



Generalized multi-SNP mediation intersection–union test

Wujuan Zhong¹ | Toni Darville² | Xiaojing Zheng^{1,2} | Jason Fine^{1,4} | Yun Li^{1,3,5}

¹ Department of Biostatistics, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

² Department of Pediatrics, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

³ Department of Genetics, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

⁴ Department of Statistics and Operations Research, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

⁵ Department of Computer Science, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

Correspondence

Xiaojing Zheng, Department of Biostatistics, Department of Pediatrics, The University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA.

Email: Xiaojinz@email.unc.edu

Jason Fine, Department of Biostatistics, The University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA.

Email: jfine@email.unc.edu

Yun Li, Department of Biostatistics, The University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA.

Email: yunli@med.unc.edu

Funding information

National Institute of Allergy and Infectious Diseases, Grant/Award Numbers: R01 AI119164, U19 AI084024, U19 AI144181; National Institutes of Health, Grant/Award Numbers: U01 DA052713, U544 HD079124; National Institute of General Medical Sciences, Grant/Award Number: R01 GM105785; National Heart, Lung, and Blood Institute, Grant/Award Numbers: R01 HL129132, R01 HL146500

Abstract

To elucidate the molecular mechanisms underlying genetic variants identified from genome-wide association studies (GWAS) for a variety of phenotypic traits encompassing binary, continuous, count, and survival outcomes, we propose a novel and flexible method to test for mediation that can simultaneously accommodate multiple genetic variants and different types of outcome variables. Specifically, we employ the intersection–union test approach combined with the likelihood ratio test to detect mediation effect of multiple genetic variants via some mediator (e.g., the expression of a neighboring gene) on outcome. We fit high-dimensional generalized linear mixed models under the mediation framework, separately under the null and alternative hypothesis. We leverage Laplace approximation to compute the marginal likelihood of outcome and use coordinate descent algorithm to estimate corresponding parameters. Our extensive simulations demonstrate the validity of our proposed methods and substantial, up to 97%, power gains over alternative methods. Applications to real data for the study of *Chlamydia trachomatis* infection further showcase advantages of our methods. We believe our proposed methods will be of value and general interest in this post-GWAS era to disentangle the potential causal mechanism from DNA to phenotype for new drug discovery and personalized medicine.

KEYWORDS

intersection–union test, mediation analysis, multiple correlated SNPs, non-Gaussian outcome

1 | INTRODUCTION

Dissection of mediation pathways underlying genetic association will enhance understanding of disease mechanisms and biomarker development. An example is *Chlamydia trachomatis* infection. Chlamydia is the leading bacterial

sexually transmitted infection in the United States (Centers for Disease Control and Prevention, 2019). Infection is often asymptomatic and after ascending to the upper genital tract may cause severe reproductive morbidities in women. Repeated infection leads to worse disease. Host genetics shapes susceptibility to chlamydia disease and/or

reinfection (Bailey *et al.*, 2009; Taylor *et al.*, 2017; Zheng *et al.*, 2018). DNA biomarkers for susceptibility to ascension or risk of reinfection are critically needed for targeted screening for women at high risk of disease and vaccine development. Genome-wide association studies (GWAS) provide candidate loci, but lack mechanistic interpretations. Although expression quantitative trait loci (eQTL) mapping can provide mechanistic hypotheses, GWAS and eQTL both only analyze two sources of data. There is a significant unmet need for simultaneous modeling of all three sources of data (namely, genetic variants, gene expression, and final outcome) by directly testing the mediation effects of multiple correlated single nucleotide polymorphisms (SNPs) via the expression of some gene (e.g., eGene associated with the eQTL SNP) on chlamydia ascension (binary outcome) and reinfection (time-to-event outcome).

Mediation analysis was first proposed by Baron and Kenny to study the association between an independent variable and an outcome by adding an intermediate variable, which is called the mediator (Baron and Kenny, 1986). In genetics and genomics studies, researchers are interested in testing mediation effects of the genetic variant(s) on the outcome through a certain mediator (e.g., the expression level of a neighboring gene). Non-Gaussian outcomes, such as binary, count, and time-to-event outcomes (e.g., disease status, time until death), are commonly present in mediation analyses but have been understudied. Huang *et al.* developed mixed model-based methods that can handle binary and time-to-event outcomes but assume *a priori* that the genetic variants under testing are eQTLs (Huang *et al.*, 2015, 2016).

We have previously proposed a method, multi-SNP mediation intersection-union test (SMUT), to assess mediation effect of high-dimensional genetic variants on any continuous outcome (Zhong *et al.*, 2019). To the best of our knowledge, none of the existing methods can jointly test mediation effects of multiple correlated SNPs (not necessarily all eQTLs) on a non-Gaussian outcome. Here, we propose a generalized multi-SNP mediation intersection-union test to evaluate mediation effects of multiple correlated SNPs on a non-Gaussian outcome without prior knowledge of eQTLs. Both SMUT and methods proposed in this work are extensions of Baron and Kenny's framework and leverage intersection-union test (IUT) (Berger and Hsu, 1996) to decompose mediation into two separate regression models. While our earlier SMUT method handles only Gaussian outcome, methods proposed here allow non-Gaussian outcomes by adopting the generalized linear mixed model (GLMM) (McCulloch *et al.*, 2008) or the mixed effects Cox proportional hazards (PH) model (Vaida and Xu, 2000; Pankratz *et al.*, 2005). More details germane to the differences between SMUT

and methods proposed here are given in Supporting Information Section 1. For presentation brevity, we hereafter refer to our method for a binary or count outcome as SMUT_GLM; while that for a time-to-event outcome as SMUT_PH.

The rest of this article is organized as follows. In Section 2, we present details of our proposed methods SMUT_GLM and SMUT_PH, followed by simulation studies and real data application in Sections 3 and 4, respectively. Finally, Section 5 concludes the article with some discussions.

2 | METHODS

2.1 | Notation

Without loss of generality, we assume that we have four types of data, namely, genotypes (as the potential causal variables), gene expression measurements (as the mediator, which can be other types of molecular measures such as metabolite levels or protein abundances), phenotypic trait (as the final outcome), and other covariates (e.g., age, gender). Let \mathbf{G} be the $n \times q$ genotype matrix, where n is the sample size, q is the number of SNPs, and G_{ij} is the number of copies of the minor allele for the i th individual at the j th SNP. Let \mathbf{X} be the $n \times p$ covariate matrix and X_{ij} denote the j th covariate variable for the i th individual. Let $\mathbf{M} = (M_1, M_2, \dots, M_n)^T$ and $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$, where M_i and Y_i denote the mediator and the outcome for the i th individual, respectively. If Y_i is a binary or count outcome, Y_i is related to the model in Equation (2); if Y_i is a time-to-event outcome, Y_i is related to the model in Equation (3) and $Y_i = (Z_i, \delta_i)$ where $Z_i = \min(T_i, C_i)$ is the observation time, T_i is the failure time, C_i is the censoring time, and $\delta_i = I(T_i \leq C_i)$ is the failure indicator; $\delta_i = 1$ indicates that the failure is observed, and $\delta_i = 0$ indicates that the response is censored. We apologize for abusing notations. Basically, we want to use the same notation Y_i to denote different types of outcomes.

2.2 | SMUT_GLM and SMUT_PH model

SMUT_GLM and SMUT_PH model the effects of SNPs on the outcome mediated by the expression level of a single gene via two models, namely a mediator model and an outcome model. We assume the expression level is continuous and consider a linear model for the mediator model (Equation 1). As for the outcome model, we fit GLMM if the outcome conditional on SNPs' effects follows an exponential family distribution (Equation 2); we fit mixed effects Cox PH model if the outcome is a time-to-event variable

(Equation 3).

$$\text{Mediator model: } M_i = \alpha_1 + \sum_{j=1}^p X_{ij}t_j^M + \sum_{j=1}^q G_{ij}\beta_j + \epsilon_i \quad (1)$$

Exponential family outcome model:

$$g\{E(Y_i|\boldsymbol{\gamma})\} = \alpha_2 + M_i\theta + \sum_{j=1}^p X_{ij}t_j + \sum_{j=1}^q G_{ij}\gamma_j \quad (2)$$

Survival outcome model:

$$\lambda(t_i) = \lambda_0(t_i) \exp\left(M_i\theta + \sum_{j=1}^p X_{ij}t_j + \sum_{j=1}^q G_{ij}\gamma_j\right) \quad (3)$$

where α_1, α_2 are fixed intercepts; fixed effects $\boldsymbol{t}^M = (t_1^M, t_2^M, \dots, t_p^M)^T$ and $\boldsymbol{t} = (t_1, t_2, \dots, t_p)^T$ are vectors of covariates' effects on the mediator and outcome, respectively; random effect $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_q)^T$ is a vector of SNPs' effects on the mediator; fixed effect θ is the mediator's effect on the outcome. The random effect $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \dots, \gamma_q)^T$ is a vector of SNPs' effects on the outcome; error terms $\epsilon_1, \epsilon_2, \dots, \epsilon_n \sim_{\text{i.i.d.}} N(0, \sigma^2)$; g is the link function; $\lambda(t_i)$ is the hazard function; $\lambda_0(t_i)$ is an unspecified baseline hazard function. We have showed that the hypotheses $H_0 : \boldsymbol{\beta}\theta = \mathbf{0}$ versus $H_1 : \boldsymbol{\beta}\theta \neq \mathbf{0}$ are valid for testing the mediation effect in Supporting Information Section 8, where $\boldsymbol{\beta}\theta \neq \mathbf{0}$ implies that SNPs exert mediation effects on the outcome. Following our previous work (Zhong *et al.*, 2019), we employ IUT to decompose the hypothesis testing $H_0 : \boldsymbol{\beta}\theta = \mathbf{0}$ versus $H_1 : \boldsymbol{\beta}\theta \neq \mathbf{0}$ into two subhypotheses $H_0^\beta : \boldsymbol{\beta} = \mathbf{0}$ versus $H_1^\beta : \boldsymbol{\beta} \neq \mathbf{0}$ and $H_0^\theta : \theta = 0$ versus $H_1^\theta : \theta \neq 0$, such that $H_0 = H_0^\beta \cup H_0^\theta$ and $H_1 = H_1^\beta \cap H_1^\theta$. Suppose the p values for testing $\boldsymbol{\beta}$ and θ being zero are p_1 and p_2 , respectively. Then the p value for testing $\boldsymbol{\beta}\theta$ being zero, using IUT, is the maximum of p_1 and p_2 . In the following sections, we provide details regarding how to separately test $\boldsymbol{\beta}$ and θ to obtain p_1 and p_2 .

2.3 | Testing $\boldsymbol{\beta}$ in the mediator model and θ in the outcome model

As in Zhong *et al.* (2019), we adopt the widely used sequence kernel association test (SKAT) method (Wu *et al.*, 2011) to test $\boldsymbol{\beta}$ in the mediator model to accommodate a potentially large number of correlated SNPs.

Our strategy for testing θ in the outcome model consists of four steps: (1) formulation of the likelihood function based on the nature of the outcome random variable

\boldsymbol{Y} , (2) Laplace approximation of the likelihood function, (3) application of the coordinate descent algorithm (Fu, 1998; Daubechies *et al.*, 2004) to estimate parameters by maximizing the approximated likelihood function, and (4) calculation of the likelihood ratio statistic. These four steps allow us to test the mediator effect θ in the outcome model.

2.4 | Likelihood function for the outcome model

To reduce the dimensionality of parameters in the outcome model, we adopted a linear mixed model for continuous outcome in our previous work (Zhong *et al.*, 2019). We assume Y_1, Y_2, \dots, Y_n are independent and identically distributed. When the outcome $Y_i (i = 1, 2, \dots, n)$ conditional on $\boldsymbol{\gamma}$ follows an exponential family distribution, we adopt the GLMM in Equation (2).

$$\begin{cases} \gamma_j \sim_{\text{i.i.d.}} N(0, \sigma_\gamma^2) \\ g(\mu_i) = \eta_i = \alpha_2 + M_i\theta + \sum_{j=1}^p X_{ij}t_j + \sum_{j=1}^q G_{ij}\gamma_j \\ E(Y_i|\boldsymbol{\gamma}) = \mu_i \\ L(\boldsymbol{y}|\boldsymbol{\gamma}) = \prod_{i=1}^n \exp\left\{\frac{y_i\tau_i - b(\tau_i)}{a(\phi)} + C(y_i, \phi)\right\} \end{cases}, \quad (4)$$

where τ_i is the canonical parameter; ϕ is the dispersion parameter; $L(\boldsymbol{y}|\boldsymbol{\gamma})$ is the likelihood function of the outcome \boldsymbol{Y} conditional on $\boldsymbol{\gamma}$. When the outcome $Y_i (i = 1, 2, \dots, n)$ is a time-to-event variable, we adopt the mixed effects Cox PH model in Equation (3).

$$\begin{cases} \gamma_j \sim_{\text{i.i.d.}} N(0, \sigma_\gamma^2) \\ \eta_i = M_i\theta + \sum_{j=1}^p X_{ij}t_j + \sum_{j=1}^q G_{ij}\gamma_j \\ \lambda(t_i) = \lambda_0(t_i) \exp \eta_i \\ PL = \prod_{i=1}^n \left(\frac{\exp \eta_i}{\sum_{k \in R_i} \exp \eta_k}\right)^{\delta_i} \end{cases}, \quad (5)$$

where $R_i = \{k : Z_k \geq Z_i\}$ is the risk set and PL is the partial likelihood function conditional on $\boldsymbol{\gamma}$. For the GLMM in (4), $\ell(\boldsymbol{y}|\boldsymbol{\gamma})$ denotes $\log\{L(\boldsymbol{y}|\boldsymbol{\gamma})\}$ and $L(\boldsymbol{y})$ denotes the likelihood function of the outcome unconditional on $\boldsymbol{\gamma}$; for the mixed effects Cox PH model in (5), $\ell(\boldsymbol{y}|\boldsymbol{\gamma})$ denotes $\log(PL)$ and $L(\boldsymbol{y})$ denotes the partial likelihood of the outcome unconditional on $\boldsymbol{\gamma}$. We again apologize for abusing notations. Our basic rationale is to employ the same notation $\ell(\boldsymbol{y}|\boldsymbol{\gamma})$ and $L(\boldsymbol{y})$ to denote different log-likelihood and likelihood functions, respectively, for different types

of outcomes. Let $f_{\gamma}(\boldsymbol{y})$ be the probability density function of \boldsymbol{y} , and $f_{\gamma}(\boldsymbol{y}) = (2\pi\sigma_{\gamma}^2)^{-\frac{q}{2}} \exp(-\frac{1}{2\sigma_{\gamma}^2}\boldsymbol{y}^T\boldsymbol{\gamma})$. Then we have the following:

$$L(\boldsymbol{y}) = \int_{R^q} \exp\{\ell(\boldsymbol{y}|\boldsymbol{\gamma})\} f_{\gamma}(\boldsymbol{y}) d\boldsymbol{\gamma} = (2\pi\sigma_{\gamma}^2)^{-\frac{q}{2}} \times \int_{R^q} \exp\{h(\boldsymbol{\gamma})\} d\boldsymbol{\gamma}, \quad (6)$$

where $h(\boldsymbol{\gamma}) = \ell(\boldsymbol{y}|\boldsymbol{\gamma}) - \frac{1}{2\sigma_{\gamma}^2}\boldsymbol{y}^T\boldsymbol{\gamma}$. Technical details are described in Supporting Information Section 2.1.

2.5 | Laplace approximation

Laplace's method is widely adopted to approximate the likelihood function (Breslow and Clayton, 1993; Raudenbush *et al.*, 2000; Pankratz *et al.*, 2005). The integral in Equation (6) can be approximated via Laplace's method by taking the Taylor expansion to the second order of $h(\boldsymbol{\gamma})$ around its maximum point $\tilde{\boldsymbol{\gamma}}$. After inserting the Taylor expansion into the integral, and taking logarithm, we have the approximated log-likelihood f .

$$\log\{L(\boldsymbol{y})\} \approx f = -\frac{q}{2} \log \sigma_{\gamma}^2 + h(\tilde{\boldsymbol{\gamma}}) - \frac{1}{2} \log |-h''(\tilde{\boldsymbol{\gamma}})|. \quad (7)$$

For the GLMM in (4), we have

$$h''(\boldsymbol{\gamma}) = \frac{\partial^2 h}{\partial \boldsymbol{\gamma} \partial \boldsymbol{\gamma}^T} = -(\boldsymbol{G}^T \boldsymbol{W} \boldsymbol{G} + \sigma_{\gamma}^{-2} \boldsymbol{I}_q), \quad (8)$$

where \boldsymbol{I}_q is a $q \times q$ identity matrix, $\boldsymbol{W} = \text{diag}(w_1, w_2, \dots, w_n)$, and w_i is recognizable as GLM (generalized linear model) iterative weight. For the mixed effects Cox PH model in (5), we have

$$h''(\boldsymbol{\gamma}) = \frac{\partial^2 h}{\partial \boldsymbol{\gamma} \partial \boldsymbol{\gamma}^T} = -(\boldsymbol{U} + \sigma_{\gamma}^{-2} \boldsymbol{I}_q), \quad (9)$$

where $\boldsymbol{U} = (u_{j_1 j_2})$, $u_{j_1 j_2} = -\frac{\partial^2 (\log PL)}{\partial \gamma_{j_1} \partial \gamma_{j_2}}$. More details of the Laplace approximation are given in Supporting Information Section 2.2.

2.6 | Coordinate descent algorithm

We apply the coordinate descent algorithm to maximize the approximated log-likelihood in Equation (7). Note that $\tilde{\boldsymbol{\gamma}}$ in Equation (7) is a function of other parameters $\boldsymbol{\xi} = (\alpha_2, \sigma_{\gamma}^2, \phi, \theta, t_1, t_2, \dots, t_p)$. Instead of taking implicit differentiation of $\tilde{\boldsymbol{\gamma}}$ (Raudenbush *et al.*, 2000), we use the approximation strategy proposed in Schelldorfer *et al.* (2014),

which regards $\tilde{\boldsymbol{\gamma}}$ as fixed when updating $\boldsymbol{\xi}$. This strategy is computationally convenient and efficient at little cost of reduced accuracy. In addition, we take further approximation when taking derivatives of the approximated log-likelihood function f . Specifically, for the GLMM in (4), we assume \boldsymbol{W} in Equation (8) varies slowly as a function of $(\mu_1, \mu_2, \dots, \mu_n)^T$ (Breslow and Clayton, 1993). For the mixed effects Cox PH model in (5), we similarly assume that \boldsymbol{U} in Equation (9) varies slowly as a function of $(\eta_1, \eta_2, \dots, \eta_n)^T$. Under the assumption, the term $-\frac{1}{2} \log |-h''(\tilde{\boldsymbol{\gamma}})|$ in Equation (7) is ignored when taking derivatives of the approximated log-likelihood function over $(\alpha_2, \phi, \theta, t_1, t_2, \dots, t_p)$. Details of the coordinate descent algorithm are given in Supporting Information Section 2.3. Finally, we employ the Newton-Raphson algorithm to sequentially update each parameter.

2.7 | Likelihood ratio test

We obtain approximated likelihood under the null and the alternative hypothesis separately, denoted by L_0 and L_1 , respectively. For GLMM, the likelihood ratio statistic $2\{\log(L_1) - \log(L_0)\}$ asymptotically follows a chi-square distribution with one degree of freedom and similarly for the partial likelihood ratio statistics for the survival outcome.

3 | SIMULATION STUDIES

3.1 | Simulation settings

To evaluate the performance of SMUT_GLM and SMUT_PH in comparison with alternative methods, we conducted extensive simulations to investigate power and type-I error. Following our previous work (Zhong *et al.*, 2019), we simulated a dataset of 10,000 pseudoindividuals measured at 2891 SNPs with minor allele frequency (MAF) $\geq 1\%$ in a 1-Mb region using the COSI coalescent model (Schaffner *et al.*, 2005) to generate realistic genetic data. The 10,000 pseudoindividuals were constructed by randomly pairing up 20,000 simulated chromosomes without replacement. To evaluate power and type-I error, we generated 500 datasets with 1000 samples each by sampling without replacement from the entire pool of 10,000 samples simulated above. We randomly selected a set of causal SNPs, which is shared across the 500 simulated datasets, from these 2891 SNPs. We then classified them into three categories: shared SNPs (*sSNPs*), mediator specific SNPs (*mSNPs*), and outcome specific SNPs (*oSNPs*). The *sSNPs* influence both the

TABLE 1 Causal SNP composition in two simulated scenarios.

Type of outcome	Sample size	Sparse or dense	Number of causal SNPs	Number of sSNPs	Number of mSNPs	Number of oSNPs	Number of Noncausal SNPs
Binary or count	1000	Sparse	10	4	3	3	890
		Dense	500	300	100	100	400
Time-to-event	200	Sparse	10	4	3	3	190
		Dense	150	90	30	30	50

The sparse(dense) scenario is to simulate datasets based on a small(large) number of causal SNPs. Causal SNPs are the union of shared SNPs, mediator specific SNPs, and outcome specific SNPs. Shared SNPs have effects on both mediator and outcome. Mediator(outcome)-specific SNPs have effects only on mediator(outcome). All these SNPs are randomly selected from the 2891 SNPs with MAF $\geq 1\%$.

mediator and the outcome, while the *mSNPs* and *oSNPs* only contribute to the mediator and outcome, respectively.

We considered two scenarios in terms of causal SNP density: sparse and dense (Table 1). For binary or count outcome, sample size is 1000 and there are 10 and 500 causal SNPs for sparse and dense scenarios, respectively. For the time-to-event outcome, sample size is 200 and there are 10 and 150 causal SNPs for sparse and dense scenarios, respectively. When we fit the model, both the causal and noncausal SNPs (Table 1) are included in the model. Thus, the distribution of coefficients of genetic variants is effectively misspecified for all the simulations. Covariates matrix \mathbf{X} consists of a continuous variable generated from $N(0, 1)$ and a binary variable generated from *Bernoulli*(0.5). We generated the mediator via $M_i = \alpha_1 + (\mathbf{G}_i^{sm})^T \boldsymbol{\beta} + (\mathbf{X}_i)^T \boldsymbol{\iota}^M + \varepsilon_i$, where \mathbf{G}_i^{sm} denotes the vector of genotype data for the *i*th individual from *sSNPs* and *mSNPs*, \mathbf{X}_i denotes the vector of the covariates for the *i*th individual, $\alpha_1 = 1$, $\boldsymbol{\iota}^M = (0.5, -0.5)^T$, $\boldsymbol{\beta} \sim c_\beta N(\mathbf{0}, \mathbf{I}_q)$, and c_β is a scalar to scale the SNPs' effects; $\varepsilon_i \sim N(0, 1)$. We generated the binary or count outcome via $g\{E(Y_i|\boldsymbol{\gamma})\} = \alpha_2 + M_i \theta + (\mathbf{G}_i^{so})^T \boldsymbol{\gamma} + (\mathbf{X}_i)^T \boldsymbol{\iota}$, where \mathbf{G}_i^{so} denotes the vector of genotype data for the *i*th individual from *sSNPs* and *oSNPs*, $\alpha_2 = 0$, $\boldsymbol{\iota} = (0.5, -0.5)^T$, $\boldsymbol{\gamma} \sim c_\gamma N(\mathbf{0}, \mathbf{I}_q)$, and $c_\gamma = 0.2$. The link function *g* was specific to the type of the outcome (Supporting Information Section 2.1). We generated the time-to-event outcome based on the Weibull baseline hazard via $t_i = \left[-\frac{\log(v)}{\lambda \exp\{M_i \theta + (\mathbf{G}_i^{so})^T \boldsymbol{\gamma} + (\mathbf{X}_i)^T \boldsymbol{\iota}\}} \right]^{\frac{1}{\rho}}$ and $c_i \sim \text{Exp}(0.001)$, where t_i is failure time and c_i is censoring time, $v \sim \text{Unif}(0, 1)$, shape $\rho = 1$, scale parameter $\lambda = 0.01$. Note that across the 500 datasets, error terms $\boldsymbol{\epsilon}$ were separately simulated for each dataset, but $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ were fixed.

In the simulations, we tested the mediation effects of these SNPs on the binary, count, or time-to-event outcome using SMUT_GLM and SMUT_PH, as well as other methods including SMUT, adapted least absolute shrinkage and selection operator (LASSO) (Tibshirani, 1996), and adapted Huang *et al.*'s method. In order to compare the performance of approximations that we adopted, we considered two versions of our method, both treat-

ing $\tilde{\boldsymbol{\gamma}}$ as fixed: (1) based on exact derivatives; (2) based on approximated derivatives. For a binary or count outcome, we refer to these two versions as SMUT_GLM exact and SMUT_GLM approxi. For a time-to-event outcome, we refer to the approximated version as SMUT_PH approxi. The exact version of SMUT_PH is not employed because it is hard to derive analytically. SMUT is naively applied to binary and count outcomes by treating them as continuous variables. The adapted LASSO approach adopts SKAT to consider all the genetic variant in the mediator model, while in the outcome model, employs LASSO for variable selection on all genetic variants as well as mediator and covariates, then refits GLM on the selected genetic variants together with mediator and covariates (latter two will be included regardless of LASSO variable selection results), and finally combines *p* values from the mediator and the refitted outcome model via IUT. The adapted Huang *et al.*'s method employs SKAT in the mediator model, adopts the original Huang *et al.*'s method in the outcome model, and then combines *p* values from the two models via IUT. We use adapted LASSO and SKAT + LASSO exchangeably. Similarly, we use adapted Huang *et al.* and SKAT + Huang *et al.* exchangeably. Details of the adapted LASSO and adapted Huang *et al.*'s method are provided in Supporting Information Section 3.

To test the robustness and generalizability of the methods, we considered two alternative situations where some assumptions are violated. The first situation is the violation of the assumption that coefficients of genetic variants follow a Gaussian distribution. The second situation is when there is an unobserved mediator that is not adjusted in the analysis. Details and results of these two simulation studies are given in Supporting Information Section 4.

3.2 | Type-I error in simulations

We evaluated the validity of SMUT_GLM and SMUT_PH along with alternative methods in simulations.

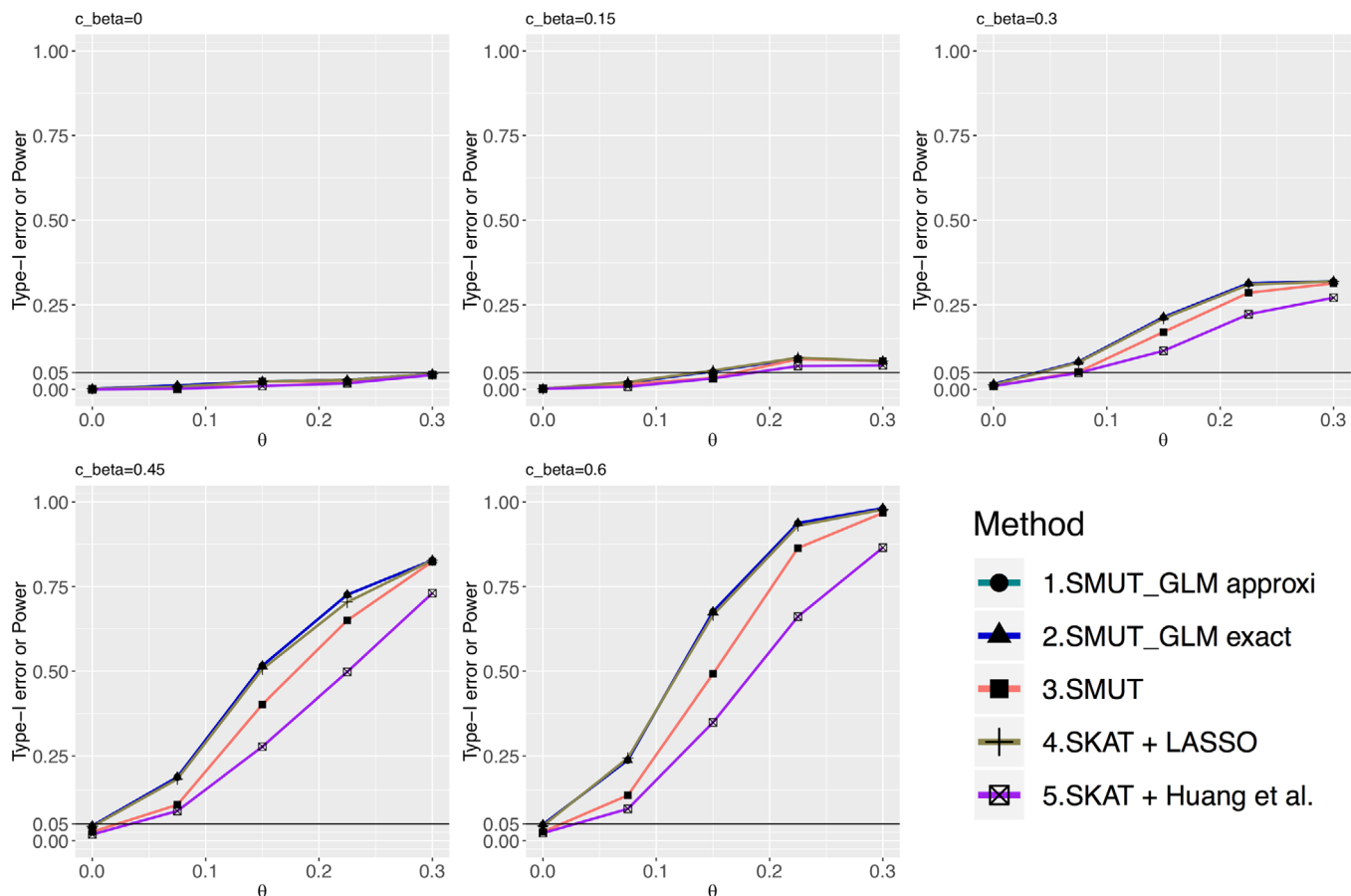


FIGURE 1 For binary outcome, power and type-I error under sparse causal SNPs scenario. The x-axis is the true mediator effect (θ) on the outcome. The y-axis is the power or type-I error. Subfigures vary in the c_β value. $c_\beta = 0$ (top-left subfigure) or $\theta = 0$ (left-most points in each subfigure) are null settings where the y-axis represents the corresponding type-I error. When $c_\beta \neq 0$ and $\theta \neq 0$, it is under alternative hypothesis and the y-axis represents the corresponding power. Line for the approximated version of SMUT_GLM is overlapped with the exact version. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

SMUT_GLM and SMUT_PH exhibited controlled type-I error rates, at $\alpha = 0.05$ level, regardless of causal SNP density and types of outcome, as shown in Figures 1 and 2 for binary outcome in sparse and dense scenarios, respectively, Figures 3 and 4 for the time-to-event outcome in sparse and dense scenarios, respectively, and Web Figures S1 and S2 for count outcome in sparse and dense scenarios, respectively. In each figure, the first panel ($c_\beta = 0$) and the leftmost point ($\theta = 0$) in other panels ($c_\beta \neq 0$) all correspond to the null of no mediation of the SNPs through the mediator. SMUT, adapted LASSO, and adapted Huang *et al.*'s method also showed protected type-I error.

3.3 | Power in simulations

SMUT_GLM and SMUT_PH demonstrated substantial power gains under both the sparse and dense scenar-

ios. We also observed that the approximated version of SMUT_GLM demonstrated very similar performance when compared with its exact counterpart. For example, for binary outcome and under the scenario of dense causal SNPs when $c_\beta = 0.6, \theta = 0.1$, exact SMUT_GLM, approximated SMUT_GLM, SMUT, adapted LASSO, and adapted Huang *et al.* had 97%, 96%, 17%, 54%, and 0% power, respectively. Thus, the power gain from the exact SMUT_GLM was 80%, 43%, and 97% compared with SMUT, adapted LASSO, and adapted Huang *et al.*, respectively. The approximated SMUT_GLM had similar power gains. For the time-to-event outcome, under the scenario of dense causal SNPs when $c_\beta = 1, \theta = 0.075$, approximated SMUT_PH and adapted LASSO had 69% and 41% power, respectively, leading to a power gain of 28%. In addition, power gains appeared more profound with increasing c_β likely because adapted LASSO and adapted Huang *et al.* becomes more conservative as the pleiotropy effect of SNPs on mediator and outcome (measured by c_β) increases.

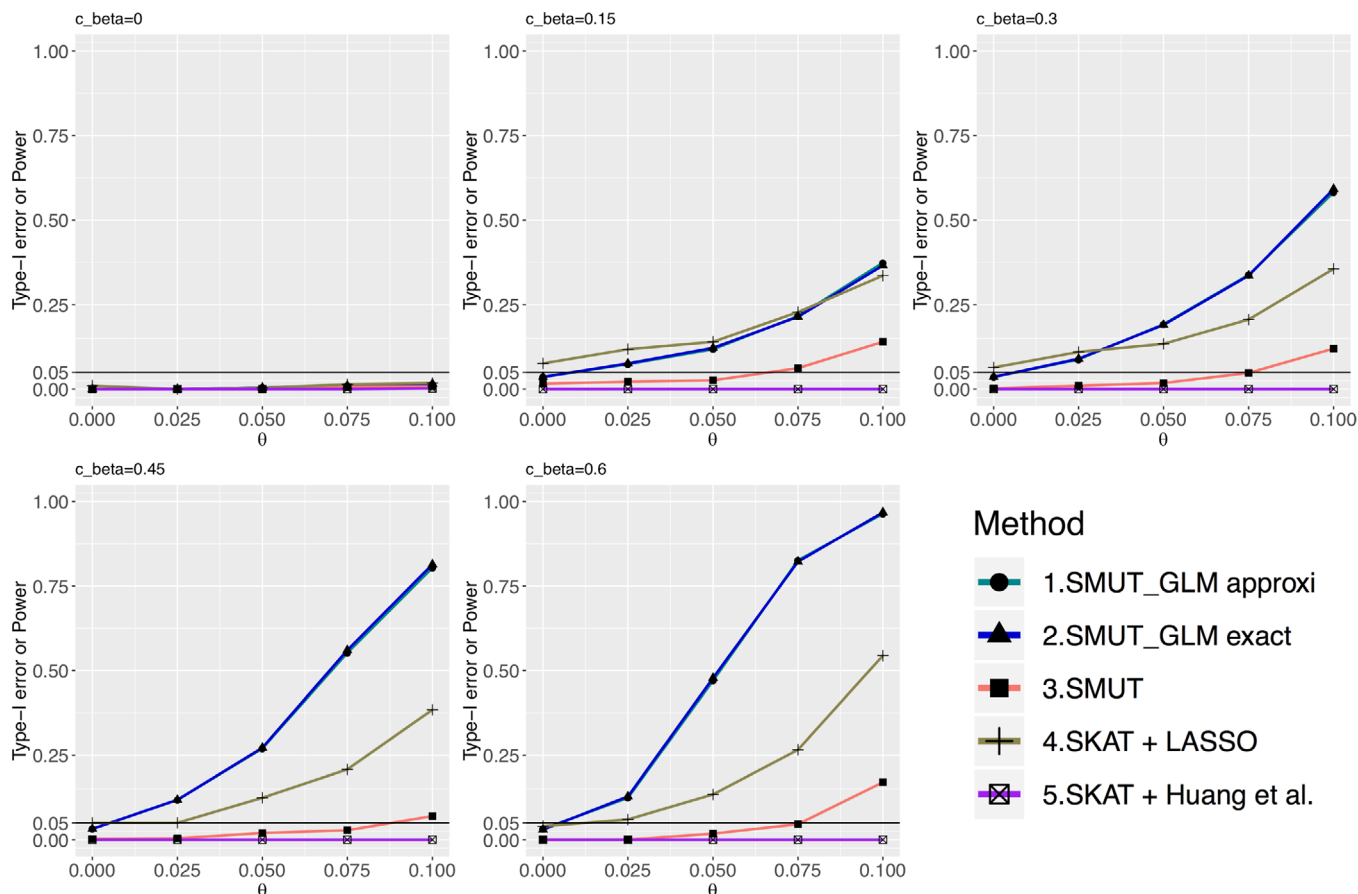


FIGURE 2 For binary outcome, power and type-I error under the dense causal SNPs scenario. The x - and y -axes are the same as in Figure 1. Line for the approximated version of SMUT_GLM is overlapped with the exact version. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

4 | REAL DATA APPLICATION

We assessed our methods and alternatives in real data from two clinical cohorts, which were designed for the study of chlamydia infection. *Chlamydia trachomatis* can ascend from the cervix to the uterus and fallopian tubes in some women, potentially resulting in pelvic inflammatory disease and severe reproductive morbidities, including infertility and ectopic pregnancy. Recurrent infection leads to a worse disease. We analyzed genotype, gene expression, and phenotype data of 200 participants combined from two cohorts, the Anaerobes and Clearance of Endometritis (ACE) cohort and the T cell Response Against Chlamydia (TRAC) cohort (Russell *et al.*, 2015). The Institutional Review Boards for Human Subject Research at the University of Pittsburgh and the University of North Carolina approved the study, and all participants provided written informed consent prior to inclusion. Descriptions of the ACE and TRAC cohorts, processing and quality control of genotype and gene expression data, and details of eQTL

analysis and mediation analysis of other genes are available in Supporting Information Section 6.

4.1 | Binary outcome

The outcome of interest is ascending chlamydia infection, among participants who had chlamydia infection at enrollment. The control group is the 71 participants who had chlamydia infection restricted to the cervix, and the case group is the 72 participants with both cervical and endometrial chlamydia infection at enrollment. We analyzed genotype, gene expression, and phenotype data from these 143 participants.

Here we presented *SOS1* and *CD151* gene, which were biologically related to the outcome, to illustrate the application of our proposed methods to a binary outcome. SOS Ras/Rac Guanine Nucleotide Exchange Factor 1 (*SOS1*) is a guanine nucleotide exchange factor that in humans is encoded by the *SOS1* gene. The importance of *SOS1* for

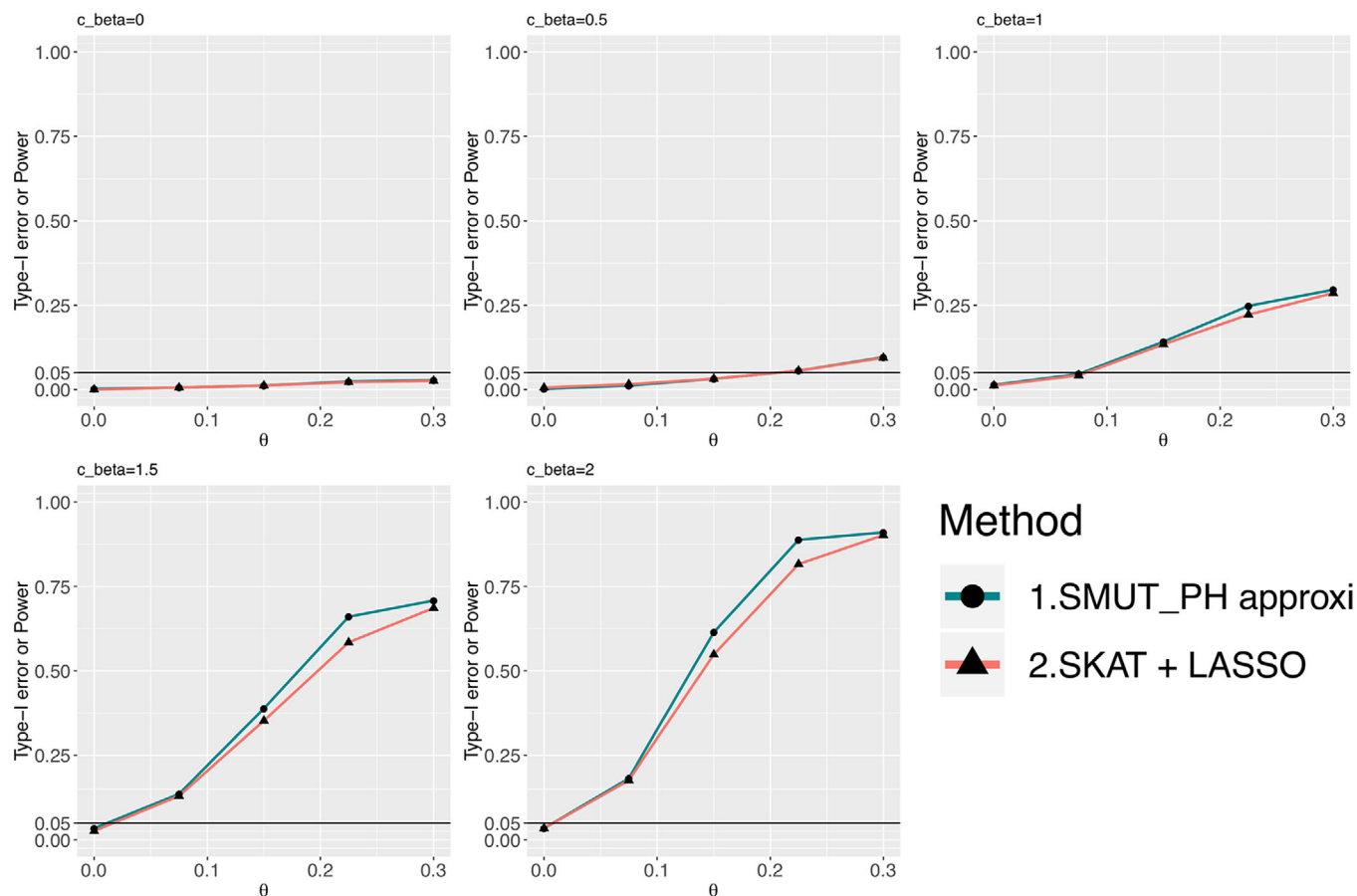


FIGURE 3 For the time-to-event outcome, power and type-I error under sparse causal SNPs scenario. The x - and y -axes are the same as in Figure 1. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

TABLE 2 Real data application to TRAC and ACE datasets.

Type of outcome	Gene	Probesets	Number of SNPs	p Values		
				SMUT_GLM	LASSO	Huang <i>et al.</i>
Binary	<i>SOS1</i>	2140519	83	0.0235	0.0691	0.0229
Binary	<i>CD151</i>	1940132	40	0.0245	0.1192	0.2289
Time-to-event	<i>BIRC3</i>	7210154	4	SMUT_PH 0.001	LASSO 0.001	Huang <i>et al.</i> 0.002

chlamydia invasion of host cells has been indicated by multiple biomedical studies (Carabeo *et al.*, 2007; Lane *et al.*, 2008; Hackstadt, 2012; Bastidas *et al.*, 2013; Mehlitz and Rudel, 2013; Elwell *et al.*, 2016). The *CD151* gene encodes a protein that is known to complex with integrins. It promotes cell adhesion and may regulate integrin trafficking and/or function. It is a member of the tetraspanin family, which is considered as the gateway for infection (Hauck and Meyer, 2003; Hemler, 2008; Hassuna *et al.*, 2009; Join-Lambert *et al.*, 2010; Monk and Partridge, 2012; Seu *et al.*, 2017). In addition, SNPs annotation database, RegulomeDB (Boyle *et al.*, 2012), demonstrates that some

SNPs in these two genes are eQTLs with experimental evidence. Thus, the presence of mediation effect via the expression of each gene is expected.

For the first gene *SOS1*, mediation testing encompassed 83 SNPs with $MAF \geq 10\%$ and significant eQTL association (with *SOS1*) at a false discovery rate (FDR) threshold of 10%, using SMUT_GLM, adapted LASSO and adapted Huang *et al.*'s method. Both SMUT_GLM and adapted Huang *et al.*'s method detected significant mediation effects, while adapted LASSO did not (Table 2). For the second gene *CD151*, our mediation (via expression of *CD151*) testing involved 40 SNPs with $MAF \geq 10\%$ and significant

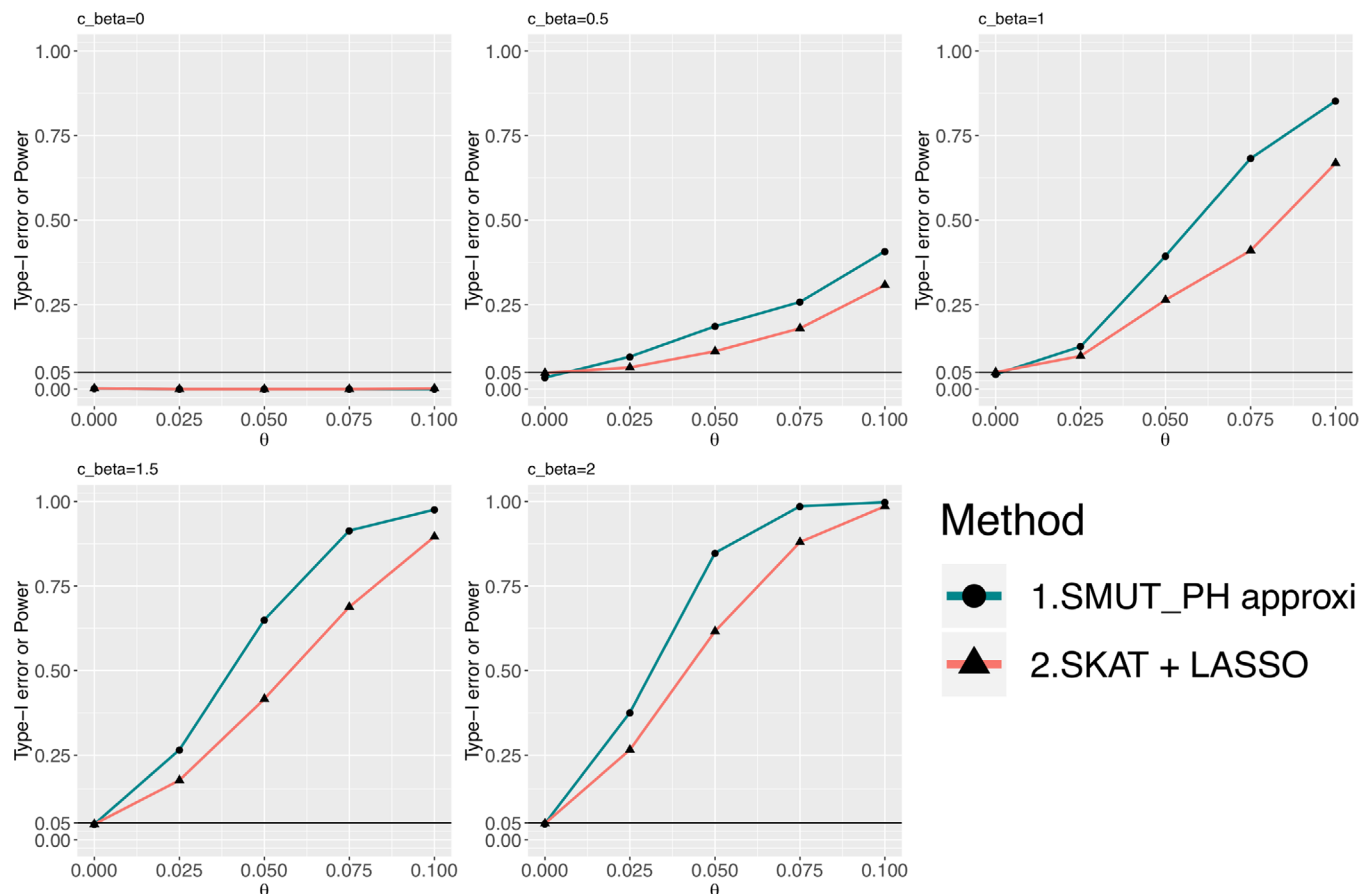


FIGURE 4 For the time-to-event outcome, power and type-I error under dense causal SNPs scenario. The x - and y -axes are the same as in Figure 1. This figure appears in color in the electronic version of this article, and any mention of color refers to that version.

eQTL (with *CDI51*) at FDR 10%. Only SMUT_GLM showed significant mediation effects of these SNPs through the expression of *CDI51* on ascending chlamydia infection (Table 2). Marginal effects of selected SNPs on *SOS1* and *CDI51* gene expression and ascending chlamydia infection are visually illustrated in Web Figures S19 and S20, respectively.

4.2 | Time-to-event outcome

TRAC participants returned for follow-up visits at 1, 4, 8, and 12 months after enrollment. The outcome of interest we evaluated here is time to the first incident chlamydia infection. We analyzed genotype, gene expression, and time-to-event data from all 181 participants in the TRAC cohort who had both genotype and gene expression data available.

Here we selected the *BIRC3* gene, which was biologically related to the outcome, to illustrate the application of our proposed methods to a time-to-event outcome. The

gene *BIRC3* encodes for Baculoviral IAP repeat containing 3, a E3 ubiquitin-protein ligase regulating NF- κ B signaling (Blankenship *et al.*, 2009; Kim *et al.*, 2010; Tan *et al.*, 2013). It acts as an important regulator of pathogen recognition receptor signaling (Bertrand *et al.*, 2009), which can have profound effects on the development of downstream adaptive immune responses (Takeda *et al.*, 2003; Palm and Medzhitov, 2009; Kumar *et al.*, 2011). In addition, biological studies suggested that *BIRC3* may protect mammalian host cells against apoptosis, leading to accommodate chlamydial growth (Bryant *et al.*, 2004; Park *et al.*, 2004; Paland *et al.*, 2006; Ying *et al.*, 2008). Therefore, the mediation effect via the expression of the *BIRC3* gene is logical. Our mediation testing involved four SNPs with MAF $\geq 10\%$ and eQTL (with *BIRC3*) at FDR 10%, using SMUT_PH, adapted LASSO and adapted Huang *et al.*'s method. All the methods showed significant mediation effects through *BIRC3* on incident chlamydia infection (Table 2). Marginal effects of selected SNPs on the *BIRC3* gene expression and time to the first incident chlamydia infection are visually illustrated in Web Figures S21.

5 | DISCUSSION

Our proposed methods, SMUT_GLM and SMUT_PH, extend our previous work (Zhong *et al.*, 2019) to test the mediation effect of multiple correlated genetic variants on a non-Gaussian outcome through a mediator. We adopt a mixed model-based approach to handle high dimension of genetic variants and do not apply any variable selection of genetic variants. Our proposed methods are statistically more powerful than alternative methods including SMUT, adapted LASSO, and adapted Huang *et al.*'s method. Analysis and discussions of possible reasons underlying alternative methods' power loss are in Supporting Information Section 5. The approximated version of SMUT_GLM and SMUT_PH are also computationally efficient (Supporting Information Section 7.2).

One limitation of our proposed methods is that we assume the effects of genetic variants follow a Gaussian distribution. This may not be correct when there are non-causal SNPs in the model, and in this case a mixture distribution might be more appropriate. It is reassuring to observe protected type-I error from our simulation studies, which included a large number of noncausal SNPs in all scenarios considered. In addition, supplementary simulation studies (Supporting Information Section 4) further demonstrate controlled type-I error when the effects of genetic variants follow a mixture of two Gaussian distributions. More properly modeling the effects of genetic variants may further increase the statistical power under the alternative hypotheses but due to modeling complexity and subsequently inevitable computational costs, we decided not to further pursue this in our current work.

Our proposed methods can be further extended to handle multiple correlated outcomes for additional power gains as well as to accommodate multiple potentially correlated mediators to jointly assess their mediation effects. Besides, we could adopt nonparametric methods to handle the mediator model and outcome model with more flexibility. Details germane to possible methodological extensions are presented in Supporting Information Section 7.1. We anticipate our proposed methods will become a powerful tool to bridge the gap in terms of molecular mechanisms between various types of phenotypes and the corresponding associated genetic variant(s) identified in the recent literature.

ACKNOWLEDGMENTS

This work was supported by the National Institutes of Health U19 AI084024, U19 AI144181, and R01 AI119164 to TD; U19 AI144181 to XZ; R01 HL129132, R01 HL146500, and U544 HD079124 to YL. YL is also partially supported

by U01 DA052713 and R01 GM105785. We thank all participants in ACE and TRAC for agreeing to take part in the studies, and all investigators in these two studies for sharing the data. We thank the editor, the associate editor, and three anonymous referees for their helpful comments and suggestions, which significantly improved the manuscript.

DATA AVAILABILITY STATEMENT

The data used in this paper to support our findings are available from the corresponding authors upon reasonable request.

ORCID

Wujuan Zhong  <https://orcid.org/0000-0002-6824-2255>

Yun Li  <https://orcid.org/0000-0002-9275-4189>

REFERENCES

- Bailey, R., Natividad-Sancho, A., Fowler, A., Peeling, R., Mabey, D., Whittle, H. et al. (2009) Host genetic contribution to the cellular immune response to *Chlamydia trachomatis*: heritability estimate from a Gambian twin study. *Drugs of Today (Barcelona, Spain: 1998)*, 45, 45–50.
- Baron, R.M. and Kenny, D.A. (1986) The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations. *Journal of Personality and Social Psychology*, 51, 1173–1182.
- Bastidas, R.J., Elwell, C.A., Engel, J.N. and Valdivia, R.H. (2013) Chlamydial intracellular survival strategies. *Cold Spring Harbor Perspectives in Medicine*, 3, a010256.
- Berger, R.L. and Hsu, J.C. (1996) Bioequivalence trials, intersection-union tests and equivalence confidence sets. *Statistical Science*, 11, 283–319.
- Bertrand, M.J., Doiron, K., Labbé, K., Korneluk, R.G., Barker, P.A. and Saleh, M. (2009) Cellular inhibitors of apoptosis cIAP1 and cIAP2 are required for innate immunity signaling by the pattern recognition receptors NOD1 and NOD2. *Immunity*, 30, 789–801.
- Blankenship, J.W., Varfolomeev, E., Goncharov, T., Fedorova, A.V., Kirkpatrick, D.S., Izrael-Tomasevic, A. et al. (2009) Ubiquitin binding modulates IAP antagonist-stimulated proteasomal degradation of c-IAP1 and c-IAP2 1. *Biochemical Journal*, 417, 149–165.
- Boyle, A.P., Hong, E.L., Hariharan, M., Cheng, Y., Schaub, M.A., Kasowski, M. et al. (2012) Annotation of functional variation in personal genomes using RegulomeDB. *Genome Research*, 22, 1790–1797.
- Breslow, N.E. and Clayton, D.G. (1993) Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, 88, 9.
- Bryant, P.A., Venter, D., Robins-Browne, R. and Curtis, N. (2004) Chips with everything: DNA microarrays in infectious diseases. *The Lancet Infectious Diseases*, 4, 100–111.
- Carabeo, R.A., Dooley, C.A., Grieshaber, S.S. and Hackstadt, T. (2007) Rac interacts with Abi-1 and WAVE2 to promote an Arp2/3-dependent actin recruitment during Chlamydial invasion. *Cellular Microbiology*, 9, 2278–2288.

- Centers for Disease Control and Prevention. (2019) *Sexually Transmitted Disease Surveillance 2018*. Technical report. Atlanta, GA: U.S. Department of Health and Human Services.
- Daubechies, I., Defrise, M. and De Mol, C. (2004) An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57, 1413–1457.
- Elwell, C., Mirrashidi, K. and Engel, J. (2016) Chlamydia cell biology and pathogenesis. *Nature Reviews Microbiology*, 14, 385.
- Fu, W.J. (1998) Penalized regressions: the bridge versus the LASSO. *Journal of Computational and Graphical Statistics*, 7, 397–416.
- Hackstadt, T. (2012). *Intracellular Pathogens I: Chlamydiales*, Volume 1. Washington, DC: American Society for Microbiology Press, pp. 126–148.
- Hassuna, N., Monk, P.N., Moseley, G.W. and Partridge, L.J. (2009) Strategies for targeting tetraspanin proteins. *BioDrugs*, 23, 341–359.
- Hauck, C.R. and Meyer, T.F. (2003) ‘Small’talk: Opa proteins as mediators of Neisseria–host-cell communication. *Current Opinion in Microbiology*, 6, 43–49.
- Hemler, M.E. (2008) Targeting of tetraspanin proteins—potential benefits and strategies. *Nature Reviews Drug Discovery*, 7, 747.
- Huang, Y.-T., Cai, T. and Kim, E. (2016) Integrative genomic testing of cancer survival using semiparametric linear transformation models. *Statistics in Medicine*, 35, 2831–344.
- Huang, Y.-T., Liang, L., Moffatt, M.F., Cookson, W.O.C.M. and Lin, X. (2015) iGWAS: integrative genome-wide association studies of genetic and genomic data for disease susceptibility using mediation analysis. *Genetic Epidemiology*, 39, 347–356.
- Join-Lambert, O., Morand, P.C., Carbonnelle, E., Coureuil, M., Bille, E., Bourdoulous, S., et al. (2010) Mechanisms of meningeal invasion by a bacterial extracellular pathogen, the example of *Neisseria meningitidis*. *Progress in Neurobiology*, 91, 130–139.
- Kim, C.W., Kim, H.K., Vo, M.-T., Lee, H.H., Kim, H.J., Min, Y.J. et al. (2010) Tristetraprolin controls the stability of cIAP2 mRNA through binding to the 3' UTR of cIAP2 mRNA. *Biochemical and Biophysical Research Communications*, 400, 46–52.
- Kumar, H., Kawai, T. and Akira, S. (2011) Pathogen recognition by the innate immune system. *International Reviews of Immunology*, 30, 16–34.
- Lane, B.J., Mutchler, C., Al Khodor, S., Grieshaber, S.S. and Carabeo, R.A. (2008) Chlamydial entry involves TARP binding of guanine nucleotide exchange factors. *PLoS Pathogens*, 4, e1000014.
- McCulloch, C.E., Searle, S.R. and Neuhaus, J.M. (2008) *Generalized, Linear, and Mixed Models*, 2nd edition. Hoboken, NJ: Wiley.
- Mehlitz, A. and Rudel, T. (2013) Modulation of host signaling and cellular responses by Chlamydia. *Cell Communication and Signaling*, 11, 90.
- Monk, P.N. and Partridge, L.J. (2012) Tetraspanins-gateways for infection. *Infectious Disorders–Drug Targets (Formerly Current Drug Targets–Infectious Disorders)*, 12, 4–17.
- Paland, N., Rajalingam, K., Machuy, N., Szczeppek, A., Wehrl, W. and Rudel, T. (2006) NF- κ B and inhibitor of apoptosis proteins are required for apoptosis resistance of epithelial cells persistently infected with *Chlamydomydia pneumoniae*. *Cellular Microbiology*, 8, 1643–1655.
- Palm, N.W. and Medzhitov, R. (2009) Pattern recognition receptors and control of adaptive immunity. *Immunological Reviews*, 227, 221–233.
- Pankratz, V.S., De Andrade, M. and Therneau, T.M. (2005) Random-effects Cox proportional hazards model: general variance components methods for time-to-event data. *Genetic Epidemiology*, 28, 97–109.
- Park, S.-M., Yoon, J.-B. and Lee, T.H. (2004) Receptor interacting protein is ubiquitinated by cellular inhibitor of apoptosis proteins (c-IAP1 and c-IAP2) in vitro. *FEBS Letters*, 566, 151–156.
- Raudenbush, S.W., Yang, M.L. and Yosef, M. (2000) Maximum likelihood for generalized linear models with nested random effects via high-order, multivariate Laplace approximation. *Journal of Computational and Graphical Statistics*, 9, 141–157.
- Russell, A.N., Zheng, X., O'Connell, C.M., Taylor, B.D., Wiesenfeld, H.C., Hillier, S.L. et al. (2015) Analysis of factors driving incident and ascending infection and the role of serum antibody in *Chlamydia trachomatis* genital tract infection. *The Journal of Infectious Diseases*, 213, 523–531.
- Schaffner, S.F., Foo, C., Gabriel, S., Reich, D., Daly, M.J. and Altshuler, D. (2005) Calibrating a coalescent simulation of human genome sequence variation. *Genome Research*, 15, 1576–1583.
- Schelldorfer, J., Meier, L. and Bühlmann, P. (2014) GLMMLasso: an algorithm for high-dimensional generalized linear mixed models using 1-penalization. *Journal of Computational and Graphical Statistics*, 23, 460–477.
- Seu, L., Tidwell, C., Timares, L., Duverger, A., Wagner, F.H., Goepfert, P.A. et al. (2017) CD151 expression is associated with a hyperproliferative T cell phenotype. *The Journal of Immunology*, 199, 3336–3347.
- Takeda, K., Kaisho, T. and Akira, S. (2003) Toll-like receptors. *Annual Review of Immunology*, 21, 335–376.
- Tan, B., Zammit, N.W., Yam, A.O., Slattery, R., Walters, S., Malle, E. and Grey, S. (2013) Baculoviral inhibitors of apoptosis repeat containing (BIRC) proteins fine-tune TNF-induced nuclear factor κ B and c-JUN N-terminal kinase signalling in mouse pancreatic beta cells. *Diabetologia*, 56, 520–532.
- Taylor, B.D., Zheng, X., Darville, T., Zhong, W., Konganti, K., Abiodun-Ojo, O. et al. (2017) Whole-exome sequencing to identify novel biological pathways associated with infertility following pelvic inflammatory disease. *Sexually Transmitted Diseases*, 44, 35.
- Tibshirani, R. (1996) Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58, 267–288.
- Vaida, F. and Xu, R. (2000) Proportional hazards model with random effects. *Statistics in Medicine*, 19, 3309–3324.
- 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. *American Journal of Human Genetics*, 89, 82–93.
- Ying, S., Christian, J.G., Paschen, S.A. and Häcker, G. (2008) *Chlamydia trachomatis* can protect host cells against apoptosis in the absence of cellular inhibitor of apoptosis proteins and Mcl-1. *Microbes and Infection*, 10, 97–101.
- Zheng, X., O'Connell, C.M., Zhong, W., Nagarajan, U.M., Tripathy, M., Russell, A.N. et al. (2018) Discovery of blood transcriptional

endotypes in women with pelvic inflammatory disease. *The Journal of Immunology*, 200, 2941–2956.

Zhong, W., Spracklen, C.N., Mohlke, K.L., Zheng, X., Fine, J. and Li, Y. (2019) Multi-SNP mediation intersection-union test. *Bioinformatics*, 35, 4724–4729.

SUPPORTING INFORMATION

Web Appendices and figures referenced in Sections 2–5 are available with this paper at the Biometrics website on Wiley Online Library. The R code for our method is the function GSMUT in the R package SMUT, which is

included in a zip file as a web supplement and also publicly available from CRAN at <https://CRAN.R-project.org/package=SMUT>.

How to cite this article: Zhong W, Darville T, Zheng X, Fine J, Li Y. Generalized multi-SNP mediation intersection–union test. *Biometrics*. 2020;1–12. <https://doi.org/10.1111/biom.13418>