

Genetics and population analysis

Re-analysis and meta-analysis of summary statistics from gene–environment interaction studies

Duy T. Pham ^{1,†}, Kenneth E. Westerman^{2,3,4,†}, Cong Pan¹, Ling Chen^{2,3}, Shylaja Srinivasan ⁵, Elvira Isganaitis ⁶, Mary Ellen Vajravelu⁷, Fida Bacha ⁸, Steve Chernausek⁹, Rose Gubitosi-Klug¹⁰, Jasmin Divers¹¹, Catherine Pihoker¹², Santica M. Marcovina¹³, Alisa K. Manning^{2,3,4}, Han Chen ^{1,*}

¹Human Genetics Center, Department of Epidemiology, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX 77030, United States

²Department of Medicine, Clinical and Translational Epidemiology Unit, Mongan Institute, Massachusetts General Hospital, Boston, MA 02114, United States

³Metabolism Program, Broad Institute of MIT and Harvard, Cambridge, MA 02142, United States

⁴Department of Medicine, Harvard Medical School, Boston, MA 02115, United States

⁵Department of Pediatrics, University of California, San Francisco, CA 94158, United States

⁶Research Division, Joslin Diabetes Center, Boston, MA 02115, United States

⁷Department of Pediatrics, University of Pittsburgh School of Medicine, Pittsburgh, PA 15224, United States

⁸Department of Pediatrics, Baylor College of Medicine, Houston, TX 77030, United States

⁹Department of Pediatrics, The University of Oklahoma College of Medicine, Oklahoma City, OK 73117, United States

¹⁰Department of Pediatrics, Case Western Reserve University, Cleveland, OH 44106, United States

¹¹Department of Foundations of Medicine, New York University, New York, NY 10016, United States

¹²Department of Pediatrics, University of Washington School of Medicine, Seattle, WA 98105, United States

¹³Northwest Lipid Metabolism and Diabetes Research Laboratories, Department of Medicine, University of Washington, Seattle, WA 98105, United States

*Corresponding author. Human Genetics Center, Department of Epidemiology, School of Public Health, The University of Texas Health Science Center at Houston, 1200 Pressler St, RAS E-517, Houston, TX 77030, United States. E-mail: han.chen.2@uth.tmc.edu

[†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Russell Schwartz

Abstract

Motivation: Summary statistics from genome-wide association studies enable many valuable downstream analyses that are more efficient than individual-level data analysis while also reducing privacy concerns. As growing sample sizes enable better-powered analysis of gene–environment interactions, there is a need for gene–environment interaction-specific methods that manipulate and use summary statistics.

Results: We introduce two tools to facilitate such analysis, with a focus on statistical models containing multiple gene–exposure and/or gene–covariate interaction terms. REGEM (RE-analysis of GEM summary statistics) uses summary statistics from a single, multi-exposure genome-wide interaction study to derive analogous sets of summary statistics with arbitrary sets of exposures and interaction covariate adjustments. METAGEM (META-analysis of GEM summary statistics) extends current fixed-effects meta-analysis models to incorporate multiple exposures from multiple studies. We demonstrate the value and efficiency of these tools by exploring alternative methods of accounting for ancestry-related population stratification in genome-wide interaction study in the UK Biobank as well as by conducting a multi-exposure genome-wide interaction study meta-analysis in cohorts from the diabetes-focused ProDIGY consortium. These programs help to maximize the value of summary statistics from diverse and complex gene–environment interaction studies.

Availability and implementation: REGEM and METAGEM are open-source projects freely available at <https://github.com/large-scale-gxe-methods/REGEM> and <https://github.com/large-scale-gxe-methods/METAGEM>.

1 Introduction

Gene–environment interaction (GEI) analysis is a key tool for understanding genetic impacts on human traits, with the potential to account for additional heritability, explain differences in genetic effects across populations, and support personalized lifestyle and therapeutic decisions. Historically, GEI studies have taken a hypothesis-driven approach, but larger cohorts (Werme *et al.* 2021), and new software programs have provided the necessary

statistical power and computational efficiency to study GEIs genome-wide (Gauderman *et al.* 2013, Bi *et al.* 2019, Kerin and Marchini 2020, Mbatchou *et al.* 2021, Westerman *et al.* 2021, Zhong *et al.* 2023). These genome-wide interaction studies (GWISs) generate summary statistics, or variant-level regression results, which have substantial value beyond locus mapping. For example, summary statistics allow for heritability analysis (Shin and Lee 2021), enrichment testing (Werme *et al.* 2021),

and genome-wide polygenic score generation (Westerman *et al.* 2020, Werme *et al.* 2021).

GEI analysis and interpretation are complicated by the densely correlated set of possible exposures that may interact with genotypes to influence human traits (the “exposome,” defined here as including demographic and physiologic traits). Two modeling implications are particularly pertinent. First, multi-exposure GEI analysis can increase statistical power by jointly testing genetic interactions with multiple exposures (Moore *et al.* 2019, Kerin and Marchini 2020). This strategy can pool signals across distinct exposures (e.g. smoking status and pollution exposure for lung function) or incorporate multiple definitions of a single-exposure category (e.g. current smoking status and pack-years of smoking). Second, proper control of confounding for GEI interaction terms requires adjustment for not just the main effects of covariates, but also their genetic interactions (Keller 2014). Inclusion of these “interaction covariates” is thus necessary to produce interpretable summary statistics.

Rigorous GEI analysis carries complexities stemming from its place at the center of traditional and genetic epidemiology. Sensitivity analyses, while commonplace in traditional epidemiology, are computationally burdensome when conducted across millions of variants genome-wide. Meanwhile, well-established meta-analysis procedures for genome-wide association study (GWAS) summary statistics become more difficult in the context of multi-exposure GEI models. Software programs do not yet exist to perform efficient meta-analysis in the context of these complex analytical designs.

We introduce methods and associated software programs to advance the field of genome-wide GEI analysis based on summary statistics. While the statistical results are general, the associated software implementations build on the results from our previously described software program for efficient GWIS, GEM (Gene-Environment interaction analysis for Millions of samples) (Westerman *et al.* 2021). Exploiting the redundancy of statistical estimates across related GEI models, we introduce the REGEM (RE-analysis of GEM summary statistics) program to derive genome-wide summary statistics corresponding to arbitrary multi-exposure and interaction covariate adjustments based on results from a single, multi-exposure GWIS. Expanding current fixed-effect meta-analysis models, we further introduce the METAGEM (META-analysis of GEM summary statistics) program to conduct efficient meta-analysis of GEI effects under complex GEI analysis models. We demonstrate the value and efficiency of these tools by exploring alternative methods of accounting for population stratification in GWIS in the UK Biobank as well as by conducting a multi-exposure GWIS meta-analysis in cohorts from the Progress in Diabetes Genetics in Youth (ProDiGY) consortium.

2 Materials and methods

2.1 GEM method

We developed two C++ software programs that use summary statistics from GEI studies. REGEM requires output from a single GEI study, while METAGEM requires output from multiple GEI studies. Both programs are designed for easy integration with output from GEM, but can use summary statistics from any GEI testing program that can output effect estimate covariances. For a single-variant test of N

unrelated individuals, GEM considers the generalized linear model:

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

for individual i , where $\mu_i = E(Y_i|X_i, G_i)$ is the conditional mean of the phenotype Y_i given p covariates X_i (including the intercept), and the genotype G_i for a single genetic variant. The interaction terms C_i and S_i are the products of G_i and c covariates and q exposures (which are disjoint subsets of X_i), respectively (Westerman *et al.* 2021). Let $Y = (Y_1 \ Y_2 \ \dots \ Y_N)^T$ be a length N vector of phenotypes, $X = (X_1^T \ X_2^T \ \dots \ X_N^T)^T$ be an $N \times p$ matrix of p covariates, $G = (G_1 \ G_2 \ \dots \ G_N)^T$ be a length N 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 gene-covariate interaction terms, and $S = (S_1^T \ S_2^T \ \dots \ S_N^T)^T$ be an $N \times q$ matrix of q gene-environment (exposure) interaction terms, we can fit a null model without any genetic effects $g(\mu_i) = X_i\beta_X$

and get a length N residual vector r . Let $\tilde{G} = G - X(X^T W X)^{-1} X^T W G$, $\tilde{C} = C - X(X^T W X)^{-1} X^T W C$, and $\tilde{S} = S - X(X^T W X)^{-1} X^T W S$ be covariate X adjusted G , C , and S , respectively, where W is a diagonal weight matrix with elements $\hat{\mu}_i(1 - \hat{\mu}_i)$ for logistic regressions ($\hat{\mu}_i$ are fitted probabilities of $Y_i = 1$ from the null model) and an identity matrix for linear regressions, GEM computes a length $(1 + c + q)$ score vector ($c \geq 0$) $U = (\tilde{G} \ \tilde{C} \ \tilde{S})^T r$, and $(1 + c + q) \times (1 + c + q)$ 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 diagonal matrix of squared residuals.

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

2.2 REGEM method

Given the full summary statistics output from GEM (version 1.3 and later), the score vector U and matrices V and Ω , can be reconstructed without access to individual-level data. Utilizing $\hat{\phi}$ and the matrices $\text{Cov}(\hat{\beta}_{G,C,S})$ and $\text{Cov}_R(\hat{\beta}_{G,C,S})$ described above, it follows that $V = \hat{\phi} \text{Cov}^{-1}(\hat{\beta}_{G,C,S})$ and $\Omega = V \text{Cov}_R(\hat{\beta}_{G,C,S}) V$. The score vector can then be recomputed as $U = V \hat{\beta}_{G,C,S}$.

REGEM supports two scenarios for re-analysis of a single GEI study. The first scenario involves the exclusion of one or more gene-covariate or GEI terms from the original model. This is achieved by filtering U to exclude the specified gene-covariate or GEI terms, resulting in the modified score vector

\hat{U} . Subsequently, the matrices V and Ω are reduced to exclude the corresponding rows and columns of the specified gene-covariate or GEI terms, denoted \hat{V} and $\hat{\Omega}$. The GEM method can then be applied to \hat{U} , \hat{V} , and $\hat{\Omega}$ to obtain new summary statistics. In the second scenario, re-analysis can be performed by conditioning on one or more GEI terms in the original GEM analysis as gene-covariate interactions or testing one or more gene-covariate interaction terms in the original GEM analysis as GEI terms of interest. In either case, the ordering of U is rearranged, denoted as \tilde{U} , to incorporate the original GEI terms into C or the original gene-covariate interaction terms into S . The rows and columns of the matrices V and Ω are also reordered and denoted as \tilde{V} and $\tilde{\Omega}$. The GEM method follows for \tilde{U} , \tilde{V} , and $\tilde{\Omega}$. Both scenarios can be applied simultaneously.

2.3 METAGEM method

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

2.4 REGEM comparison and benchmark

To demonstrate the computational benefits of REGEM, we test and compare four variations of the waist-hip ratio (WHR) model originally described by [Westerman et al. \(2021\)](#). The original model is defined as follows (excluding the array covariate and PC6-PC10):

$$\text{WHR} \sim G + \text{sex} + \text{age} + \text{age}^2 + \text{BMI} + \text{PC1} + \dots + \text{PC5} + G \times \text{sex} + G \times \text{BMI}, \quad (2)$$

where WHR is the phenotype, sex is the primary exposure of interest, BMI is the interaction covariate, and age, age², and PC1-PC5 are the covariates. Here, we retrieved PCs calculated as part of the Pan-UKBB project ([Pan-UKBB team 2020](#); <https://pan.ukbb.broadinstitute.org>). All terms in the model were centered. First, we performed a genome-wide analysis of the original model using GEM (version 1.5) using 362 449 unrelated European-ancestry participants, and filtered variants with minor allele frequency (MAF) <0.001, leaving 16 539 280 variants for re-analysis. Next, we derived associated genome-wide summary statistics corresponding to variations of the original model using REGEM, comparing their results and runtimes to simply re-running that same model genome-wide using GEM. [Supplementary Table S1](#) summarizes the variations of the original models, including the original model. These variations involve the joint testing of $G \times \text{sex}$ and $G \times \text{BMI}$ (M1), testing for $G \times \text{BMI}$ while adjusting for $G \times \text{sex}$ (M2), testing for $G \times \text{sex}$ while removing the $G \times \text{BMI}$ term (M3), and testing for $G \times \text{BMI}$ while removing the $G \times \text{sex}$ term (M4). All analyses were performed on the DNAnexus platform using the *mem1_ssd1_v2_x16* instance type, and we reported the runtime and memory usage of each run. The GEM and REGEM summary statistic comparisons were visualized using the scattermore and ggplot2 R packages.

2.5 METAGEM comparison and benchmark

To evaluate the computational efficiency of METAGEM, we conducted a simulation study using phenotype and genotype data from the Pan-UKBB ([Pan-UKBB team 2020](#)). We randomly split the phenotype data, which comprised 362 449 samples, into 11 datasets: one with 100 000 samples, two with 50 000 samples, seven with 10 000 samples, and one with 92 449 samples. For each dataset, we conducted a genome-wide gene-sex interaction test and filtered out variants with a MAF < 0.001, resulting in 15.46 to 16.85 million variants per dataset, and a total of 17 993 341 unique variants across all datasets. We then performed a gene-sex interaction meta-analysis using METAGEM and the METAL software (version 2010-02-08) ([Willer et al. 2010](#)) with the joint meta-analysis patch ([Manning et al. 2012](#)) and compared the results. Additionally, we conducted a genome-wide joint gene-sex and gene-BMI interaction test for each dataset and performed a meta-analysis using METAGEM to evaluate its performance in the presence of multiple interaction terms. All analyses were conducted on the DNAnexus platform using a single core and the *mem1_ssd1_v2_x16* instance type. We reported the CPU time and memory usage for each analysis. We used the scattermore and ggplot2 R packages to visualize the comparison of summary statistics between METAGEM and METAL.

2.6 Multi-exposure interactions influencing WHR in the UK Biobank

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

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

2.7 ProDiGY dataset

ProDiGY is a multi-ethnic resource including three studies: Treatment Options for Type 2 Diabetes in Adolescents and Youth ([TODAY Study Group et al. 2007](#)), SEARCH for Diabetes in Youth ([SEARCH Study Group 2004](#)), and T2D-

GENES. In total, the dataset contains 2820 youth and 4858 adult cases with T2D, and 656 diabetes-free youth and 4934 adult controls after removing individuals with maturity-onset diabetes of the young and type I diabetes. Samples were genotyped on the Infinium GWAS array by the Genetic Analysis Platform at the Broad Institute of MIT and Harvard. Details on quality control procedures for the genotype data have been previously described (Srinivasan *et al.* 2021). Genotype data were imputed on the TOPMed Imputation Server using the TOPMed v2 reference panel. Variants passing an imputation quality threshold (R^2) of 0.5 were retained for analysis. Genetic ancestry groups were assigned to ProDiGY samples based on genetic principal components analysis after merging with the 1000 Genomes dataset.

2.8 Application multi-interaction to T2D in ProDiGY

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

$$T2D \sim G + \text{sex} + \text{age} + PC1 + \dots + PC10 + G \times \text{sex} + G \times \text{age}. \quad (3)$$

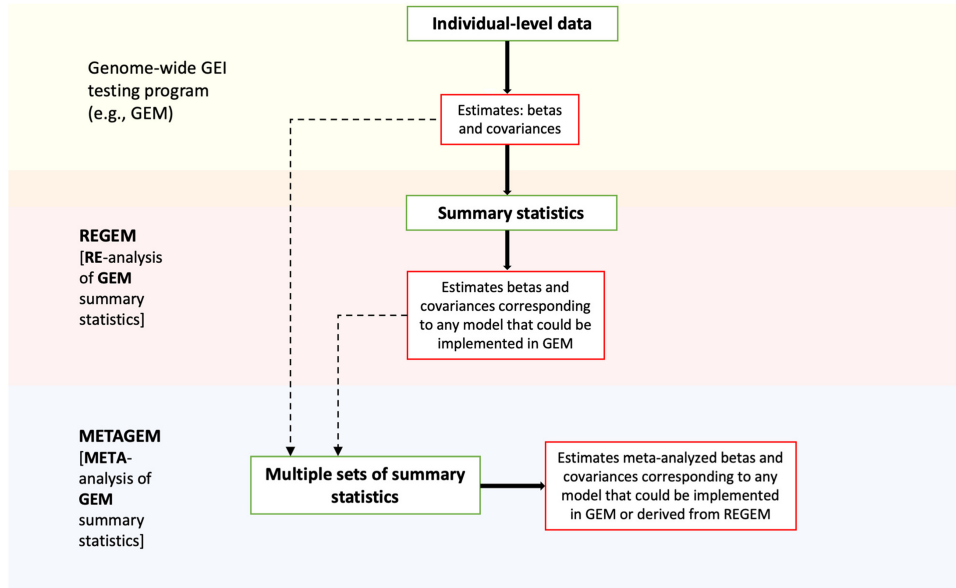
Using the full output from GEM, we performed cross-ancestry meta-analysis using METAGEM in both youth group and adult group analyses. We also conducted equivalent single-exposure GWIS with sex and age separately for comparison with the multi-exposure scan. Meta-analysis for these single-exposure tests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

3 Results

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

3.1 REGEM computational performance

We compared results obtained from genome-wide interactions tests using the REGEM and GEM methods across four distinct GEI models. The benchmark results are presented in



Software	Exposures	Interaction Covariates	Summary Statistics ignored	Associated (implicit) regression model	Interaction test H_0
GEM (or other)	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	$\beta_{gE_1}^* = 0$

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

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

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).

Figure 1. Large-scale GEI methods software suite and connections in the context of an analysis workflow. GEM (previously published) conducts GWISs for single datasets. Given multi-exposure summary statistics from GEM (version 1.3 or later) or an alternative GEI testing program, REGEM can estimate genome-wide summary statistics from an associated model that re-partitions any subset of exposures into interaction covariates and simple main effect adjustments without interaction. Given multiple sets of summary statistics from GEM or an alternative GEI testing program and/or REGEM, METAGEM conducts meta-analysis for any number of jointly tested exposures and interaction covariates.

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

Benchmark	GEM				REGEM			
	M1	M2	M3	M4	M1	M2	M3	M4
CPU time (min)	13 972.17	13 618.44	10 959.33	10 994.26	5.22	5.20	4.43	4.06
Memory (MB)	2325.37	2342.48	2188.14	2188.14	13.66	13.64	11.43	11.63

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

Benchmark	METAL	METAGEM	
	1 – exposure	1 – exposure	2 – exposures
CPU time (min)	16.38	14.38	19.55
Memory (GB)	7.10	6.11	6.96

Table 1. For each model, REGEM completed a genome-wide run in <6 min, while GEM required several CPU days to achieve the same outcome. Additionally, re-analyses for multiple interactions (M1 and M2) using REGEM took only about a minute of additional CPU time compared to single-exposure re-analyses (M3 and M4). Overall, REGEM saved considerable time, ranging from hours to days of computation time. Moreover, the memory requirements for REGEM were minimal, primarily depending on the number of GEI terms, which are usually small. Finally, the effect and variance estimates from REGEM were consistent with those obtained from GEM for each of the four models (M1–M4) as shown in [Supplementary Figs S1–S4](#).

3.2 METAGEM computational performance

Genome-wide meta-analysis runs of ~17.99 million variants, derived from 11 simulated UKB datasets, were carried out using the METAGEM and METAL methods with a single core. [Table 2](#) summarizes the CPU time and memory usage of the runs. For a single-exposure meta-analysis, METAGEM showed a modest improvement in performance compared to METAL, completing the run ~2 min faster and using ~1 GB less memory. We note that METAGEM meta-analyzed all 17 993 341 variants, while METAL skipped 25 670 multi-allelic variants that contained duplicate variant identifiers. However, the impact of the skipped variants on the benchmark results was negligible. Model-based and robust meta-analysis results from METAGEM and METAL are compared in [Supplementary Fig. S5](#). As expected, the summary statistics and joint *P*-values were consistent between the two methods. To test the performance of METAGEM in conducting meta-analysis with multiple interactions, we performed genome-wide joint meta-analysis with gene–sex and gene–BMI as the interactions using METAGEM. As shown in [Table 2](#), METAGEM efficiently completed the run in an additional ~6 min of CPU time and <1 GB of additional memory compared to the single-exposure meta-analysis.

3.3 Accounting for ancestry in pooled analysis of WHR

In order to test the functionality of REGEM and METAGEM on real datasets, we further explored the expanded WHR GWIS model used for benchmarking. The primary analysis tested

genetic interactions with two exposures (sex and BMI) in a pooled dataset containing six ancestry groups. Without additional adjustments, this pooled dataset produced highly inflated summary statistics (genomic inflation $\lambda = 5.35$), but after inclusion of interaction covariates (gene–ancestry and exposure–ancestry interaction terms), this inflation was reduced to a level identical to that of a European ancestry-only analysis ($\lambda = 1.18$ for both; [Fig. 2a](#)). This λ value is reasonable for a highly polygenic trait with known sex dimorphism, comparable to our prior observations, and consistent with expectations for analogous main effect GWAS after accounting for differences in statistical power ([Westerman et al. 2021](#)). This properly-adjusted pooled analysis uncovered 55 independent loci using a standard genome-wide significance threshold of 5×10^{-8} . Using REGEM to produce equivalent single-exposure interaction tests (sex or BMI), we saw that the sex-only GWIS revealed a highly overlapping set of loci (57 loci in total, 47 of which overlapped loci from the multi-exposure test), while the BMI-only GWIS revealed many fewer (six loci in total, five of which overlapped loci from the multi-exposure test; [Fig. 2b](#)).

Using METAGEM, we then conducted a meta-analysis of six ancestry-specific GWIS, finding 54 total loci, all of which overlapped loci from the primary ancestry-pooled analysis ([Fig. 2c](#)). This high concordance reinforces two conclusions. First, proper adjustment for interaction covariates can allow rigorous pooled-ancestry GWIS and avoid the need for stratification. Second, in situations where pooled analysis is not possible for logistical or analytical reasons, the ability to adjust for interaction covariates and possibly include multiple exposures in conducting GWIS meta-analysis can be critical for proper interpretation and control of inflation.

3.4 Sex and age interaction effects on T2D in the ProDiGY dataset

We performed a genome-wide, multi-exposure test of sex and age interactions affecting T2D analysis in the ProDiGY dataset, separately in the youth (youth cases versus youth controls) and adult (adult cases versus adult controls) subsets. After cross-ancestry meta-analysis, we did not detect any significant signals using the interaction test, but using the joint test found eight independent loci passed the genome-wide significance threshold in the youth group ([Supplementary Table S2](#)) and three loci in the adult group ([Supplementary Table S3](#)). Of the eight loci in the youth group, two were known associations, at *TCF7L2* ($p_{\text{joint}} = 1.30 \times 10^{-9}$) and *MC4R* ($p_{\text{joint}} = 9.22 \times 10^{-9}$). Only one, rs7903146 at *TCF7L2*, showed a significant effect in the marginal genetic effect test (excluding interaction effects). Six of the eight signals were not reported in previous T2D GWAS studies (as per the Common Metabolic Disease Knowledge Portal). One variant, rs114578532, upstream of *FGF6*, passed the genome-wide significance threshold in the marginal test ($p_{\text{marginal}} = 2.18 \times 10^{-8}$), but not joint test ($p_{\text{joint}} = 7.25 \times 10^{-7}$). These

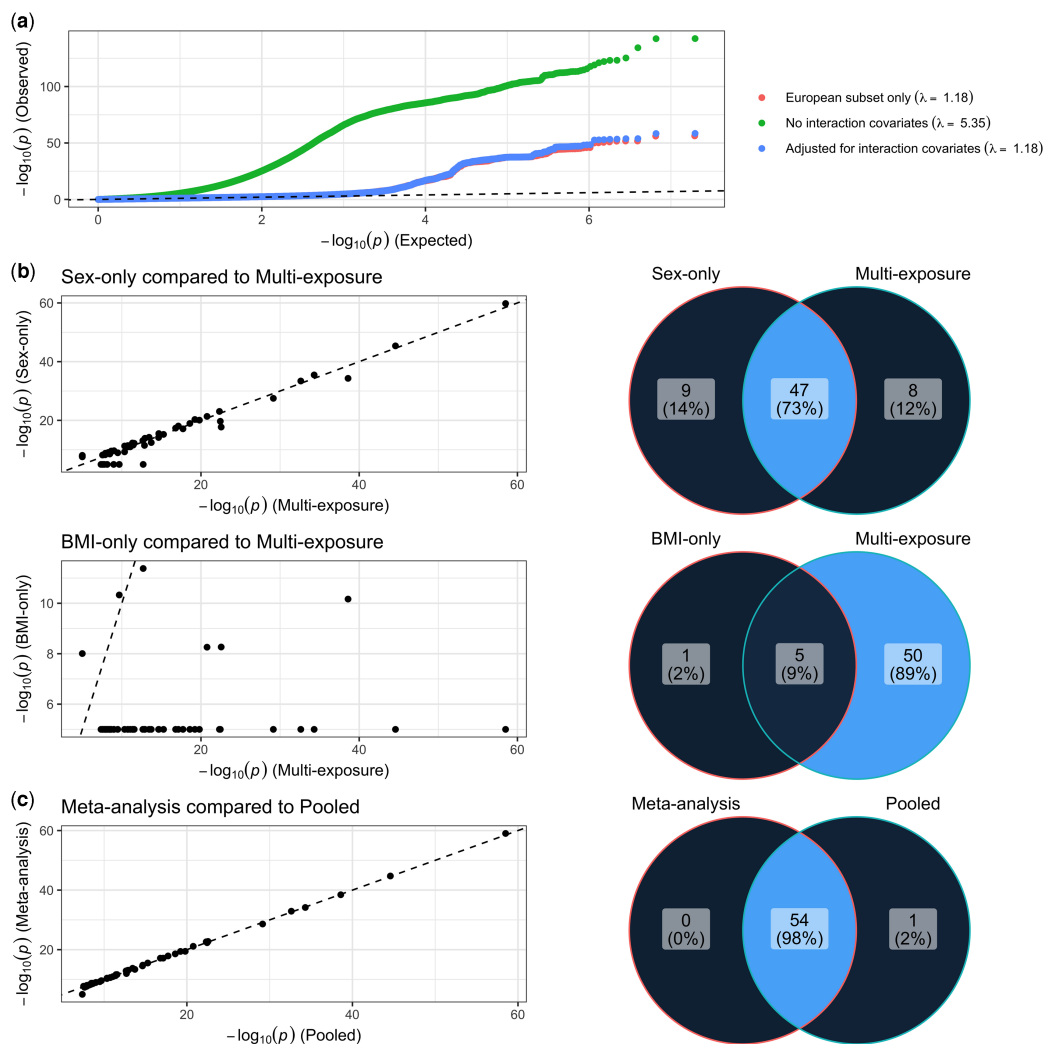


Figure 2. Results from multi-exposure, multi-ancestry GWIS for WHR. (a) Quantile–quantile plots display observed versus expected P -values for selected analyses. (b) Results from REGEM-derived, single-exposure GWIS results for sex (top panel) and BMI (bottom panel). Scatter plots compare P -values between single- and multi-exposure interaction tests and Venn diagrams display the overlap in independent loci discovered using single- and multi-exposure interaction tests. (c) As in (b), but replacing REGEM-derived, single-exposure results with METAGEM-derived, multi-ancestry meta-analysis results.

signals, with the exception of *TCF7L2*, did not show strong effects in the adult group analysis. In the adult cases versus adult controls comparison, out of three signals, two were known to be associated with T2D and also showed statistical significance in the marginal test (rs35198068 at *TCF7L2* and rs2237892 at *KCNQ1*). The third locus, with lead variant rs62287662 within an intron of *KCNAB1*, has not been previously associated with T2D ($p_{\text{joint}} = 1.79 \times 10^{-8}$; $p_{\text{interaction}} = 6.27 \times 10^{-8}$). *KCNAB1* encodes a protein involved in diverse functions including heart rate and insulin secretion. This locus did not show meaningful association in the youth group analysis.

To evaluate the added value of multi-exposure analysis, we ran analogous single-exposure meta-analyses, separately for sex and age. Of eight multi-exposure signals in the youth group joint test, we found that five reached significance in the sex-only analysis (plus two additional signals) and three in the age-only analysis (plus one additional signal) (Fig. 3). In the adult group, two of three loci were found in all three models, with the third found in both the multi-exposure and age-only tests but not the sex-only test (Supplementary Fig. S6).

4 Discussion

GEI studies are becoming increasingly challenging due to complex structured models involving multiple interaction terms. Here, we introduce two software programs, REGEM and METAGEM, to enable further downstream analysis of such studies using only summary statistics. They integrate easily with the GEM program, but are designed to allow for input from any alternative GEI testing program that can produce sufficiently detailed summary statistics. We show that both programs are much more computationally efficient than the corresponding individual-level data analyses and validate their results in comparison to existing software options. Additionally, we demonstrate how REGEM and METAGEM can be applied to improve GEI studies related to anthropometric traits in the UK Biobank and diabetes in the ProDiGY resource.

REGEM is a powerful tool that exploits the GEM methodology to enable rapid estimation of genome-wide summary statistics for any re-partition of a set of exposures and interaction covariates. One potential application of REGEM is in sensitivity analyses, a common epidemiological tool used to assess genetic confounding. In our analysis, we demonstrate that proper

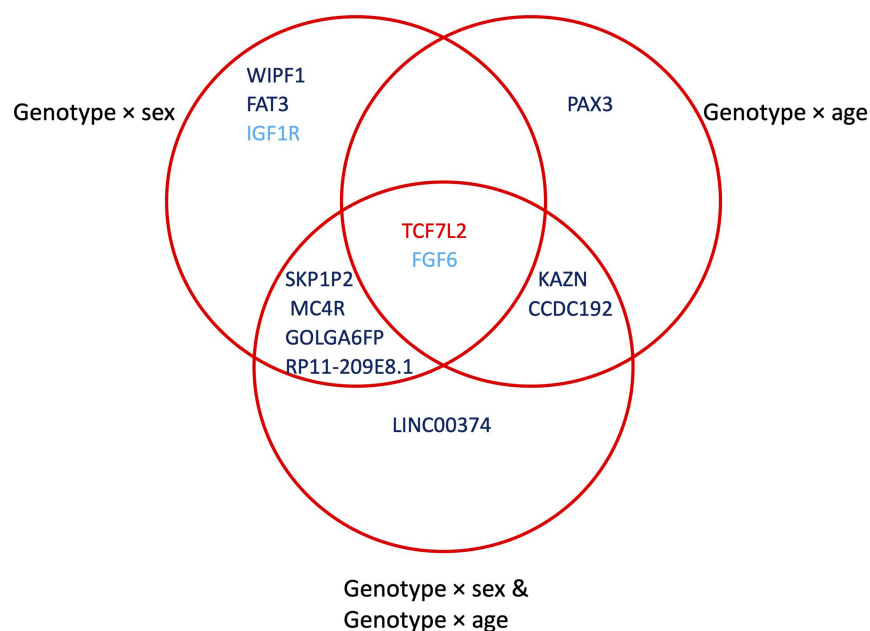


Figure 3. Results from multi-exposure GWIS for incident T2D in the ProDiGY youth cohort. Venn diagram displays overlap between loci discovered at genome-wide significance using the joint test of genetic and interaction effects ($p_{\text{joint}} = 5 \times 10^{-8}$), from each of: sex-only, age-only, and multi-exposure (sex and age) analyses. Variants are labeled according to the closest gene, and colors correspond to the test(s) in which significance was achieved: marginal genetic effect (light blue), joint genetic effect (dark blue), or both joint and marginal genetic effects (red).

adjustment for interaction covariates can significantly reduce highly inflated summary statistics and increase the discovery of genetic loci. Such discoveries could have been missed due to the computational expense of repeated genome-wide calculations on individual-level data. While recent algorithms have enabled multi-threading capabilities (Chang *et al.* 2015, Westerman *et al.* 2021), high-performance computing, and cloud environments enable parallel genome-wide analysis, the pre-processing time required to set up these environments may add additional computational time and financial cost to individual-level genome-wide analysis. In our REGEM benchmark study, we show that by avoiding repeated computation on individual-level data, a genome-wide re-analysis can be completed within minutes, requiring minimal computation resources while still producing valid summary statistic results. REGEM is lightweight and can be run on local machines, greatly reducing runtime and cost compared to an equivalent individual-level data analysis.

Additionally, REGEM can also serve as a valuable pre-processing tool to harmonize summary statistics results from multiple GEI studies for downstream meta-analysis. This is particularly valuable in situations where different studies may test different combinations of exposure and interaction covariates. For instance, one study may jointly test $G \times \text{sex}$ and $G \times \text{BMI}$, while another may only test $G \times \text{sex}$. By applying REGEM to the first study, summary statistics from a model testing only $G \times \text{sex}$ can be obtained without having to re-analyze individual-level genotypes in that study. The resulting summary statistics from both studies can then be combined for meta-analysis without sharing individual-level data. Traditionally, harmonizing data from multiple GEI studies has been challenging due to lack of data sharing, privacy protection issues and logistics in data transportation and storage of individual-level data (Reales and Wallace 2023). Summary statistics-based algorithms help bypass such restrictions to facilitate collaborative research, and REGEM helps extend this family of tools to the GEI space.

Various GEI software programs can fit models with multiple interaction terms (Lin *et al.* 2014, Chang *et al.* 2015, Westerman *et al.* 2021). However, limited statistical power remains a challenge, requiring larger study cohorts, especially in underrepresented populations (Laville *et al.* 2022). By enabling more flexible summary statistic-based meta-analysis, METAGEM provides an alternative strategy toward increasing overall sample size and statistical power for such analyses. For a single-exposure meta-analysis without gene-by-covariate interactions, existing software options, such as the popular METAL program, are adequate. However, a nuanced set of considerations are required to determine whether it is appropriate to include additional terms in meta-analysis, whether related to additional exposure terms, gene-by-covariate interactions (Keller 2014), or genetic main effects (Laville *et al.* 2022). For multiple interaction meta-analysis, METAGEM demonstrated efficient CPU time, though large memory space is required for larger numbers of interaction terms and unique variants across studies.

By facilitating more comprehensive, genome-wide analyses and meta-analyses involving interactions using only summary statistics, REGEM and METAGEM enable researchers to maximize the value of GWISs while minimizing computational time. A few limitations should be noted. Firstly, the GEM model corrects for standard covariates by removing them from the genotype and interaction matrices in a single projection step. While this approach improves computational performance of the primary GWIS considerably, it also takes away the possibility of modifying covariate main effect adjustments in subsequent re-analysis. Any such modification (e.g. seeking an interaction effect while completely removing a covariate main effect from the statistical model) would require a new analysis using individual-level data. Additionally, while REGEM has been shown to produce results that are consistent with those of GEM, improper GEI analysis using GEM, particularly in the case of rare variants, can lead to spurious summary statistics results, and may invalidate re-analysis

results. Therefore, researchers must ensure valid summary statistics (e.g. well-controlled genomic inflation) are generated from GEI methods before performing a re-analysis. In this vein, it is also important that study-specific interaction terms to be meta-analyzed have equivalent interpretations; e.g. METAGEM cannot conduct valid meta-analysis when there are discrepant study-specific variable coding choices in terms of exposure (and covariate) centering. Finally, the current METAGEM implementation does not allow for extensions, such as aggregate tests of rare variants. Meta-analysis of such aggregate tests for GEI is available in the MAGEE tool (Wang *et al.* 2022), and further power and computational efficiency may be gained by leveraging genomic annotations and sparse weighted linkage disequilibrium (LD) matrices, following approaches, such as MetaSTAAR (Li *et al.* 2023) and using latest LD references provided by TOP-LD (Huang *et al.* 2022).

In summary, we have introduced REGEM and METAGEM for further complex downstream analysis of GEI studies. REGEM and METAGEM, along with our GEM tool for genome-wide interaction analysis and corresponding workflows for reproducible and scalable deployment in cloud computing environments, are publicly available at (<https://github.com/large-scale-gxe-methods>). The suite of tools, including GEM, REGEM, and METAGEM, provides key software infrastructure for maximizing the utility of summary statistics from diverse and complex GEI studies.

Acknowledgements

This research was conducted using the UK Biobank Resource under Application Numbers 27892 and 42646. ProDiGY acknowledgements are included in the [Supplemental Material](#).

Author contributions

D.T.P. and H.C. developed the METAGEM and REGEM algorithms. D.T.P., H.C., and C.P. implemented the METAGEM and REGEM software programs. D.T.P. and K.E.W. implemented software programs as cloud workflows. D.T.P. and H.C. designed the benchmark simulation study and carried out the analyses. K.E.W., L.C., and A.K.M. carried out the real-data analyses. 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 related to analysis of the ProDiGY dataset. K.E.W., D.T.P., H.C., and A.K.M. wrote the manuscript. All authors critically read the manuscript.

Supplementary data

[Supplementary data](#) are available at *Bioinformatics* online.

Conflict of interest

None declared.

Funding

This work was supported by the National Institutes of Health [grant numbers R01 HL145025, K01 DK133637 to K.E.W.]. ProDiGY funding sources are included in the [Supplemental Material](#).

Data availability

Data from UK Biobank is accessible to researchers by applying to: <https://www.ukbiobank.ac.uk/enable-your-research/apply-for-access>. The ProDiGY dataset analyzed in the current study is available at dbGaP (dbGaP Study Accession: phs001511.v1.p1, https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001511.v1.p1).

References

- Bi W, Zhao Z, Dey R *et al.* A fast and accurate method for genome-wide scale phenome-wide $G \times E$ analysis and its application to UK Biobank. *Am J Hum Genet* 2019;105:1182–92.
- Chang CC, Chow CC, Tellier LC *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 2015;4:7.
- Gauderman WJ, Zhang P, Morrison JL *et al.* Finding novel genes by testing $G \times E$ interactions in a genome-wide association study. *Genet Epidemiol* 2013;37:603–13.
- Huang L, Rosen JD, Sun Q *et al.*; NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium. TOP-LD: a tool to explore linkage disequilibrium with TOPMed whole-genome sequence data. *Am J Hum Genet* 2022;109:1175–81.
- Keller MC. Gene \times environment interaction studies have not properly controlled for potential confounders: The problem and the (simple) solution. *Biol Psychiatry* 2014;75:18–24.
- Kerin M, Marchini J. Inferring gene-by-environment interactions with a Bayesian whole-genome regression model. *Am J Hum Genet* 2020;107:698–713.
- Kim J, Ziyatdinov A, Laville V *et al.* Joint analysis of multiple interaction parameters in genetic association studies. *Genetics* 2019;211:483–94.
- Laville V, Majarian T, Sung YJ *et al.*; CHARGE Gene-Lifestyle Interactions Working Group. Gene-lifestyle interactions in the genomics of human complex traits. *Eur J Hum Genet* 2022;30:730–9.
- Li X, Quick C, Zhou H *et al.*; NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, TOPMed Lipids Working Group. Powerful, scalable and resource-efficient meta-analysis of rare variant associations in large whole genome sequencing studies. *Nat Genet* 2023;55:154–64.
- Lin D-Y, Tao R, Kalsbeek WD *et al.* Genetic association analysis under complex survey sampling: the Hispanic community health study/study of Latinos. *Am J Hum Genet* 2014;95:675–88.
- Manning AK, Hivert M-F, Scott RA *et al.*; Multiple Tissue Human Expression Resource (MUTHER) Consortium. A genome-wide approach accounting for body mass index identifies genetic variants influencing fasting glycemic traits and insulin resistance. *Nat Genet* 2012;44:659–69.
- Mbatchou J, Barnard L, Backman J *et al.* Computationally efficient whole-genome regression for quantitative and binary traits. *Nat Genet* 2021;53:1097–103.
- Moore R, Casale FP, Jan Bonder M *et al.*; BIOS Consortium. A linear mixed-model approach to study multivariate gene-environment interactions. *Nat Genet* 2019;51:180–6.
- Pan-UKB team. 2020. <https://pan.ukbb.broadinstitute.org>.
- Reales G, Wallace C. Sharing GWAS summary statistics results in more citations. *Commun Biol* 2023;6:116.
- SEARCH Study Group. SEARCH for diabetes in youth: a multicenter study of the prevalence, incidence and classification of diabetes mellitus in youth. *Control Clin Trials* 2004;25:458–71.
- Shin J, Lee SH. GxEsum: a novel approach to estimate the phenotypic variance explained by genome-wide $G \times E$ interaction based on GWAS summary statistics for biobank-scale data. *Genome Biol* 2021;22:183.
- Srinivasan S, Chen L, Todd J *et al.*; ProDiGY Consortium. The first genome-wide association study for type 2 diabetes in youth: the Progress in Diabetes Genetics in Youth (ProDiGY) consortium. *Diabetes* 2021;70:996–1005.

- TODAY Study Group; Zeitler P, Epstein L, Grey M *et al.* Treatment options for type 2 diabetes in adolescents and youth: a study of the comparative efficacy of metformin alone or in combination with rosiglitazone or lifestyle intervention in adolescents with type 2 diabetes. *Pediatr Diabetes* 2007;8:74–87.
- Wang X, Pham DT, Westerman KE *et al.* Genomic summary statistics and meta-analysis for set-based gene-environment interaction tests in large-scale sequencing studies. medRxiv, 2022, 2022.05.08. 22274819.
- Werme J, van der Sluis S, Posthuma D *et al.* Genome-wide gene-environment interactions in neuroticism: an exploratory study across 25 environments. *Transl Psychiatry* 2021;11:180.
- Westerman K, Liu Q, Liu S *et al.* A gene-diet interaction-based score predicts response to dietary fat in the women’s health initiative. *Am J Clin Nutr* 2020;111:893–902.
- Westerman KE, Pham DT, Hong L *et al.* GEM: scalable and flexible gene-environment interaction analysis in millions of samples. *Bioinformatics* 2021;37:3514–20.
- Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010; 26:2190–1.
- Zhong W, Chhibber A, Luo L *et al.* A fast and powerful linear mixed model approach for genotype-environment interaction tests in large-scale GWAS. *Brief Bioinform* 2023;24:bbac547.