

1 **Summary statistics from large-scale gene-environment**

2 **interaction studies for re-analysis and meta-analysis**

3 Duy T. Pham,^{1,15} Kenneth E. Westerman,^{2,3,4,15} Cong Pan,¹ Ling Chen,^{2,3} Shylaja Srinivasan,⁵
4 Elvira Isganaitis,⁶ Mary Ellen Vajravelu,⁷ Fida Bacha,⁸ Steve Chernausek,⁹ Rose Gubitosi-
5 Klug,¹⁰ Jasmin Divers,¹¹ Catherine Pihoker,¹² Santica M. Marcovina,¹³ Alisa K. Manning,^{2,3,4} Han
6 Chen^{1,14,*}

7 ¹ Human Genetics Center, Department of Epidemiology, Human Genetics and Environmental
8 Sciences, School of Public Health, The University of Texas Health Science Center at Houston,
9 Houston, Texas 77030, USA

10 ² Department of Medicine, Clinical and Translational Epidemiology Unit, Mongan Institute,
11 Massachusetts General Hospital, Boston, MA 02114, USA

12 ³ Metabolism Program, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

13 ⁴ Department of Medicine, Harvard Medical School, Boston, MA 02115, USA

14 ⁵ Department of Pediatrics, University of California, San Francisco, CA 94158, USA

15 ⁶ Joslin Diabetes Center, Boston, MA, 02115, USA

16 ⁷ Department of Pediatrics, University of Pittsburgh School of Medicine, Pittsburgh, PA 15224,
17 USA

18 ⁸ Department of Pediatrics, Baylor College of Medicine, Houston, TX 77030, USA

19 ⁹ Department of Pediatrics, The University of Oklahoma College of Medicine, Oklahoma City,
20 OK 73117, USA

21 ¹⁰ Department of Pediatrics, Case Western Reserve University, Cleveland OH, 44106, USA

22 ¹¹ Department of Foundations of Medicine, New York University, New York, NY 10016, USA

23 ¹² Department of Pediatrics, University of Washington School of Medicine, Seattle, WA, 98105,
24 USA

25 ¹³ Northwest Lipid Metabolism and Diabetes Research Laboratories, University of Washington,
26 Seattle, WA, 98105, USA

27 ¹⁴ Center for Precision Health, McWilliams School of Biomedical Informatics, The University of
28 Texas Health Science Center at Houston, Houston, Texas 77030, USA

29 ¹⁵ These authors contributed equally.

30 *Correspondence: han.chen.2@uth.tmc.edu

31

32

33

34

35

36

37

38

39 **Abstract**

40 Summary statistics from genome-wide association studies enable many valuable downstream
41 analyses that are more efficient than individual-level data analysis while also reducing privacy
42 concerns. As growing sample sizes enable better-powered analysis of gene-environment
43 interactions (GEIs), there is a need for GEI-specific methods that manipulate and use summary
44 statistics. We introduce two tools to facilitate such analysis, with a focus on statistical models
45 containing multiple gene-exposure and/or gene-covariate interaction terms. REGEM (RE-
46 analysis of GEM summary statistics) uses summary statistics from a single, multi-exposure
47 genome-wide interaction study (GWIS) to derive analogous sets of summary statistics with
48 arbitrary sets of exposures and interaction covariate adjustments. METAGEM (META-analysis
49 of GEM summary statistics) extends current fixed-effects meta-analysis models to incorporate
50 multiple exposures from multiple studies. We demonstrate the value and efficiency of these
51 tools by exploring alternative methods of accounting for ancestry-related population stratification
52 in GWIS in the UK Biobank as well as by conducting a multi-exposure GWIS meta-analysis in
53 cohorts from the diabetes-focused ProDiGY consortium. These programs help to maximize the
54 value of summary statistics from diverse and complex GEI studies.

55

56

57

58

59

60

61 **Introduction**

62 Gene-environment interaction (GEI) analysis is a key tool for understanding genetic impacts on
63 human traits, with the potential to account for additional heritability, explain differences in
64 genetic effects across populations, and support personalized lifestyle and therapeutic decisions.
65 Historically, GEI studies have taken a hypothesis-driven approach, but larger cohorts,¹ and new
66 software programs have provided the necessary statistical power and computational efficiency
67 to study GEIs genome-wide.^{2,3,4,5,6,7} These genome-wide interaction studies (GWIS) generate
68 summary statistics, or variant-level regression results, which have substantial value beyond
69 locus mapping. For example, summary statistics allow for heritability analysis,⁸ enrichment
70 testing,¹ and genome-wide polygenic score generation.^{1,9}

71 GEI analysis and interpretation are complicated by the densely correlated set of possible
72 exposures that may interact with genotypes to influence human traits (the “exposome”, defined
73 here as including demographic and physiologic traits). Two modeling implications are
74 particularly pertinent. First, multi-exposure GEI analysis can increase statistical power by jointly
75 testing genetic interactions with multiple exposures.^{5,10,11} This strategy can pool signals across
76 distinct exposures (e.g., smoking status and pollution exposure for lung function) or incorporate
77 multiple definitions of a single exposure category (e.g., current smoking status and pack-years
78 of smoking). Second, proper control of confounding for GEI interaction terms requires
79 adjustment for not just the main effects of covariates, but also their genetic interactions.¹²
80 Inclusion of these “interaction covariates” is thus necessary to produce interpretable summary
81 statistics.

82 Rigorous GEI analysis carries complexities stemming from its place at the center of traditional
83 and genetic epidemiology. Sensitivity analyses, while commonplace in traditional epidemiology,
84 are computationally burdensome when conducted across millions of variants genome-wide.

85 Meanwhile, well-established meta-analysis procedures for genome-wide association study
86 (GWAS) summary statistics become more difficult in the context of multi-exposure GEI models.
87 Software programs do not yet exist to perform efficient meta-analysis in the context of these
88 complex analytical designs.

89 We introduce methods and associated software programs to advance the field of genome-wide
90 GEI analysis based on summary statistics. While the statistical results are general, the
91 associated software implementations build on the results from our previously described software
92 program for efficient GWIS, GEM.² Exploiting the redundancy of statistical estimates across
93 related GEI models, we introduce the REGEM (RE-analysis of GEM summary statistics)
94 program to derive genome-wide summary statistics corresponding to arbitrary multi-exposure
95 and interaction covariate adjustments based on results from a single, multi-exposure GWIS.
96 Expanding current fixed-effect meta-analysis models, we further introduce the METAGEM
97 (META-analysis of GEM summary statistics) program to conduct efficient meta-analysis of GEI
98 effects under complex GEI analysis models. We demonstrate the value and efficiency of these
99 tools by exploring alternative methods of accounting for population stratification in GWIS in the
100 UK Biobank as well as by conducting a multi-exposure GWIS meta-analysis in cohorts from the
101 ProDiGY consortium.

102 **Material and methods**

103 **GEM method**

104 We developed two C++ software programs that use summary statistics from GEI studies.
105 REGEM requires output from a single GEI study, while METAGEM requires output from multiple
106 GEI studies. Both programs are designed for easy integration with output from GEM. Here we
107 summarize the GEM methodology. For a single-variant test of N unrelated individuals, GEM
108 considers the generalized linear model:

$$g(\mu_i) = X_i\beta_X + G_i\beta_G + C_i\beta_C + S_i\beta_S \#(1)$$

109 for individual i , where $\mu_i = E(Y_i|X_i, G_i)$ is the conditional mean of the phenotype Y_i given p
 110 covariates X_i (including the intercept), and the genotype G_i for a single genetic variant. The
 111 interaction terms C_i and S_i are the products of G_i and c covariates and q exposures (which are
 112 disjoint subsets of X_i), respectively.² Let $Y = (Y_1 Y_2 \dots Y_N)^T$ be a length N vector of phenotypes,
 113 $X = (X_1^T X_2^T \dots X_N^T)^T$ be an $N \times p$ matrix of p covariates, $G = (G_1 G_2 \dots G_N)^T$ be a length N
 114 vector of genotypes for this single genetic variant, $C = (C_1^T C_2^T \dots C_N^T)^T$ be an $N \times c$ matrix of c
 115 gene-covariate interaction terms, $S = (S_1^T S_2^T \dots S_N^T)^T$ be an $N \times q$ matrix of q gene-
 116 environment (exposure) interaction terms, we can fit a null model without any genetic effects
 117 $g(\mu_i) = X_i\beta_X$ and get a length N residual vector r . Let $\tilde{G} = G - X(X^T W X)^{-1} X^T W G$, $\tilde{C} = C -$
 118 $X(X^T W X)^{-1} X^T W C$ and $\tilde{S} = S - X(X^T W X)^{-1} X^T W S$ be covariate X adjusted G , C and S ,
 119 respectively, where W is a diagonal weight matrix with elements $\hat{\mu}_i(1 - \hat{\mu}_i)$ for logistic
 120 regressions ($\hat{\mu}_i$ are fitted probabilities of $Y_i = 1$ from the null model) and an identity matrix for
 121 linear regressions, GEM computes a length $(1 + c + q)$ score vector ($c \geq 0$) $U = (\tilde{G} \tilde{C} \tilde{S})^T r$, and
 122 $(1 + c + q) \times (1 + c + q)$ matrices $V = (\tilde{G} \tilde{C} \tilde{S})^T W (\tilde{G} \tilde{C} \tilde{S})$, $\Omega = (\tilde{G} \tilde{C} \tilde{S})^T D (\tilde{G} \tilde{C} \tilde{S})$, where D is a
 123 diagonal matrix of squared residuals.

124 For M variants in a genome-wide scan, we retrieve the dispersion parameter estimate, $\hat{\phi}$ (which
 125 is fixed at 1 for logistic regressions and the residual variance estimate from the null model for
 126 linear regressions), the genetic main effect, gene-covariate interaction effects and gene-
 127 environment (exposure) interaction effects, as well as both model-based and robust standard
 128 errors and covariances for G , C and S . The effect estimates are computed as $\hat{\beta}_{G,C,S} = V^{-1}U$.
 129 The full $(1 + c + q) \times (1 + c + q)$ model-based and robust variance-covariance matrices are
 130 computed as $Cov(\hat{\beta}_{G,C,S}) = \hat{\phi}V^{-1}$ and $Cov_R(\hat{\beta}_{G,C,S}) = V^{-1}\Omega V^{-1}$, respectively. In the full output,
 131 GEM (version 1.3 and later) reports the model-based and robust standard errors of effect

132 estimates, which are the square root of the diagonal elements of $Cov(\hat{\beta}_{G,C,S})$ and $Cov_R(\hat{\beta}_{G,C,S})$,
133 as well as the model-based and robust covariances for these effect estimates (the off-diagonal
134 elements of $Cov(\hat{\beta}_{G,C,S})$ and $Cov_R(\hat{\beta}_{G,C,S})$).

135 **REGEM Method**

136 Given the full summary statistics output from GEM (version 1.3 and later), the score vector U
137 and matrices V and Ω , can be reconstructed without access to individual-level data. Utilizing $\hat{\phi}$
138 and the matrices $Cov(\hat{\beta}_{G,C,S})$ and $Cov_R(\hat{\beta}_{G,C,S})$ described above, it follows that

139 $V = \hat{\phi}Cov^{-1}(\hat{\beta}_{G,C,S})$ and $\Omega = VCov_R(\hat{\beta}_{G,C,S})V$. The score vector can then be recomputed as

$$140 \quad U = V\hat{\beta}_{G,C,S}.$$

141 REGEM supports two scenarios for re-analysis of a single GEI study. The first scenario involves
142 the exclusion of one or more gene-covariate or gene-environment interaction terms from the
143 original model. This is achieved by filtering U to exclude the specified gene-covariate or gene-
144 environment interaction terms, resulting in the modified score vector \dot{U} . Subsequently, the
145 matrices V and Ω are reduced to exclude the corresponding rows and columns of the specified
146 gene-covariate or gene-environment interaction terms, denoted \dot{V} and $\dot{\Omega}$. The GEM method can
147 then be applied to \dot{U} , \dot{V} , and $\dot{\Omega}$ to obtain new summary statistics. In the second scenario, re-
148 analysis can be performed by conditioning on one or more gene-environment interaction terms
149 in the original GEM analysis as gene-covariate interactions or testing one or more gene-
150 covariate interaction terms in the original GEM analysis as gene-environment interaction terms
151 of interests. In either case, the ordering of U is rearranged, denoted as \ddot{U} , to incorporate the
152 original gene-environment interaction terms into C or the original gene-covariate interaction
153 terms into S . The rows and columns of the matrices V and Ω are also reordered and denoted as

154 \hat{U} and $\hat{\Omega}$. The GEM method follows for \hat{U} , \hat{V} , and $\hat{\Omega}$. Both scenarios can be applied
155 simultaneously.

156 **METAGEM method**

157 METAGEM combines summary statistics from K independent studies using the inverse-
158 variance weighted approach. For individual studies $k = 1, 2, \dots, K$, with effect estimates $\hat{\beta}_k$ and
159 the variance-covariance matrix Cov_k from the GEM output (model-based or robust), the
160 summary effect estimates are computed as $\hat{\beta} = (\sum_{k=1}^K Cov_k^{-1})^{-1} (\sum_{k=1}^K Cov_k^{-1} \hat{\beta}_k)$, with the
161 model-based or robust variance-covariance matrix $Cov = (\sum_{k=1}^K Cov_k^{-1})^{-1}$.

162 **REGEM Comparison and Benchmark**

163 To demonstrate the computational benefits of REGEM, we test and compare four variations of
164 the waist-hip ratio (WHR) model originally described by Westerman et al. The original model is
165 defined as follows (excluding the array covariate and PC6 - PC10):

$$WHR \sim G + sex + age + age^2 + BMI + PC1 + \dots + PC5 + G \times sex + G \times BMI \#(2)$$

166 where WHR is the phenotype, sex is the primary exposure of interest, BMI is the interaction
167 covariate, and age, age², and PC1-PC5 are the covariates. Here, we retrieved PCs calculated
168 as part of the Pan-UKBB project.¹³ All terms in the model were centered. First, we performed a
169 genome-wide analysis of the original model using GEM (version 1.5) using 362,449 unrelated
170 European-ancestry participants, and filtered variants with minor allele frequency (MAF) < 0.001,
171 leaving 16,539,280 variants for re-analysis. Next, we derived associated genome-wide summary
172 statistics corresponding to variations of the original model using REGEM, comparing their
173 results and runtimes to simply re-running that same model genome-wide using GEM. Table S1
174 summarizes the variations of the original models, including the original model. These variations

175 involve the joint testing of G x sex and G x BMI (M1), testing for G x BMI while adjusting for G x
176 sex (M2), testing for G x sex while removing the G x BMI term (M3), and testing for G x BMI
177 while removing the G x sex term (M4). All analyses were performed on the DNAnexus platform
178 using the *mem1_ssd1_v2_x16* instance type, and we reported the runtime and memory usage
179 of each run. The GEM and REGEM summary statistic comparisons were visualized using the
180 scattermore and ggplot2 R packages.

181 **METAGEM Comparison and Benchmark**

182 To evaluate the computational efficiency of METAGEM, we conducted a simulation study using
183 phenotype and genotype data from the Pan-UKBB.¹³ We randomly split the phenotype data,
184 which comprised 362,449 samples, into 11 datasets: one with 100,000 samples, two with
185 50,000 samples, seven with 10,000 samples, and one with 92,449 samples. For each dataset,
186 we conducted a genome-wide gene-sex interaction test and filtered out variants with a MAF <
187 0.001, resulting in 15.46 to 16.85 million variants per dataset, and a total of 17,993,341 unique
188 variants across all datasets. We then performed a gene-sex interaction meta-analysis using
189 METAGEM and the METAL software (version 2010-02-08),¹⁴ with the joint meta-analysis
190 patch,¹⁵ and compared the results. Additionally, we conducted a genome-wide joint gene-sex
191 and gene-BMI interaction test for each dataset and performed a meta-analysis using
192 METAGEM to evaluate its performance in the presence of multiple interaction terms. All
193 analyses were conducted on the DNAnexus platform using a single core and the
194 *mem1_ssd1_v2_x16* instance type. We reported the CPU time and memory usage for each
195 analysis. We used the scattermore and ggplot2 R packages to visualize the comparison of
196 summary statistics between METAGEM and METAL.

197 **Multi-exposure interactions influencing waist-hip ratio in the UK Biobank**

198 Expanding the WHR analyses described above, we performed multiple GWIS, with downstream
199 analysis using REGEM and METAGEM, to investigate genetic interactions with sex and BMI
200 across multiple ancestries. The primary model, run using GEM, was conducted in unrelated
201 individuals from multiple ancestries (N = 379,092) and followed model (2) above with the
202 addition of gene-ancestry interaction covariates. Ancestry labels (AFR, AMR, CSA, EAS, EUR,
203 and MID) were retrieved from the Pan-UKBB effort and were coded using five indicator
204 variables, with EUR as the reference group. Using REGEM, we then derived summary statistics
205 corresponding to equivalent single-exposure GWIS in the pooled-ancestry sample (testing only
206 gene-sex or only gene-BMI interactions, while adjusting for only the main effect of the other).
207 Additionally, we ran ancestry-stratified, multi-exposure analyses (using the same model but
208 removing all covariate and interaction covariate terms containing ancestry labels). These
209 ancestry-stratified analyses were then combined using METAGEM to generate meta-analyzed,
210 multi-exposure interaction tests for comparison to the results from the ancestry-pooled analysis.

211 To compare locus discoveries across analysis strategies (e.g., ancestry-pooled vs. cross-
212 ancestry meta-analysis), we first independently clumped summary statistics from each analysis
213 using a distance-based method that grouped variants within 500kb of each lead variant. We
214 then concatenated the clumped results from the two analyses and performed a secondary
215 clumping using the same strategy, such that clumped loci in this second stage were considered
216 to represent the same locus.

217 **Progress in Diabetes Genetics in Youth (ProDiGY) dataset**

218 ProDiGY is a multi-ethnic resource including three studies: Treatment Options for Type 2
219 Diabetes in Adolescents and Youth (TODAY),¹⁶ SEARCH for Diabetes in Youth (SEARCH),¹⁷
220 and T2D-GENES. In total, the dataset contains 2,820 youth and 4,858 adult cases with T2D,
221 and 656 diabetes-free youth and 4,934 adult controls after removing individuals with maturity-

222 onset diabetes of the young (MODY) and type I diabetes. Samples were genotyped on the
223 Infinium GWAS array by the Genetic Analysis Platform at the Broad Institute of MIT and
224 Harvard. Details on quality control procedures for the genotype data have been previously
225 described.¹⁸ Genotype data were imputed on the TOPMed Imputation Server using the
226 TOPMed v2 reference panel. Variants passing an imputation quality threshold (R^2) of 0.5 were
227 retained for analysis. Genetic ancestry groups were assigned to ProDiGY samples based on
228 genetic principal components analysis after merging with the 1000 Genomes dataset.

229 **Application multi-interaction to T2D in ProDiGY**

230 To show the performance of METAGEM in the multi-gene-environment interactions with a real
231 and genome-wide study, we first used GEM to conduct a multi-exposure gene-sex and gene-
232 age interaction analysis for incident T2D, separately within each genetic ancestry group in two
233 different comparisons: youth cases vs. youth controls (youth group) and adult cases vs. adult
234 controls (adult group). Sex and age were both used as exposures and tested jointly for
235 interaction using robust standard errors. Covariates included age, sex and 10 genetic principal
236 components.

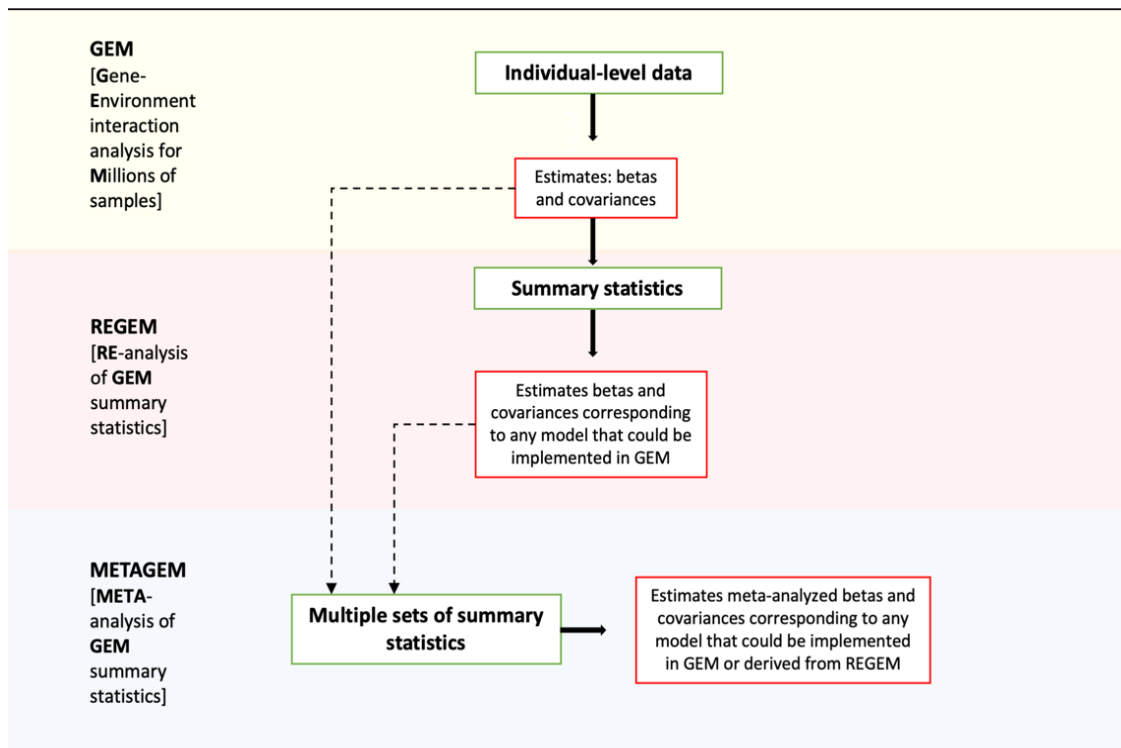
$$T2D \sim G + sex + age + PC1 + \dots + PC10 + G \times sex + G \times age \#(3)$$

237 Using the full output from GEM, we performed cross-ancestry meta-analysis using METAGEM
238 in both youth group and adult group analyses. We also conducted equivalent single-exposure
239 GWIS with sex and age separately for comparison with the multi-exposure scan. Meta-analysis
240 for these single-exposure tests was conducted using METAL, for both the joint (genetic plus
241 interaction effect) test (patched version 2010-02-08; the only version for which the patch is
242 available) and marginal test (version 2011-03-25) to conduct the marginal meta-analysis test
243 across genetic ancestry groups. A threshold of $p < 5 \times 10^{-8}$ was used to define genome-wide
244 significance.

245 **RESULTS**

246 Figure 1 shows the suite of software tools described here in the context of an analysis workflow,
 247 along with an example set of associated statistical models.

248



Software	Exposures	Interaction Covariates	Summary Statistics ignored	Associated (implicit) regression model	Interaction test H_0
GEM	E1, E2, E3	None	None	Model 1	$\beta_{gE_1} = \beta_{gE_2} = \beta_{gE_3} = 0$
REGEM	E1, E2	E3	None	Model 1	$\beta_{gE_1} = \beta_{gE_2} = 0$
REGEM	E1	None	E2, E3	Model 2	$\beta_{gE_1}^* = 0$

Model 1: $E(Y|g, E_1, E_2, E_3, C) = \beta_0 + \beta_g g + \beta_{E_1} E_1 + \beta_{E_2} E_2 + \beta_{E_3} E_3 + \beta_{gE_1} g E_1 + \beta_{gE_2} g E_2 + \beta_{gE_3} g E_3 + \beta_C C$

Model 2: $E(Y|g, E_1, E_2, E_3, C) = \beta_0^* + \beta_g^* g + \beta_{E_1}^* E_1 + \beta_{E_2}^* E_2 + \beta_{E_3}^* E_3 + \beta_{gE_1}^* g E_1 + \beta_C^* C$

Note: Main effect betas for covariates, intercept, and E main effects are the same in Models 1 & 2 in GEM's implementation (though not in the classical model).

249

250 **Figure 1:** Large-scale GxE methods software suite and connections in the context of an
251 analysis workflow. GEM (previously published) conducts genome-wide interaction studies for
252 single datasets. Given multi-exposure summary statistics from GEM (version 1.3 and later),
253 REGEM can estimate genome-wide summary statistics from an associated model that re-
254 partitions any subset of exposures into interaction covariates and simple main effect
255 adjustments without interaction. Given multiple sets of summary statistics from GEM and/or
256 REGEM, METAGEM conducts meta-analysis for any number of jointly-tested exposures and
257 interaction covariates.

258

259 **REGEM computational performance**

260 We compared results obtained from genome-wide interactions tests using the REGEM and
261 GEM methods across four distinct GEI models. The benchmark results, presented in Table 1,
262 indicate that REGEM significantly reduces CPU time by eliminating the need for computation on
263 individual-level data. For each model, REGEM completed a genome-wide run in less than 6
264 minutes, while GEM required several CPU days to achieve the same outcome. Additionally, re-
265 analyses for multiple interactions (M1 and M2) using REGEM took only about a minute of
266 additional CPU time compared to single exposure re-analyses (M3 and M4). Overall, REGEM
267 saved considerable time, ranging from hours to days of computation time. Moreover, the
268 memory requirements for REGEM were minimal, primarily depending on the number of gene-
269 environment interaction terms, which are usually small. Finally, the effect and variance
270 estimates from REGEM were consistent with those obtained from GEM for each of the four
271 models (M1-M4) as shown in Figures S1-S4.

272

Benchmark	GEM				REGEM			
	M1	M2	M3	M4	M1	M2	M3	M4
CPU time (Mins)	13,972.17	13,618.44	10,959.33	10,994.26	5.22	5.20	4.43	4.06
Memory (MB)	2,325.37	2,342.48	2,188.14	2,188.14	13.66	13.64	11.43	11.63

273 **Table 1.** Genome-wide re-analysis benchmark comparison between GEM and REGEM.

274

275 **METAGEM computational performance**

276 Genome-wide meta-analysis runs of ~17.99 million variants, derived from 11 simulated UKB
277 datasets, were carried out using the METAGEM and METAL methods with a single core. Table
278 2 summarizes the CPU time and memory usage of the runs. For a single exposure meta-
279 analysis, METAGEM showed a modest improvement in performance compared to METAL,
280 completing the run approximately 2 minutes faster and using approximately 1 GB less memory.
281 We note that METAGEM meta-analyzed all 17,993,341 variants, while METAL skipped 25,670
282 multi-allelic variants that contained duplicate variant identifiers. However, the impact of the
283 skipped variants on the benchmark results was negligible. Model-based and robust meta-
284 analysis results from METAGEM and METAL are compared in Figure S5. As expected, the
285 summary statistics and joint *P*-values were consistent between the two methods. To test the
286 performance of METAGEM in conducting meta-analysis with multiple interactions, we performed
287 genome-wide joint meta-analysis with gene-sex and gene-BMI as the interactions using
288 METAGEM. As shown in Table 2, METAGEM efficiently completed the run in an additional ~6
289 minutes of CPU time and less than 1 GB of additional memory compared to the single exposure
290 meta-analysis.

291

Benchmark	METAL	METAGEM	
	1 - Exposure	1 - Exposure	2 - Exposures
CPU time (Mins)	16.38	14.38	19.55
Memory (GB)	7.10	6.11	6.96

292 **Table 2.** Genome-wide meta-analysis benchmark between METAL and METAGEM for
293 17,993,341 variants using a single core.

294

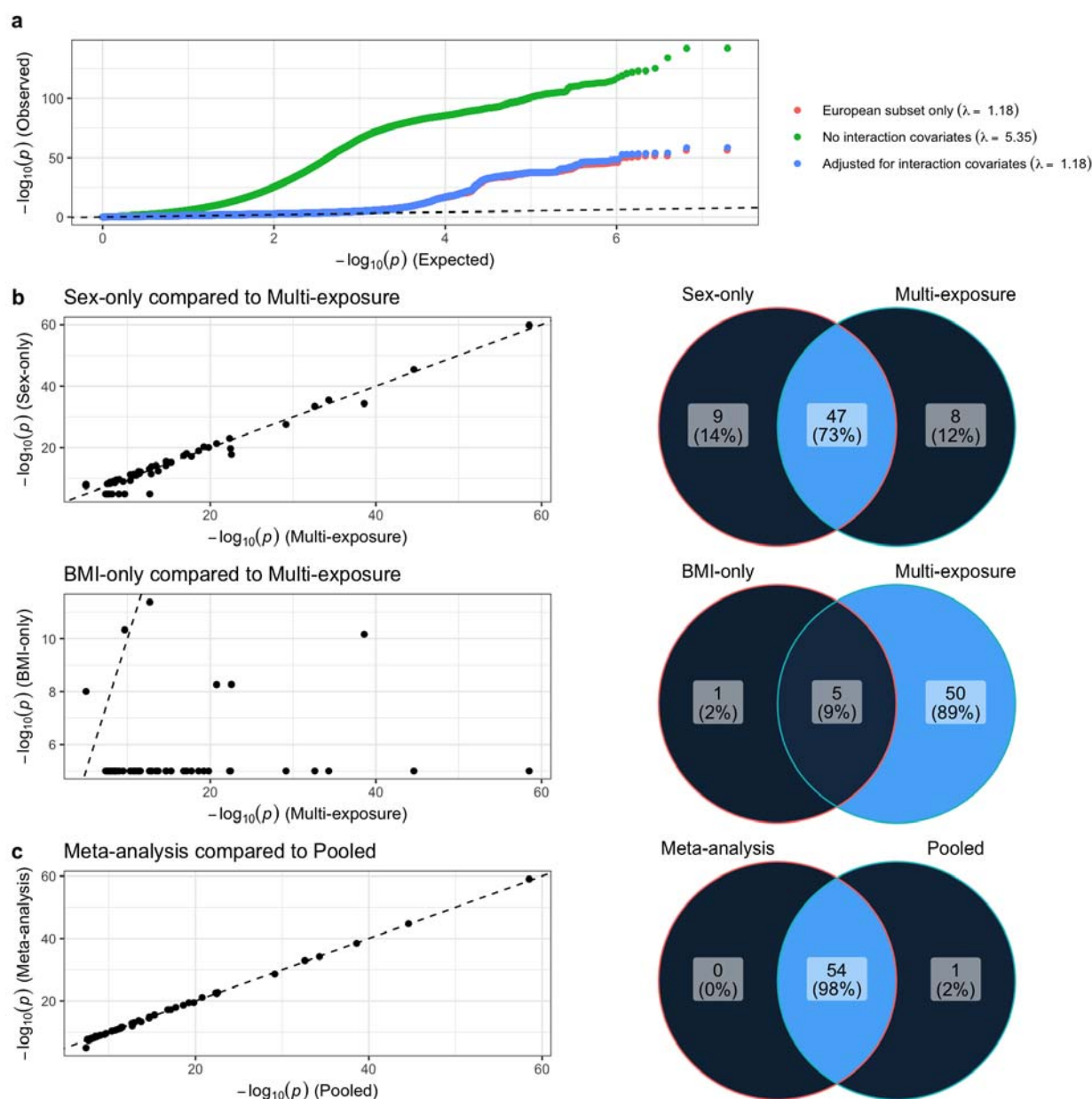
295 **Accounting for ancestry in pooled analysis of waist-hip ratio**

296 In order to test the functionality of REGEM and METAGEM on real datasets, we further explored
297 the expanded WHR GWIS model used for benchmarking. The primary analysis tested genetic
298 interactions with two exposures (sex and BMI) in a pooled dataset containing six ancestry
299 groups. Without additional adjustments, this pooled dataset produced highly inflated summary
300 statistics (genomic inflation $\lambda = 5.35$), but after inclusion of interaction covariates (gene-
301 ancestry and exposure-ancestry interaction terms), this inflation was reduced to a level identical
302 to that of a European ancestry-only analysis ($\lambda = 1.18$ for both; Figure 2a). This properly-
303 adjusted pooled analysis uncovered 55 independent loci using a standard genome-wide
304 significance threshold of 5×10^{-8} . Using REGEM to produce equivalent single-exposure
305 interaction tests (sex or BMI), we saw that the sex-only GWIS revealed a highly overlapping set
306 of loci (57 loci in total, 47 of which overlapped loci from the multi-exposure test), while the BMI-
307 only GWIS revealed many fewer (6 loci in total, 5 of which overlapped loci from the multi-
308 exposure test; Figure 2b).

309 Using METAGEM, we then conducted a meta-analysis of six ancestry-specific GWIS, finding 54
310 total loci, all of which overlapped loci from the primary ancestry-pooled analysis (Figure 2c).

311 This high concordance reinforces two conclusions. First, proper adjustment for interaction
 312 covariates can allow rigorous pooled-ancestry GWIS and avoid the need for stratification.
 313 Second, in situations where pooled analysis is not possible for logistical or analytical reasons,
 314 the ability to adjust for interaction covariates and possibly include multiple exposures in
 315 conducting GWIS meta-analysis can be critical for proper interpretation and control of inflation.

316



317

318 **Figure 2:** Results from multi-exposure, multi-ancestry GWIS for waist-hip ratio. a) Quantile-
319 Quantile plots display observed vs. expected p-values for selected analyses. b) Results from
320 REGEM-derived, single-exposure GWIS results for sex (top panel) and BMI (bottom panel).
321 Scatter plots compare p-values between single- and multi-exposure interaction tests and Venn
322 diagrams display the overlap in independent loci discovered using single- and multi-exposure
323 interaction tests. c) As in (b), but replacing REGEM-derived, single-exposure results with
324 METAGEM-derived, multi-ancestry meta-analysis results.

325

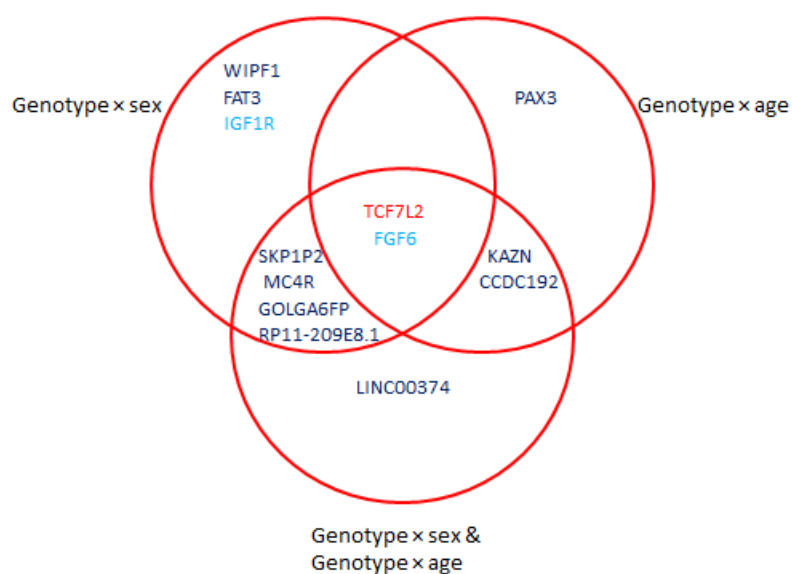
326 **Sex and age interaction effects on T2D in the ProDiGY dataset**

327 We performed a genome-wide, multi-exposure test of sex and age interactions affecting T2D
328 analysis in the ProDiGY dataset, separately in the youth (youth cases vs youth controls) and
329 adult (adult cases vs. adult controls) subsets. After cross-ancestry meta-analysis, we did not
330 detect any significant signals using the interaction test, but using the joint test found 8
331 independent loci passed the genome-wide significance threshold in the youth group (Table S2)
332 and 3 loci in the adult group (Table S3). Of the 8 loci in the youth group, two were known
333 associations, at *TCF7L2* ($p_{joint} = 1.30 \times 10^{-9}$) and *MC4R* ($p_{joint} = 9.22 \times 10^{-9}$). Only one, rs7903146
334 at *TCF7L2*, showed a significant effect in the marginal genetic effect test (excluding interaction
335 effects). Six of the 8 signals were not reported in previous T2D GWAS studies (as per the
336 Common Metabolic Disease Knowledge Portal). One variant, rs114578532, upstream of *FGF6*,
337 passed the genome-wide significance threshold in the marginal test ($p_{marginal} = 2.18 \times 10^{-8}$), but
338 not joint test ($p_{joint} = 7.25 \times 10^{-7}$). These signals, with the exception of *TCF7L2*, did not show
339 strong effects in the adult group analysis. In the adult cases vs. adult controls comparison, out
340 of three signals, two were known to be associated with T2D and also showed statistical
341 significance in the marginal test (rs35198068 at *TCF7L2* and rs2237892 at *KCNQ1*). The third

342 locus, with lead variant rs62287662 within an intron of *KCNAB1*, has not been previously
343 associated with T2D ($p_{joint} = 1.79 \times 10^{-8}$; $p_{interaction} = 6.27 \times 10^{-8}$). *KCNAB1* encodes a protein
344 involved in diverse functions including heart rate and insulin secretion. This locus did not show
345 meaningful association in the youth group analysis.

346 To evaluate the added value of multi-exposure analysis, we ran analogous single-exposure
347 meta-analyses, separately for sex and age. Of 8 multi-exposure signals in the youth group joint
348 test, we found that 5 reached significance in the sex-only analysis (plus 2 additional signals) and
349 3 in the age-only analysis (plus 1 additional signal) (Figure 3). In the adult group, 2 of 3 loci
350 were found in all three models, with the third found in both the multi-exposure and age-only
351 tests but not the sex-only test (Figure S6).

352



353

354 **Figure 3:** Results from multi-exposure GWIS for incident T2D in the ProDiGY youth cohort.
355 Venn diagram displays overlap between loci discovered at genome-wide significance using the
356 joint test of genetic and interaction effects ($p_{joint} = 5 \times 10^{-8}$), from each of: sex-only, age-only, and
357 multi-exposure (sex and age) analyses. Variants are labeled according to the closest gene, and

358 colors correspond to the test(s) in which significance was achieved: marginal genetic effect
359 (light blue), joint genetic effect (dark blue), or both joint and marginal genetic effects (red).

360

361 **DISCUSSION**

362 GEI studies are becoming increasingly challenging due to complex structured models involving
363 multiple interaction terms. Here we introduce two software programs, REGEM and METAGEM,
364 to enable further downstream analysis of such studies using only summary statistics. We show
365 that both programs are much more computationally efficient than the corresponding individual-
366 level data analyses and validate their results in comparison to existing software options.
367 Additionally, we demonstrate how REGEM and METAGEM can be applied to improve GEI
368 studies related to anthropometric traits in the UK Biobank and diabetes in the ProDiGY
369 resource.

370 REGEM is a powerful tool that exploits the GEM methodology to enable rapid estimation of
371 genome-wide summary statistics for any re-partition of a set of exposures and interaction
372 covariates. One potential application of REGEM is in sensitivity analyses, a common
373 epidemiological tool used to assess genetic confounding. In our analysis, we demonstrate that
374 proper adjustment for interaction covariates can significantly reduce highly inflated summary
375 statistics and increase the discovery of genetic loci. Such discoveries could have been missed
376 due to the computational expense of repeated genome-wide calculations on individual-level
377 data. While recent algorithms have enabled multi-threading capabilities,^{2,19} high-performance
378 computing, and cloud environments enable parallel genome-wide analysis, the pre-processing
379 time required to set up these environments may add additional computational time and financial
380 cost to individual-level genome-wide analysis. In our REGEM benchmark study, we show that
381 by avoiding repeated computation on individual-level data, a genome-wide re-analysis can be

382 completed within minutes, requiring minimal computation resources while still producing valid
383 summary statistic results. REGEM is lightweight and can be run on local machines, greatly
384 reducing runtime and cost compared to an equivalent individual-level data analysis.

385 Additionally, REGEM can also serve as a valuable pre-processing tool to harmonize summary
386 statistics results from multiple GEI studies for downstream meta-analysis. This is particularly
387 valuable in situations where different studies may test different combinations of exposure and
388 interaction covariates. For instance, one study may jointly test G x sex and G x BMI, while
389 another may only test G x sex. By applying REGEM to the first study, summary statistics from a
390 model testing only G x sex can be obtained without having to re-analyze individual-level
391 genotypes in that study. The resulting summary statistics from both studies can then be
392 combined for meta-analysis without sharing individual-level data. Traditionally, harmonizing data
393 from multiple GEI studies has been challenging due to lack of data sharing, privacy protection
394 issues and logistics in data transportation and storage of individual-level data.²⁰ Summary
395 statistics-based algorithms help bypass such restrictions to facilitate collaborative research, and
396 REGEM helps extend this family of tools to the GEI space.

397 Various GEI software programs can fit models with multiple interaction terms.^{2,19,21} However,
398 limited statistical power remains a challenge, requiring larger study cohorts, especially in
399 underrepresented populations.²² By enabling more flexible summary statistic-based meta-
400 analysis, METAGEM provides an alternative strategy towards increasing overall sample size
401 and statistical power for such analyses. For a single exposure meta-analysis without gene-by-
402 covariate interactions, existing software options, such as the popular METAL program, are
403 adequate. However, a nuanced set of considerations are required to determine whether it is
404 appropriate to include additional terms in meta-analysis, whether related to additional exposure
405 terms,¹⁰ gene-by-covariate interactions,¹² or genetic main effects.²² For multiple interaction

406 meta-analysis, METAGEM demonstrated efficient CPU time, though large memory space is
407 required for larger numbers of interaction terms and unique variants across studies.

408 By facilitating more comprehensive, genome-wide analyses and meta-analyses involving
409 interactions using only summary statistics, REGEM and METAGEM enable researchers to
410 maximize the value of genome-wide interaction studies while minimizing computational time. A
411 few limitations should be noted. Firstly, the GEM model corrects for standard covariates by
412 removing them from the genotype and interaction matrices in a single projection step. While this
413 approach improves computational performance of the primary GWIS considerably, it also takes
414 away the possibility of modifying covariate main effect adjustments in subsequent re-analysis.
415 Any such modification (e.g., seeking an interaction effect while completely removing a covariate
416 main effect from the statistical model) would require a new analysis using individual-level data.
417 Additionally, while REGEM has been shown to produce results that are consistent with those of
418 GEM, improper GEI analysis using GEM, particularly in the case of rare variants, can lead to
419 spurious summary statistics results, and may invalidate re-analysis results. Therefore,
420 researchers must ensure valid summary statistics (for example, well-controlled genomic
421 inflation) are generated from GEI methods before performing a re-analysis. In this vein, it is also
422 important that study-specific interaction terms to be meta-analyzed have equivalent
423 interpretations; for example, METAGEM cannot conduct valid meta-analysis when there are
424 discrepant study-specific variable coding choices in terms of exposure (and covariate)
425 centering.

426 In summary, we have introduced REGEM and METAGEM for further complex downstream
427 analysis of GEI studies. REGEM and METAGEM, along with our GEM tool for genome-wide
428 interaction analysis and corresponding workflows for reproducible and scalable deployment in
429 cloud computing environments, are publicly available at ([https://github.com/large-scale-gxe-](https://github.com/large-scale-gxe-methods)
430 [methods](https://github.com/large-scale-gxe-methods)). The suite of tools, including GEM, REGEM and METAGEM, provides key software

431 infrastructure for maximizing the utility of summary statistics from diverse and complex GEI
432 studies.

433 **Declaration of interests**

434 The authors declare no competing interests.

435 **Acknowledgements**

436 This research was conducted using the UK Biobank Resource under Application Numbers
437 27892 and 42646. This work was supported by NIH grant R01 HL145025. KEW was supported
438 by NIH grant K01 DK133637. ProDiGY acknowledgements and funding sources are included in
439 the Supplemental Material.

440 **Author contributions**

441 D.T.P. and H.C. developed the METAGEM and REGEM algorithms. D.T.P., H.C., and C.P.
442 implemented the METAGEM and REGEM software programs. D.T.P. and K.E.W. implemented
443 software programs as cloud workflows. D.T.P. and H.C. designed the benchmark simulation
444 study and carried out the analyses. K.E.W., L.C., and A.K.M. carried out the real-data analyses.
445 S.S., E.I., M.E.V., F.B., S.C., R.G.-K., J.D., C.P., and S.M.M. provided guidance and input
446 related to analysis of the ProDiGY dataset. K.E.W., D.T.P., H.C., and A.K.M. wrote the
447 manuscript. All authors critically read the manuscript.

448 **Web resources**

449 GEM, <https://github.com/large-scale-gxe-methods/GEM>

450 GEM Workflow, <https://github.com/large-scale-gxe-methods/gem-workflow>

451 METAGEM, <https://github.com/large-scale-gxe-methods/METAGEM>

452 METAGEM Workflow, <https://github.com/large-scale-gxe-methods/metagem-workflow>

453 REGEM, <https://github.com/large-scale-gxe-methods/REGEM>

454 REGEM Workflow, <https://github.com/large-scale-gxe-methods/regem-workflow>

455 **Data and code availability**

456 METAGEM and REGEM are both open source projects freely available at

457 <https://github.com/large-scale-gxe-methods/METAGEM> and [458 \[methods/REGEM\]\(https://github.com/large-scale-gxe-methods/REGEM\). Workflows for both programs are also available at \[459 \\[scale-gxe-methods/metagem-workflow\\]\\(https://github.com/large-scale-gxe-methods/metagem-workflow\\) and \\[460 \\\[workflow\\\]\\\(https://github.com/large-scale-gxe-methods/regem-workflow\\\).\\]\\(https://github.com/large-scale-gxe-methods/regem-</p></div><div data-bbox=\\)\]\(https://github.com/large-</p></div><div data-bbox=\)](https://github.com/large-scale-gxe-</p></div><div data-bbox=)

461 **References**

462 1. Werme, J., van der Sluis, S., Posthuma, D., and de Leeuw, C.A. (2021). Genome-wide
463 gene-environment interactions in neuroticism: an exploratory study across 25
464 environments. *Transl. Psychiatry* 11, 180. [10.1038/s41398-021-01288-9](https://doi.org/10.1038/s41398-021-01288-9).

465 2. Westerman, K.E., Pham, D.T., Hong, L., Chen, Y., Sevilla-González, M., Sung, Y.J., Sun,
466 Y.V., Morrison, A.C., Chen, H., and Manning, A.K. (2021). GEM: scalable and flexible
467 gene-environment interaction analysis in millions of samples. *Bioinformatics* 37, 3514–
468 3520. [10.1093/bioinformatics/btab223](https://doi.org/10.1093/bioinformatics/btab223).

469 3. Bi, W., Zhao, Z., Dey, R., Fritsche, L.G., Mukherjee, B., and Lee, S. (2019). A Fast and
470 Accurate Method for Genome-wide Scale Phenome-wide G × E Analysis and Its
471 Application to UK Biobank. *Am. J. Hum. Genet.* 105, 1182–1192.
472 [10.1016/j.ajhg.2019.10.008](https://doi.org/10.1016/j.ajhg.2019.10.008).

- 473 4. Gauderman, W.J., Zhang, P., Morrison, J.L., and Lewinger, J.P. (2013). Finding novel
474 genes by testing G × E interactions in a genome-wide association study. *Genet. Epidemiol.*
475 *37*, 603–613. [10.1002/gepi.21748](https://doi.org/10.1002/gepi.21748).
- 476 5. Kerin, M., and Marchini, J. (2020). Inferring Gene-by-Environment Interactions with a
477 Bayesian Whole-Genome Regression Model. *Am. J. Hum. Genet.* *107*, 698–713.
478 [10.1016/j.ajhg.2020.08.009](https://doi.org/10.1016/j.ajhg.2020.08.009).
- 479 6. Mbatchou, J., Barnard, L., Backman, J., Marcketta, A., Kosmicki, J.A., Ziyatdinov, A.,
480 Benner, C., O'Dushlaine, C., Barber, M., Boutkov, B., et al. (2021). Computationally
481 efficient whole-genome regression for quantitative and binary traits. *Nat. Genet.* *53*, 1097–
482 1103. [10.1038/s41588-021-00870-7](https://doi.org/10.1038/s41588-021-00870-7).
- 483 7. Zhong, W., Chhibber, A., Luo, L., Mehrotra, D.V., and Shen, J. (2023). A fast and powerful
484 linear mixed model approach for genotype-environment interaction tests in large-scale
485 GWAS. *Brief. Bioinform.* *24*. [10.1093/bib/bbac547](https://doi.org/10.1093/bib/bbac547).
- 486 8. Shin, J., and Lee, S.H. (2021). GxEsum: a novel approach to estimate the phenotypic
487 variance explained by genome-wide GxE interaction based on GWAS summary statistics
488 for biobank-scale data. *Genome Biol.* *22*, 183. [10.1186/s13059-021-02403-1](https://doi.org/10.1186/s13059-021-02403-1).
- 489 9. Westerman, K., Liu, Q., Liu, S., Parnell, L.D., Sebastiani, P., Jacques, P., DeMeo, D.L., and
490 Ordovás, J.M. (2020). A gene-diet interaction-based score predicts response to dietary fat
491 in the Women's Health Initiative. *Am. J. Clin. Nutr.* *111*, 893–902. [10.1093/ajcn/nqaa037](https://doi.org/10.1093/ajcn/nqaa037).
- 492 10. Kim, J., Ziyatdinov, A., Laville, V., Hu, F.B., Rimm, E., Kraft, P., and Aschard, H. (2019).
493 Joint Analysis of Multiple Interaction Parameters in Genetic Association Studies. *Genetics*
494 *211*, 483–494. [10.1534/genetics.118.301394](https://doi.org/10.1534/genetics.118.301394).

- 495 11. Moore, R., Casale, F.P., Jan Bonder, M., Horta, D., BIOS Consortium, Franke, L., Barroso,
496 I., and Stegle, O. (2019). A linear mixed-model approach to study multivariate gene-
497 environment interactions. *Nat. Genet.* 51, 180–186. 10.1038/s41588-018-0271-0.
- 498 12. Keller, M.C. (2014). Gene x environment interaction studies have not properly controlled for
499 potential confounders: the problem and the (simple) solution. *Biol. Psychiatry* 75, 18–24.
500 10.1016/j.biopsych.2013.09.006.
- 501 13. Pan-UKB team. <https://pan.ukbb.broadinstitute.org>. 2020.
- 502 14. Willer, C.J., Li, Y., and Abecasis, G.R. (2010). METAL: fast and efficient meta-analysis of
503 genomewide association scans. *Bioinformatics* 26, 2190–2191.
504 10.1093/bioinformatics/btq340.
- 505 15. Manning, A.K., Hivert, M.-F., Scott, R.A., Grimsby, J.L., Bouatia-Naji, N., Chen, H., Rybin,
506 D., Liu, C.-T., Bielak, L.F., Prokopenko, I., et al. (2012). A genome-wide approach
507 accounting for body mass index identifies genetic variants influencing fasting glycaemic traits
508 and insulin resistance. *Nat. Genet.* 44, 659–669. 10.1038/ng.2274.
- 509 16. TODAY Study Group, Zeitler, P., Epstein, L., Grey, M., Hirst, K., Kaufman, F., Tamborlane,
510 W., and Wilfley, D. (2007). Treatment options for type 2 diabetes in adolescents and youth:
511 a study of the comparative efficacy of metformin alone or in combination with rosiglitazone
512 or lifestyle intervention in adolescents with type 2 diabetes. *Pediatr. Diabetes* 8, 74–87.
513 10.1111/j.1399-5448.2007.00237.x.
- 514 17. SEARCH Study Group (2004). SEARCH for Diabetes in Youth: a multicenter study of the
515 prevalence, incidence and classification of diabetes mellitus in youth. *Control. Clin. Trials*
516 25, 458–471. 10.1016/j.cct.2004.08.002.

- 517 18. Srinivasan, S., Chen, L., Todd, J., Divers, J., Gidding, S., Chernausek, S., Gubitosi-Klug,
518 R.A., Kelsey, M.M., Shah, R., Black, M.H., et al. (2021). The First Genome-Wide
519 Association Study for Type 2 Diabetes in Youth: The Progress in Diabetes Genetics in
520 Youth (ProDiGY) Consortium. *Diabetes* 70, 996–1005. [10.2337/db20-0443](https://doi.org/10.2337/db20-0443).
- 521 19. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015).
522 Second-generation PLINK: rising to the challenge of larger and richer datasets.
523 *Gigascience* 4, 7. [10.1186/s13742-015-0047-8](https://doi.org/10.1186/s13742-015-0047-8).
- 524 20. Reales, G., and Wallace, C. (2023). Sharing GWAS summary statistics results in more
525 citations. *Commun Biol* 6, 116. [10.1038/s42003-023-04497-8](https://doi.org/10.1038/s42003-023-04497-8).
- 526 21. Lin, D.-Y., Tao, R., Kalsbeek, W.D., Zeng, D., Gonzalez, F., 2nd, Fernández-Rhodes, L.,
527 Graff, M., Koch, G.G., North, K.E., and Heiss, G. (2014). Genetic association analysis
528 under complex survey sampling: the Hispanic Community Health Study/Study of Latinos.
529 *Am. J. Hum. Genet.* 95, 675–688. [10.1016/j.ajhg.2014.11.005](https://doi.org/10.1016/j.ajhg.2014.11.005).
- 530 22. Laville, V., Majarian, T., Sung, Y.J., Schwander, K., Feitosa, M.F., Chasman, D.I., Bentley,
531 A.R., Rotimi, C.N., Cupples, L.A., de Vries, P.S., et al. (2022). Gene-lifestyle interactions in
532 the genomics of human complex traits. *Eur. J. Hum. Genet.* 30, 730–739. [10.1038/s41431-](https://doi.org/10.1038/s41431-022-01045-6)
533 [022-01045-6](https://doi.org/10.1038/s41431-022-01045-6).