements and challenges. *Eur. J. Hum. Genet. EJHG* **16**, 1164–1172 (2008).
3. Castaldi, P. J. *et al.* Screening for interaction effects in gene expression data. *PloS One* **12**, e0173847 (2017).
4. Kim, J. *et al.* Joint Analysis of Multiple Interaction Parameters in Genetic Association Studies. *Genetics* **211**, 483–494 (2019).
5. Dai, J. Y. *et al.* Simultaneously testing for marginal genetic association and gene-environment interaction. *Am. J. Epidemiol.* **176**, 164–173 (2012).
6. Patel, C. J., Chen, R., Kodama, K., Ioannidis, J. P. A. & Butte, A. J. Systematic identification of interaction effects between genome- and environment-wide associations in type 2 diabetes mellitus. *Hum. Genet.* **132**, 495–508 (2013).
7. Almasy, L. & Blangero, J. Variance component methods for analysis of complex phenotypes. *Cold Spring Harb. Protoc.* **2010**, pdb.top77 (2010).
8. Veerman, J. R., Leday, G. G. R. & Wiel, M. A. van de. Estimation of variance components, heritability and the ridge penalty in high-dimensional generalized linear models. *Commun. Stat. - Simul. Comput.* **0**, 1–19 (2019).
9. Speed, D., Hemani, G., Johnson, M. R. & Balding, D. J. Improved Heritability Estimation from Genome-wide SNPs. *Am. J. Hum. Genet.* **91**, 1011–1021 (2012).



10. Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.* **49**, 986–992 (2017).
11. Speed, D. & Balding, D. J. SumHer better estimates the SNP heritability of complex traits from summary statistics. *Nat. Genet.* **51**, 277–284 (2019).
12. Evans, L. M. *et al.* Comparison of methods that use whole genome data to estimate the heritability and genetic architecture of complex traits. *Nat. Genet.* **50**, 737–745 (2018).
13. Gazal, S., Marquez-Luna, C., Finucane, H. K. & Price, A. L. *Reconciling S-LDSC and LDAK models and functional enrichment estimates.*  
<http://biorxiv.org/lookup/doi/10.1101/256412> (2018) doi:10.1101/256412.
14. Yang, J. *et al.* Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
15. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
16. Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
17. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
18. Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits shows action of negative selection. *Nat. Genet.* **49**, 1421–1427 (2017).
19. Hou, K. *et al.* Accurate estimation of SNP-heritability from biobank-scale data irrespective of genetic architecture. *Nat. Genet.* **51**, 1244–1251 (2019).

20. de Los Campos, G., Hickey, J. M., Pong-Wong, R., Daetwyler, H. D. & Calus, M. P. L. Whole-genome regression and prediction methods applied to plant and animal breeding. *Genetics* **193**, 327–345 (2013).
21. Mayhew, A. J. & Meyre, D. Assessing the Heritability of Complex Traits in Humans: Methodological Challenges and Opportunities. *Curr. Genomics* **18**, 332–340 (2017).
22. Browning, S. R. & Browning, B. L. Population structure can inflate SNP-based heritability estimates. *Am. J. Hum. Genet.* **89**, 191–193; author reply 193-195 (2011).
23. Shewchuk, J. R. *An Introduction to the Conjugate Gradient Method Without the Agonizing Pain.* (1994).
24. Nogueira, B. & Pinheiro, R. G. S. A GPU based local search algorithm for the unweighted and weighted maximum s-plex problems. *Ann. Oper. Res.* **284**, 367–400 (2020).
25. Poppitt, S. D. *et al.* Long-term effects of ad libitum low-fat, high-carbohydrate diets on body weight and serum lipids in overweight subjects with metabolic syndrome. *Am. J. Clin. Nutr.* **75**, 11–20 (2002).
26. Moore, R. *et al.* A linear mixed model approach to study multivariate gene-environment interactions. *Nat. Genet.* **51**, 180–186 (2019).
27. Robinson, M. R. *et al.* Genotype-covariate interaction effects and the heritability of adult body mass index. *Nat. Genet.* **49**, 1174–1181 (2017).
28. Dahl, A. *et al.* A Robust Method Uncovers Significant Context-Specific Heritability in Diverse Complex Traits. *Am. J. Hum. Genet.* **106**, 71–91 (2020).
29. Sulc, J. *et al.* Quantification of the overall contribution of gene-environment interaction for obesity-related traits. *Nat. Commun.* **11**, (2020).

30. Kerin, M. & Marchini, J. Inferring Gene-by-Environment Interactions with a Bayesian Whole-Genome Regression Model. *Am. J. Hum. Genet.* **107**, 698–713 (2020).
31. Sniderman, A. D. *et al.* Apolipoprotein B Particles and Cardiovascular Disease: A Narrative Review. *JAMA Cardiol.* **4**, 1287–1295 (2019).
32. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLOS Med.* **12**, e1001779 (2015).
33. De La Vega, F. M. & Bustamante, C. D. Polygenic risk scores: a biased prediction? *Genome Med.* **10**, 100 (2018).
34. Graf, R. G. & Alf, E. F. Correlations Redux: Asymptotic Confidence Limits for Partial and Squared Multiple Correlations. *Appl. Psychol. Meas.* **23**, 116–119 (1999).

## Figure Legends

### **Figure 1 | Summary of Gene-by-Environment (GxE) analysis conducted with MonsterLM.**

Initial simulation studies were conducted to verify the properties of MonsterLM; simulated phenotypes with known values for variance explained were regressed under varying SNP partitioning and interaction structure conditions to ensure robust estimations (**blue panel**). Real trait analyses were conducted with UK Biobank data (**grey panels**). Genome-wide SNP heritability estimates with and without waist-hip-ratio (WHR) interactions revealed significant interaction effects for five of eight biomarkers and were further assessed with a directionality of effects and stratification analysis (**bottom left panel**). The model was further explored by recovering genotype and interaction variance explained through partitioning SNPs based on genotype and interaction univariate regressions thus providing insights into the model's

architecture (**bottom middle panel**). Lastly, sequential incorporation of subsets of SNPs with significant  $P_{GE}$  derived from univariate interaction regressions of the genotype SNPs on their respective traits revealed modest improvements of polygenic scores ( $PS_G, PS_{GE}$ ) in four of the five biomarkers tested (**bottom right panel**).

**Figure 2 | Estimation of variance explained by GxE interactions for 100 simulated**

**phenotypes.** Estimation of variance explained by GxE interactions under three simulation scenarios. (+) indicates that the presence of a specific condition, while (-) indicates the absence of a condition. The “E dependent on G” condition denotes the case where environment effect SNPs are a subset of the same genetic effect SNPs. The “SNP ( $\beta_{G,i} \neq 0, \beta_{GE,i} \neq 0$ )” condition denotes the case where a single SNP has both non-zero genetic and interaction effects. Dashed blue lines denote the true variance set by simulations. **a**, Estimation of variance explained by GxE interactions using a single block in chromosome 18 in three scenarios. **b**, Estimation of variance explained by GxE interactions under the three multi-block simulation scenarios for chromosome 11 (3 blocks). 95% CIs were calculated for simulations as described in the methods. *P*-values were derived via *Z*-test.

**Figure 3 | Estimates of genetic, interaction, and environment (WHR)  $\overline{R^2}$  for eight**

**biomarkers and associated directionality of effects.** Studied biomarkers were residualized for age, sex, WHR and the first 20 genetic principal components. Phenotypes were quantile normalized and mean imputed as per methods. 95% CIs were calculated for each estimate as described in the online methods. **a**, Genetic, interaction, and environment (WHR) variance estimated  $\overline{R^2}$  for each biomarker using the MonsterLM protocol. **b**, The directionality of effects

for derived interaction estimates. SNPs were filtered based on univariate  $P_G$ ,  $P_{GE}$  and LD ( $r^2 < 0.1$ ) for each biomarker. Directionality is concordant when  $\hat{\beta}_G$  and  $\hat{\beta}_{GE}$  have the same sign (+/, -/-) and discordant when they have opposite signs (+/-, -/+). Two-proportion Z-tests were used to compare each directionality result with a null value of 0.5.

**Figure 4 | Proportion of  $\overline{R^2}_{G-val}$  and  $\overline{R^2}_{GWE-val}$  as a function of  $P_G$ .** **a**, The proportion of total  $\overline{R^2}_{G-val}$  recovered in the validation set at each discovery sample  $P_G$  for the five biomarkers with significant interaction variance. **b**, The proportion of total interaction  $\overline{R^2}_{GWE-val}$  recovered in the validation set at each discovery sample  $P_G$  threshold for the same biomarkers. 95% CI were derived based on the upper and lower bounds of each estimate in proportion to either total  $\overline{R^2}_{G-val}$  or  $\overline{R^2}_{GWE-val}$ .

**Figure 5 | Proportion of  $\overline{R^2}_{GWE-val}$  estimates as a function of  $P_{GE}$  thresholds.** Proportion of total  $\overline{R^2}_{GWE-val}$  recovered in the validation set at each univariate  $P_{GE}$  threshold for the five biomarkers with significant interaction variance. 95% CI were derived based on the upper and lower bounds of each estimate in proportion to total  $\overline{R^2}_{GWE-val}$ .

**Figure 6 | Polygenic score prediction  $\overline{R^2}$  with and without incorporation of interaction effects.** For each biomarker, there are 20 different conditions based on discovery sample  $P_G$  and  $P_{GE}$  thresholds. The polygenic score  $\overline{R^2}$  was estimated in the validation sample based on discovery sample  $\hat{\beta}_G$ ,  $\hat{\beta}_{GE}$  values.

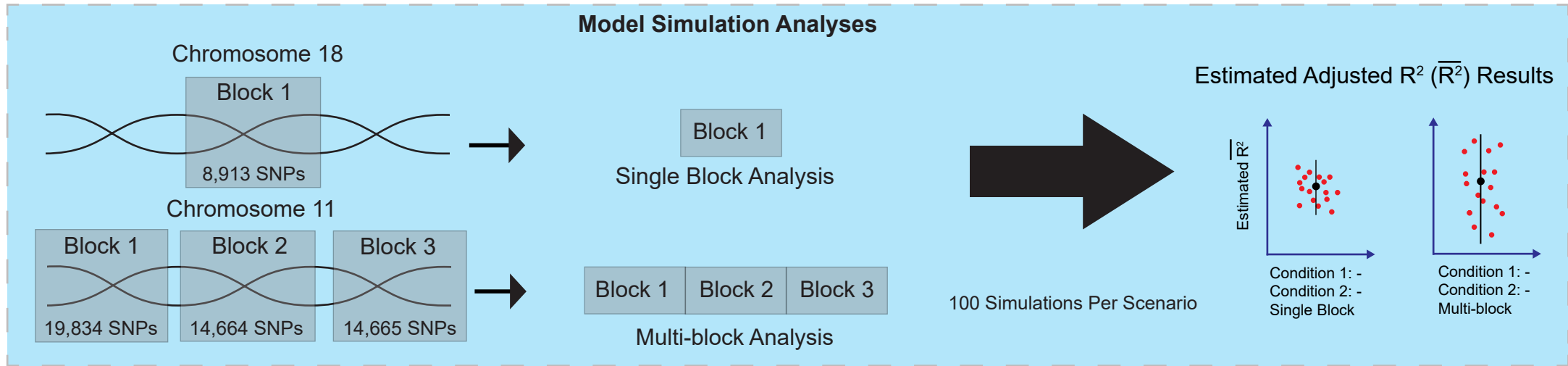
## Tables

Table 1 - Comparison of current methods estimating GxE contributions to MonsterLM

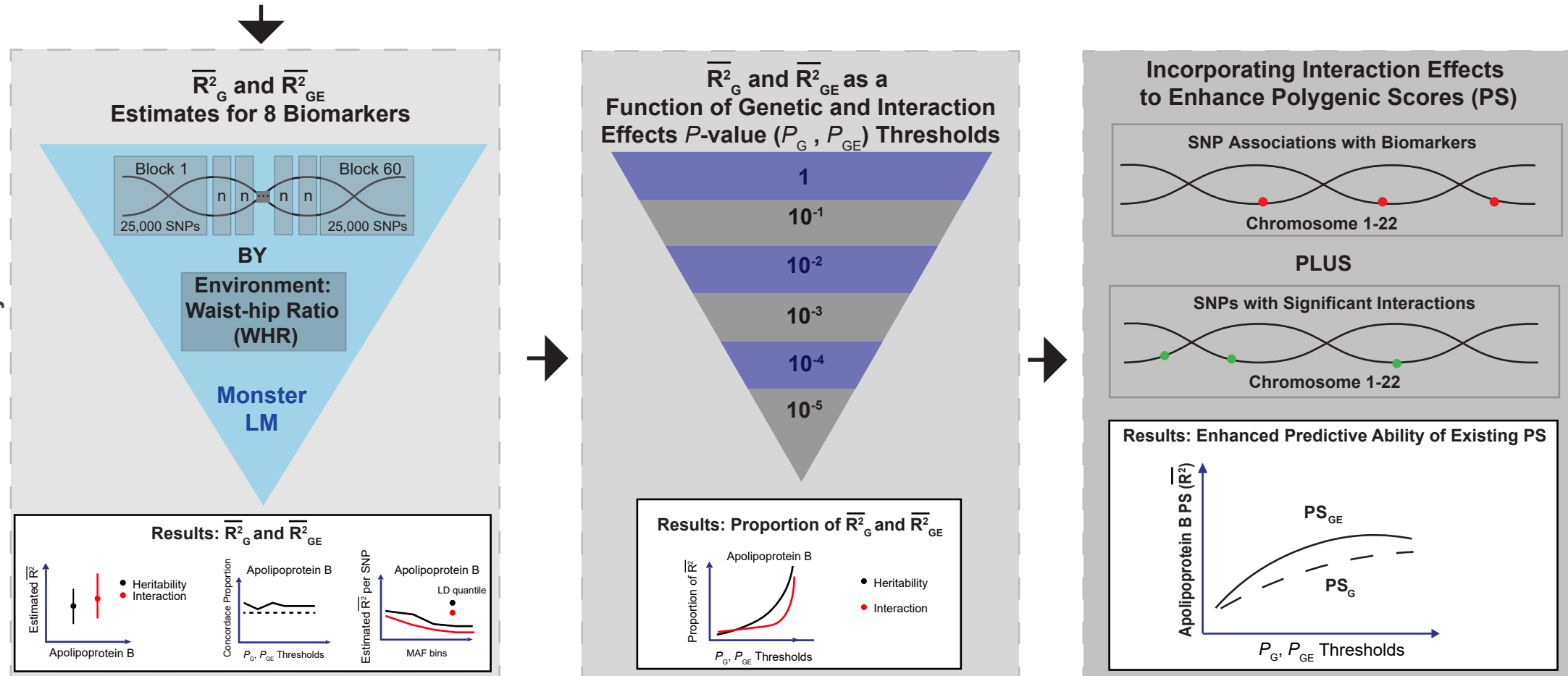
| <b>Method</b>                    | <b>Description</b>   | <b>Advantages</b>   | <b>Disadvantages</b>   | <b>MonsterLM</b>  |
|----------------------------------|--|---|--|---|
| StructLMM <sup>26</sup>          | Evaluates interaction variance for multiple environmental factors with a single SNP.                                     | Fast, robust model for a variety of different environmental exposures.  | Limited to interaction effects of only a single SNP or genotype.   | Analyzes variance explained by interactions genome-wide (after LD-pruning).   |
| CGI-GREML <sup>27</sup>          | Uses a mix of parametrized models and restricted likelihood methods to estimate variance explained by GxE.               | Well-structured for identifying GxE interactions with categorical exposures.  | >250 Likelihood Ratio Tests; Slow; Cannot use continuous traits.   | Can analyze continuous traits without categorizing them; Quick, efficient Wald-test/CI for R2 of interaction effects. |
| GxE <sup>EMM</sup> <sup>28</sup> | Linear mixed model method to detect GxE interactions across the genome and a single exposure.                            | Multiple parametrizations available to efficiently model GxE interaction effects.                                     | Small sample size only; Minimal number of SNPs.  | Can analyze a large sample size with many SNPs through partition of genotype matrix & Conjugate Gradient method.      |
| GRSxE <sup>29</sup>              | Method to detect total GxE interactions with a Gene-Risk Score.  | Estimates the GxE contribution for all possible environmental factors with SNPs.                                      | Assumes each SNP interacts equally with E.   | Accounts for the unique interaction effect of each SNP with E.  |
| LEMMA <sup>30</sup>              | Linear mixed model method to detect GxE interactions across the genome and an estimated linear combination of exposures. | Considers the impact of overlapping environmental exposures when computing total GxE contributions across the genome. | Requires parametrization and is dependent on model specification; Uses an estimated linear combination of exposures, assuming all E's interact with the same SNPs. | Tests for specific interaction with E rather than a linear combination of E.  |

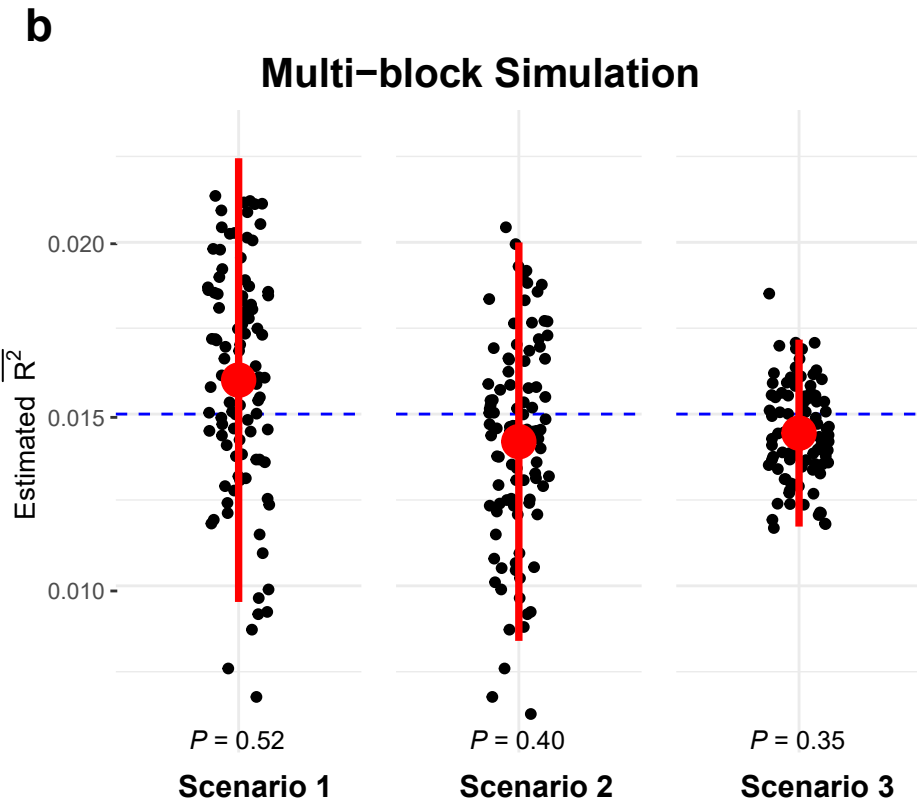
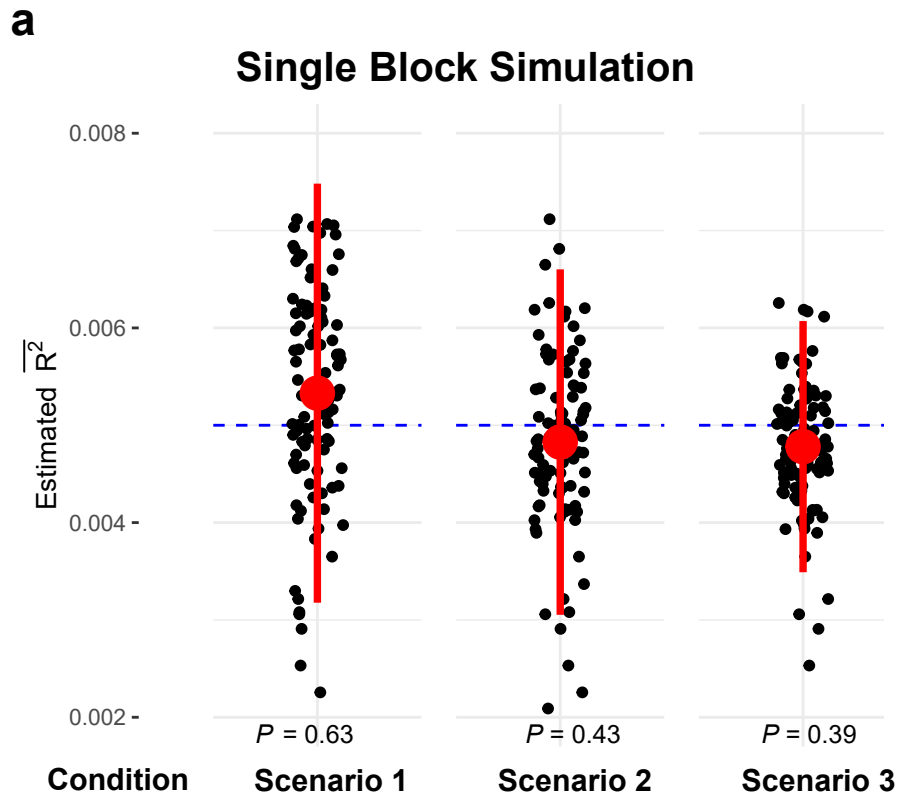
# CENTRAL ILLUSTRATION: GxE Interaction Variance Estimation (MonsterLM)

Model



Real Data Analyses





$E$  Dependent on  $G$

—

+

+

—

+

+

SNP ( $\beta_G \neq 0, \beta_{GE} \neq 0$ )

—

—

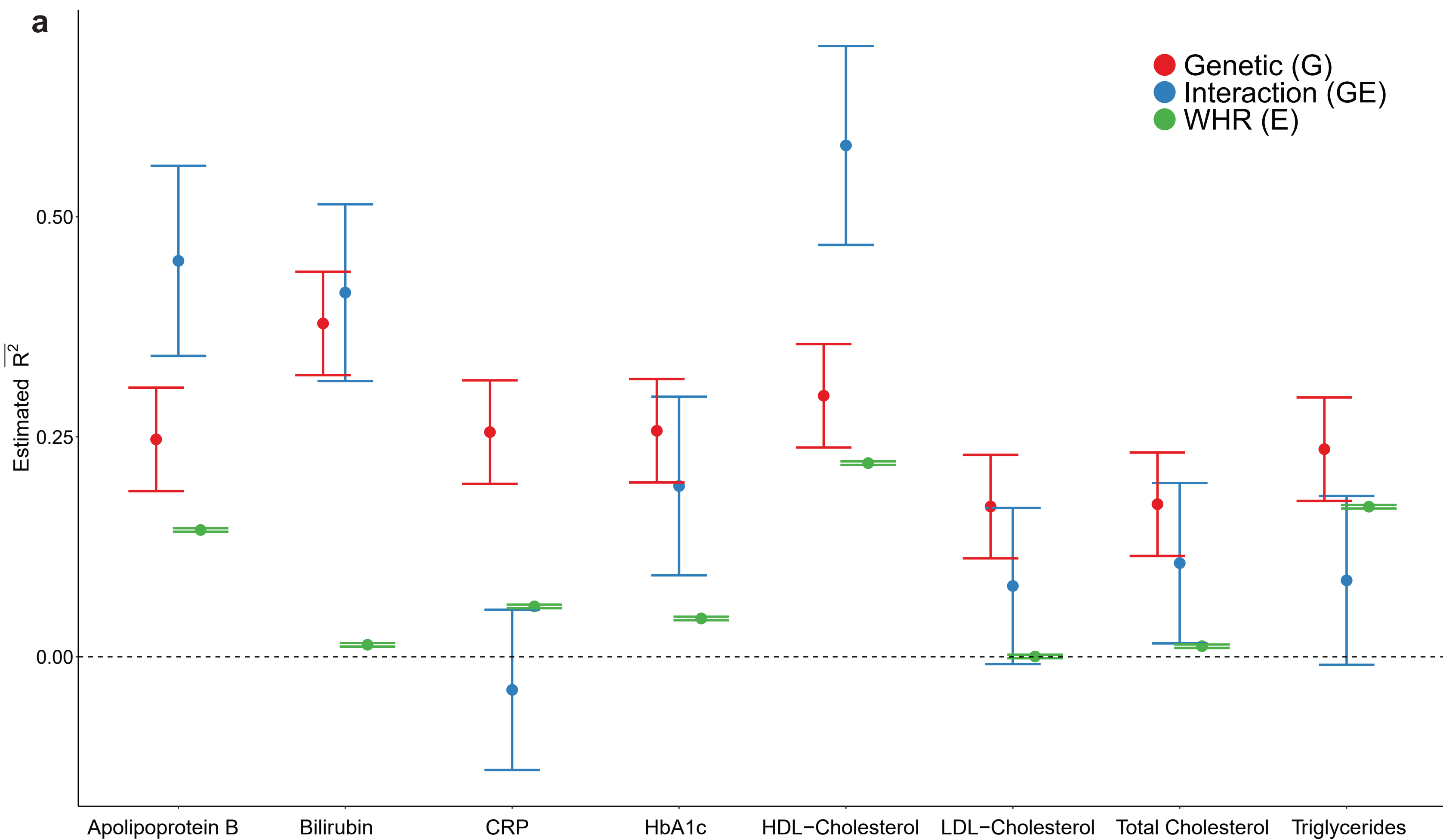
+

—

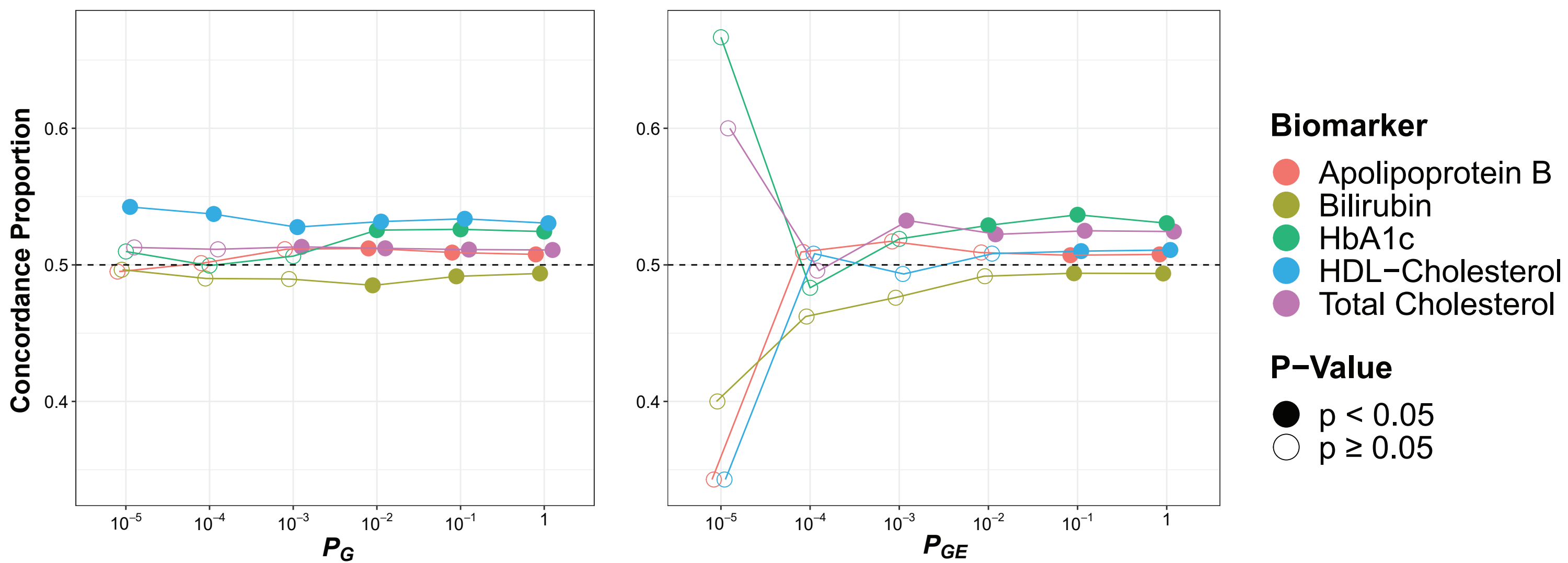
—

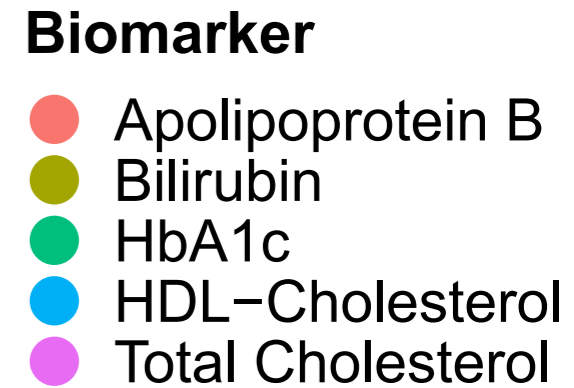
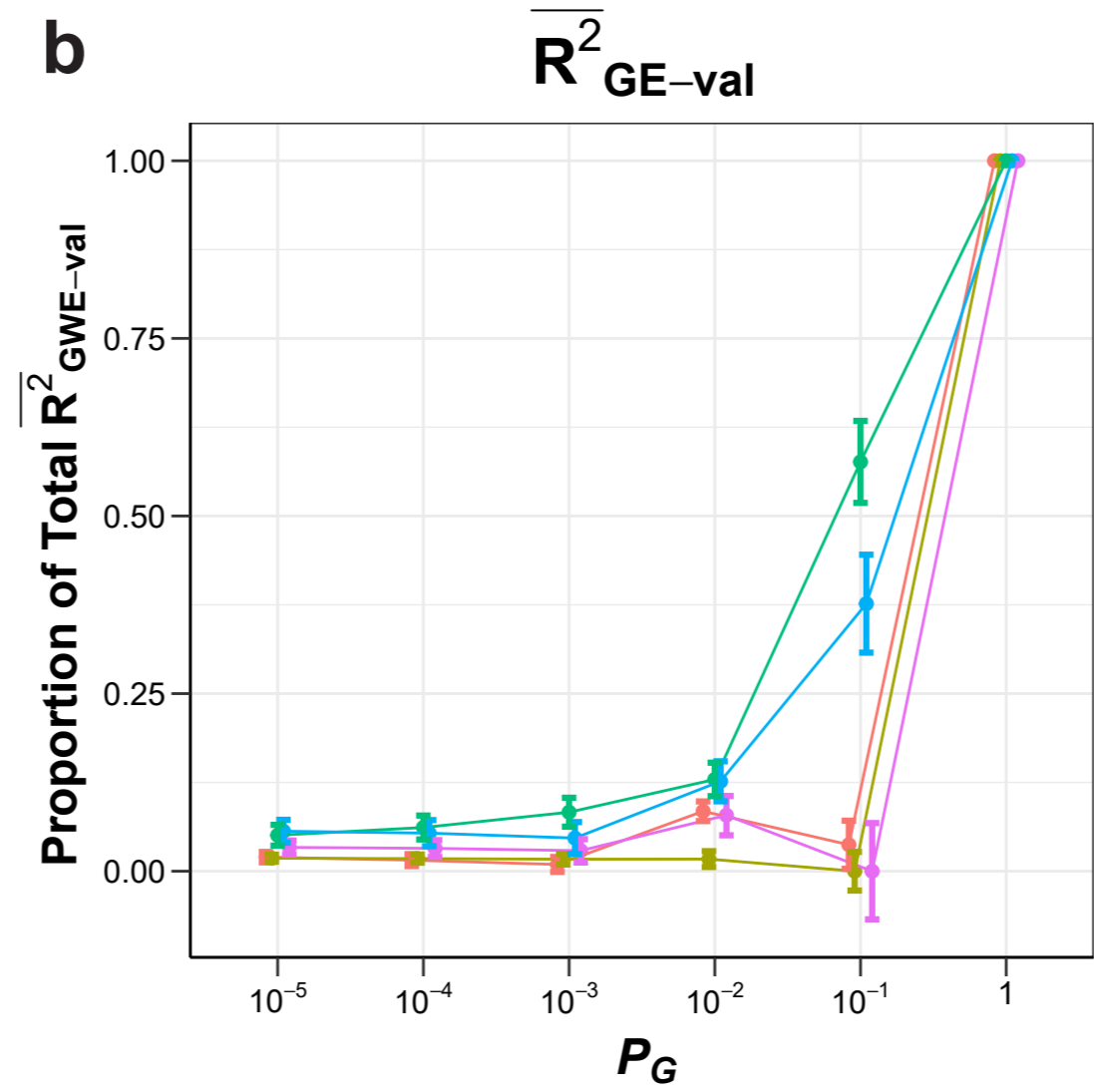
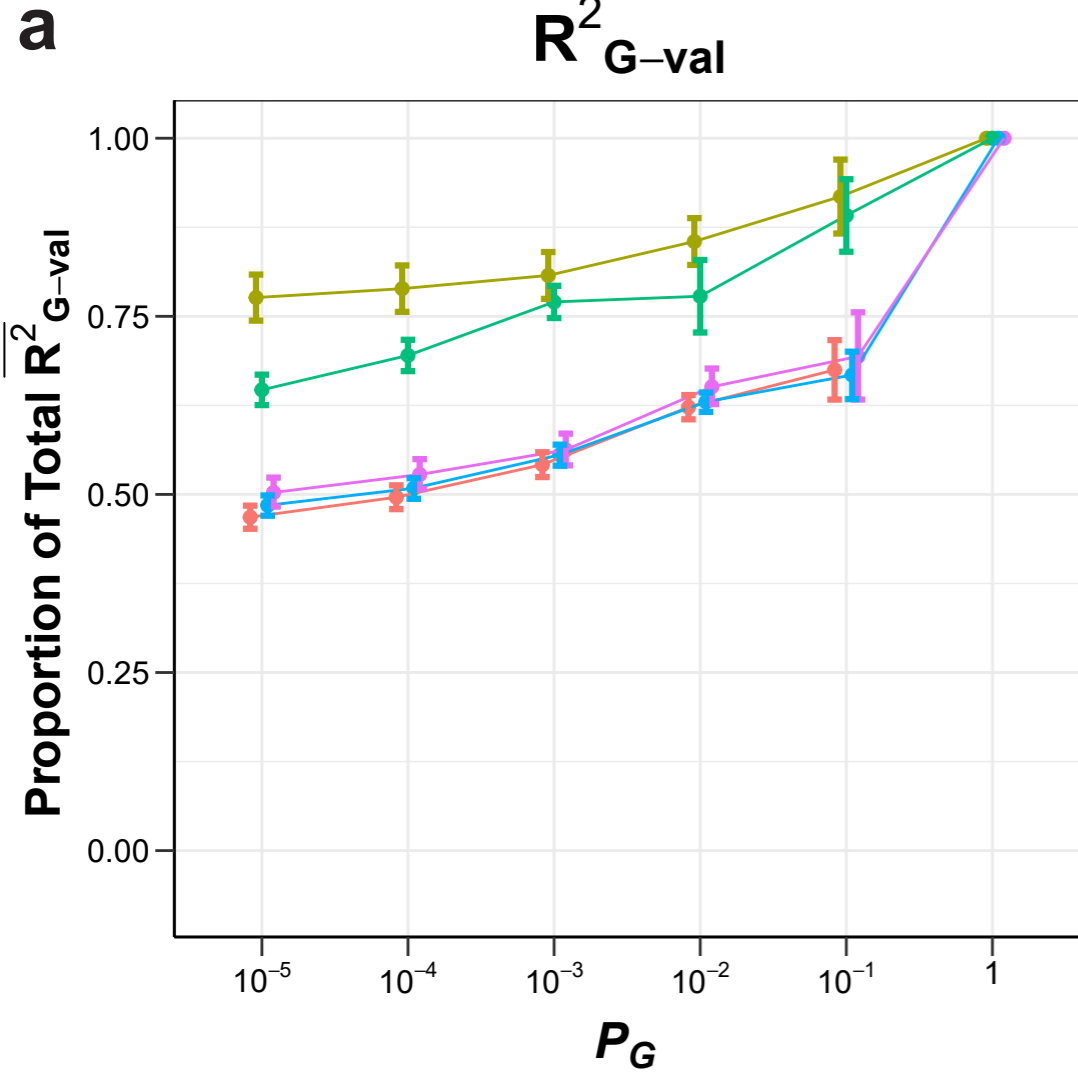
+

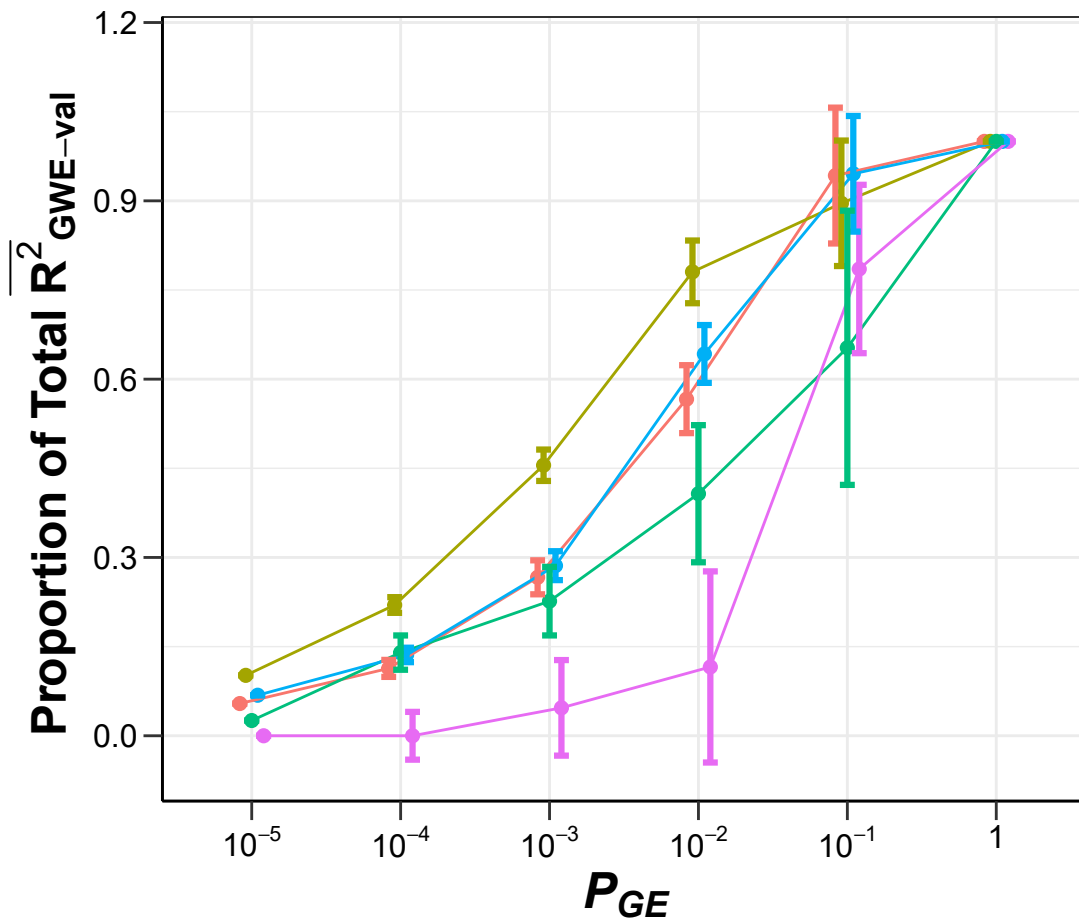




**b**      **Directionality of Effects Analysis**



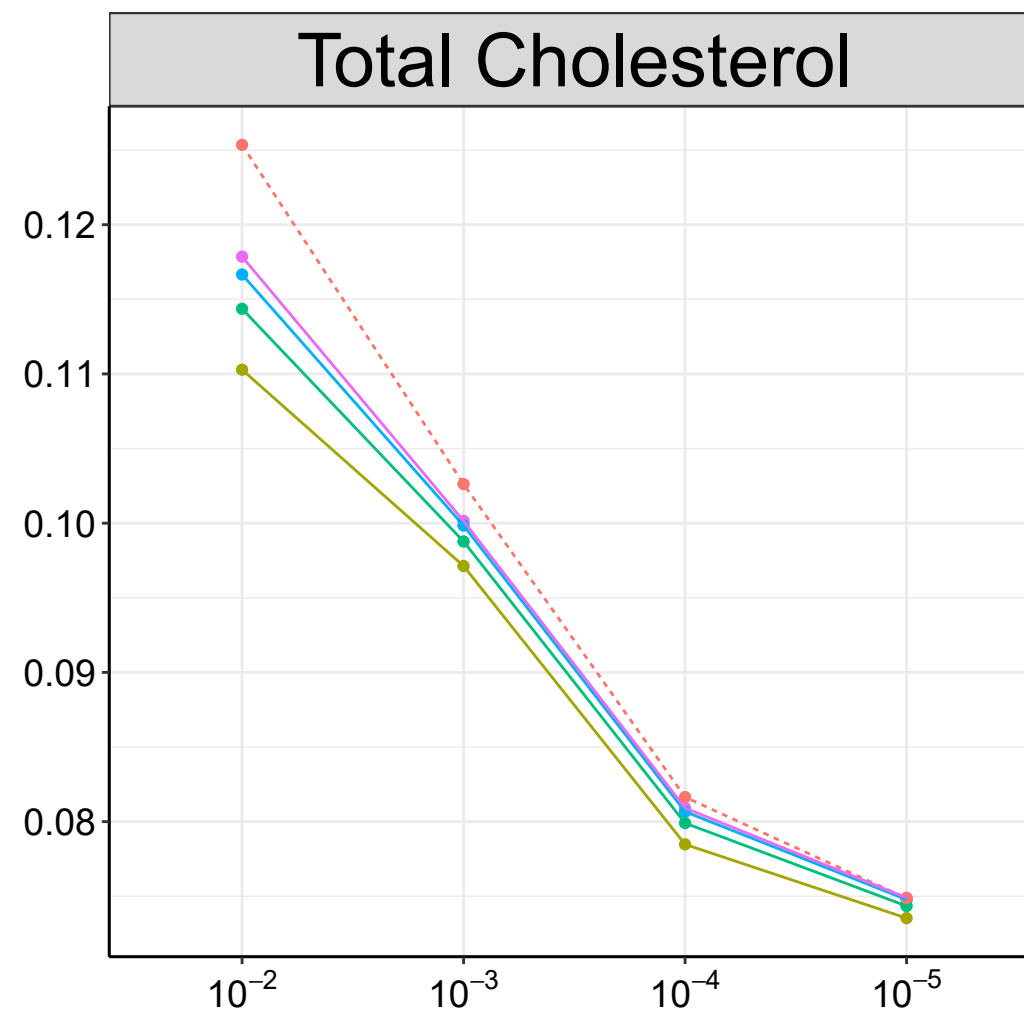
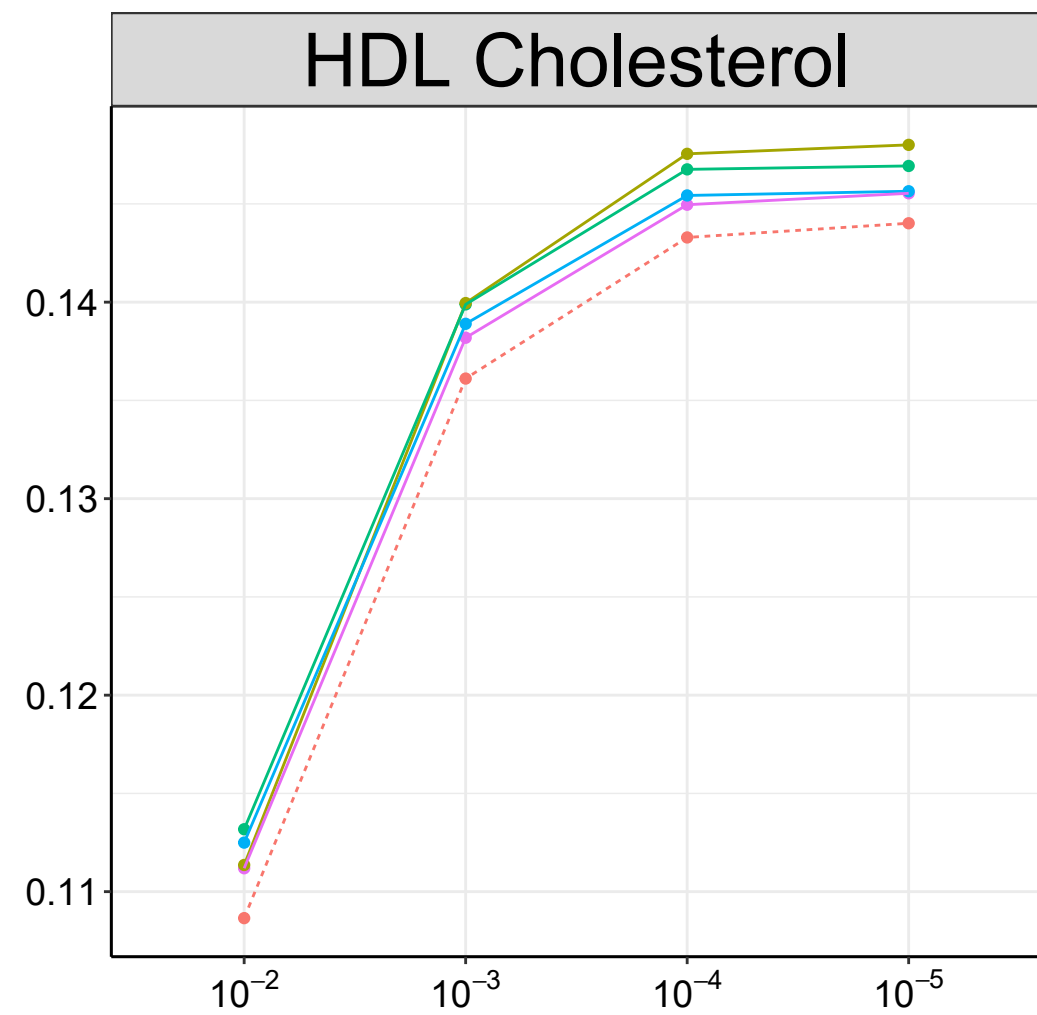
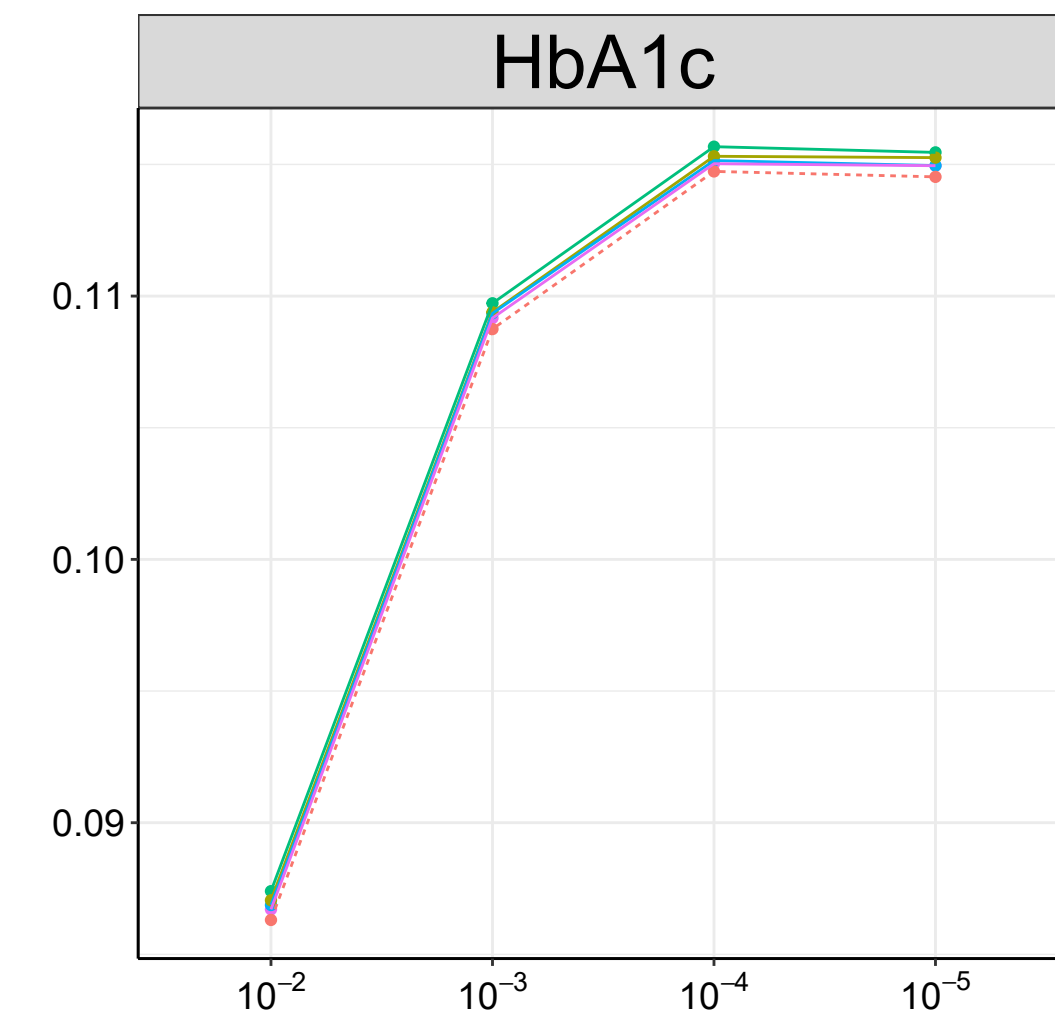
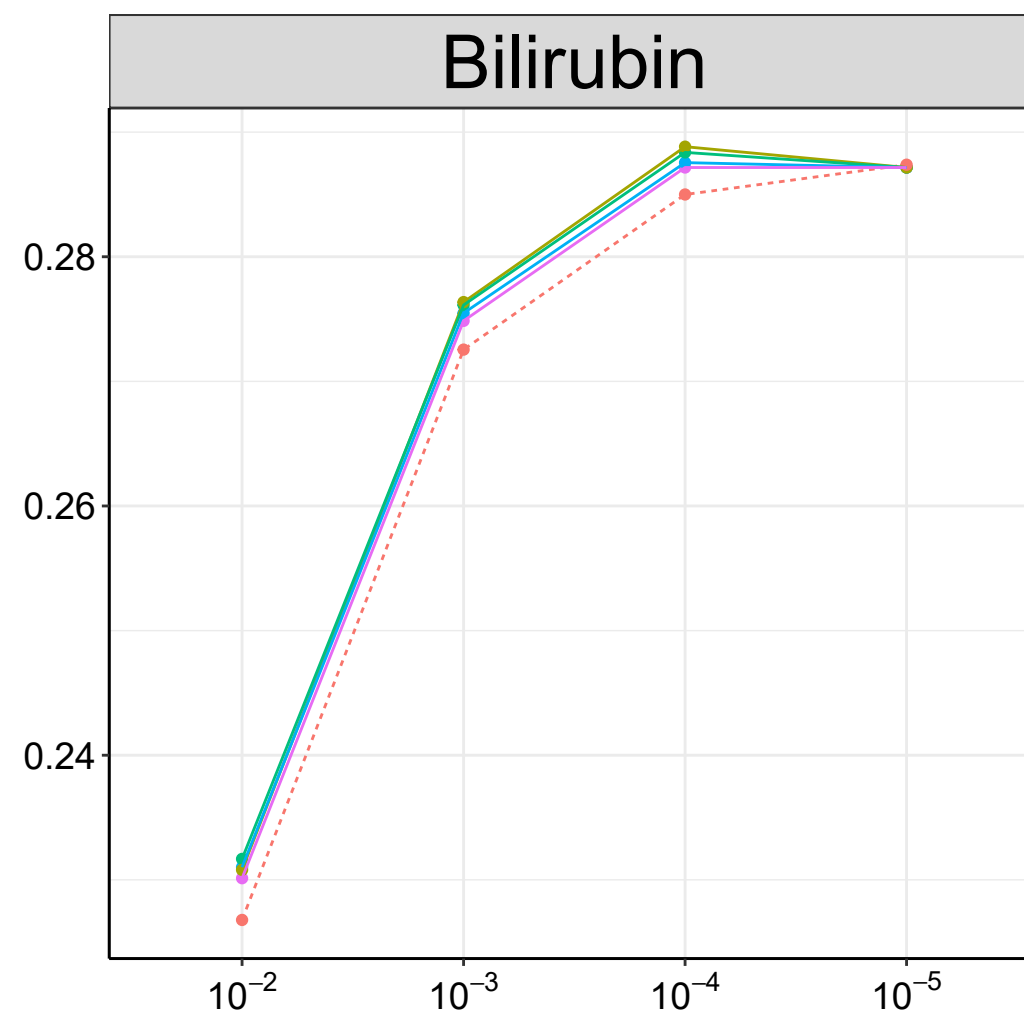
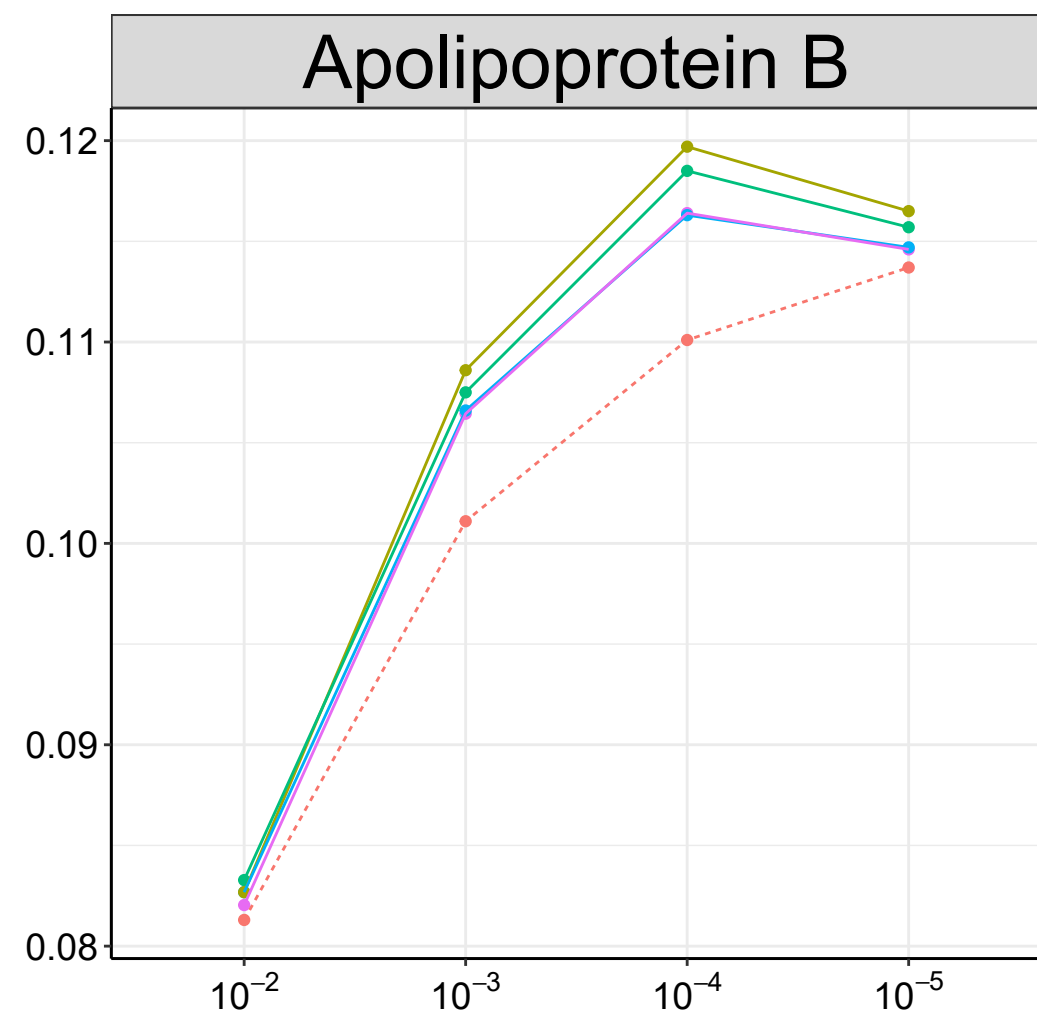


$\overline{R^2}_{GE-val}$ 

### Biomarker

- Apolipoprotein B
- Bilirubin
- HbA1c
- HDL-Cholesterol
- Total Cholesterol

Polygenic Score  $R^2_{val}$



$P_{GE}$

- No Interaction Term
- $10^{-2}$
- $10^{-3}$
- $10^{-4}$
- $10^{-5}$

$P_G$