

1 **GxEsum: a novel approach to estimate the phenotypic variance explained by**  
2 **genome-wide GxE interaction based on GWAS summary statistics for biobank-**  
3 **scale data**

4

5

6 **Jisu Shin<sup>1,2</sup>, S Hong Lee<sup>1,2\*</sup>**

7

8 <sup>1</sup>Australian Centre for Precision Health, University of South Australia Cancer  
9 Research Institute, University of South Australia, Adelaide, SA 5000, Australia

10 <sup>2</sup>UniSA: Allied Health and Human Performance, University of South Australia,  
11 Adelaide, SA 5000, Australia

12 \*Correspondence [Hong.Lee@unisa.edu.au](mailto:Hong.Lee@unisa.edu.au)

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

1 **Abstract**

2 Genetic variation in response to the environment is fundamental in the biology of  
3 complex traits and diseases, i.e. genotype-by-environment interaction (GxE).  
4 However, existing methods are computationally demanding and infeasible to handle  
5 biobank-scale data. Here we introduce GxEsum, a method for estimating the  
6 phenotypic variance explained by genome-wide GxE based on GWAS summary  
7 statistics. Through comprehensive simulations and analysis of UK Biobank with  
8 288,837 individuals, we show that GxEsum can handle a large-scale biobank dataset  
9 with controlled type I error rates and unbiased GxE estimates, and its computational  
10 efficiency can be hundreds of times higher than existing GxE methods.

11

12

13

14 **Keywords**

15 GxE interaction, whole-genome approach, Biobank-scale data, reaction norm model,  
16 LDSC

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

## 1 Background

2

3 The success of the human genome project has led to a paradigm-shift in the  
4 complex trait analysis that focuses on the genome-wide association studies (GWAS)  
5 [1]. GWAS have been incredibly successful at identifying genome-wide significant  
6 single nucleotide polymorphisms (SNPs) that are associated with causal variants  
7 underlying complex traits [2, 3]. Moreover, whole-genome approaches, using all  
8 common SNPs across the genome, have been useful to dissect the genetic  
9 architecture of complex traits, e.g. SNP-based heritability and genetic correlation [4].  
10 However, the analytical modelling used in GWAS and whole-genome approaches  
11 usually assumes that there is no genotype-environment interaction (GxE), which can  
12 be often violated against the true genetic architecture of complex traits. Indeed,  
13 interaction is fundamental in biology and there has been increasing interest in  
14 estimating GxE, using genome-wide SNPs [5-7].

15 Current state-of-the-art whole genome methods for estimating GxE include  
16 genotype-covariate interaction genomic restricted maximum likelihood (GREML) and  
17 random regression GREML [8]. Recently, a multivariate reaction norm model (RNM)  
18 has been introduced [9], which can disentangle GxE from genotype-environmental  
19 correlation, providing more reliable GxE estimations. These methods typically  
20 employ the GREML approach that requires individual level genotypes and is  
21 computationally intensive. Especially when using biobank-scale data, the approach  
22 becomes computationally intractable.

23 To reduce the computational limitation of GREML, linkage disequilibrium score  
24 regression (LDSC) was introduced to estimate SNP-based heritability and genetic  
25 correlation [10]. LDSC is computationally efficient and requires no individual-level  
26 genotypes. Instead, it uses GWAS summary statistics, regressing the association  
27 test statistics of SNPs on their LD score. However, existing LDSC methods are  
28 limited to additive models only [11-14].

29 In this study, we propose a novel approach to estimate the phenotypic variance  
30 explained by genome-wide GxE based on GWAS summary statistics (GxEsum) for a  
31 large-scale biobank dataset, correctly accounting for genotype-environment  
32 correlation and scale effects. In simulated and real data analyses, we show that the  
33 computational efficiency of the proposed approach is substantially higher than RNM,  
34 an existing GREML-based method, while the estimates are reasonably accurate and  
35 precise. Because of this computational advantage, GxEsum may be an efficient tool  
36 to estimate GxE that can be applied to large-scale data across multiple complex  
37 traits.

38

39

40

41

## 1 Results

### 2 Method Overview

3 We propose a method to estimate the phenotypic variance explained by the whole-  
4 genome GxE, based on GWAS summary statistics, referred to as GxEsum. GxEsum  
5 can be a computationally efficient RNM using an extension of the LDSC approach.  
6 While the existing LDSC approach is designed to use estimated additive SNP effects  
7 in GWAS summary statistics (Supplementary Note 1), GxEsum requires summary  
8 statistics of SNP-by-environment interaction effects. For SNP effects modulated by  
9 an environment, the expected chi-square statistic ( $\chi_j^2$ ) is

$$E[\chi_j^2 | \ell_j] = \frac{N\sigma_{g_1}^2}{M} * \ell_j + 1 + 2(\sigma_{g_1}^2 + \sigma_{\tau_1}^2)$$

10 where N is the number of individuals, M is the number of SNPs,  $\sigma_{g_1}^2$  is the variance  
11 due to GxE,  $\sigma_{\tau_1}^2$  is the variance due to residual heterogeneity or scale effects caused  
12 by residual-environment interaction (RxE) and  $\ell_j$  is the LD score at the variant  $j$  that  
13 can be estimated from a reference panel (please see Methods for a full derivation of  
14 this equation). The  $\chi_j^2$  test statistics corresponds to the regression coefficient for the  
15 interaction between the  $j$ th SNP and the environmental covariate (E). The outcome  
16 trait is pre-adjusted for confounders and the main effects of E and then a regression  
17 model with the main and interaction effects is run by SNP-by-SNP. If chi-square  
18 statistics from GWAS are regressed on LD scores, non-genetic interaction effects  
19 ( $\sigma_{\tau_1}^2$ ) are captured by the intercept, from which GxE ( $\sigma_{g_1}^2$ ) can be disentangled.  
20 Consequently, GxE effects estimated by GxEsum are equivalent to that adjusted for  
21 RxE when using RNM [9].

22 To validate the proposed model, i.e. GxEsum, we used various simulations based on  
23 real genotype data (see Supplementary Note 2 for a full description of the simulation  
24 models). In simulations with and without GxE, we assessed the type I error rates and  
25 the accuracy of estimated GxE. We deliberately generated confounding effects such  
26 as genotype-environment (G-E) correlation, RxE and residual-environment (R-E)  
27 correlation to see if the type I error rate and the accuracy of GxEsum were affected  
28 by these confounding factors.

29 In the real data analysis, we used the UK Biobank data with 288,837 unrelated  
30 individuals after stringent quality control. Subsets of the data with various sample  
31 sizes were analysed to compare the precision (i.e. power) and the computational  
32 efficiency of GxEsum and GREML-based GxE model (i.e. RNM).

33 Finally, we show how the genetic effects of a complex trait (e.g. BMI, hypertension or  
34 type 2 diabetes) are modulated by environment (e.g. neuroticism score, alcohol  
35 intake frequency, physical activity or age) by using the proposed method.

36

37

38

## 1 Simulations

2 For a continuous trait, under the null (no GxE), whether or not there were  
 3 confounding effects (RxE and G-E and R-E correlations), the type I error rate of  
 4 GxEsum was not significantly inflated (Table 1). Note that the use of 500 replicates  
 5 for each simulation scenario can detect a type I error of greater than 0.07 or less  
 6 than 0.03 as significantly different from 0.05, using the binomial distribution theory  
 7 [15, 16]. Even with larger confounding effects (Supplementary Table 1), there was no  
 8 inflation for the type I error rate of GxEsum.

9 **Table 1. Type I error rates of GxEsum to detect GxE at a significance threshold of p-**  
 10 **value < 0.05.**

Scenarios	Type I Error rate
Var(GxE <sup>a</sup> ) = 0, var(RxE <sup>b</sup> ) = 0	0.066
Var(GxE) = 0, var(RxE) = 0, G-E correlation <sup>c</sup> = 0.1	0.064
Var(GxE) = 0, var(RxE) = 0, R-E correlation <sup>d</sup> = 0.1	0.044
Var(GxE) = 0, var(RxE) = 0, G-E correlation=0.1, R-E correlation=0.1	0.056
Var(GxE <sup>a</sup> ) = 0, var(RxE <sup>b</sup> ) = 0.1	0.044
Var(GxE) = 0, var(RxE) = 0.1, G-E correlation <sup>c</sup> = 0.1	0.034
Var(GxE) = 0, var(RxE) = 0.1, R-E correlation <sup>d</sup> = 0.1	0.028
Var(GxE) = 0, var(RxE) = 0.1, G-E correlation=0.1, R-E correlation=0.1	0.054
<b>Average</b>	<b>0.049</b>

11

12 <sup>a</sup>GxE: Genotype-Environment interaction, <sup>b</sup>RxE: Residual-Environment interaction, <sup>c</sup>G-E  
 13 correlation: Genotype-Environment correlation, <sup>d</sup>R-E correlation: Residual-Environment  
 14 correlation. We simulated phenotypic data based on a real genotypic dataset (ARIC GWAS)  
 15 including 7,263 participants with 583,085 SNPs, using various scenarios. The phenotypes  
 16 were standardised such that the phenotypic mean was 0 and the phenotypic variance was 1.  
 17 Type I error rate (i.e. false-positive) was estimated from 500 replicates for each scenario.

18

19 In simulation with non-zero interactions, estimated GxE (g1) was not remarkably  
 20 different from the true values whether there were significant G-E and R-E  
 21 correlations or not (see supplementary Figure 1). It was noted that RxE component  
 22 was correctly captured by the intercept and not confounded with GxE estimates even  
 23 when using non-normal environmental variables (Supplementary Note 3 and  
 24 Supplementary Tables 2 and 3). In the absence of RxE, estimated GxE was also  
 25 unbiased (Supplementary Figure 2). The estimated GxE seemed robust to different  
 26 values of G-E and R-E correlations ranging from 0.05 to 0.2, respectively  
 27 (Supplementary Figures 3 and 4).

28 On the other hand, estimated main genetic variance (g0) was slightly biased  
 29 especially when using a large G-E or R-E correlation (Supplementary Figure 4). This  
 30 is probably because of the fact that the main genetic effects are over-adjusted for the

1 environment due to the large correlations (between the trait and environment) in the  
 2 model.

3

4 **Table 2. Type I error rates of GxEsum when using binary disease traits with various**  
 5 **population prevalence.**

Scenarios	Population prevalence (k)	Type I error rate
<b>Var(GxE) = 0, Var(RxE) = 0</b>	0.025	0.052
	0.05	0.042
	0.1	0.076
	0.5	0.054
<b>Var(GxE) = 0, Var(RxE) = 0.1</b> (on the liability scale)	0.025	0.044
	0.05	0.036
	0.1	0.050
	0.5	0.052
<b>Average</b>		<b>0.050</b>

6

7 We simulated quantitative phenotypic data based on a real genotypic dataset (ARIC GWAS)  
 8 including 7,263 individuals with 583,085 SNPs. The phenotypes were standardised such that  
 9 the mean was 0 and variance was 1, for which we applied the liability threshold model to  
 10 generate affected or unaffected disease status for each individual, using various values for  
 11 the population prevalence ( $k = 0.025, 0.05, 0.1$  or  $0.5$ ). Type I error rate at a significance  
 12 threshold of  $p\text{-value} < 0.05$  was estimated from 500 replicates for each scenario and  
 13 population prevalence.

14

15 We also validated that there was no inflation for the type I error rate when applying  
 16 GxEsum to binary (disease) traits (Table 2 and Supplementary Table 4), showing  
 17 that GxEsum appears to be robust to false positives in the scenarios of various  
 18 confounders. In addition, we estimated the variance component of GxE on the  
 19 observed scale and transformed to that on the liability scale, using Robertson  
 20 transformation [17]. As shown in Supplementary Figure 5, the transformed estimates  
 21 were close to the true simulated values on the liability scale, although the precision  
 22 of estimates (represented as 95% CI) was shown to be decreased when the  
 23 population prevalence approached an extreme (e.g.  $k=0.025$ ). GxE estimates were  
 24 biased when simulating a large effect size of GxE (e.g. 10% of phenotypic variance  
 25 explained by GxE) in the case of  $k=0.025$  (Supplementary Figure 6) although they  
 26 were mostly unbiased in the case of  $k=0.1$  (Supplementary Figure 7). The level of  
 27 biasedness appeared to be increased when there were RxE effects (Supplementary  
 28 Figure 6). Finally, caution should be given in interpreting GxE estimates when there  
 29 are large confounding effects such as substantial G-E and R-E correlations  
 30 (Supplementary Figure 8). The inflated GxE estimates were probably due to the fact  
 31 that the phenotypes were over-adjusted for the environment in the model because of  
 32 the correlation between the main trait and environment (G-E and R-E correlations).  
 33 This resulted in a reduced phenotypic variance (Supplementary Figure 9), hence an

1 inflated GxE estimates that is the ratio of the estimated GxE variance to the  
2 phenotypic variance.

3 Nevertheless, those confounders including RxE interaction and G-E/R-E correlations  
4 would not produce false positives whether using continuous quantitative or binary  
5 responses as shown above (also see Supplementary Tables 1 – 4). We additionally  
6 tested if the type I error rate of GxEsum was controlled when there is collider bias,  
7 which is a concern especially when using a self-report study (e.g. UK Biobank data)  
8 [18]. In simulations with collider bias, although estimated SNP-heritability was  
9 substantially (and unrealistically) underestimated (Supplementary Figures 10 and 11),  
10 the type I error rate of GxEsum was well controlled whether using continuous or  
11 binary responses (Supplementary Tables 5 and 6).

12 The estimated variance of the main genetic effects was mostly unbiased when using  
13 binary disease traits without G-E/R-E correlations (Supplementary Figure 12). When  
14 there were significant G-E and /or R-E correlations, the estimated variance of the  
15 main genetic effects appeared to be underestimated especially when there was RxE  
16 interaction (Supplementary Figure 13), which confirmed the fact that the main  
17 genetic effects are over-adjusted for the environment due to correlations between  
18 the trait and environment (Supplementary Figure 9).

19

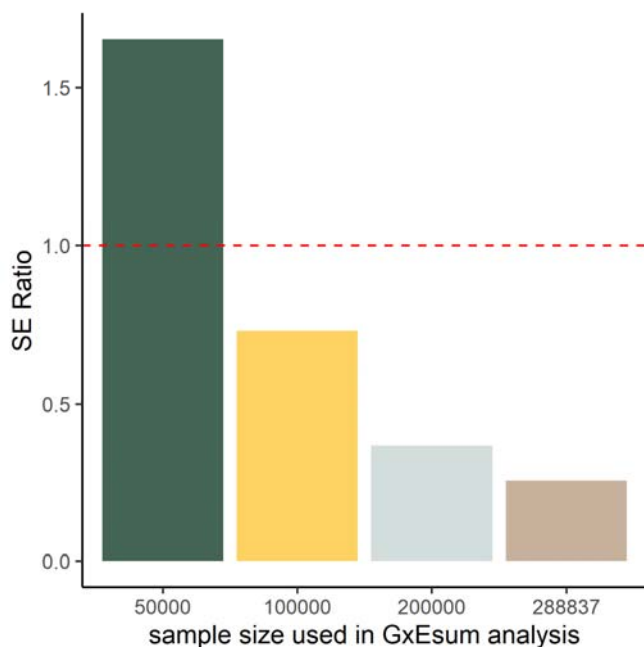
## 20 **Precision and Computational Efficiency**

21 The precision was assessed by comparing the standard error (SE) of GxEsum and  
22 RNM estimates. The SE of GxEsum was obtained from LDSC software (using a  
23 jackknife method). The SE of RNM for GxE component can be obtained from the  
24 information matrix [19] or from well-established theory [20] (see Supplementary  
25 Table 7). Figure 1 shows that the SE of GxEsum was 1.65 times higher than that of  
26 RNM when using the same sample size of 50,000. However, when sample size  
27 increased for GxEsum up to 288,837, for which RNM estimation is infeasible, the  
28 ratio reduced to 0.2. GxEsum can use a larger sample size (e.g. > 1,000,000), for  
29 which the ratio is expected to be further decreased, although the largest sample size  
30 tested in this study was 288,837 (Supplementary Figure 14).

31 While the precision of GxEsum is competitive with that of RNM, the computational  
32 efficiency is dramatically different between two methods (Figure 2 and  
33 supplementary Table 8). For example, when using a sample size of 50,000, the  
34 computing time for RNM was taken more than a thousand times than GxEsum. Even  
35 for GxEsum with a sample size of 288,837, its computational efficiency was still  
36 substantially higher than RNM with a sample size of 50,000 (Supplementary Figure  
37 14 and Supplementary Table 7). This justifies that GxEsum is a computationally  
38 efficient tool that can be applied to biobank scale data for multiple complex traits and  
39 diseases. It is noted that we assumed that preliminary analyses for each method  
40 were already done (e.g. GRM for RNM, and LD scores and GWAS for GxEsum)  
41 (Supplementary Table 8).

- 1 GxEsum uses the Wald test to get a p-value for the null hypothesis, i.e. absence of
- 2 GxE interaction, using an estimated GxE variance and its standard error. Therefore,
- 3 the power of the method is closely related to the precision.



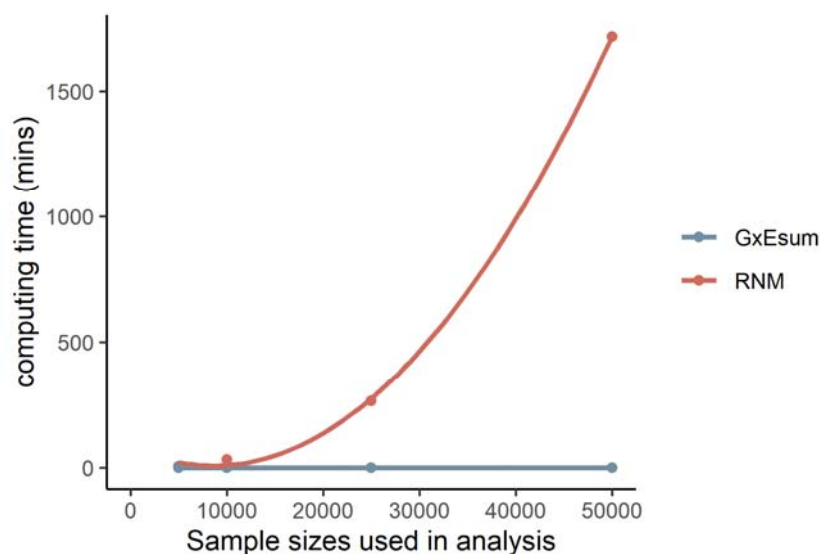


1

2 **Figure 1. The ratio of standard error (SE) from GxEsum to that from RNM using**  
3 **UK Biobank data.** The SEs of GxE variance estimated from GxEsum with various  
4 sample sizes ranging from 50,000 to 288,837 were obtained, and they were  
5 compared to the SE of GxE variance estimated from RNM with a sample size of  
6 50,000. The dashed horizontal line represents the ratio as 1.

7

8



9

10 **Figure 2. Computing time with various sample sizes used in GxEsum and RNM**  
11 **analyses.** As the sample size increases, the computing time of RNM (red) increases  
12 exponentially, while that of GxEsum (blue) is almost invariant (less than a minute).

## 1 **Real data analyses**

2 We applied GxEsum to estimate genetic effects of body mass index (BMI) that were  
3 modulated by an environment such as age, alcohol intake frequency, neuroticism  
4 scores or physical activity. The significant GxE was observed from the analyses  
5 using neuroticism scores. On the other hand, we did not find significant GxE when  
6 using age, alcohol intake frequency and physical activity after Bonferroni correction  
7 (Bonferroni p-value =  $0.05/10 = 0.005$  since there were 10 significance tests in this  
8 study) (Table 3). The GxEsum approach applied to a binary disease was conducted  
9 using hypertension or type 2 diabetes as the main trait, and BMI, waist-hip ratio  
10 (WHR), body fat percentage (BFP), or systolic/diastolic blood pressures (BP) as an  
11 environmental variable. Table 3 shows that the genetic effects of hypertension and  
12 type 2 diabetes were significantly modulated by BMI, but not by other environmental  
13 variables. For a comparison, LDSC estimates (i.e. from the null model without GxE)  
14 are also shown in Supplementary Table 9.

15 Because not all variables were without missing, we imputed missing phenotypes  
16 using the mean value for each variable in the analyses, in order to maximise the  
17 sample size. In this real data analysis, there was no remarkable difference in the  
18 results whether using phenotypic imputation or not although some variables  
19 improved their significance, e.g. NEU (Supplementary Tables 10 -12). It is noted that  
20 our main phenotypes had a small proportion of missing values, i.e. 0.3%, 7.9% and  
21 0.2% for BMI, hypertension and type 2 diabetes, respectively. If missing rate is  
22 substantially high, we recommend to exclude missing values from the analysis, or a  
23 better phenotypic imputation method [21] should be used.

24

25

26

27

28

29

1 **Table 3. Estimates obtained from GxEsum analysis using real data.**

Main trait	Environmental variable	Main additive genetic variance ( $\sigma_{g_0}^2$ )	GxE interaction variance ( $\sigma_{g_1}^2$ )	p-value for GxE
<b>BMI</b>	Age	0.216 (0.007)	0.004 (0.002)	1.86E-02
	NEU <sup>a</sup>	0.216 (0.007)	0.007 (0.002)	1.61E-05
	PA <sup>b</sup>	0.218 (0.007)	0.003 (0.001)	2.57E-02
	ALC <sup>c</sup>	0.216 (0.007)	0.003 (0.002)	5.98E-02
<b>Hypertension</b>	BMI	0.152 (0.008)	0.006 (0.002)	2.09E-03
	WHR <sup>d</sup>	0.154 (0.008)	0.005 (0.002)	3.21E-02
	BFP <sup>e</sup>	0.151 (0.008)	0.008 (0.003)	2.66E-02
<b>Type 2 Diabetes</b>	BMI	0.141 (0.014)	0.085 (0.022)	1.58E-04
	Diastolic BP <sup>f</sup>	0.198 (0.014)	-0.004 (0.006)	5.38E-01
	Systolic BP	0.204 (0.014)	-0.006 (0.006)	3.17E-01

2

3 We used a quantitative trait (BMI) and binary disease traits (hypertension and type 2 diabetes) because BMI is known to be  
 4 modulated by age/lifestyle such as NEU, ALC, PA [8, 22, 23], and hypertension and type 2 diabetes are known to be caused by  
 5 obese traits such as BMI and WHR [24, 25]. The p-value is from a Wald test for the estimated GxE variance not being different from  
 6 zero. The estimates on the observed scale for the binary trait, hypertension, were transformed to those on the liability scale using  
 7 Robertson transformation[17, 26]. All estimates were from the GxEsum model.

8 <sup>a</sup>NEU: Neuroticism Score

9 <sup>b</sup>PA: Physical Activity

10 <sup>c</sup>ALC: Alcohol intake frequency

11 <sup>d</sup>WHR: Waist-Hip Ratio

12 <sup>e</sup>BFP: Body Fat Percentage

13 <sup>f</sup>BP: Blood pressure

## 1 Discussion

2 In this study, we propose GxEsum, a novel whole-genome GxE method, of which the  
3 computational efficiency is a thousand times higher than existing methods. The  
4 estimation of GxE using GWAS summary statistics has great flexibility in the  
5 application of the method to multiple complex traits and diseases. The proposed  
6 method and theory have been explicitly verified using comprehensive simulations  
7 that were carried out for both quantitative trait and binary disease. Moreover, we  
8 showed that the type I error rate of the proposed method was not inflated by  
9 moderate to severe collider bias [18] that caused a substantial underestimation of  
10 heritability shown in our simulation (Supplementary Figures 10 and 11).

11 In the real data analysis, we show that the genetic effects of BMI were significantly  
12 modulated by NEU, which agrees with previous studies [9]. It is noted that the  
13 significance of GxE was improved because we used a larger sample size, compared  
14 with the previous studies. Our result agrees with Robinson et al. (2017) who found  
15 no significant GxE evidence for age when analysing BMI using the UK Biobank in  
16 which the participants aged 40-69 at the recruitment. However, a dataset with a  
17 wider range of ages is desirable, which would increase the power to detect GxE on  
18 age. For example, a significant GxE was found in a BMI-age analysis using a dataset  
19 including samples aged 18-80 at the recruitment [8]. For hypertension, its causal  
20 relationship with BMI has been reported by a number of studies using Mendelian  
21 randomization [24, 27]. However, it was not clear if the causal relationship is due to  
22 GxE or something else, e.g. unknown non-genetic effects of the disease modulated  
23 by BMI status. We show that the causal relationship between hypertension and BMI,  
24 and type 2 diabetes and BMI [28] reported in the previous studies [24, 27] may be  
25 partly due to genome-wide GxE interaction effects. Interestingly, there is no  
26 significant evidence of genome-wide GxE for hypertension-WHR or hypertension-  
27 BFP causal relationship that was observed in Mendelian randomization studies [29,  
28 30].

29 The estimated intercept from GxEsum should be interpreted with caution. We show  
30 that estimated intercepts were unbiased from the theoretically predicted values when  
31 using the simulation of quantitative traits, as a proof of concept, i.e. the phenotypic  
32 variance explained by RxE effects ( $h_{\tau_1}^2$ ) can be obtained as  $h_{\tau_1}^2 = (\text{intercept} - 1 - 2h_{g_1}^2)$   
33 / 2 from eq. (4), or more generally,  $h_{\tau_1}^2 = (\text{intercept} - 1 - (k_{urtosis} - 1)h_{g_1}^2) / (k_{urtosis} - 1)$  from  
34 eq. (5). However, in real data analyses, there may be additional confounding effects  
35 such as scale effects, residual heteroscedasticity or/and sample heterogeneity that  
36 are often attributed to unknown factors. Moreover, when using binary traits,  
37 substantial scale effects can be generated (statistical RxE effects) because only  
38 affected and unaffected status are observed and individual differences within  
39 affected or unaffected group are ignored. These additional confounding effects and  
40 statistical scale effects are captured and estimated as an intercept in GxEsum [10],  
41 resulting in unreliable RxE estimates. It is noted that RxE estimation is not the main  
42 interest of GxEsum and can be more reliably estimated in RNM that is designed to  
43 model both GxE and RxE.

1 The existing GxE methods require individual-level genotype data which often has a  
2 restriction to share, and their computational burden is typically high. Moreover, it is  
3 not clear how they perform when the representativeness of the samples is limited,  
4 e.g. selection bias due to a collider in the UK Biobank samples. On the contrary, the  
5 proposed approach, GxEsum, is computationally efficient and can detect GxE  
6 interaction correctly for both quantitative and binary disease traits even when there is  
7 moderate to server collider bias. If GWAS summary statistics of estimated main  
8 additive and interaction effects can be made publicly available, a meta-analysis  
9 across multiple cohorts can be possible for an ever-large GxE study (like the context  
10 of LDSC SNP-heritability meta-analysis). There are some issues that the measure of  
11 environmental variable may not be standardised across study cohorts, and the  
12 environmental variable may be even unavailable in some cohorts. However, these  
13 issues can be remedied when the information of exposome that is the standardised  
14 measure of all exposures for individuals, complemented to the genome, is available.

15 There is a GxE method that can use GWAS summary statistics, i.e. VarExp, which is  
16 recently published. While VarExp benefits computationally from using GWAS  
17 summary statistics, it needs to invert the correlation matrix between SNPs, which  
18 prevents from using a large number of SNPs [31]. Furthermore, the theoretical  
19 frameworks of GxEsum and VarExp are fundamentally different in that the latter  
20 does not account for confounding effects such as scale effects, residual  
21 heterogeneity or RxE that can be captured by the estimated intercept of GxEsum.  
22 Finally, the performance of VarExp has been verified with a limited magnitude of  
23 interaction effects up to 1.5% and 0.25% of the phenotypic variance for quantitative  
24 and binary traits, respectively [31].

25 Like RNM, GxEsum can fit environmental exposures such that the genetic effects of  
26 a trait can be modelled as a nonlinear function of a continuous environmental  
27 gradient. The potential modifier of the genetic effects is not limited to environmental  
28 exposures but can be extended to novel variables from multi-omics data such as  
29 gene-expression, protein expression and methylation data [32, 33]. Polygenic risk  
30 scores [34, 35] can also be considered as an environmental variable in the model.  
31 This novel approach may allow dissecting a latent biological architecture of a  
32 complex trait in a future application of GxEsum.

33 In the analysis of binary disease traits, estimates on the liability scale, transformed  
34 from those on the observed scale using Robertson transformation, should be  
35 interpreted with caution. Biased estimates on the liability scale are likely due to the  
36 violation of the normality assumption that is essentially required for the Robertson  
37 transformation, i.e. large interaction effects can cause a non-normal phenotypic  
38 distribution. It is also known that if the transformation involves substantial non-  
39 additive effects, it can give biased estimates on the liability scale [17, 26, 36].  
40 However, when non-additive effects are small, the transformation can give  
41 reasonably accurate estimates on the liability scale, which is also evidenced by our  
42 simulations with small interaction effects. As shown in the real data analysis, the  
43 magnitude of genome-wide GxE is not large (< 10% of the phenotypic variance),  
44 showing that the bias of transformation due to the assumption violation may not be  
45 substantial in general. Nevertheless, it is required to develop a better transformation

1 method for large interaction effects in a further study, e.g. when using multiple  
2 environmental variables simultaneously, the interaction effects are aggregated and  
3 can be substantially large.

4 There are a number of limitations in our study. First, like RNM, GxEsum does not  
5 determine the causal direction between variables, which can be provided from  
6 previous studies or other epidemiologic methods, e.g. Mendelian randomisation, as  
7 prior information. Second, we only modelled the first order of random regression  
8 coefficients with a single environmental variable, and there may be significant  
9 additional effects when modelling a higher-order interaction or multiple  
10 environmental variables. It is possible to extend GxEsum model to fit additional  
11 quadratic and polynomial terms or multiple environmental variables simultaneously.  
12 However, assessing the performance of these advanced models is a formidable task,  
13 requiring a further study. Third, the estimation for the main genetic effects can be  
14 biased when there are large G-E and/or R-E correlations. Because of such  
15 correlations, the main genetic effects are over-adjusted when the phenotypes of the  
16 main trait are adjusted for the environmental variable in the model. Therefore, a  
17 careful interpretation of the estimated main genetic effects is required when using  
18 GxEsum. Fourth, we did not investigate the performance of GxEsum for ascertained  
19 case-control studies in which cases are over-sampled. A further study is required to  
20 extend the method to non-random case-control samples so that it can be applied to  
21 consortium data with multiple case-control studies. Lastly, when using the same  
22 sample size, the precision of GxEsum is not better than GREML-based GxE  
23 methods, implying that the former is only useful when using a large sample size that  
24 the latter cannot handle.

## 25 **Conclusions**

26 Despite these caveats, GxEsum can be a useful tool to estimate whole-genome GxE  
27 as it can achieve a higher precision (i.e. power) from a larger sample size, compared  
28 to existing GxE methods. Especially when the scale of available resources increases,  
29 GxEsum may be a unique method that can be applied to large-scale data across  
30 multiple complex traits and diseases in the context of GxE.

31

32

33

## 34 **Methods**

35

### 36 **GxEsum**

37 Following Ni et al. (2019), RNM can be written as

$$\mathbf{y} = \mathbf{b} + \mathbf{g} + \boldsymbol{\tau} = \mathbf{b} + \mathbf{g}_0 + \mathbf{g}_1 \times \mathbf{E} + \boldsymbol{\tau}_0 + \boldsymbol{\tau}_1 \times \mathbf{E}$$

38 where  $\mathbf{y}$  is  $N$  vector of phenotypic observations,  $\mathbf{b}$  is a vector of fixed effects,  $\mathbf{g}$  is  $N$   
39 individual genetic effects, which can be decomposed into the first and second order

1 of genetic random regression coefficients,  $\mathbf{g}_0$  and  $\mathbf{g}_1$ ,  $\boldsymbol{\tau}$  is residual effects,  
 2 decomposed into the first and second order of residual random regression  
 3 coefficients,  $\boldsymbol{\tau}_0$  and  $\boldsymbol{\tau}_1$ , and  $\mathbf{E}$  is an  $N$  vector of environmental variable. Note that  $\mathbf{E}$   
 4 can be also any covariate variable (e.g. smoking, alcohol intake frequency).

5 Assuming that the phenotypes ( $\mathbf{y}$ ) are pre-adjusted for the main genetic effects ( $\mathbf{g}_0$ ),  
 6 environmental or covariate variable ( $\mathbf{E}$ ) and other fixed effects ( $\mathbf{b}$ ), the model can be  
 7 rewritten as

$$\mathbf{y} = \mathbf{g}_1 \times \mathbf{E} + \boldsymbol{\tau}_0 + \boldsymbol{\tau}_1 \times \mathbf{E} = \mathbf{X}\boldsymbol{\beta}_1 \times \mathbf{E} + \boldsymbol{\tau}_0 + \boldsymbol{\tau}_1 \times \mathbf{E}$$

8 where  $\mathbf{X}$  is an  $N \times M$  standardised genotype matrix for  $M$  SNPs,  $\boldsymbol{\beta}_1$  is an  $M$  vector of  
 9 SNP interaction effects modulated by the environment (i.e. GxE SNP effects). It is  
 10 noted that  $\boldsymbol{\tau}_0$  is residual effects that are consistent across environment whereas  $\boldsymbol{\tau}_1$   
 11 captures heterogeneous residual effects across environment (i.e. RxE).

12 Following Bulik-Sullivan et al. (2015), assuming  $\mathbb{E}[\mathbf{g}_1] = \mathbb{E}[\boldsymbol{\beta}_1] = 0$ , the expected chi-  
 13 square statistics of variant  $j$  for the GxE is

$$\mathbb{E}[\chi_j^2] = N \cdot \text{Var}(\widehat{\beta}_{1j}) \quad \text{eq. (1)}$$

14 Using the law of total variance,  $\text{Var}[\widehat{\beta}_{1j}]$  can be obtained as

$$\begin{aligned} \text{Var}(\widehat{\beta}_{1j}) &= \mathbb{E}[\text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X})] + \text{Var}[\mathbb{E}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X})] \\ &= \mathbb{E}[\text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X})] \end{aligned}$$

15  
 16 where  $\mathbf{E}\mathbf{X}$  is an  $N \times M$  matrix with each column having the Hadamard product  
 17 between  $\mathbf{E}$  and  $\mathbf{X}_j$  (standardised genotypes at the  $j$ th SNP) and the conditional  
 18 expectation of  $\widehat{\beta}_{1j}$  is  $\mathbb{E}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X}) = 0$ .

19 Noting that the least-square estimate of  $\widehat{\beta}_{1j}$  can be obtained as  $\widehat{\beta}_{1j} = (\mathbf{E}\mathbf{X}_j)' \mathbf{y} / N$ ,  
 20  $\text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X})$  can be rearranged as

21

$$\begin{aligned} \text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X}) &= \text{Var}[(\mathbf{E}\mathbf{X}_j)' \mathbf{y} / N | \mathbf{E}\mathbf{X}] \\ &= \frac{1}{N^2} \text{Var}[(\mathbf{E}\mathbf{X}_j)' \mathbf{y} | \mathbf{E}\mathbf{X}] \\ &= \frac{1}{N^2} (\mathbf{E}\mathbf{X}_j)' \text{Var}(\mathbf{y} | \mathbf{E}\mathbf{X}) (\mathbf{E}\mathbf{X}_j) \\ &= \frac{1}{N^2} (\mathbf{E}\mathbf{X}_j)' \text{Var}[(\mathbf{E}\mathbf{X})\boldsymbol{\beta}_1 + \boldsymbol{\tau}_0 + \mathbf{E}\boldsymbol{\tau}_1 | \mathbf{E}\mathbf{X}] (\mathbf{E}\mathbf{X}_j) \\ &= \frac{1}{N^2} (\mathbf{E}\mathbf{X}_j)' \text{Var}[(\mathbf{E}\mathbf{X})\boldsymbol{\beta}_1 | \mathbf{E}\mathbf{X}] (\mathbf{E}\mathbf{X}_j) \\ &\quad + \frac{1}{N^2} (\mathbf{E}\mathbf{X}_j)' (\mathbf{E}\mathbf{X}_j) \text{Var}(\boldsymbol{\tau}_0) + \frac{1}{N^2} (\mathbf{E}\mathbf{X}_j)' (\mathbf{E})(\mathbf{E})' (\mathbf{E}\mathbf{X}_j) \text{Var}(\boldsymbol{\tau}_1) \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{N^2} (\mathbf{EX}_j)' (\mathbf{EX}) (\mathbf{EX})' (\mathbf{EX}_j) \text{Var}(\boldsymbol{\beta}_1 | \mathbf{EX}) \\
 &+ \frac{1}{N^2} [N \text{Var}(\boldsymbol{\tau}_0) + (\mathbf{EX}_j)' (\mathbf{E}) (\mathbf{E})' (\mathbf{EX}_j) \text{Var}(\boldsymbol{\tau}_1)] \\
 &= \frac{1}{N^2} \frac{h_{g_1}^2}{M} (\mathbf{EX}_j)' (\mathbf{EX}) (\mathbf{EX})' (\mathbf{EX}_j) \\
 &+ \frac{1}{N^2} [N(1 - h_{g_1}^2 - h_{\tau_1}^2) + h_{\tau_1}^2 (\mathbf{EX}_j)' (\mathbf{E}) (\mathbf{E})' (\mathbf{EX}_j)]
 \end{aligned}$$

- 1 where  $h_{g_1}^2$  and  $h_{\tau_1}^2$  is the proportion of phenotypic variance explained by GxE and
- 2 RxE, respectively.
- 3 Therefore,  $\mathbb{E}[\chi_j^2]$  in eq. (1) can be written as

$$\begin{aligned}
 \mathbb{E}[\chi_j^2] &= N \cdot \text{Var}(\widehat{\beta}_{1j} | \mathbf{EX}) \\
 &= \frac{1}{N} \left[ \frac{h_{g_1}^2}{M} (\mathbf{EX}_j)' (\mathbf{EX}) (\mathbf{EX})' (\mathbf{EX}_j) + N(1 - h_{g_1}^2 - h_{\tau_1}^2) + h_{\tau_1}^2 (\mathbf{EX}_j)' (\mathbf{E}) (\mathbf{E})' (\mathbf{EX}_j) \right] \quad \text{eq. (2)}
 \end{aligned}$$

4

- 5 According to Bulik-Sullivan et al. (2015), the products of the standardised genotypes
- 6 at variant  $j$  and other variants can be expressed as a function of LD scores, i.e.

$$\begin{aligned}
 \frac{1}{N^2} (\mathbf{X}_j)' (\mathbf{X}) (\mathbf{X})' (\mathbf{X}_j) &= \frac{1}{N^2} (N^2 + [N * (1 - \tilde{r}_j^2) + \tilde{r}_j^2 * N^2] * (M - 1)) \\
 &= \mathbf{1} + \left[ \frac{1 - \tilde{r}_j^2}{N} + \tilde{r}_j^2 \right] * (M - 1) \\
 &= \ell_j + \left[ \frac{(M - 1)(1 - \tilde{r}_j^2)}{N} \right] \\
 &\approx \ell_j + \frac{M - \ell_j}{N}
 \end{aligned}$$

- 7 where  $\tilde{r}_j^2$  defined as the expected sample correlation between genotypes at the  $j$ th
- 8 variant and the other  $(M-1)$  variants, and  $\ell_j = 1 + \tilde{r}_j^2(M - 1)$  is the LD scores of the
- 9  $j$ th SNP.

- 10 According to the central moment theory of standard normal distribution of three
- 11 independent random variables ( $\mathbf{X}_1$ ,  $\mathbf{X}_2$  and  $\mathbf{E}$ ), each with an  $N$  vector, useful
- 12 equations are

$$\begin{aligned}
 \mathbb{E}[(\mathbf{X}_1)' (\mathbf{X}_1)] &= N \\
 \mathbb{E}[(\mathbf{X}_1)' (\mathbf{X}_2) (\mathbf{X}_2)' (\mathbf{X}_1)] &= N \\
 \mathbb{E}[(\mathbf{X}_1)' (\mathbf{X}_1) (\mathbf{X}_1)' (\mathbf{X}_1)] &= N^2 \\
 \mathbb{E}[(\mathbf{EX}_1)' (\mathbf{E}) (\mathbf{E})' (\mathbf{EX}_1)] &= \mathbb{E}[(\mathbf{EX}_1)' (\mathbf{EX}_2) (\mathbf{EX}_2)' (\mathbf{EX}_1)] = 3N
 \end{aligned}$$



1 and

$$\mathbb{E}[(\mathbf{E}\mathbf{X}_1)'(\mathbf{E}\mathbf{X}_1)(\mathbf{E}\mathbf{X}_1)'(\mathbf{E}\mathbf{X}_1)] = N^2.$$

2

3 Therefore, assuming that  $\mathbf{E}$  and  $\mathbf{X}_j$  have negligible correlation for a polygenic trait (i.e.  
4 a tiny proportion of the phenotypic variance of  $\mathbf{E}$  can be explained by a single SNP,  
5  $\mathbf{X}_j$ ), the term  $(\mathbf{E}\mathbf{X}_j)'(\mathbf{E})(\mathbf{E}\mathbf{X}_j)'(\mathbf{E}\mathbf{X}_j)$  can be expressed as a function of LD scores as

$$\begin{aligned} \frac{1}{N^2}(\mathbf{E}\mathbf{X}_j)'(\mathbf{E})(\mathbf{E}\mathbf{X}_j)'(\mathbf{E}\mathbf{X}_j) &= \frac{1}{N^2} \left( N^2 + [3N(1 - \tilde{r}_j^2) + \tilde{r}_j^2 * N^2] * (M - 1) \right) \\ &= 1 + \left[ \frac{3(1 - \tilde{r}_j^2)}{N} + \tilde{r}_j^2 \right] * (M - 1) \\ &= \ell_j + \frac{3(M - \ell_j)}{N} \end{aligned}$$

6 Thus, a part in eq. (2) can be rearranged as

$$\begin{aligned} &\frac{1}{N} \left[ \frac{h_{g1}^2}{M} (\mathbf{E}\mathbf{X}_j)'(\mathbf{E})(\mathbf{E}\mathbf{X}_j)'(\mathbf{E}\mathbf{X}_j) + N(1 - h_{g1}^2) \right] \\ &= \frac{1}{N} \left[ \frac{N^2 h_{g1}^2}{M} \left( \ell_j + \frac{3(M - \ell_j)}{N} \right) + N(1 - h_{g1}^2) \right] \\ &= \frac{N h_{g1}^2}{M} \left( \ell_j + \frac{3(M - \ell_j)}{N} \right) + 1 - h_{g1}^2 \\ &= \frac{(N - 3) h_{g1}^2}{M} \ell_j + 1 + 2h_{g1}^2 \\ &= \frac{N(1 - 3/N) h_{g1}^2}{M} \ell_j + 1 + 2h_{g1}^2 \\ &= \frac{N h_{g1}^2}{M} \ell_j + 1 + 2h_{g1}^2 \end{aligned} \tag{3}$$

7

8 The term,  $1 - 3/N$ , in eq. (3) can be approximated as 1 in the analysis using biobank  
9 scale data, which contains over  $10^5$  samples.

10 The remaining part in eq. (2) can be rearranged as

$$\frac{1}{N} [N(-h_{\tau_1}^2) + h_{\tau_1}^2 (\mathbf{E}\mathbf{X}_j)'(\mathbf{E})(\mathbf{E}\mathbf{X}_j)] = 2h_{\tau_1}^2$$

11 where  $\mathbb{E}[(\mathbf{E}\mathbf{X}_1)'(\mathbf{E})(\mathbf{E}\mathbf{X}_1)] = 3N$  according to the central moment theory of  
12 standard normal distribution (see above), assuming that  $\mathbf{E}$  and each column of  $\mathbf{X}$   
13 have a negligible correlation, which satisfies if  $\mathbf{E}$  is an environmental variable or a  
14 polygenic trait.

1 Therefore,

$$2 \quad \mathbb{E}[\chi_j^2] = N \cdot \text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X}) = \frac{N h_{g1}^2}{M} \ell_j + 1 + 2h_{g1}^2 + 2h_{\tau_1}^2 \quad \text{eq. (4)}$$

3 where  $1 + 2h_{g1}^2 + 2h_{\tau_1}^2$  can be obtained as the intercept of the outcome by fitting to  
4 the proposed model (GxEsum). It is noted Eq. (4) is valid when  $\mathbf{X}_j$  is not strictly  
5 normally distributed, i.e. the centred and standardised genotypes of  $j$ th SNP, which is  
6 already shown in Bulik-Sullivan et al. [10]. When the environmental variable ( $\mathbf{E}$ ) is  
7 non-normal, the general form of the fourth central moment term can be expressed as  
8  $\mathbb{E}[(\mathbf{E}\mathbf{X}_1)'(\mathbf{E})(\mathbf{E})'(\mathbf{E}\mathbf{X}_1)] = k_{urtosis} \cdot N$  where  $k_{urtosis}$  is the kurtosis of  $\mathbf{E}$ . And, only the  
9 intercept part of the Eq. (4) is slightly modified as

$$10 \quad \mathbb{E}[\chi_j^2] = N \cdot \text{Var}(\widehat{\beta}_{1j} | \mathbf{E}\mathbf{X}) = \frac{N h_{g1}^2}{M} \ell_j + 1 + (k_{urtosis} - 1)h_{g1}^2 + (k_{urtosis} - 1)h_{\tau_1}^2 \quad \text{eq. (5)}$$

11 Eq. (4) and (5) are verified using simulations (see Supplementary Note 3 and  
12 Supplementary Table 2).

13 To validate the proposed model in general, we used comprehensive phenotypic  
14 simulations that were based on real genotype data (see Supplementary Note 4).

## 15 Real data

16 UK Biobank data were used, which contains 0.5 million individuals aged between 40-  
17 69 years. The data consists of health-related information for each participant who  
18 was recruited in 2006-2010, and their imputed genomic data (~92 million SNPs) has  
19 been distributed through European Genome-phenome Archive. A stringent quality  
20 control process for individuals was set as followings: 1) who were reported as non-  
21 white British, 2) who were having mismatched gender between the reported and the  
22 inferred by the genotypic data, 3) who were having missing rate over 0.05, 4) who  
23 were having putative sex chromosome aneuploidy. In addition, only HapMap3 SNPs  
24 were used which were passed from the stringent quality controls for SNPs. The filter  
25 for SNPs is set as followings: 1) which were having INFO score less than 0.6, 2)  
26 which were having a MAF less than 1%, 3) which were having Hardy-Weinberg  
27 Equilibrium (HWE) P-value less than 1E-4, 4) one of which from the duplicated SNPs.  
28 From those passing the tough procedures, we additionally excluded one of pair of  
29 samples who were having the genomic relationship higher than 0.05. After quality  
30 control, 288,837 individuals and 1,133,273 SNPs were remained. We estimated LD  
31 scores using the genotypic data of UK Biobank after these quality control processes.

32 Among trait phenotypes available in the UK biobank, we arbitrarily selected BMI (a  
33 quantitative trait), hypertension and type 2 diabetes (binary disease traits) and tested  
34 if the genetic effects of the complex traits were significantly modulated by an  
35 environmental variable, i.e. NEU, ALC, PA or age (for testing BMI), BMI, WHR or  
36 BFP (for hypertension), and BMI, diastolic BP or systolic BP (for type 2 diabetes).  
37 The number cases for hypertension and type 2 diabetes was 134,499 (population  
38 prevalence is 0.51) and 11,694 (population prevalence is 0.04), respectively. The  
39 phenotypes of the main trait were adjusted for potential confounders such as age,  
40 gender, year of birth, assessment centre, Townsend Deprivation Index, genetic batch,

1 household income, educational qualification, the first 10 principal components, and  
2 the environmental variable. For any phenotypic missing value for each variable, we  
3 used the mean of the phenotypes of the variable, i.e. phenotypic imputation with the  
4 mean. A better phenotypic imputation method [21] can be used, which is likely to  
5 improve the significance of GxE. Further details of the variables used on this study  
6 are in Supplementary Note 4.

7 In GWAS, we used a linear model for quantitative traits as well as for binary  
8 responses. The use of a linear model applied to binary responses is because it has  
9 been reported that a logistic regression may generate biased estimates in some  
10 instances [37] and our simulations (Supplementary Note 2) were based on a probit  
11 model (i.e. a linear transformation of the inverse standard normal distribution) that  
12 can be well approximated by a linear model [38].

13

## 14 **Abbreviations**

15 GxE: genotype-by-environment interaction  
16 RNM: reaction norm model  
17 GWAS: genome-wide association studies  
18 SNPs: single nucleotide polymorphisms  
19 GREML: genomic restricted maximum likelihood  
20 LDSC: linkage disequilibrium score regression  
21 RxE: residual-environment interaction  
22 G-E correlation: genotype-environment correlation  
23 R-E correlation: residual-environment correlation  
24 SE: standard error  
25 BMI: body mass index  
26 WHR: waist-hip ratio  
27 BFP: body fat percentage  
28 NEU: neuroticism score  
29 PA: physical activity  
30 ALC: alcohol intake frequency

31

## 32 **Acknowledgements**

33 The HPC resources were provided by the Australian Government through Gadi under the National  
34 Computational Merit Allocation Scheme (NCMAS) and by the university of South Australia through  
35 Tango 2.0. We would like to thank staff and participants of the ARIC study and the UK Biobank for  
36 their valuable contributions. The Atherosclerosis Risk in Communities Study is carried out as a  
37 collaborative study supported by National Heart, Lung, and Blood Institute contracts  
38 (HHSN268201100005C, HHSN268201100006C, HHSN268201100007C, HHSN268201100008C,  
39 HHSN268201100009C, HHSN268201100010C, HHSN268201100011C, and  
40 HHSN268201100012C). We thank the staff and participants of the ARIC study for their  
41 important contributions. Funding for GENEVA was provided by National Human Genome  
42 Research Institute grant U01HG004402 (E. Boerwinkle).

1

## 2 **Authors' contribution**

3 S.H.L. conceived the idea and directed the study. S.H.L. and J.S. derived and  
4 verified the theory. J.S performed the analyses. J.S. and S.H.L. drafted the  
5 manuscript.

6

## 7 **Funding**

8 This research is supported by Australian Research Council (DP 190100766, FT  
9 160100229).

## 10 **Availability of Data and Materials**

### 11 **Data:**

12 The UK Biobank data are accessed via <https://www.ukbiobank.ac.uk/>

13 The ARIC study data are accessed via dbGaP (<https://dbgap.ncbi.nlm.nih.gov>) and  
14 its accession code is phs000280.v7.p1.

### 15 **Software:**

16 GxEsum model is implemented in the script that are publicly available [at](#)  
17 <https://github.com/honglee0707/GxEsum>, and the demonstration of GxEsum  
18 software is described in Supplementary Note 5. The version of source code used in  
19 the manuscript is deposited with DOI: 10.5281/zenodo.4659681 at  
20 <https://zenodo.org/record/4659681#.YGkZXc9xeUk>.

21 LDSC can be download from <https://github.com/bulik/ldsc>

22 PLINK version 1.9 can be download from <https://www.cog-genomics.org/plink/1.9/>

23 MTG2 version 2.15 can be download from

24 <https://sites.google.com/site/honglee0707/mtg2>

## 25 **Ethics approval and consent to participant**

26 The current study was approved by the University of South Australia Human  
27 Research Ethics Committee. The ARIC Study was approved by the institutional  
28 review boards of all participating institutions, including the University of Minnesota,  
29 Johns Hopkins University, University of North Carolina, University of Mississippi  
30 Medical Centre, and Wake Forest University. The UK Biobank was approved by the  
31 North West Multi-centre Research Ethics Committee (11/NW/0382). the reference  
32 number approved by the UK Biobank is 14575. All UK Biobank and ARIC Study  
33 participants gave written informed consent.

## 34 **Competing interests**

35 The authors declare that they have no competing interests.

1

## 2 **References**

- 3 1. Visscher PM, Brown MA, McCarthy MI, Yang J: **Five years of GWAS discovery.** *Am J Hum*  
4 *Genet* 2012, **90**:7-24.
- 5 2. Sud A, Kinnersley B, Houlston RS: **Genome-wide association studies of cancer: current**  
6 **insights and future perspectives.** *Nature Reviews Cancer* 2017, **17**:692-704.
- 7 3. Pasaniuc B, Price AL: **Dissecting the genetics of complex traits using summary association**  
8 **statistics.** *Nature Reviews Genetics* 2017, **18**:117-127.
- 9 4. Yang J, Lee SH, Goddard ME, Visscher PM: **GCTA: a tool for genome-wide complex trait**  
10 **analysis.** *Am J Hum Genet* 2011, **88**:76-82.
- 11 5. Arnau-Soler A, Macdonald-Dunlop E, Adams MJ, Clarke T-K, MacIntyre DJ, Milburn K,  
12 Navrady L, Hayward C, McIntosh AM, Thomson PA, et al: **Genome-wide by environment**  
13 **interaction studies of depressive symptoms and psychosocial stress in UK Biobank and**  
14 **Generation Scotland.** *Translational Psychiatry* 2019, **9**:14.
- 15 6. Gong J, Hutter CM, Newcomb PA, Ulrich CM, Bien SA, Campbell PT, Baron JA, Berndt SI,  
16 Bezieau S, Brenner H, et al: **Genome-Wide Interaction Analyses between Genetic Variants**  
17 **and Alcohol Consumption and Smoking for Risk of Colorectal Cancer.** *PLoS Genet* 2016,  
18 **12**:e1006296.
- 19 7. Manning AK, Hivert M-F, Scott RA, Grimsby JL, Bouatia-Naji N, Chen H, Rybin D, Liu C-T,  
20 Bielak LF, Prokopenko I, et al: **A genome-wide approach accounting for body mass index**  
21 **identifies genetic variants influencing fasting glycemic traits and insulin resistance.** *Nature*  
22 *Genetics* 2012, **44**:659-669.
- 23 8. Robinson MR, English G, Moser G, Lloyd-Jones LR, Triplett MA, Zhu Z, Nolte IM, van Vliet-  
24 Ostaptchouk JV, Snieder H, LifeLines Cohort S, et al: **Genotype-covariate interaction effects**  
25 **and the heritability of adult body mass index.** *Nat Genet* 2017, **49**:1174-1181.
- 26 9. Ni G, van der Werf J, Zhou X, Hypponen E, Wray NR, Lee SH: **Genotype-covariate correlation**  
27 **and interaction disentangled by a whole-genome multivariate reaction norm model.** *Nat*  
28 *Commun* 2019, **10**:2239.
- 29 10. Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the  
30 Psychiatric Genomics C, Patterson N, Daly MJ, Price AL, Neale BM: **LD Score regression**  
31 **distinguishes confounding from polygenicity in genome-wide association studies.** *Nat*  
32 *Genet* 2015, **47**:291-295.
- 33 11. Finucane HK, Bulik-Sullivan B, Gusev A, Trynka G, Reshef Y, Loh PR, Anttila V, Xu H, Zang C,  
34 Farh K, et al: **Partitioning heritability by functional annotation using genome-wide**  
35 **association summary statistics.** *Nat Genet* 2015, **47**:1228-1235.
- 36 12. Gazal S, Marquez-Luna C, Finucane HK, Price AL: **Reconciling S-LDSC and LDK functional**  
37 **enrichment estimates.** *Nat Genet* 2019, **51**:1202-1204.
- 38 13. Hou K, Burch KS, Majumdar A, Shi H, Mancuso N, Wu Y, Sankararaman S, Pasaniuc B:  
39 **Accurate estimation of SNP-heritability from biobank-scale data irrespective of genetic**  
40 **architecture.** *Nat Genet* 2019, **51**:1244-1251.
- 41 14. Ni G, Moser G, Schizophrenia Working Group of the Psychiatric Genomics C, Wray NR, Lee  
42 SH: **Estimation of Genetic Correlation via Linkage Disequilibrium Score Regression and**  
43 **Genomic Restricted Maximum Likelihood.** *Am J Hum Genet* 2018, **102**:1185-1194.
- 44 15. Rosner B: *Fundamentals of biostatistics.* Nelson Education; 2015.
- 45 16. Austin PC: **Type I error rates, coverage of confidence intervals, and variance estimation in**  
46 **propensity-score matched analyses.** *Int J Biostat* 2009, **5**:Article 13.
- 47 17. Lee SH, Wray NR, Goddard ME, Visscher PM: **Estimating missing heritability for disease**  
48 **from genome-wide association studies.** *American journal of human genetics* 2011, **88**:294-  
49 305.

- 1 18. Munafò MR, Tilling K, Taylor AE, Evans DM, Davey Smith G: **Collider scope: when selection**  
2 **bias can substantially influence observed associations.** *International Journal of*  
3 *Epidemiology* 2017, **47**:226-235.
- 4 19. Lynch M, Walsh B: *Genetics and analysis of quantitative traits*. Sinauer Sunderland, MA;  
5 1998.
- 6 20. Visscher PM, Hemani G, Vinkhuyzen AA, Chen GB, Lee SH, Wray NR, Goddard ME, Yang J:  
7 **Statistical power to detect genetic (co)variance of complex traits using SNP data in**  
8 **unrelated samples.** *PLoS Genet* 2014, **10**:e1004269.
- 9 21. Hormozdiari F, Kang Eun Y, Bilow M, Ben-David E, Vulpe C, McLachlan S, Lusi Aldons J, Han  
10 B, Eskin E: **Imputing Phenotypes for Genome-wide Association Studies.** *The American*  
11 *Journal of Human Genetics* 2016, **99**:89-103.
- 12 22. Sutin AR, Ferrucci L, Zonderman AB, Terracciano A: **Personality and obesity across the adult**  
13 **life span.** *Journal of personality and social psychology* 2011, **101**:579-592.
- 14 23. Rask-Andersen M, Karlsson T, Ek WE, Johansson Å: **Gene-environment interaction study for**  
15 **BMI reveals interactions between genetic factors and physical activity, alcohol**  
16 **consumption and socioeconomic status.** *PLoS genetics* 2017, **13**:e1006977-e1006977.
- 17 24. Hyppönen E, Mulugeta A, Zhou A, Santhanakrishnan VK: **A data-driven approach for**  
18 **studying the role of body mass in multiple diseases: a phenome-wide registry-based case-**  
19 **control study in the UK Biobank.** *The Lancet Digital Health* 2019, **1**:e116-e126.
- 20 25. Lee M-R, Lim Y-H, Hong Y-C: **Causal association of body mass index with hypertension using**  
21 **a Mendelian randomization design.** *Medicine* 2018, **97**:e11252-e11252.
- 22 26. Dempster ER, Lerner IM: **Heritability of Threshold Characters.** *Genetics* 1950, **35**:212-236.
- 23 27. Larsson SC, Bäck M, Rees JMB, Mason AM, Burgess S: **Body mass index and body**  
24 **composition in relation to 14 cardiovascular conditions in UK Biobank: a Mendelian**  
25 **randomization study.** *European heart journal* 2020, **41**:221-226.
- 26 28. Morrison J, Knoblauch N, Marcus JH, Stephens M, He X: **Mendelian randomization**  
27 **accounting for correlated and uncorrelated pleiotropic effects using genome-wide**  
28 **summary statistics.** *Nature Genetics* 2020, **52**:740-747.
- 29 29. Si S, Tewara MA, Li Y, Li W, Chen X, Yuan T, Liu C, Li J, Wang B, Li H, et al: **Causal Pathways**  
30 **from Body Components and Regional Fat to Extensive Metabolic Phenotypes: A Mendelian**  
31 **Randomization Study.** *Obesity* 2020, **28**:1536-1549.
- 32 30. Zanetti D, Tikkanen E, Gustafsson S, Priest James R, Burgess S, Ingelsson E: **Birthweight, Type**  
33 **2 Diabetes Mellitus, and Cardiovascular Disease.** *Circulation: Genomic and Precision*  
34 *Medicine* 2018, **11**:e002054.
- 35 31. Lavielle V, Bentley AR, Privé F, Zhu X, Gauderman J, Winkler TW, Province M, Rao DC, Aschard  
36 H: **VarExp: estimating variance explained by genome-wide GxE summary statistics.**  
37 *Bioinformatics* 2018, **34**:3412-3414.
- 38 32. Zhou X, Im HK, Lee SH: **CORE GREML for estimating covariance between random effects in**  
39 **linear mixed models for complex trait analyses.** *Nature Communications* 2020, **11**:4208.
- 40 33. Zhou X, Lee SH: **An integrative analysis of genomic and exposomic data for complex traits**  
41 **and phenotypic prediction.** *bioRxiv* 2020:2020.2011.2009.373704.
- 42 34. Khera AV, Chaffin M, Aragam KG, Haas ME, Roselli C, Choi SH, Natarajan P, Lander ES, Lubitz  
43 SA, Ellinor PT, Kathiresan S: **Genome-wide polygenic scores for common diseases identify**  
44 **individuals with risk equivalent to monogenic mutations.** *Nature Genetics* 2018, **50**:1219-  
45 1224.
- 46 35. Truong B, Zhou X, Shin J, Li J, van der Werf JHJ, Le TD, Lee SH: **Efficient polygenic risk scores**  
47 **for biobank scale data by exploiting phenotypes from inferred relatives.** *Nature*  
48 *Communications* 2020, **11**:3074.
- 49 36. Van Vleck LD: **Estimation of Heritability of Threshold Characters.** *Journal of Dairy Science*  
50 1972, **55**:218-225.

- 1 37. Sun R, Carroll RJ, Christiani DC, Lin X: **Testing for gene-environment interaction under**
- 2 **exposure misspecification.** *Biometrics* 2018, **74**:653-662.
- 3 38. Lee SH, Goddard ME, Wray NR, Visscher PM: **A Better Coefficient of Determination for**
- 4 **Genetic Profile Analysis.** *Genetic Epidemiology* 2012, **36**:214-224.

5