

# Efficient Variant Set Mixed Model Association Tests for Continuous and Binary Traits in Large-Scale Whole-Genome Sequencing Studies

Han Chen,<sup>1,2</sup> Jennifer E. Huffman,<sup>3</sup> Jennifer A. Brody,<sup>4</sup> Chaolong Wang,<sup>5</sup> Seunggeun Lee,<sup>6</sup> Zilin Li,<sup>7</sup> Stephanie M. Gogarten,<sup>8</sup> Tamar Sofer,<sup>9,10</sup> Lawrence F. Bielak,<sup>11</sup> Joshua C. Bis,<sup>4</sup> John Blangero,<sup>12</sup> Russell P. Bowler,<sup>13</sup> Brian E. Cade,<sup>9,10</sup> Michael H. Cho,<sup>14,15</sup> Adolfo Correa,<sup>16</sup> Joanne E. Curran,<sup>12</sup> Paul S. de Vries,<sup>1</sup> David C. Glahn,<sup>17,18</sup> Xiuqing Guo,<sup>19</sup> Andrew D. Johnson,<sup>20</sup> Sharon Kardia,<sup>11</sup> Charles Kooperberg,<sup>21</sup> Joshua P. Lewis,<sup>22</sup> Xiaoming Liu,<sup>23</sup> Rasika A. Mathias,<sup>24</sup> Braxton D. Mitchell,<sup>22,25</sup> Jeffrey R. O'Connell,<sup>22</sup> Patricia A. Peyser,<sup>11</sup> Wendy S. Post,<sup>26</sup> Alex P. Reiner,<sup>21</sup> Stephen S. Rich,<sup>27</sup>

(Author list continued on next page)

With advances in whole-genome sequencing (WGS) technology, more advanced statistical methods for testing genetic association with rare variants are being developed. Methods in which variants are grouped for analysis are also known as variant-set, gene-based, and aggregate unit tests. The burden test and sequence kernel association test (SKAT) are two widely used variant-set tests, which were originally developed for samples of unrelated individuals and later have been extended to family data with known pedigree structures. However, computationally efficient and powerful variant-set tests are needed to make analyses tractable in large-scale WGS studies with complex study samples. In this paper, we propose the variant-set mixed model association tests (SMMAT) for continuous and binary traits using the generalized linear mixed model framework. These tests can be applied to large-scale WGS studies involving samples with population structure and relatedness, such as in the National Heart, Lung, and Blood Institute's Trans-Omics for Precision Medicine (TOPMed) program. SMMATs share the same null model for different variant sets, and a virtue of this null model, which includes covariates only, is that it needs to be fit only once for all tests in each genome-wide analysis. Simulation studies show that all the proposed SMMATs correctly control type I error rates for both continuous and binary traits in the presence of population structure and relatedness. We also illustrate our tests in a real data example of analysis of plasma fibrinogen levels in the TOPMed program ( $n = 23,763$ ), using the Analysis Commons, a cloud-based computing platform.

## Introduction

In recent years, massive DNA sequence data have been generated. Large-scale whole-genome sequencing projects, such as the National Heart, Lung, and Blood Institute's (NHLBI) Trans-Omics for Precision Medicine (TOPMed) program and the National Human Genome Research Institute's (NHGRI) Genome Sequencing Project (GSP),

have produced whole-genome sequences from more than 120,000 samples. The designs of the studies from which participants are drawn need not be uniform or simple; for example, TOPMed includes population-based cohorts, family studies, and case-control studies, some of which are conducted in recently admixed populations, and some of which involve large pedigrees of closely related participants.

<sup>1</sup>Human Genetics Center, Department of Epidemiology, Human Genetics and Environmental Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA; <sup>2</sup>Center for Precision Health, School of Public Health and School of Biomedical Informatics, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA; <sup>3</sup>Center for Population Genomics, VA Boston Healthcare System, Jamaica Plain, MA 02130, USA; <sup>4</sup>Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA 98101, USA; <sup>5</sup>Department of Epidemiology and Biostatistics, School of Public Health, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, Hubei 430030, China; <sup>6</sup>Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA; <sup>7</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA; <sup>8</sup>Department of Biostatistics, University of Washington, Seattle, WA 98195, USA; <sup>9</sup>Division of Sleep and Circadian Disorders, Brigham and Women's Hospital, Boston, MA 02115, USA; <sup>10</sup>Division of Sleep Medicine, Harvard Medical School, Boston, MA 02115, USA; <sup>11</sup>Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI 48109, USA; <sup>12</sup>Department of Human Genetics and South Texas Diabetes and Obesity Institute, School of Medicine, The University of Texas Rio Grande Valley, Brownsville, TX 78520, USA; <sup>13</sup>Division of Pulmonary Medicine, Department of Medicine, National Jewish Health, Denver, CO 80206, USA; <sup>14</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA 02115, USA; <sup>15</sup>Division of Pulmonary and Critical Care Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115, USA; <sup>16</sup>Jackson Heart Study, Department of Medicine, University of Mississippi Medical Center, Jackson, MS 39216, USA; <sup>17</sup>Department of Psychiatry, Yale University School of Medicine, New Haven, CT 06510, USA; <sup>18</sup>Olin Neuropsychiatric Research Center, Institute of Living, Hartford Hospital, Hartford, CT 06106, USA; <sup>19</sup>The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, Los Angeles Biomedical Research Institute at Harbor-UCLA Medical Center, Torrance, CA 90502, USA; <sup>20</sup>Framingham Heart Study, National Heart, Lung, and Blood Institute and Boston University, Framingham, MA 01702, USA; <sup>21</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA; <sup>22</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD 21201, USA; <sup>23</sup>USF Genomics, College of Public Health, University of South Florida, Tampa, FL 33612, USA; <sup>24</sup>Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD 21287, USA; <sup>25</sup>Geriatrics Research and Education Clinical Center, Baltimore VA Medical Center, Baltimore, MD 21201, USA; <sup>26</sup>Division of Cardiology, Johns Hopkins University, Baltimore, MD 21287, USA; <sup>27</sup>Center for Public Health Genomics, University of Virginia, Charlottesville, VA 22908, USA; <sup>28</sup>Sections of Preventive Medicine and Epidemiology, and of Cardiology, Department of Medicine, Boston University School of Medicine, Boston, MA 02118, USA;

(Affiliations continued on next page)



Jerome I. Rotter,<sup>19</sup> Edwin K. Silverman,<sup>14,15</sup> Jennifer A. Smith,<sup>11</sup> Ramachandran S. Vasan,<sup>20,28,29</sup> James G. Wilson,<sup>30</sup> Lisa R. Yanek,<sup>24</sup> NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, TOPMed Hematology and Hemostasis Working Group, Susan Redline,<sup>9,10,31</sup> Nicholas L. Smith,<sup>4,32,33,34</sup> Eric Boerwinkle,<sup>1,35</sup> Ingrid B. Borecki,<sup>8</sup> L. Adrienne Cupples,<sup>20,36</sup> Cathy C. Laurie,<sup>8</sup> Alanna C. Morrison,<sup>1</sup> Kenneth M. Rice,<sup>8</sup> and Xihong Lin<sup>7,37,\*</sup>

In population-based cohorts and case-control studies, population stratification and cryptic relatedness are major sources of confounding that need to be accounted for in association tests. For common single-variant analysis, linear mixed models that use an estimated genetic relationship matrix (GRM) to account for both population stratification and cryptic relatedness have been widely applied in genome-wide association studies (GWASs) to analyze structured and related samples.<sup>1–6</sup> For binary traits, however, we previously showed that linear mixed models may not be appropriate in the presence of population stratification due to misspecified mean-variance relationships. Therefore, we instead proposed a computationally efficient method GMMAT<sup>7</sup> to perform common single-variant tests in GWASs by fitting generalized linear mixed models (GLMMs),<sup>8</sup> which simultaneously account for population structure, cryptic relatedness, and shared environmental effects, using multiple variance components and/or random effects.

Hundreds of millions of genetic variants, mostly with a low and extremely rare minor allele frequency (MAF), are being analyzed in large-scale sequencing projects such as TOPMed and GSP. Yet, single-variant tests that have been widely used in GWASs are generally underpowered for analyzing rare genetic variants from sequencing studies. To circumvent this problem, statistical tests such as the burden test,<sup>9–12</sup> sequence kernel association test (SKAT),<sup>13</sup> and their various combinations<sup>14–16</sup> have been proposed. These tests analyze multiple genetic variants in sets, grouped by genes, genomic regions, or other bioinformatic aggregation units. Most of these tests were originally developed to analyze samples from unrelated individuals, as well as extensions to analyze family data with known pedigree structures in the parametric mixed model and semiparametric generalized estimating equation frameworks.<sup>17–23</sup>

Linear mixed models using a single random effect with the GRM covariance matrix to account for population structure have been developed and implemented in software programs for sequencing data analysis, such as EPACTS and Rvtests.<sup>24</sup> Meta-analysis methods for family data have been developed and implemented in seqMeta and RAREMETAL,<sup>25,26</sup> but only for continuous traits in the linear mixed model framework. Moreover, these existing methods do not account for cryptic relatedness and between-subject relatedness from multiple sources and

have not been applied to large-scale whole-genome sequencing studies with complex study samples, due to statistical and computational challenges.

One challenge is that among traditional variant set tests such as burden tests and SKAT, no single approach is uniformly most powerful. Another challenge is that existing hybrid tests that combine burden tests and SKAT, such as SKAT-O,<sup>14</sup> MiST,<sup>15</sup> and aSPU,<sup>16</sup> are powerful but are subject to much greater computational loads than either the burden test or SKAT alone in the GLMM framework. Of note, SKAT-O is slower than SKAT because it searches on a grid for the optimal linear combination of the burden test and SKAT statistics. MiST requires adjusting for the genetic burden as a covariate in the SKAT model and hence needs to fit a burden model for each variant set. In large samples of possibly related individuals, extension of MiST is not as practical as in unrelated samples, since fitting a mixed effects model using the burden score for each variant set (or each test unit) is computationally intensive across the genome. Finally, aSPU uses a permutation or Monte Carlo simulation procedure to compute the p values, which can also be challenging in the context of large-scale whole-genome sequencing studies with both population structure and relatedness. Therefore, there is a pressing need to develop powerful and computationally efficient statistical methods for large-scale whole-genome sequencing studies.

To address these statistical and computational challenges, we develop the variant set mixed model association tests (SMMATs), computationally efficient variant set tests for both continuous and binary traits, which are applicable to structured and related samples with potential multiple sources of correlations, from large-scale whole-genome sequencing studies. We include four tests in the SMMAT framework: the burden test (SMMAT-B), SKAT (SMMAT-S), SKAT-O (SMMAT-O), and an efficient hybrid test to combine the burden test and SKAT (SMMAT-E), with power improvements over mixed model-based burden test, SKAT and SKAT-O. All four SMMATs share the same reduced model under the null hypothesis, i.e., the GLMM with only covariates, which needs to be fit only once for all genetic variant sets in an analysis. We show that all of these tests can be constructed using shared single-variant scores and their covariance matrices, thus further improving the computational efficiency in practice compared to

<sup>29</sup>Department of Epidemiology, Boston University School of Public Health, Boston, MA 02118, USA; <sup>30</sup>Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS 39216, USA; <sup>31</sup>Division of Pulmonary, Critical Care, and Sleep Medicine, Beth Israel Deaconess Medical Center, Boston, MA 02115, USA; <sup>32</sup>Kaiser Permanente Washington Health Research Institute, Seattle, WA 98101, USA; <sup>33</sup>Seattle Epidemiologic Research and Information Center, Department of Veterans Affairs Office of Research and Development, Seattle, WA 98108, USA; <sup>34</sup>Department of Epidemiology, University of Washington, Seattle, WA 98195, USA; <sup>35</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA; <sup>36</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA; <sup>37</sup>Department of Statistics, Harvard University, Cambridge, MA 02138, USA

\*Correspondence: [xlin@hsph.harvard.edu](mailto:xlin@hsph.harvard.edu)  
<https://doi.org/10.1016/j.ajhg.2018.12.012>.

performing these tests separately. Moreover, it has been shown that single-variant scores and their covariance matrices can also be used in the meta-analysis of variant set tests,<sup>25,27</sup> and thus SMMAT has been implemented to be directly applicable to combining multi-cohort studies ranging from unstructured independent samples to structured and related samples. Finally, we develop a unified analysis pipeline in our software package GMMAT that implements SMMAT variant set tests in both single study (pooled analysis) and meta-analysis contexts to facilitate research on rare genetic variants from large-scale sequencing studies. We demonstrate the application of our method to the analysis of fibrinogen levels in the TOPMed study.

## Material and Methods

### Generalized Linear Mixed Models (GLMMs)

We formulate the SMMATs (SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E) from the same GLMM

$$g(\mu_i) = \mathbf{X}_i \boldsymbol{\alpha} + \mathbf{G}_i \boldsymbol{\beta} + b_i, \quad (\text{Equation 1})$$

where  $g(\cdot)$  is a monotonic “link” function that connects the mean of phenotype  $y_i$ , denoted by  $\mu_i = E(y_i | \mathbf{X}_i, \mathbf{G}_i, b_i)$ , for subject  $i$  of  $n$  samples, to the covariate row vector  $\mathbf{X}_i$ , the genotype row vector  $\mathbf{G}_i$  for  $q$  genetic variants in a set, and the random effects  $b_i$  that accounts for population structure and relatedness. The phenotypes  $y_i$  follow a distribution in the exponential family. For continuous traits, we usually assume that  $y_i$  follow a normal distribution and use an identity link function; for binary traits, we assume  $y_i$  follow a Bernoulli distribution and use a logit link function. In Equation 1,  $\boldsymbol{\alpha}$  is a  $p \times 1$  vector of fixed covariate effects including an intercept, and the genotype effects  $\boldsymbol{\beta}$  are assumed to be a  $q \times 1$  vector whose distribution has mean  $\mathbf{W} \mathbf{1}_q \beta_0$  and covariance  $\theta \mathbf{W}^2$ , where  $\mathbf{W} = \text{diag}\{w_j\}$  is a pre-specified  $q \times q$  matrix assigning weights to each variant,  $\theta$  is a variance component parameter, and  $\mathbf{1}_q$  is a column vector of length  $q$  with all elements 1. We assume that  $\mathbf{b} \sim N(0, \sum_{k=1}^K \nu_k \boldsymbol{\Phi}_k)$  is an  $n \times 1$  vector of random effects with each entry  $b_i$ ,  $K$  variance component parameters  $\nu_k$ , and known  $n \times n$  relatedness matrices  $\boldsymbol{\Phi}_k$  ( $1 \leq k \leq K$ ). We allow for multiple random effects to account for complex sampling designs such as hierarchical designs, shared environmental effects, and repeated-measures from longitudinal studies.

### SMMAT-B, SMMAT-S, and SMMAT-O

In Equation 1, testing the genotype effects of  $q$  variants  $H_0 : \boldsymbol{\beta} = 0$  is equivalent to testing the null hypothesis that  $H_0 : \beta_0 = 0$  and  $\theta = 0$ . The reduced GLMM under this null hypothesis specifies that

$$g(\mu_{0i}) = \mathbf{X}_i \boldsymbol{\alpha} + b_i, \quad (\text{Equation 2})$$

where  $\mu_{0i} = E(y_i | \mathbf{X}_i, b_i)$ . If we test  $H_0 : \beta_0 = 0$  under the assumption that  $\theta = 0$ , a burden score test SMMAT-B can be constructed as

$$T_B = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}_0)^T \mathbf{G} \mathbf{W} \mathbf{1}_q \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0)}{\hat{\phi}^2},$$

where  $\mathbf{y} = (y_1 \ y_2 \ \dots \ y_n)^T$  is an  $n \times 1$  vector of phenotypes  $y_i$ ,  $\hat{\boldsymbol{\mu}}_0$  is a vector of fitted mean values under the model in Equation 2,

$\mathbf{G} = (\mathbf{G}_1^T \ \mathbf{G}_2^T \ \dots \ \mathbf{G}_n^T)^T$  is an  $n \times q$  genotype matrix of the variant set in the test, and  $\hat{\phi}$  is an estimate of the dispersion parameter (or the residual variance)  $\phi$ . Under  $H_0 : \beta_0 = 0$ , the statistic  $T_B$  asymptotically follows  $\xi_B \chi_1^2$ , where the scalar  $\xi_B = \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q$ ,  $\chi_1^2$  is a chi-square distribution with 1 df, and  $\hat{\mathbf{P}} = \hat{\boldsymbol{\Sigma}}^{-1} - \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{X} (\mathbf{X}^T \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\boldsymbol{\Sigma}}^{-1}$  is the  $n \times n$  projection matrix of the null GLMM (Equation 2),  $\mathbf{X} = (\mathbf{X}_1^T \ \mathbf{X}_2^T \ \dots \ \mathbf{X}_n^T)^T$  is an  $n \times p$  covariate matrix,  $\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{V}} + \sum_{k=1}^K \hat{\nu}_k \boldsymbol{\Phi}_k$  with  $\hat{\mathbf{V}} = \hat{\phi} \mathbf{I}_n$  for continuous traits in linear mixed models, and  $\hat{\mathbf{V}} = \text{diag}\{1/(\hat{\mu}_{0i}(1 - \hat{\mu}_{0i}))\}$  for binary traits in logistic mixed models (where the dispersion parameter  $\phi$  is known to be 1).

On the other hand, if we test  $H_0 : \theta = 0$  under the assumption  $\beta_0 = 0$ , a variance component score-type test SMMAT-S can be constructed as

$$T_S = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}_0)^T \mathbf{G} \mathbf{W} \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0)}{\hat{\phi}^2}.$$

Under  $H_0 : \theta = 0$ ,  $T_S$  asymptotically follows  $\sum_{j=1}^q \xi_{Sj} \chi_{1,j}^2$ , where  $\chi_{1,j}^2$  are independent chi-square distributions with 1 df, and  $\xi_{Sj}$  are the eigenvalues of  $\boldsymbol{\Xi}_S = \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}$ .

If one assumes  $\beta_0$  has mean 0 and variance  $\gamma$ ,  $\boldsymbol{\beta}$  then follows a distribution 0 and covariance  $\tau \mathbf{W}\{(1 - \rho) \mathbf{I}_q + \rho \mathbf{1}_q \mathbf{1}_q^T\} \mathbf{W}$ , where  $\tau = \gamma + \theta$  and  $\rho = \gamma/(\gamma + \theta)$ , which takes values between 0 and 1. The joint null hypothesis  $H_0 : \beta_0 = 0$  and  $\theta = 0$  is equivalent to  $H_0 : \tau = 0$ . Given  $\rho$ , a variance component score-type test can be constructed as

$$T_\rho = \rho T_B + (1 - \rho) T_S.$$

If  $\rho = 1$ ,  $T_\rho$  becomes the SMMAT-B burden statistic  $T_B$ , which assumes  $\boldsymbol{\beta}$  are the same for all  $q$  variants after weighting. If  $\rho = 0$ ,  $T_\rho$  becomes the SMMAT-S SKAT statistic  $T_S$ . If an optimal  $\rho$  is obtained by minimizing the p value of  $T_\rho$ , then SMMAT-O can be constructed, with its p value calculated using a one-dimensional numerical integration, following SKAT-O.<sup>14</sup> A key advantage of SMMAT-O is that it maximizes the power by using the optimal linear combination of the mixed model burden test SMMAT-B and the mixed model SKAT SMMAT-S. As it requires a grid search over  $\rho$ , it is computationally considerably more expensive than SMMAT-B and SMMAT-S. We propose in the next section a computationally much more efficient method to combine SMMAT-B and SMMAT-S.

### SMMAT-E

An alternative joint test to SMMAT-O for  $H_0 : \beta_0 = 0$  and  $\theta = 0$  can be constructed using two asymptotically independent tests: a test for  $H_0 : \beta_0 = 0$  versus  $H_1 : \beta_0 \neq 0$  under the constraint  $\theta = 0$  and a test for  $H_0 : \theta = 0$  versus  $H_1 : \theta > 0$  with  $\beta_0$  as a nuisance parameter that is estimated under  $H_0 : \theta = 0$ . In unrelated samples, this testing strategy is a special case of MiST adjusting for the genotype burden score as a single fixed-effects covariate,<sup>15</sup> which requires the burden model to be fit for each SNP set. We note that the first test is SMMAT-B  $T_B$  in the SMMAT framework, and the second test  $T_\theta$  can be constructed from the null burden GLMM

$$g(\mu_{B_i}) = \mathbf{X}_i \boldsymbol{\alpha} + \mathbf{G}_i \mathbf{W} \mathbf{1}_q \beta_0 + b_i, \quad (\text{Equation 3})$$

where  $\mu_{B_i} = E(y_i | \mathbf{X}_i, \mathbf{G}_i \mathbf{W} \mathbf{1}_q, b_i)$  is the mean of  $y_i$  in the burden GLMM. We can construct a SKAT-type statistic adjusting for the genetic burden

$$T_\theta = \frac{(\mathbf{y} - \hat{\boldsymbol{\mu}}_B)^T \mathbf{G} \mathbf{W} \mathbf{W}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_B)}{\hat{\phi}^2},$$

where  $\hat{\boldsymbol{\mu}}_B$  is a vector of fitted values  $\hat{\mu}_{B_i}$  using the burden GLMM in Equation 3 for a given variant set. However, fitting this burden GLMM separately for each variant set is computationally expensive in large-scale whole-genome association studies.

Therefore, we propose a different computationally efficient strategy by assuming that the mean of genetic effects  $\beta_0$  is not large, a reasonable assumption for most genomic regions and most complex human diseases. Then we can construct  $T_\theta$  efficiently without refitting the burden GLMMs in Equation 3 for each variant set across the genome. We show in Appendix A that  $T_\theta$  can be approximated by

$$T_\theta \approx \hat{\phi}^{-2} (\mathbf{y} - \hat{\boldsymbol{\mu}}_0)^T \mathbf{G} \mathbf{W} \left\{ \mathbf{I}_q - \mathbf{1}_q \left( \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q \right)^{-1} \right. \\ \left. \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \right\} \left\{ \mathbf{I}_q - \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q \right. \\ \left. \left( \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q \right)^{-1} \mathbf{1}_q^T \right\} \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0).$$

Therefore, under  $H_0 : \theta = 0$ ,  $T_\theta$  asymptotically approximately follows  $\sum_{j=1}^q \xi_{\theta j} \chi_{1,j}^2$ , where  $\chi_{1,j}^2$  are independent chi-square distributions with 1 df, and  $\xi_{\theta j}$  are the eigenvalues of  $\Xi_\theta = \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} - \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q (\mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q)^{-1} \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}$ . By the central limit theorem, both  $\mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_B) / \hat{\phi}$  and  $\mathbf{1}_q^T \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0) / \hat{\phi}$  are asymptotically normal, and their covariance matrix is

$$\text{Cov} \left( \frac{\mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_B)}{\hat{\phi}}, \frac{\mathbf{1}_q^T \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0)}{\hat{\phi}} \right) \\ \approx \left\{ \mathbf{I}_q - \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q \left( \mathbf{1}_q^T \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q \right)^{-1} \mathbf{1}_q^T \right\} \\ \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \mathbf{1}_q = \mathbf{0}.$$

Therefore,  $T_\theta$  and  $T_B$  are approximately asymptotically independent. Let  $p_\theta$  and  $p_B$  be the p value of the two tests, respectively, then SMMAT-E p value  $p_E$  is computed using Fisher's method with a chi-square distribution with 4 df as  $p_E = P(\chi_4^2 > -2 \log(p_\theta p_B))$ .

## Meta-analysis

SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E can all be conducted in the meta-analysis context. Assuming the single-variant scores  $\mathbf{S} = \mathbf{G}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}_0) / \hat{\phi}$  and their covariance matrix  $\boldsymbol{\Psi} = \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G}$  are computed for each variant set in each study, we can reconstruct  $T_B = \mathbf{S}^T \mathbf{W} \mathbf{1}_q \mathbf{1}_q^T \mathbf{W} \mathbf{S}$  with  $\xi_B = \mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q$ ;  $T_S = \mathbf{S}^T \mathbf{W} \mathbf{W} \mathbf{S}$  with  $\xi_S = \mathbf{W} \boldsymbol{\Psi} \mathbf{W}$ ;  $T_\rho = \rho T_B + (1 - \rho) T_S$  and  $T_\theta = \mathbf{S}^T \mathbf{W} \{ \mathbf{I}_q - \mathbf{1}_q (\mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q)^{-1} \mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \} \{ \mathbf{I}_q - \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q (\mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q)^{-1} \mathbf{1}_q^T \} \mathbf{W} \mathbf{S}$  with  $\Xi_\theta = \mathbf{W} \boldsymbol{\Psi} \mathbf{W} - \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q (\mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W} \mathbf{1}_q)^{-1} \mathbf{1}_q^T \mathbf{W} \boldsymbol{\Psi} \mathbf{W}$ .

For each variant set, let  $m = 1, 2, \dots, M$  be the index of studies and  $\mathbf{S}_m$  and  $\boldsymbol{\Psi}_m$  be the single-variant scores and covariance matrix from study  $m$ . In testing the “weak” null hypothesis<sup>28</sup> of summary genetic effects  $H_0 : \beta = 0$ ,<sup>25,27</sup> we can compute meta summary statistics  $\mathbf{S} = \sum_{m=1}^M \mathbf{S}_m$  and  $\boldsymbol{\Psi} = \sum_{m=1}^M \boldsymbol{\Psi}_m$  and use them in SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E. If a genetic variant is monomorphic in a study, its single-variant score statistic and the corresponding row and column in the covariance matrix will be set to 0 for that study. When combining studies with very different sample characteristics, testing the “strong” null hypothesis<sup>28</sup> that genetic effects in all studies are 0 is sometimes desired. In

the general case, we may choose to group studies that are similar and test whether the summary genetic effects in all groups are 0, for example, in the meta-analysis of multi-ethnic samples. Let  $c = 1, 2, \dots, C$  be a partition of  $M$  studies ( $C \leq M$ ), where  $C$  is the number of ethnicities,  $\mathbf{S}_{c,m}$  and  $\boldsymbol{\Psi}_{c,m}$  be the single-variant scores and covariance matrix from study  $m$  in partition  $c$  ( $m = 1, 2, \dots, M_c$  in partition  $c$ , and  $\sum_{c=1}^C M_c = M$ ), such that genetic effects for the same variant are summarized within each partition  $c$  but heterogeneous across partitions,<sup>27</sup> we can also compute summary statistics  $\mathbf{S} = \left( \sum_{m=1}^{M_1} \mathbf{S}_{1,m}^T \quad \sum_{m=1}^{M_2} \mathbf{S}_{2,m}^T \quad \dots \quad \sum_{m=1}^{M_C} \mathbf{S}_{C,m}^T \right)^T$  and  $\boldsymbol{\Psi} = \text{diag} \{ \sum_{m=1}^{M_c} \boldsymbol{\Psi}_{c,m} \}$ . Note that  $\mathbf{S}$  is now a vector of length  $Cq$  and  $\boldsymbol{\Psi}$  is a block-diagonal matrix with  $C$  blocks of  $q \times q$  matrices, one for each partition of studies (with total dimension  $Cq \times Cq$ ), so we should replace  $\mathbf{W}$ ,  $\mathbf{1}_q$ , and  $\mathbf{I}_q$  by  $\mathbf{I}_C \otimes \mathbf{W}$  (where  $\otimes$  denotes the Kronecker product),  $\mathbf{1}_{Cq}$ , and  $\mathbf{I}_{Cq}$ , respectively, in the above expressions for  $T_B$ ,  $T_\rho$ ,  $T_S$ , and  $T_\theta$  for meta-analysis.

## Simulation Studies

### Type I Error in Single-Cohort Studies

We performed coalescent simulations to generate sequence data with 100 genetic variants in each set, and 10,000 independent sets for 8,000 individuals from a  $20 \times 20$  grid of spatially continuous populations with migration rate between adjacent cells  $M = 10$  (Figure 1A). Within each cell, we paired 20 individuals into 10 families and simulated 2 children for each family using gene dropping,<sup>29</sup> and in total we had 4,000 families and 16,000 individuals. For continuous traits, in each simulation replicate, we simulated the phenotype  $y_{ij}$  for individual  $j$  in family  $i$  under the null hypothesis of no genetic association from

$$y_{ij} = \alpha_1 Z_i + b_{ij} + \varepsilon_{ij}, \quad (\text{Equation 4})$$

where the “population effect”  $\alpha_1 = 1$  and the population indicator  $Z_i = 1$  if family  $i$  was from a  $10 \times 10$  grid in the top left of the map (population 1) and  $Z_i = 0$  otherwise (population 2). The familial random effects were simulated as

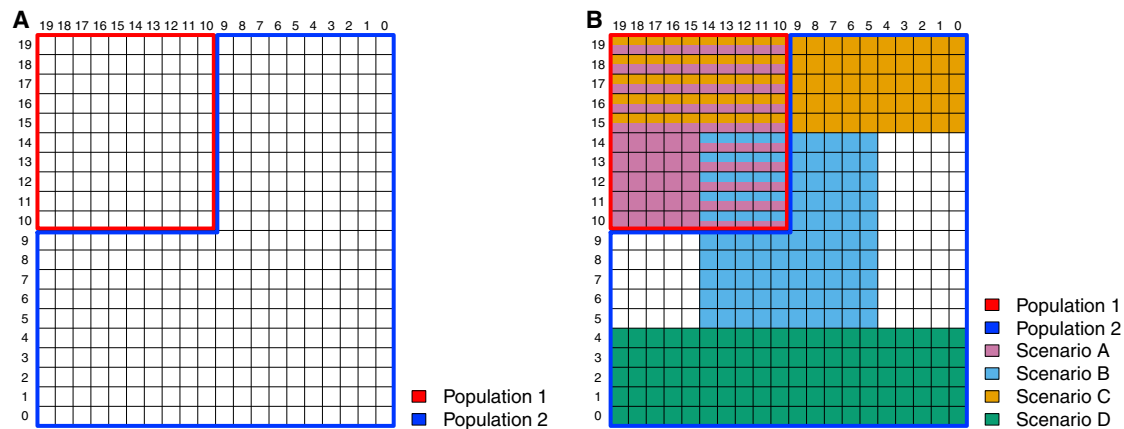
$$\mathbf{b}_i = \begin{pmatrix} b_{i1} \\ b_{i2} \\ b_{i3} \\ b_{i4} \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.5 & 0 & 0.25 & 0.25 \\ 0 & 0.5 & 0.25 & 0.25 \\ 0.25 & 0.25 & 0.5 & 0.25 \\ 0.25 & 0.25 & 0.25 & 0.5 \end{pmatrix} \right), \quad (\text{Equation 5})$$

and the random error  $\varepsilon_{ij} \sim N(0, 1)$  for each individual  $j$  in family  $i$ . Then we randomly sampled 3,500 individuals from the  $10 \times 10$  grid in the top left and 6,500 individuals from the rest of the map. The family identifier was removed for all individuals in the analysis, so that there were both population structure and cryptic relatedness in the sample. We compared SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in analyzing 10,000 independent variant sets based on a linear mixed model using our GMMAT package, including random effects with their covariance matrix proportional to the GRM, and adjusted for the first ten principal components (PCs) of ancestry. We repeated this 4,000 times to get p values combined from 40 million independent genetic variant sets for each test.

For binary traits, in each simulation replicate, we simulated the phenotype  $y_{ij}$  for individual  $j$  in family  $i$  under the null hypothesis of no genetic association from

$$\log \left( \frac{P(y_{ij} = 1)}{1 - P(y_{ij} = 1)} \right) = \alpha_0 + b_{ij}, \quad (\text{Equation 6})$$





**Figure 1. Map of Spatially Continuous Populations from Which Genotypes Were Simulated Based on the Coalescent Model**

(A) Map for a single-cohort simulation study: the top left 10 × 10 grid formed population 1, and the rest formed population 2.

(B) Map for a meta-analysis simulation study: scenario A studies were unrelated individuals sampled from population 1 only; scenario B studies were related individuals sampled from specific regions in population 1 and population 2; scenario C studies were unrelated individuals sampled from specific regions in population 1 and population 2; and scenario D studies were related individuals sampled from specific regions in population 2 only.

where  $\alpha_0$  was chosen such that the disease prevalence was 0.01 in all populations, and the familial random effects  $b_{ij}$  were simulated in the same way as for continuous traits. Then we randomly sampled 2,500 case subjects and 1,000 control subjects from the 10 × 10 grid in the top left (population 1), and 2,500 case subjects and 4,000 control subjects from the rest of the map (population 2) to form a hypothetical study with balanced case and control subjects in combined populations. Therefore, there was confounding by population structure resulting from unequal sampling, even though the disease prevalence was the same. We removed the family identifier, compared SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in analyzing 10,000 independent variant sets based on a logistic mixed model using our GMMAT package, similarly as described above, and repeated this 4,000 times to get p values combined from 40 million independent genetic variant sets for each test.

#### Type I Error in Meta-analysis

We also conducted simulation studies in the meta-analysis context to evaluate the type I error rates. We considered four scenarios: unrelated individuals, without confounding by population structure (scenario A studies); related individuals, with confounding by population structure (scenario B studies); unrelated individuals, with confounding by population structure (scenario C studies); and related individuals, without confounding by population structure (scenario D studies).

For scenario A studies, we simulated 16 unrelated individuals in each cell from the 10 × 10 grid in the top left of the map (Figure 1B). For continuous traits, we simulated the phenotype  $y_{ij}$  from Equation 4, with  $\alpha_1 = 0$  and  $b_{ij} = 0$  and randomly sampled 1,000 individuals. For binary traits, we simulated  $y_{ij}$  from Equation 6, with  $b_{ij} = 0$ , and randomly sampled 500 case subjects and 500 control subjects.

For scenario B studies, we simulated eight unrelated individuals, paired them into four families, and simulated two children for each family in each cell from the 10 × 10 grid in the center of the map (Figure 1B). For continuous traits, we simulated the phenotype  $y_{ij}$  from Equation 4, with  $\alpha_1 = 1$  and the population indicator  $Z_i = 1$  if family  $i$  was from population 1, and  $Z_i = 0$  if from population 2. Familial random effects  $b_{ij}$  were simulated using Equation 5, and we randomly sampled 350 individuals from population 1 and 650 individuals from population 2. For binary traits,

we simulated  $y_{ij}$  from Equation 6, with  $b_{ij}$  from Equation 5, and randomly sampled 250 case subjects and 100 control subjects from population 1, and 250 case subjects and 400 control subjects from population 2.

For scenario C studies, we simulated 16 unrelated individuals in each cell from the 20 × 5 grid in the top of the map (Figure 1B). For continuous traits, we simulated the phenotype  $y_{ij}$  from Equation 4, with  $\alpha_1 = 1$ , the population indicator  $Z_i = 1$  if family  $i$  was from population 1 and  $Z_i = 0$  if from population 2, and  $b_{ij} = 0$ , and we randomly sampled 350 individuals from population 1 and 650 individuals from population 2. For binary traits, we simulated  $y_{ij}$  from Equation 6, with  $b_{ij} = 0$ , and randomly sampled 250 case subjects and 100 control subjects from population 1 and 250 case subjects and 400 control subjects from population 2.

For scenario D studies, we simulated 8 unrelated individuals, paired them into 4 families and simulated 2 children for each family in each cell from the 20 × 5 grid in the bottom of the map (Figure 1B). For continuous traits, we simulated the phenotype  $y_{ij}$  from Equation 4, with  $\alpha_1 = 0$ , familial random effects  $b_{ij}$  simulated using Equation 5, and we randomly sampled 1,000 individuals. For binary traits, we simulated  $y_{ij}$  from Equation 6, with  $b_{ij}$  from Equation 5, and randomly sampled 500 case subjects and 500 control subjects.

In each simulation replicate, we simulated 3 studies from each scenario, totaling 12 studies with a combined sample size of 12,000 (6,000 case subjects and 6,000 control subjects for binary traits). We compared SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E using two meta-analysis strategies: all studies in the same group, and scenario A, B, C, and D studies in four separate groups. In the latter case, three studies from the same scenario were grouped in the same partition with shared genetic effects, while studies from different scenarios were allowed to have heterogeneous genetic effects. Variants are included in the meta-analysis as long as they are polymorphic in at least one of the 12 studies. We repeated 4,000 simulation replicates to get p values from 40 million independent genetic variant sets.

#### Power

We used the same genotype data as in the single-cohort type I error simulations and evaluated the empirical power of SMMAT-B, SMMAT-S, SMMAT-O, SMMAT-E, and the GLMM extension of MiST (GLMM-MiST) that combines the p value of SMMAT-B

**Table 1. Empirical Type I Error Rates of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in Single-Cohort Simulation Studies at Significance Levels of 0.05, 0.0001, and  $2.5 \times 10^{-6}$**

Level	Continuous Traits			Binary Traits		
	0.05	0.0001	$2.5 \times 10^{-6}$	0.05	0.0001	$2.5 \times 10^{-6}$
SMMAT-B	0.047	$8.7 \times 10^{-5}$	$2.0 \times 10^{-6}$	0.049	$9.6 \times 10^{-5}$	$2.0 \times 10^{-6}$
SMMAT-S	0.048	$8.7 \times 10^{-5}$	$2.0 \times 10^{-6}$	0.049	$9.5 \times 10^{-5}$	$2.3 \times 10^{-6}$
SMMAT-O	0.050	$1.1 \times 10^{-4}$	$3.0 \times 10^{-6}$	0.052	$1.2 \times 10^{-4}$	$3.0 \times 10^{-6}$
SMMAT-E	0.050	$1.0 \times 10^{-4}$	$3.0 \times 10^{-6}$	0.050	$9.9 \times 10^{-5}$	$2.0 \times 10^{-6}$

The total sample size was 10,000, and results from 4,000 simulation replicates were combined to get 40 million genetic variant sets.

(Equation 2) and the p value of SMMAT-S (Equation 3) using Fisher's method. All tests were performed using weights equal to a beta distribution density function with parameters 1 and 25 on the MAF of each variant.<sup>13</sup> We considered 9 scenarios, with the proportion of causal variants in a test unit changing from 10% to 20% to 50%, and the proportion of variants with negative effects out of causal variants changing from 100% to 80% to 50%. For continuous traits, we simulated the phenotype  $y_{ij}$  for individual  $j$  in family  $i$  from

$$y_{ij} = \alpha_1 Z_i + \sum_l G_{ijl} \beta_l + b_{ij} + \varepsilon_{ij},$$

where  $\alpha_1 = 1$ , the population indicator  $Z_i = 1$  if family  $i$  was from population 1 and  $Z_i = 0$  if from population 2,  $g_{ijl}$  was the centered genotype for causal variant  $l$  of individual  $j$  in family  $i$ , the causal effect size was  $|\beta_l| = c |\log_{10} \text{MAF}_l|$  for variant  $l$  with  $\text{MAF}_l$ , where the constant  $c$  was set to 0.2, 0.1, and 0.05 when the proportion of causal variants was 10%, 20%, and 50%, the familial random effects  $b_{ij}$  were simulated using Equation 5, and the random error  $\varepsilon_{ij} \sim N(0, 1)$ . We randomly sampled 35% individuals from population 1 and 65% individuals from population 2.

For binary traits, we simulated the phenotype  $y_{ij}$  for individual  $j$  in family  $i$  from

$$\log \left( \frac{P(y_{ij} = 1)}{1 - P(y_{ij} = 1)} \right) = \alpha_0 + \sum_l G_{ijl} \beta_l + b_{ij},$$

where  $\alpha_0$  was chosen such that the disease prevalence was 0.01 in all populations,  $G_{ijl}$  was the centered genotype for causal variant  $l$  of individual  $j$  in family  $i$ , the causal effect size was  $|\beta_l| = c |\log_{10} \text{MAF}_l|$  for variant  $l$  with  $\text{MAF}_l$ , where the constant  $c$  was set to 0.3, 0.2, and 0.1 when the proportion of causal variants was 10%, 20%, and 50%, the familial random effects  $b_{ij}$  were simulated using Equation 5. We randomly sampled 35% individuals (with 25% case subjects and 10% control subjects out of the total sample size) from population 1, and 65% individuals (with 25% case subjects and 40% control subjects out of the total sample size) from population 2 to form a hypothetical study with balanced case and control subjects in combined populations.

For both continuous and binary traits, we varied the total sample size from 2,000 to 5,000 to 10,000, repeated 1,000 simulation replicates for each scenario under the alternative hypothesis, and compared the empirical power at the significance level of  $2.5 \times 10^{-6}$ .

### TOPMed Example Involving Fibrinogen Levels

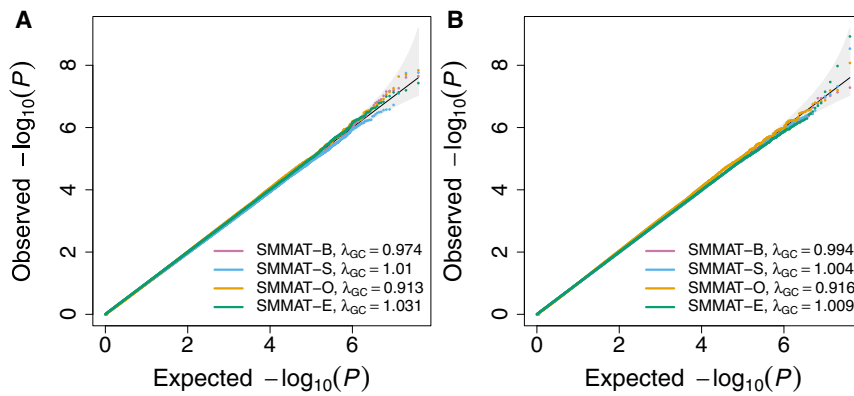
Samples with both plasma fibrinogen measures and whole-genome sequence data (Freeze 5b) from the following 11 TOPMed studies

were included in the analysis: the Old Order Amish Study (Amish), Cleveland Family Study (CFS), Genetic Epidemiology of COPD Study (COPDGene), Framingham Heart Study (FHS), Jackson Heart Study (JHS), San Antonio Family Study (SAFS), the Atherosclerosis Risk in Communities (ARIC) Study, Genetic Studies of Atherosclerosis Risk (GeneSTAR), Genetic Epidemiology Network of Arteriopathy (GENOA), the Multi-Ethnic Study of Atherosclerosis (MESA), and Women's Health Initiative (WHI). The TOPMed studies were approved by institutional review boards at participating institutions, and informed consent was obtained from all study participants. Amish, CFS, FHS, JHS, and SAFS are family-based studies with differing degrees of relatedness. The total sample size was 23,763. Within each study and each ethnicity, measured fibrinogen levels were adjusted for age, sex, and study-specific covariates, and the residuals were rank normalized and rescaled by multiplying by the original standard deviation, so that the transformed phenotype data have the same variances as on the original scale. The transformed phenotype data were pooled together in the analysis, using a heteroscedastic linear mixed model<sup>30</sup> allowing for different residual variances in each study/ethnicity, adjusting for study, ethnicity, sequence center, and top ten ancestry PCs<sup>31</sup> as fixed-effects covariates, and including a GRM calculated by mixed model analysis for pedigrees and populations (MMAP) to model the random effects for relatedness. Rare and low-frequency genetic variants on chromosome 4 with MAF less than 5%, including all singletons and extremely rare variants, were included in our rare variant association analysis of fibrinogen levels using the sliding window method<sup>32</sup> with 4 kb non-overlapping windows, using SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E with weights equal to a beta distribution density function with parameters 1 and 25 on the MAF of each variant.<sup>13</sup> As sensitivity analyses, we also included 1 kb, 10 kb, and 40 kb non-overlapping sliding windows, as well as an analysis using 4 kb windows with no ancestry PC adjustment. The analyses were performed using the GMMAT App (v.0.9.3), which includes the implementation of the SMMAT method, with 32 parallel threads on a single computing node with 240 GB total memory in the Analysis Commons.<sup>33</sup> To benchmark the computational speed in running SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E, we also ran re-analyses to perform each test separately, using summary statistics from the sliding window analysis and a single thread on a computing node with 15 GB total memory in the Analysis Commons.

## Results

### Simulation Studies

Table 1 shows the empirical type I error rates of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E at significance levels of 0.05, 0.0001, and  $2.5 \times 10^{-6}$  in the variant set analyses



**Figure 2.** Quantile-Quantile Plots of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in the Analysis of 10,000 Samples in Single-Cohort Studies with Both Population Structure and Cryptic Relatedness, under the Null Hypothesis of No Genetic Association

(A) Continuous traits in linear mixed models.

(B) Binary traits in logistic mixed models.

of continuous and binary traits in single-cohort simulation studies. All four tests have well-controlled type I error rates at these significance levels, suggesting that GLMMs can be effective in adjusting for population structure and cryptic relatedness in complex study samples. This is also consistent with the quantile-quantile (QQ) plots in Figure 2, which show neither inflation nor deflation in the tail.

Table 2 and Figure 3 show simulation results of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E assuming all studies in the same group (hom) or in four separate groups (het) in meta-analyses for combining four types of studies: with and without confounding by population structure, with and without cryptic relatedness. We note that SMMAT-B statistic  $T_B$  has the same form in these two meta-analysis strategies,<sup>27</sup> so we included seven tests in the simulation studies. In het SMMAT-S, SMMAT-O, and SMMAT-E, studies from the same scenario were grouped together to assume shared genetic effects. Under the null hypothesis of no genetic associations, hom SMMAT-O shows very mild inflation in our simulation settings, but all other six tests in the SMMAT framework control type I error rates well at significance levels of 0.05, 0.0001, and  $2.5 \times 10^{-6}$  and have well-calibrated tail probabilities, for both continuous and binary traits.

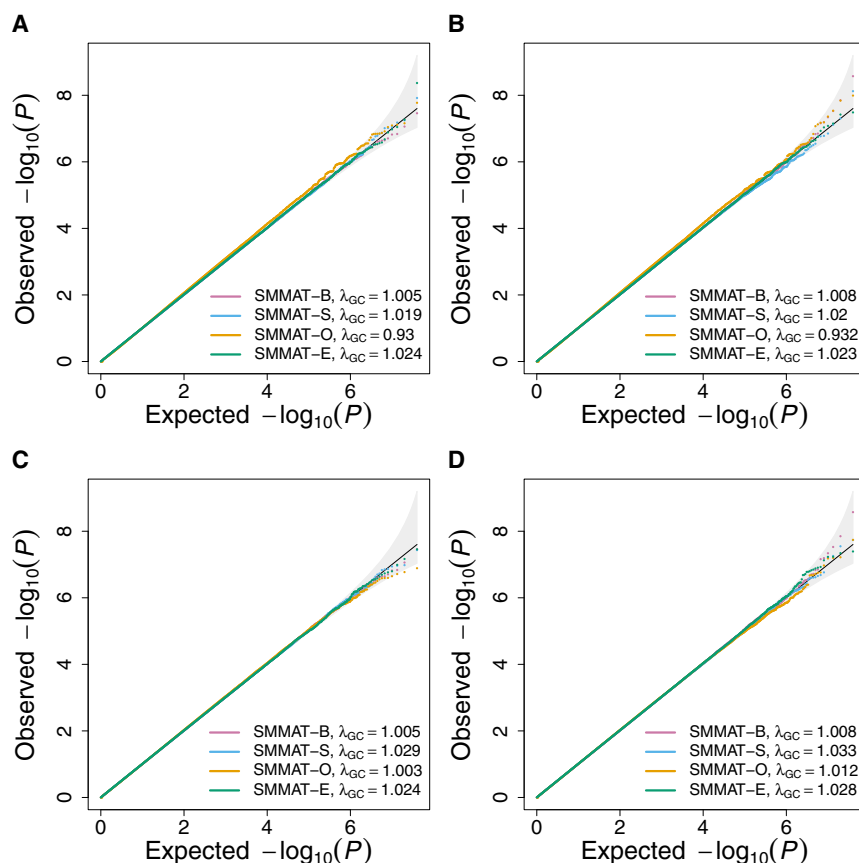
Figures 4 and 5 present the empirical power for causal variant sets at the significance level of  $2.5 \times 10^{-6}$  for continuous and binary traits, respectively. The power increases with the sample size. As the proportion of causal variants with effects in the same direction drops from 100% to 80% to 50% in each row, the power drops for all tests, but most substantially for the burden test SMMAT-B. When the sample size is large (i.e., 10,000 samples), SMMAT-E and GLMM-MiST have the highest power, for both continuous and binary traits in all nine simulation scenarios. SMMAT-E and GLMM-MiST have almost the same power in all these settings, while GLMM-MiST requires fitting a separate GLMM for each variant set. When all genetic variants in a test unit are causal with large effects in the same direction (a simulation scenario in favor of SMMAT-B, see Supplemental Material and Methods for details), SMMAT-B has the highest power, followed by SMMAT-O and SMMAT-E or GLMM-MiST (Figures S1A and S1B). On the log scale, SMMAT-E and GLMM-MiST p values are very close (Figures S1C and S1D).

In the presence of genetic relatedness from multiple sources (see Supplemental Material and Methods for details), linear and logistic mixed models with single GRM random effects and multiple random effects all control type I errors for continuous and binary traits (Figure S2). The multiple random effects model is more powerful than the single GRM random effects models for

**Table 2.** Empirical Type I Error Rates of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E Assuming All Studies in the Same Group (hom) and Scenario A, B, C, and D Studies in Four Separate Groups (het), in Meta-analysis Simulation Studies at Significance Levels of 0.05, 0.0001, and  $2.5 \times 10^{-6}$

Level	Continuous Traits			Binary Traits		
	0.05	0.0001	$2.5 \times 10^{-6}$	0.05	0.0001	$2.5 \times 10^{-6}$
SMMAT-B	0.051	$1.0 \times 10^{-4}$	$2.6 \times 10^{-6}$	0.051	$1.1 \times 10^{-4}$	$2.5 \times 10^{-6}$
Hom SMMAT-S	0.051	$1.0 \times 10^{-4}$	$2.6 \times 10^{-6}$	0.051	$1.1 \times 10^{-4}$	$2.1 \times 10^{-6}$
Het SMMAT-S	0.051	$1.0 \times 10^{-4}$	$2.8 \times 10^{-6}$	0.052	$1.0 \times 10^{-4}$	$2.4 \times 10^{-6}$
Hom SMMAT-O	0.053	$1.3 \times 10^{-4}$	$4.0 \times 10^{-6}$	0.053	$1.4 \times 10^{-4}$	$3.4 \times 10^{-6}$
Het SMMAT-O	0.052	$1.1 \times 10^{-4}$	$2.6 \times 10^{-6}$	0.052	$1.1 \times 10^{-4}$	$2.2 \times 10^{-6}$
Hom SMMAT-E	0.051	$1.0 \times 10^{-4}$	$2.5 \times 10^{-6}$	0.051	$1.1 \times 10^{-4}$	$2.6 \times 10^{-6}$
Het SMMAT-E	0.051	$1.0 \times 10^{-4}$	$2.8 \times 10^{-6}$	0.052	$1.1 \times 10^{-4}$	$3.0 \times 10^{-6}$

The total sample size was 12,000 from 12 studies, and results from 4,000 simulation replicates were combined to get 40 million genetic variant sets.



**Figure 3. Quantile-Quantile Plots of SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in the Meta-analysis of 12 Studies with a Total Sample Size of 12,000, under the Null Hypothesis of No Genetic Association**

(A) Continuous traits in linear mixed models, all studies in the same group.

(B) Binary traits in logistic mixed models, all studies in the same group.

(C) Continuous traits in linear mixed models, scenario A, B, C, and D studies in four separate groups.

(D) Binary traits in logistic mixed models, scenario A, B, C, and D studies in four separate groups.

continuous traits in our simulation settings, although the single GRM random effects models with and without ancestry PC adjustment almost have the same power (Figure S3). For binary traits, compared to the single GRM random effects model adjusting for ten ancestry PCs as fixed effects, the multiple random effects model is slightly more powerful and the single GRM random effects model with no ancestry PC adjustment is generally slightly less powerful, in our simulation settings (Figure S4).

#### TOPMed Example Involving Fibrinogen Levels

We compared the results from SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E in an analysis of fibrinogen levels, using chromosome 4 (including the genomic region that encodes the fibrinogen protein, *FGB*) whole-genome sequence data from 11 TOPMed studies. Previous studies have reported two rare variants within *FGB* on chromosome 4, rs6054 (hg38 position 154,568,456) and rs201909029 (hg38 position 154,567,636) associated with lower fibrinogen levels, with similar effect sizes in all ancestry groups.<sup>34</sup> In the sliding window analysis, we grouped low-frequency and rare genetic variants with MAF less than 5% into 46,859 non-overlapping 4 kb windows containing at least one variant. The number of variants in each window passing the MAF filter ranged from 1 to 1,290, with a median of 351 (25% quartile 326 and 75% quartile 380). The QQ plot (Figure 6A) shows that all four tests have well-calibrated tail probabilities. Table 3 summarizes heteroscedastic

linear mixed model-based SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E p values in *FGB* and flanking regions. SMMAT-S, SMMAT-O, and SMMAT-E give the most significant results in the 4 kb window 154,554–154,558 kb, with p values  $1.6 \times 10^{-17}$ ,  $8.9 \times 10^{-17}$ , and  $6.2 \times 10^{-19}$ , respectively, while SMMAT-B p value is much larger ( $6.9 \times 10^{-5}$ ). In the 4 kb window that covers both known association rare variants rs6054 and rs201909029 (window 154,566–

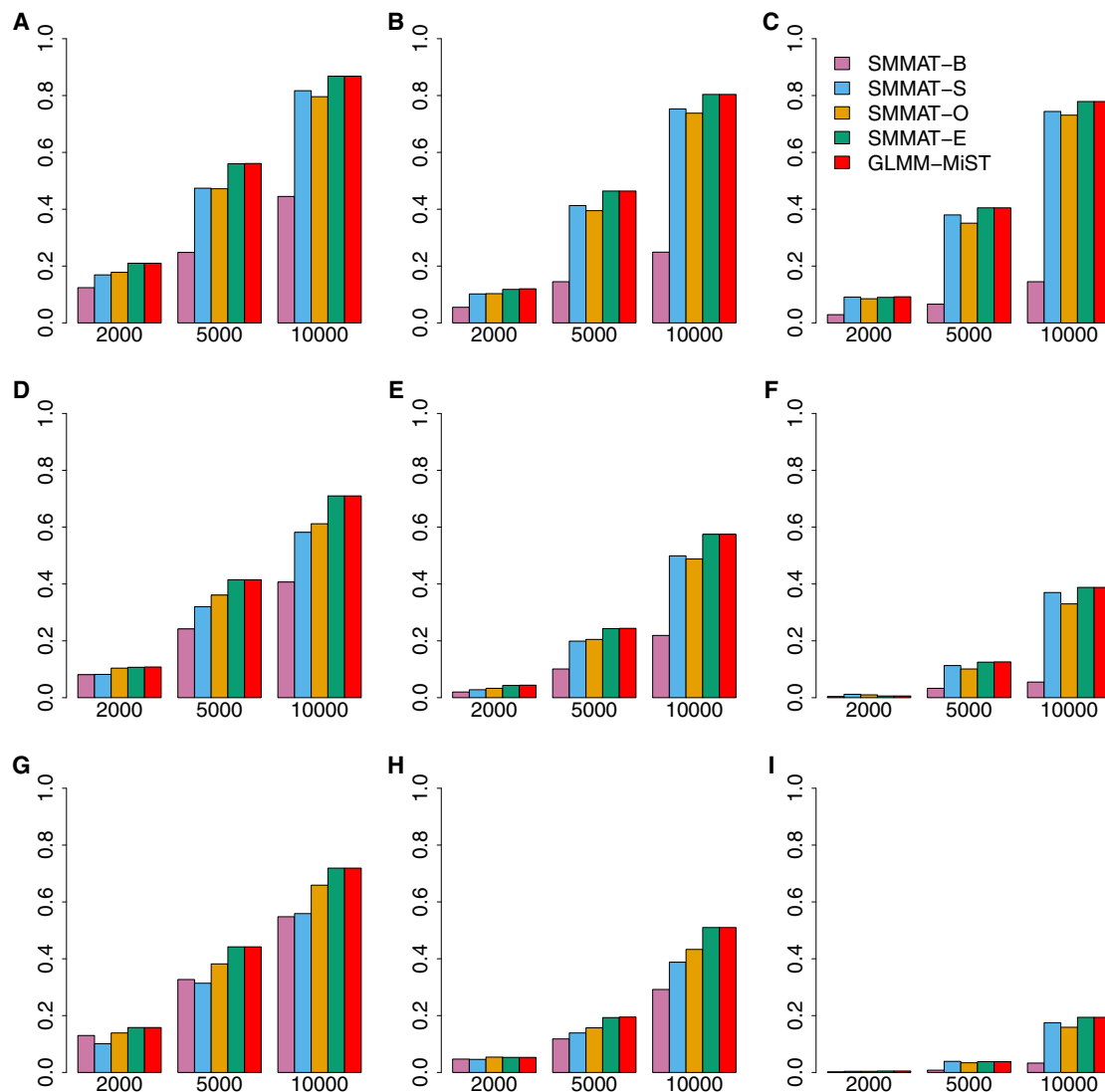
154,570 kb), SMMAT-E gives the smallest p value ( $3.1 \times 10^{-17}$ ), followed by SMMAT-S (p value  $9.7 \times 10^{-17}$ ), SMMAT-O (p value  $3.3 \times 10^{-16}$ ), and SMMAT-B (p value  $1.6 \times 10^{-8}$ ).

In this TOPMed data example, linear mixed models with and without adjusting for ten ancestry PCs as fixed-effects covariates gave very close p values (Figure S5). When we changed the window size from 4 kb to 1 kb (Figure S6), 10 kb (Figure S7), and 40 kb (Figure S8), the QQ plots showed that the analyses were well calibrated and the same association was identified. Regardless of the window size, SMMAT-E almost always gave the smallest p values, except for the 1 kb window 154,567–154,568 kb, which covers rs201909029. For this 1 kb window, none of the tests gave significant p values after adjusting for multiple testing, indicating potential lack of power, since rs201909029 has only 33 minor allele counts in our TOPMed samples (Table S1).

#### Computation Time

Table 4 shows the CPU time for running the sliding window analysis for 23,763 individuals with TOPMed whole-genome sequence data and fibrinogen levels, using summary statistics from 46,859 non-overlapping 4 kb windows on chromosome 4. The GMMAT App (v.0.9.3) in the Analysis Commons cloud computing platform has implemented SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E, with the option of running one or more tests





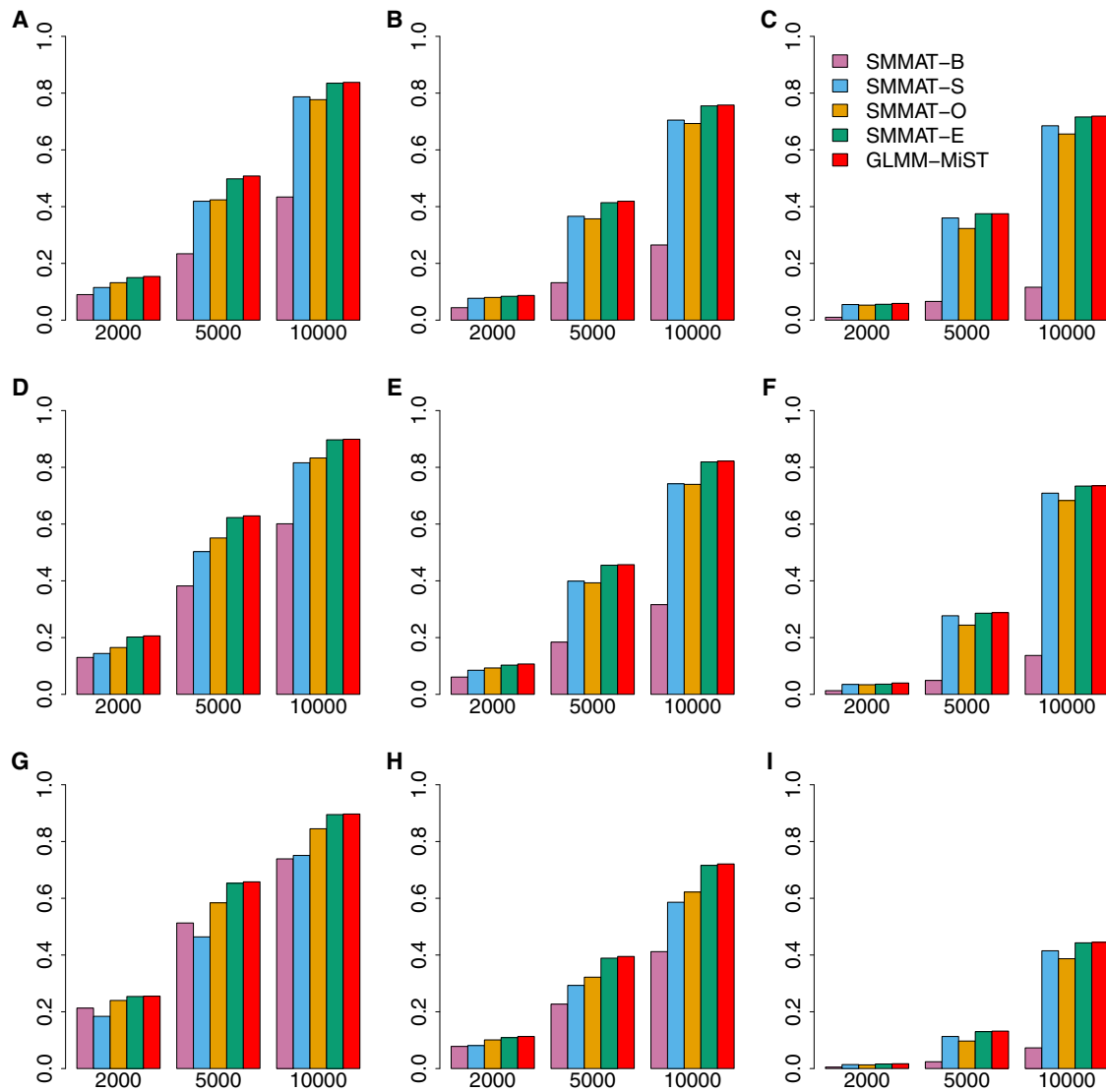
**Figure 4. Empirical Power of Linear Mixed Model-Based SMMAT-B, SMMAT-S, SMMAT-O, SMMAT-E, and GLMM-MiST in Continuous Trait Analysis of 2,000, 5,000, and 10,000 Samples**  
 (A–C) 10% causal variants with 100% (A), 80% (B), or 50% (C) negative effects.  
 (D–F) 20% causal variants with 100% (D), 80% (E), or 50% (F) negative effects.  
 (G–I) 50% causal variants with 100% (G), 80% (H), or 50% (I) negative effects.  
 Effect sizes were simulated using the same parameter in each row, but different across rows.

in an analysis. SMMAT-B results are automatically included when running SMMAT-O or SMMAT-E, and SMMAT-S p values will also be output when running SMMAT-O. Of the four tests in Table 4, SMMAT-B takes shortest time as the p value calculation does not involve any eigen-decomposition of covariance matrices. SMMAT-S takes only about 10 min longer than SMMAT-B for the eigen-decomposition of 46,859 covariance matrices. SMMAT-E takes about 12 min longer than SMMAT-S and gives both SMMAT-B and SMMAT-E p values. SMMAT-O takes 175 min longer than SMMAT-S, as more eigen-decompositions are performed in SMMAT-O when it searches for the optimal combination of SMMAT-B and SMMAT-S on a grid of  $\rho$  values. We did not include GLMM-MiST in the analysis, because it took

159 min CPU time to fit a GLMM for this TOPMed sample. By extrapolation, it would take more than 14 years CPU time for analyzing 23,763 related individuals with 46,859 windows using GLMM-MiST.

## Discussion

We have developed and implemented SMMAT, a family of computationally efficient variant set mixed model association tests for continuous and binary traits in large-scale whole-genome sequencing studies. This framework includes extensions of three widely used variant set tests for unrelated individuals to complex study samples with population structure and cryptic relatedness: the burden test (SMMAT-B), SKAT (SMMAT-S), and SKAT-O



**Figure 5. Empirical Power of Logistic Mixed Model-Based SMMAT-B, SMMAT-S, SMMAT-O, SMMAT-E, and GLMM-MiST in Binary Trait Analysis of 2,000, 5,000, and 10,000 Samples**

(A–C) 10% causal variants with 100% (A), 80% (B), or 50% (C) negative effects.

(D–F) 20% causal variants with 100% (D), 80% (E), or 50% (F) negative effects.

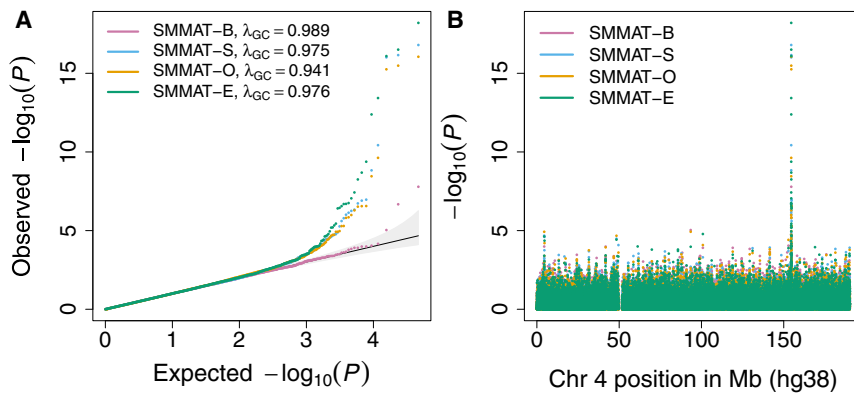
(G–I) 50% causal variants with 100% (G), 80% (H), or 50% (I) negative effects.

Effect sizes were simulated using the same parameter in each row, but different across rows.

(SMMAT-O), as well as a new efficient hybrid test that combines the mixed model burden and SKAT tests (SMMAT-E). Specifically, SMMAT-E is constructed by combining the burden test and an adjusted mixed model SKAT statistic that is approximately asymptotically independent from the mixed model burden test statistic, in a similar spirit to MiST in non-mixed model setting,<sup>15</sup> but that differs from MiST in that it does not require fitting separate mixed effect burden models for each variant set with the set genetic burden as a fixed-effects covariate. Instead, we use matrix projections to approximate the adjusted SKAT statistic from a global null model without any fixed effects for the variant set-specific genetic burden. Of note, this global null model needs to be fit only once in a whole-genome analysis, which greatly reduces the computational

cost. The approximation is highly accurate, even in the presence of large genetic effects. We show in simulation studies and the TOPMed fibrinogen example that SMMAT-E is more powerful than the other three tests in large samples, at the computational cost almost on the same scale of SMMAT-B and SMMAT-S. Therefore, SMMAT-E is recommended in the analysis of large-scale whole-genome sequencing studies.

In the SMMAT framework, different weighting strategies can be used. One can use a function of the MAF,<sup>11,13</sup> or external measures based on functional annotation such as CADD,<sup>35</sup> Eigen,<sup>36</sup> FATHMM-XF,<sup>37</sup> or tissue-specific annotations, such as GENOSKYLINE,<sup>38</sup> as the weight for each variant in a set. In the analysis of fibrinogen levels in TOPMed, we used MAF-based weights. Recently, unified



**Figure 6. TOPMed Fibrinogen Level SMMAT Analysis Results via a Heteroscedastic Linear Mixed Model on Rare Variants with MAF < 5% in Non-overlapping 4 kb Sliding Windows on Chromosome 4 (n = 23,763)**

(A) Quantile-quantile plot.

(B) p values on the log scale versus physical positions of the windows on chromosome 4 (build hg38).

variant set tests allowing for multiple functional annotations have been developed,<sup>39</sup> and the SMMAT framework can possibly be extended to accommodate multiple weights. Nevertheless, the optimal weighting strategy in rare variant analysis remains an open question and an active field of research.

As SMMAT-E combines the burden test p value  $p_B$  with an asymptotically independent adjusted SKAT p value  $p_\theta$  using Fisher's method in our SMMAT implementation in the GMMAT App, we note that other forms of combinations may also be applied.<sup>40</sup> For example, previous studies have shown that Tippett's procedure based on the minimum of  $p_\theta$  and  $p_B$  might be more powerful than Fisher's method in MiST when only one of the p values is small.<sup>15</sup> Alternatively, instead of combining the p values, weighted linear combinations of chi-square statistics have been proposed<sup>41–43</sup> and they can also be applied to combine the burden test statistic  $T_B$  and the asymptotically independent SKAT statistic  $T_\theta$  in the SMMAT framework.

SMMAT also has some limitations. SMMAT p values are computed based on asymptotic distributions, which may not be accurate in small samples, especially for binary traits and heavily skewed continuous traits. For continuous traits, small-sample inference procedures have been proposed for SKAT,<sup>44,45</sup> and the same methodology can be applied to SMMAT. For ultra-rare genetic variants with very low minor allele counts, the single-variant scores used to construct SMMAT-B, SMMAT-S, SMMAT-O, and SMMAT-E may not be close to a normal distribution, even if the total sample size is large. If there are only ultra-rare variants (e.g., singletons, doubletons) in a test region and the number of variants is small, SMMAT-B might be the best analysis strategy as its asymptotic property depends on the cumulative minor allele counts. Moreover, the asymptotic issue of single-variant scores also exists for binary traits with highly unbalanced case-control ratios, and a saddlepoint approximation approach has been proposed to match the cumulant generating function of the single-variant scores,<sup>46</sup> and it has recently been extended to GLMMs.<sup>47</sup>

Fitting GLMMs with a GRM has  $O(n^3)$  complexity in general, where  $n$  is the sample size. We have overcome this computational challenge by fitting only one GLMM

SMMAT. In large-scale whole-genome sequencing studies, solutions to other computational challenges are being proposed. For example, when the number of variants  $q$  in SKAT is very large, eigendecomposition of the covariance matrix, which has  $O(\min(n, q)^3)$  complexity, could be computationally expensive. Recently, the fastSKAT approach has been proposed to efficiently approximate the null distribution of SKAT when  $q$  is very large,<sup>48</sup> and the same strategy can be applied to speed up SMMAT p value calculation for very large  $q$ . On the other hand, as the sample size in ongoing large-scale sequencing projects such as TOPMed eventually expands to hundreds of thousands, using a full  $n \times n$  GRM would not be computationally practical in pooled analyses, as it may take several weeks to fit even only one GLMM with  $O(n^3)$  complexity and  $O(n^2)$  memory footprint. Meta-analyses may be a more appealing analysis strategy in that situation by combining summary statistics from study-specific or ancestry-specific analyses. Essentially equivalently, in pooled analyses, using a sparse and/or block-diagonal GRM with each block corresponding to an individual study in meta-analyses, will help reduce the computational cost in fitting GLMMs, providing one uses specialized routines for manipulation of sparse matrices.<sup>49</sup> Although whole-genome sequencing studies have not yet been conducted in large biobanks with sample sizes on the scale of millions of individuals, it is expected that calculating the GRM itself would become a major computational bottleneck. Recently, GRM-free mixed effects models such as BOLT-LMM<sup>6,50</sup> and SAIGE<sup>47</sup> have been developed for single variant tests, and we note that extension of these methods to the SMMAT framework will further reduce the computational cost in biobank-scale whole-genome sequencing studies in the future.

In summary, SMMAT provides a flexible and practical statistical framework for large-scale whole-genome sequencing studies with complex study samples, with balanced power and computational performance. With continuing advances in technology, lowering cost and development of new analytical methods, large-scale whole-genome sequencing studies will facilitate human genetic research and enhance our understandings of complex diseases and traits.

**Table 3. TOPMed Fibrinogen-Level SMMAT p Values in Known Association Gene *FGB* and Flanking Regions on Chromosome 4, using a Heteroscedastic Linear Mixed Model on Rare Variants with MAF < 5% (n = 23,763)**

Start (kb)	End (kb)	No. of Variants	SMMAT-B	SMMAT-S	SMMAT-O	SMMAT-E
154,554	154,558	348	$6.9 \times 10^{-5}$	$1.6 \times 10^{-17}$	$8.9 \times 10^{-17}$	$6.2 \times 10^{-19}$
154,558	154,562	370	0.078	$3.7 \times 10^{-11}$	$2.4 \times 10^{-10}$	$3.7 \times 10^{-14}$
154,562	154,566	326	0.76	$1.5 \times 10^{-9}$	$3.5 \times 10^{-9}$	$4.2 \times 10^{-10}$
154,566	154,570	309	$1.6 \times 10^{-8}$	$9.7 \times 10^{-17}$	$3.3 \times 10^{-16}$	$3.1 \times 10^{-17}$
154,570	154,574	332	0.030	$1.9 \times 10^{-7}$	$5.2 \times 10^{-7}$	$8.9 \times 10^{-8}$
154,574	154,578	349	$2.1 \times 10^{-7}$	$7.3 \times 10^{-7}$	$2.8 \times 10^{-7}$	$4.1 \times 10^{-13}$
154,578	154,582	342	$1.7 \times 10^{-4}$	$2.7 \times 10^{-5}$	$2.8 \times 10^{-5}$	$2.1 \times 10^{-9}$

Physical positions of each window are on build hg38.

## Appendix A: Approximations in SMMAT-E

Here we derive the approximations used in SMMAT-E to construct the SKAT-type statistic adjusting for the genetic burden

$$T_\theta = \frac{(\mathbf{y} - \tilde{\mu}_B)^T \mathbf{G} \mathbf{W} \mathbf{W}^T (\mathbf{y} - \tilde{\mu}_B)}{\tilde{\phi}^2}.$$

Let  $\tilde{\phi}$ ,  $\tilde{\alpha}$ ,  $\tilde{\beta}_0$ ,  $\tilde{b}_i$ ,  $\tilde{\mathbf{V}}$ , and  $\tilde{\Sigma}$  be estimates for  $\phi$ ,  $\alpha$ ,  $\beta_0$ ,  $b_i$ ,  $\mathbf{V}$ , and  $\Sigma$ , respectively, from the burden GLMM (Equation 3). We define  $\tilde{\mathbf{Y}} = \mathbf{y}$  as the phenotype vector for continuous traits, and the “working vector” with components  $\tilde{Y}_i = \mathbf{X}_i \tilde{\alpha} + \mathbf{G}_i \mathbf{W}_i \mathbf{1}_q \tilde{\beta}_0 + \tilde{b}_i + \{\tilde{\mu}_{B_i}(1 - \tilde{\mu}_{B_i})\}^{-1}(y_i - \tilde{\mu}_{B_i})$  at convergence of the logistic burden mixed model for binary traits (Equation 3), where  $\tilde{\alpha}$ ,  $\tilde{\beta}_0$ ,  $\tilde{b}_i$  are fixed-effects and random-effects estimates from the burden GLMM. We have

$$\begin{aligned} \frac{\mathbf{y} - \tilde{\mu}_B}{\tilde{\phi}} &= \tilde{\mathbf{V}}^{-1} (\tilde{\mathbf{Y}} - \mathbf{X} \tilde{\alpha} - \mathbf{G} \mathbf{W}_q \mathbf{1}_q \tilde{\beta}_0 - \tilde{\mathbf{b}}) \\ &= \tilde{\Sigma}^{-1} (\tilde{\mathbf{Y}} - \mathbf{X} \tilde{\alpha} - \mathbf{G} \mathbf{W}_q \mathbf{1}_q \tilde{\beta}_0) \\ &= \tilde{\Sigma}^{-1} \left\{ \tilde{\mathbf{Y}} - (\mathbf{X} \quad \mathbf{G} \mathbf{W}_q \mathbf{1}_q) \right. \\ &\quad \times \begin{pmatrix} \mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X} & \mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \\ \mathbf{1}_q^T \mathbf{W}_q^T \tilde{\Sigma}^{-1} \mathbf{X} & \mathbf{1}_q^T \mathbf{W}_q^T \tilde{\Sigma}^{-1} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \end{pmatrix}^{-1} \\ &\quad \times \begin{pmatrix} \mathbf{X}^T \tilde{\Sigma}^{-1} \\ \mathbf{1}_q^T \mathbf{W}_q^T \tilde{\Sigma}^{-1} \end{pmatrix} \tilde{\mathbf{Y}} \left. \right\} \\ &= \left\{ \tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \tilde{\Sigma}^{-1} \right\} \tilde{\mathbf{Y}} \\ &\quad - \left\{ \tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \tilde{\Sigma}^{-1} \right\} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \\ &\quad \left[ \mathbf{1}_q^T \mathbf{W}_q^T \left\{ \tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \tilde{\Sigma}^{-1} \right\} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right]^{-1} \\ &\quad \times \mathbf{1}_q^T \mathbf{W}_q^T \left\{ \tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \tilde{\Sigma}^{-1} \right\} \tilde{\mathbf{Y}}. \end{aligned}$$

Note that  $\tilde{\phi} = 1$  for binary traits. Moreover, since the true value of  $\beta_0$  is small, assuming including the genetic burden

$\mathbf{G}_i \mathbf{W}_i \mathbf{1}_q$  in the second term in Equation 3 does not dramatically change the variance component estimates for  $\nu_k$  and  $\phi$  (and for binary traits, also the “working vector”  $\tilde{\mathbf{Y}}$  at convergence of the model from Equation 2), we have the approximation  $\tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \tilde{\Sigma}^{-1} \approx \hat{\mathbf{P}}$  and  $(\mathbf{y} - \hat{\mu}_0)/\hat{\phi} \approx \hat{\mathbf{P}} \tilde{\mathbf{Y}}$ , then

$$\begin{aligned} \frac{\mathbf{W} \mathbf{G}^T (\mathbf{y} - \tilde{\mu}_B)}{\tilde{\phi}} &\approx \mathbf{W} \mathbf{G}^T \left\{ \hat{\mathbf{P}} \tilde{\mathbf{Y}} - \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right. \\ &\quad \times \left( \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right)^{-1} \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \tilde{\mathbf{Y}} \left. \right\} \\ &\approx \left\{ \mathbf{I}_q - \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \left( \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right)^{-1} \mathbf{1}_q^T \right\} \\ &\quad \times \frac{\mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\mu}_0)}{\hat{\phi}}. \end{aligned}$$

Therefore,

$$\begin{aligned} T_\theta &= \frac{(\mathbf{y} - \tilde{\mu}_B)^T \mathbf{G} \mathbf{W} \mathbf{W}^T (\mathbf{y} - \tilde{\mu}_B)}{\tilde{\phi}^2} \\ &\approx \hat{\phi}^{-2} (\mathbf{y} - \hat{\mu}_0)^T \mathbf{G} \mathbf{W} \left\{ \mathbf{I}_q - \mathbf{1}_q \left( \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right)^{-1} \right. \\ &\quad \times \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W} \left. \right\} \left\{ \mathbf{I}_q - \mathbf{W} \mathbf{G}^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right. \\ &\quad \times \left( \mathbf{1}_q^T \mathbf{W}_q^T \hat{\mathbf{P}} \mathbf{G} \mathbf{W}_q \mathbf{1}_q \right)^{-1} \mathbf{1}_q^T \left. \right\} \mathbf{W} \mathbf{G}^T (\mathbf{y} - \hat{\mu}_0). \end{aligned}$$

**Table 4. CPU Time in the TOPMed Fibrinogen Level SMMAT using Summary Statistics from a Sliding Window Analysis using Non-overlapping 4 kb Windows on Chromosome 4 (n = 23,763)**

Test	Time (min)
SMMAT-B	81
SMMAT-S	91
SMMAT-O	266
SMMAT-E	103

Tests were performed using the GMMAT App (v.0.9.3) with one single thread on a computing node with 15 GB total memory in the Analysis Commons.



## Supplemental Data

Supplemental Data include eight figures, one table, Supplemental Material and Methods, Supplemental Acknowledgments, and the full authorship list with affiliations of the Trans-Omics for Precision Medicine (TOPMed) Consortium and can be found with this article online at <https://doi.org/10.1016/j.ajhg.2018.12.012>.

## Acknowledgments

This work was supported by National Institutes of Health grants R00 HL130593 (to H.C.), U01 HL120393 (to H.C. and J.E.H.), and R35 CA197449, P01-CA134294, U01-HG009088, U19-CA203654, and R01-HL113338 (to X. Lin). The authors acknowledge the Texas Advanced Computing Center (TACC, <https://www.tacc.utexas.edu>) at The University of Texas at Austin for providing high performance computing (HPC) resources that have contributed to the research results reported within this paper. Whole-genome sequence analysis of fibrinogen levels in TOPMed was performed in the Analysis Commons on DNAnexus, a hosting platform that uses Amazon Web Services (AWS) to provide a cloud data management and computing environment for large genomic data projects. The Analysis Commons was funded by NIH R01 HL131136. Phenotype harmonization and aggregation of the fibrinogen levels across TOPMed studies were supported in part by NIH R01 HL139553. Detailed TOPMed and study-specific acknowledgments can be found in [Supplemental Acknowledgments](#). The views expressed in this manuscript are those of the authors and do not necessarily represent the views of the National Heart, Lung, and Blood Institute, the National Institutes of Health, or the U.S. Department of Health and Human Services. We thank the referees for their helpful comments that have helped improve the paper.

## Declaration of Interests

In the past three years, E.K.S. received honoraria from Novartis for Continuing Medical Education Seminars and grant and travel support from GlaxoSmithKline. M.H.C. has received grant support from GlaxoSmithKline.

Received: August 18, 2018

Accepted: December 17, 2018

Published: January 10, 2019

## Web Resources

Analysis Commons, <http://analysiscommons.com/>  
DNAnexus, <https://www.dnanexus.com/>  
EPACTS, <https://github.com/statgen/EPACTS>  
GMMAT (including implementation of SMMAT), <https://github.com/hanchenphd/GMMAT>  
MMA, <https://github.com/MMA>  
RAREMETAL, [https://genome.sph.umich.edu/wiki/RAREMETAL\\_Documentation](https://genome.sph.umich.edu/wiki/RAREMETAL_Documentation)  
Rvtests, <http://zhanxw.github.io/rvtests/>  
seqMeta, <https://cran.r-project.org/packages/seqMeta/index.html>  
Xihong Lin's group (including GMMAT), <https://content.sph.harvard.edu/xlin/software.html>

## References

1. Kang, H.M., Zaitlen, N.A., Wade, C.M., Kirby, A., Heckerman, D., Daly, M.J., and Eskin, E. (2008). Efficient control of population structure in model organism association mapping. *Genetics* 178, 1709–1723.
2. Kang, H.M., Sul, J.H., Service, S.K., Zaitlen, N.A., Kong, S.Y., Freimer, N.B., Sabatti, C., and Eskin, E. (2010). Variance component model to account for sample structure in genome-wide association studies. *Nat. Genet.* 42, 348–354.
3. Lippert, C., Listgarten, J., Liu, Y., Kadie, C.M., Davidson, R.I., and Heckerman, D. (2011). FaST linear mixed models for genome-wide association studies. *Nat. Methods* 8, 833–835.
4. Zhou, X., and Stephens, M. (2012). Genome-wide efficient mixed-model analysis for association studies. *Nat. Genet.* 44, 821–824.
5. Pirinen, M., Donnelly, P., and Spencer, C.C.A. (2013). Efficient computation with a linear mixed model on large-scale data sets with applications to genetic studies. *Ann. Appl. Stat.* 7, 369–390.
6. Loh, P.R., Tucker, G., Bulik-Sullivan, B.K., Vilhjálmsson, B.J., Finucane, H.K., Salem, R.M., Chasman, D.I., Ridker, P.M., Neale, B.M., Berger, B., et al. (2015). Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* 47, 284–290.
7. Chen, H., Wang, C., Conomos, M.P., Stilp, A.M., Li, Z., Sofer, T., Szpiro, A.A., Chen, W., Brehm, J.M., Celedón, J.C., et al. (2016). Control for population structure and relatedness for binary traits in genetic association studies via logistic mixed Mmodels. *Am. J. Hum. Genet.* 98, 653–666.
8. Breslow, N.E., and Clayton, D.G. (1993). Approximate inference in generalized linear mixed models. *J. Am. Stat. Assoc.* 88, 9–25.
9. Morgenthaler, S., and Thilly, W.G. (2007). A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutat. Res.* 615, 28–56.
10. Li, B., and Leal, S.M. (2008). Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *Am. J. Hum. Genet.* 83, 311–321.
11. Madsen, B.E., and Browning, S.R. (2009). A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet.* 5, e1000384.
12. Morris, A.P., and Zeggini, E. (2010). An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genet. Epidemiol.* 34, 188–193.
13. Wu, M.C., Lee, S., Cai, T., Li, Y., Boehnke, M., and Lin, X. (2011). Rare-variant association testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.* 89, 82–93.
14. Lee, S., Wu, M.C., and Lin, X. (2012). Optimal tests for rare variant effects in sequencing association studies. *Biostatistics* 13, 762–775.
15. Sun, J., Zheng, Y., and Hsu, L. (2013). A unified mixed-effects model for rare-variant association in sequencing studies. *Genet. Epidemiol.* 37, 334–344.
16. Pan, W., Kim, J., Zhang, Y., Shen, X., and Wei, P. (2014). A powerful and adaptive association test for rare variants. *Genetics* 197, 1081–1095.
17. Schifano, E.D., Epstein, M.P., Bielak, L.F., Jhun, M.A., Kardia, S.L., Peyser, P.A., and Lin, X. (2012). SNP set association analysis for familial data. *Genet. Epidemiol.* 36, 797–810.

18. Chen, H., Meigs, J.B., and Dupuis, J. (2013). Sequence kernel association test for quantitative traits in family samples. *Genet. Epidemiol.* 37, 196–204.
19. Ouallache, K., Dastani, Z., Li, R., Cingolani, P.E., Spector, T.D., Hammond, C.J., Richards, J.B., Ciampi, A., and Greenwood, C.M. (2013). Adjusted sequence kernel association test for rare variants controlling for cryptic and family relatedness. *Genet. Epidemiol.* 37, 366–376.
20. Wang, X., Lee, S., Zhu, X., Redline, S., and Lin, X. (2013). GEE-based SNP set association test for continuous and discrete traits in family-based association studies. *Genet. Epidemiol.* 37, 778–786.
21. Jiang, D., and McPeck, M.S. (2014). Robust rare variant association testing for quantitative traits in samples with related individuals. *Genet. Epidemiol.* 38, 10–20.
22. Yan, Q., Tiwari, H.K., Yi, N., Gao, G., Zhang, K., Lin, W.Y., Lou, X.Y., Cui, X., and Liu, N. (2015). A sequence kernel association test for dichotomous traits in family samples under a generalized linear mixed model. *Hum. Hered.* 79, 60–68.
23. Park, J.Y., Wu, C., Basu, S., McGue, M., and Pan, W. (2018). Adaptive SNP-Set association testing in generalized linear mixed models with application to family studies. *Behav. Genet.* 48, 55–66.
24. Zhan, X., Hu, Y., Li, B., Abecasis, G.R., and Liu, D.J. (2016). RVTESTS: an efficient and comprehensive tool for rare variant association analysis using sequence data. *Bioinformatics* 32, 1423–1426.
25. Liu, D.J., Peloso, G.M., Zhan, X., Holmen, O.L., Zawistowski, M., Feng, S., Nikpay, M., Auer, P.L., Goel, A., Zhang, H., et al. (2014). Meta-analysis of gene-level tests for rare variant association. *Nat. Genet.* 46, 200–204.
26. Feng, S., Pistis, G., Zhang, H., Zawistowski, M., Mulas, A., Zoledziwska, M., Holmen, O.L., Busonero, F., Sanna, S., Hveem, K., et al. (2015). Methods for association analysis and meta-analysis of rare variants in families. *Genet. Epidemiol.* 39, 227–238.
27. Lee, S., Teslovich, T.M., Boehnke, M., and Lin, X. (2013). General framework for meta-analysis of rare variants in sequencing association studies. *Am. J. Hum. Genet.* 93, 42–53.
28. Rice, K., Higgins, J.P., and Lumley, T. (2017). A re-evaluation of fixed effect(s) meta-analysis. *J. R. Stat. Soc. A* 181, 205–227.
29. MacCluer, J.W., VandeBerg, J.L., Read, B., and Ryder, O.A. (1986). Pedigree analysis by computer simulation. *Zoo Biol.* 5, 147–160.
30. Conomos, M.P., Laurie, C.A., Stilp, A.M., Gogarten, S.M., McHugh, C.P., Nelson, S.C., Sofer, T., Fernández-Rhodes, L., Justice, A.E., Graff, M., et al. (2016). Genetic diversity and association studies in US Hispanic/Latino populations: Applications in the Hispanic Community Health Study/Study of Latinos. *Am. J. Hum. Genet.* 98, 165–184.
31. Conomos, M.P., Miller, M.B., and Thornton, T.A. (2015). Robust inference of population structure for ancestry prediction and correction of stratification in the presence of relatedness. *Genet. Epidemiol.* 39, 276–293.
32. Morrison, A.C., Huang, Z., Yu, B., Metcalf, G., Liu, X., Ballantyne, C., Coresh, J., Yu, F., Muzny, D., Feofanova, E., et al. (2017). Practical approaches for whole-genome sequence analysis of heart- and blood-related traits. *Am. J. Hum. Genet.* 100, 205–215.
33. Brody, J.A., Morrison, A.C., Bis, J.C., O'Connell, J.R., Brown, M.R., Huffman, J.E., Ames, D.C., Carroll, A., Conomos, M.P., Gabriel, S., et al.; NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium; Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) Consortium; TOPMed Hematology and Hemostasis Working Group; and CHARGE Analysis and Bioinformatics Working Group (2017). Analysis commons, a team approach to discovery in a big-data environment for genetic epidemiology. *Nat. Genet.* 49, 1560–1563.
34. Huffman, J.E., de Vries, P.S., Morrison, A.C., Sabater-Lleal, M., Kacprowski, T., Auer, P.L., Brody, J.A., Chasman, D.I., Chen, M.H., Guo, X., et al. (2015). Rare and low-frequency variants and their association with plasma levels of fibrinogen, FVII, FVIII, and vWF. *Blood* 126, e19–e29.
35. Kircher, M., Witten, D.M., Jain, P., O'Roak, B.J., Cooper, G.M., and Shendure, J. (2014). A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* 46, 310–315.
36. Ionita-Laza, I., McCallum, K., Xu, B., and Buxbaum, J.D. (2016). A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat. Genet.* 48, 214–220.
37. Rogers, M.F., Shihab, H.A., Mort, M., Cooper, D.N., Gaunt, T.R., and Campbell, C. (2018). FATHMM-XF: accurate prediction of pathogenic point mutations via extended features. *Bioinformatics* 34, 511–513.
38. Lu, Q., Powles, R.L., Wang, Q., He, B.J., and Zhao, H. (2016). Integrative tissue-specific functional annotations in the human genome provide novel insights on many complex traits and improve signal prioritization in genome wide association studies. *PLoS Genet.* 12, e1005947.
39. He, Z., Xu, B., Lee, S., and Ionita-Laza, I. (2017). Unified sequence-based association tests allowing for multiple functional annotations and meta-analysis of noncoding variation in Metabochip data. *Am. J. Hum. Genet.* 101, 340–352.
40. Koziol, J.A., and Perlman, M.D. (1978). Combining independent chi-squared tests. *J. Am. Stat. Assoc.* 73, 753–763.
41. Wu, M.C., Maity, A., Lee, S., Simmons, E.M., Harmon, Q.E., Lin, X., Engel, S.M., Mouldrem, J.J., and Armistead, P.M. (2013). Kernel machine SNP-set testing under multiple candidate kernels. *Genet. Epidemiol.* 37, 267–275.
42. Ionita-Laza, I., Lee, S., Makarov, V., Buxbaum, J.D., and Lin, X. (2013). Sequence kernel association tests for the combined effect of rare and common variants. *Am. J. Hum. Genet.* 92, 841–853.
43. Su, Y.R., Di, C., Bien, S., Huang, L., Dong, X., Abecasis, G., Berndt, S., Bezieau, S., Brenner, H., Caan, B., et al. (2018). A mixed-effects model for powerful association tests in integrative functional genomics. *Am. J. Hum. Genet.* 102, 904–919.
44. Chen, J., Chen, W., Zhao, N., Wu, M.C., and Schaid, D.J. (2016). Small sample kernel association tests for human genetic and microbiome association studies. *Genet. Epidemiol.* 40, 5–19.
45. Zhou, J.J., Hu, T., Qiao, D., Cho, M.H., and Zhou, H. (2016). Boosting gene mapping power and efficiency with efficient exact variance component tests of single nucleotide polymorphism sets. *Genetics* 204, 921–931.
46. Dey, R., Schmidt, E.M., Abecasis, G.R., and Lee, S. (2017). A fast and accurate algorithm to test for binary phenotypes and its application to PheWAS. *Am. J. Hum. Genet.* 101, 37–49.
47. Zhou, W., Nielsen, J.B., Fritsche, L.G., Dey, R., Gabrielsen, M.E., Wolford, B.N., LeFaive, J., VandeHaar, P., Gagliano,

- S.A., Gifford, A., et al. (2018). Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat. Genet.* 50, 1335–1341.
48. Lumley, T., Brody, J., Peloso, G., Morrison, A., and Rice, K. (2018). FastSKAT: Sequence kernel association tests for very large sets of markers. *Genet. Epidemiol.* 42, 516–527.
  49. Bates, D., Maechler, M., Davis, T.A., Oehlschlägel, J., Riedy, J.; and R Core Team. (2018). Matrix: Sparse and Dense Matrix Classes and Methods. R package Version 1.2-14, <https://CRAN.R-project.org/package=Matrix>.
  50. Loh, P.R., Kichaev, G., Gazal, S., Schoech, A.P., and Price, A.L. (2018). Mixed-model association for biobank-scale datasets. *Nat. Genet.* 50, 906–908.