

Estimation of Direct and Indirect Polygenic Effects and Gene-Environment Interactions using Polygenic Scores in Case-Parent Trio Studies

Ziqiao Wang¹, Luke Grosvenor^{2,3}, Debashree Ray^{1,4}, Ingo Ruczinski¹, Terri H. Beaty⁴, Heather Volk^{3,4}, Christine Ladd-Acosta^{3,4}, Nilanjan Chatterjee^{1,5,*}

1. Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, United States of America 21205;
2. Division of Research, Kaiser Permanente Northern California, Pleasanton, CA, United States of America 94588;
3. Department of Mental Health, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, United States of America 21205;
4. Department of Epidemiology, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, United States of America 21205;
5. Department of Oncology, School of Medicine, Johns Hopkins University, Baltimore, MD, United States of America 21205.

*Correspondence to: Nilanjan Chatterjee (nilanjan@jhu.edu), Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, 615 N. Wolfe Street, Baltimore, MD 21205.

Abstract

Family-based studies provide a unique opportunity to characterize genetic risks of diseases in the presence of population structure, assortative mating, and indirect genetic effects. We propose a novel framework, PGS-TRI, for the analysis of polygenic scores (PGS) in case-parent trio studies for estimation of the risk of an index condition associated with direct effects of inherited PGS, indirect effects of parental PGS, and gene-environment interactions. Extensive simulation studies demonstrate the robustness of PGS-TRI in the presence of complex population structure and assortative mating compared to alternative methods. We apply PGS-TRI to multi-ancestry trio studies of autism spectrum disorders ($N_{\text{trio}} = 1,517$) and orofacial clefts ($N_{\text{trio}} = 1,904$) to establish the first transmission-based estimates of risk associated with pre-defined PGS for these conditions and other related traits. For both conditions, we further explored offspring risk associated with polygenic gene-environment interactions, and direct and indirect effects of genetically predicted levels of gene expression and metabolite traits.

Introduction

Large genome-wide association studies (GWAS) of unrelated individuals have been widely used to derive polygenic scores (PGS) or polygenic risk scores (PRS) for complex traits. While it is standard practice to account for population stratification using genetic principal components, recent studies¹⁻⁵ have demonstrated the potential for overestimating genetic effects due to residual confounding with geographical variations and/or assortative mating in the population. This complicates the translational applications and interpretations for PGS across various analyses, including risk predictions and Mendelian randomization analysis.⁶ Family-based association studies,⁷ which involve estimating genetic effects through within-family comparisons, can protect against such biases when assessing the effects of individual genetic variants as well as PGS. Further, family-based studies with parental data can uniquely be used to estimate the indirect effects⁸⁻¹¹ of parental genetic variation on offspring outcomes, possibly through parental environmental factors. Thus, these PGS could be used to obtain genetic-based evidence of the effect of parental exposures on children's health outcomes.

To date, there have been limited studies on methodologies estimating PGS effects associated with disease risks through family-based studies. One study⁷ described methods for separating within and between family effects of PGS based on random-effect models to account for family-specific effects. Another study¹¹ pioneered the use of parental genotype data on probands to separate the direct and indirect effects of PGS on traits within families. These methods were primarily developed for the analysis of quantitative traits in randomly sampled families and are not suitable for other important study designs including case-parent trios, mother-child dyads, or other family-based study designs that ascertain participants based on affected probands. One recent study¹² introduced the polygenic transmission disequilibrium test (pTDT) based on case-parent trio designs and detected evidence of polygenic risk of autism spectrum disorders (ASD), irrespective of the presence of high *de novo* variants in probands. The method, however, does not provide estimates of effect-sizes in a suitable risk-scale, a critically important task for many purposes such as for comparisons of risk estimates from population-based

studies, conducting Mendelian randomization studies, and defining concepts of interactions.

To address these limitations and meet the current needs of the field for the analysis of PGS in family-based studies, we introduce PGS-TRI. This method can be applied to case-parent trio study designs to estimate the risk of a condition in offspring associated with direct (inherited) effects of PGS and its interaction with environmental exposures (PGSx E), and indirect effects of parental PGS.^{8,10} The method allows disease risk and PGS distribution to vary across families in a flexible manner, making it highly resilient to population stratification and assortative mating. We show that under our modeling framework, the PGS distribution in ascertained families can be derived in a compact form and can be conveniently partitioned into transmission and parental components. Based on this factorization, we present novel methods for estimating the direct effect of PGS and the effect of PGSx E interactions on offspring outcomes using the transmission component, while using a key scale parameter estimate from the parental PGS distribution. Additionally, we show that parental PGS data can be used to derive a simple and highly robust estimator for the difference in indirect effects of maternal and paternal PGS on offspring's outcomes. We conduct extensive simulation studies to demonstrate the validity and power of PGS-TRI for detecting direct and indirect effects, and PGSx E interactions under the presence of complex population structures and assortative mating.

To illustrate varied applications of the proposed method, we use multi-ancestry case-parent trio studies of two distinct developmental conditions: autism spectrum disorders (ASD) and orofacial clefts (OFCs). Both of these conditions are known to be highly heritable, have been associated with specific maternal risk factors, and have pre-defined PGS from recent large GWAS. We obtain transmission-based estimates of effect sizes for PGS underlying these two traits and show that they are comparable to corresponding effect sizes reported in prior studies, primarily conducted based on unrelated cases and controls. For ASD, we further investigate transmission-based estimates of risk associated with a number of other psychiatric and cognitive traits. For both ASD and OFCs, we examine PGSx E interactions for several known maternal risk factors. Finally, we showcase the effectiveness of PGS-TRI as a discovery tool by

analyzing PGS for gene expression and metabolite traits, obtained from the OMICSPRED study,¹³ to investigate their potential direct or indirect effects on the risk of the two conditions.

Material and Methods

Modeling Direct Genetic Effects and Gene-Environment Interactions

An overview of the method is presented in Fig.1. We assume that PGS for the index condition and other traits of interest can be evaluated for family trio participants, using published meta-data from prior association studies. Our goal is to investigate the association of these PGS with an index condition, such as ASD, using case-parent trios. In a trio, let PGS_C , PGS_M , and PGS_F denote the PGS values of the child, mother, and father respectively. Additionally, let D_C denote the binary outcome status of the child. We assume $i = 1, \dots, N$ families have been sampled following the case-parent design. We also assume a set of environmental exposures, denoted by \mathbf{E}_{ic} , has been ascertained for each child across the different families. We assume the prospective disease risk of the child in the population follows a log-linear model:

$$\begin{aligned} pr(D_{ic}|PGS_{ic}, PGS_{iF}, PGS_{iM}, \mathbf{E}_{ic}) &= pr(D_{ic}|PGS_{ic}, \mathbf{E}_{ic}) \\ &= \exp\{\alpha_i + \beta_G PGS_{ic} + \boldsymbol{\beta}_E^T \mathbf{E}_{ic} + \boldsymbol{\beta}_{GE}^T PGS_{ic} \times \mathbf{E}_{ic}\}. \end{aligned} \quad (1)$$

here, β_G is the direct genetic effect (DE). In (1), we assume no indirect genetic effects, i.e., within each family, the parental PGS values only affect the child's disease risk mediated through the child's own PGS value, but we will relax this assumption later (see subsequent sections and Supplemental Notes). Additionally, this model incorporates family-specific intercept terms α_i without any further assumption about their distribution, thereby allowing disease risks to vary arbitrarily across families.

We assume that the joint distributions of PGS across families in the underlying population follow tri-variate normal distributions of the form:

$$(PGS_{ic}, PGS_{iM}, PGS_{iF})^T \sim MVN_3 \left\{ \mu_i \mathbf{1}_3, \sigma_i^2 \begin{pmatrix} 1 & 0.5 & 0.5 \\ 0.5 & 1 & 0 \\ 0.5 & 0 & 1 \end{pmatrix} \right\}. \quad (2)$$

where the correlation of 0.5 between PGS values of individual parents and children follows from Mendel’s law of inheritance. We allow family-specific mean (μ_i) and variance (σ_i^2) terms without imposing any assumptions regarding their distributions. We assume that, within a family, PGS for the parents are independently distributed.

We explain here how a model of form (2) can flexibly account for both population structure and assortative mating. We conceptualize a sampling mechanism where for each family sampled, we assume it belongs to a unique subpopulation of "homogeneous" families with distinct genetic ancestry and trait characteristics that influence assortative mating. The number of these subpopulations can be arbitrarily large, making each subpopulation at an extremely fine level. It is assumed that random mating occurs within these homogeneous fine-level subpopulations, implying the PGS correlation between partners is 0 within these subpopulations, but not necessarily across them. By allowing arbitrary family-specific parameters for the PGS distribution, the model can accommodate fine-level population structure and the presence of any population-level correlation in PGS values between parents due to assortative mating. Finally, we assume $E_{iC} \perp (PGS_{iC}, PGS_{iM}, PGS_{iF})$, but allow the distribution E_{iC} to remain unspecified. From a population perspective, this assumption can again be viewed as “gene-environment” independence *within* highly homogenous subpopulations, but the model can still accommodate gene-environment correlation at the population-level that may arise due to population substructure and assortative mating.

Retrospective Likelihood and Parameter Estimation

The retrospective likelihood for the case-parent trio data in each family i can be decomposed into an offspring’s (L_{iC}) and a parents’ (L_{iP}) component as

$$\begin{aligned} L_i &= pr(PGS_{iC}, PGS_{iM}, PGS_{iF} | \mathbf{E}_{iC}, D_{iC} = 1) \\ &= pr(PGS_{iC} | PGS_{iM}, PGS_{iF}, \mathbf{E}_{iC}, D_{iC} = 1) \times pr(PGS_{iM}, PGS_{iF} | \mathbf{E}_{iC}, D_{iC} = 1) \\ &\stackrel{\text{def}}{=} L_{iC} \times L_{iP}. \end{aligned}$$

Under the above model, the likelihood components associated with children (L_{iC}) and parental PGS (L_{iP}) data for each family can be derived in terms of the following normal distributions (see Supplemental Notes for detailed derivation):

$$[PGS_{iC} | PGS_{iM}, PGS_{iF}, \mathbf{E}_{iC}, D_{iC} = 1] \sim N(\mu_{iC} + 0.5\sigma_i^2(\beta_G + \boldsymbol{\beta}_{GE}^T \mathbf{E}_{iC}), 0.5\sigma_i^2),$$

and

$$[PGS_{iM/iF} | \mathbf{E}_{iC}, D_{iC} = 1] \sim N(\mu_i + 0.5\sigma_i^2(\beta_G + \boldsymbol{\beta}_{GE}^T \mathbf{E}_{iC}), \sigma_i^2).$$

Maximum-likelihood estimation based on $L = \prod_{i=1}^N L_i$ can be complex due to the presence of large dimensional nuisance parameters $(\mu_i, \sigma_i)_{i=1}^N$. Instead, we propose a combination of likelihood- and moment-based estimation. First, we observe that the likelihood $L_C = \prod_{i=1}^N L_{iC}$ is informative for the estimation of β_G and $\boldsymbol{\beta}_{GE}$, but one complication is that it requires estimates of family-specific variance parameters σ_i^2 . We now observe that under the above model, the PGS values for two parents within each family are expected to have identical distribution, and thus $\hat{\sigma}_i^2 = 0.5(PGS_{iM} - PGS_{iF})^2$ provides an unbiased estimator for $\sigma_i^2, i = 1, \dots, N$. We can now obtain estimates of β_G and $\boldsymbol{\beta}_{GE}$ based on the likelihood L_C with plugged-in values for $\hat{\sigma}_i^2$. Under the above framework, we show that the final estimator can be derived in an analytic form as a solution to a weighted least-square problem as

$$\hat{\boldsymbol{\beta}} = (\mathbf{E}^T \widehat{\mathbf{W}} \mathbf{E})^{-1} \mathbf{E}^T \mathbf{Z}, \quad (3)$$

where $\boldsymbol{\beta} = (\beta_G, \boldsymbol{\beta}_{GE}^T)^T$, $\mathbf{E} = (\mathbf{E}_1^T, \dots, \mathbf{E}_N^T)^T$, $\mathbf{E}_i = (1, \mathbf{E}_{iC}^T)^T$, $\widehat{\mathbf{W}} = \mathbf{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_N^2)$, and $\mathbf{Z} = 2(PGS_{1C} - \mu_{1C}, \dots, PGS_{NC} - \mu_{NC})^T$, and $\mu_{iC} = 0.5(PGS_{iM} + PGS_{iF})$.

In the special case when no gene-environment interaction terms are incorporated, the estimate of the DE of PGS takes the simple form

$$\hat{\beta}_G = \frac{2 \sum_{i=1}^N (PGS_{iC} - 0.5(PGS_{iM} + PGS_{iF})) / N}{\sum_{i=1}^N \hat{\sigma}_i^2 / N}. \quad (4)$$

It is noteworthy that the numerator of (4) forms the basis of the pTDT test.¹² In pTDT, the numerator is normalized by the variance of average PGS values of the parents across families. However, our derivation of (4) suggests that obtaining an unbiased estimate of effect size for PGS in TDT-type analysis requires normalization of the transmission disequilibrium statistics, i.e., the numerator, by an estimate of *within-family* variance.

Incorporating Indirect Parental Genetic Effects

Next, we extend model (1) to incorporate indirect parental genetic effects as

$$\begin{aligned} & pr(D_{iC} | PGS_{iC}, PGS_{iF}, PGS_{iM}, \mathbf{E}_{iC}) \\ &= \exp\{\alpha_i + \beta_G PGS_{iC} + \beta_M PGS_{iM} + \beta_F PGS_{iF} + \boldsymbol{\beta}_E^T \mathbf{E}_{iC} + \boldsymbol{\beta}_{GE}^T PGS_{iC} \times \mathbf{E}_{iC}\}, \end{aligned} \quad (5)$$

where β_M, β_F capture indirect effects of parental PGS on the disease risk of the children that are not mediated through the children's genotypes.

The likelihood components associated with children (L_{iC}) remain unchanged after incorporating the indirect genetic effects. We, however, show that the conditional distribution of PGS value in the mother/father in the i -th ascertained family, when parental effects are incorporated, needs to be updated as $L_{iP} = [PGS_{iM/iF} | \mathbf{E}_{iC}, D_{iC} = 1] \sim N(\mu_i + \sigma_i^2 [\beta_{M/F} + 0.5(\beta_G + \boldsymbol{\beta}_{GE}^T \mathbf{E}_{iC})], \sigma_i^2)$. Now we observe that, unlike the previous setting, here the two parents within a family could have asymmetric distribution depending on the difference in the magnitude of their indirect effects (β_M and β_F). We can exploit this parental asymmetry in PGS distribution to derive an estimator for the difference of parental indirect genetic effect (δ -IDE) as

$$\hat{\delta}\text{-IDE} = \hat{\beta}_M - \hat{\beta}_F = \frac{\sum_{i=1}^N (PGS_{iM} - PGS_{iF})/N}{\sigma_{sum}^2/N}.$$

We further show that in the presence of indirect effects, an approximately unbiased estimator for $\sigma_{sum}^2/N = \sum_{i=1}^n \sigma_i^2/N$ can be derived as

$$\hat{\sigma}_{sum}^2/N = \sum_{i=1}^N \hat{\sigma}_i^2/N \approx \frac{1}{2(N-1)} \sum_{i=1}^N \left[(PGS_{iM} - PGS_{iF}) - \sum_{i=1}^N (PGS_{iM} - PGS_{iF})/N \right]^2.$$

Further, when the indirect effects of parental PGS is incorporated, we observe the form of the estimates of direct effect parameters ($\boldsymbol{\beta}$) as shown in (3) remains unchanged, but the estimates $\hat{\sigma}_i^2, i = 1, \dots, N$ in defining the weight matrix $\widehat{\mathbf{W}} = \mathbf{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_N^2)$ needs to be modified as $\hat{\sigma}_i^2 = 0.5[(PGS_{iM} - PGS_{iF}) - \sum_{i=1}^N (PGS_{iM} - PGS_{iF})/N]^2$ to obtain more accurate estimator of the quantity $\mathbf{E}^T \mathbf{W} \mathbf{E}$.

Detailed derivations of all mathematical results and asymptotic variance estimators are presented in Supplemental Notes.

Simulation Studies

We conducted two types of simulation studies to evaluate the performance of the proposed method under complex population substructure and assortative mating. In the first setting, we simulated PGS values in over 1 million randomly sampled trios from trivariate normal distribution based on Equation (2) and disease status in children under the full model Equation (5), and then further selected families with affected children ($D = 1$). We varied $\rho_G = \text{cor}(\alpha_i, \mu_i)$ to create different scenarios of population-stratification bias, with a value of 0 indicating no relationship between variation in disease risk and PGS distribution across the underlying substructure – a scenario where we do not expect any bias in population-based association studies. For the investigation of gene-environment interactions, we incorporate a binary (E_1) and a continuous (E_2) variable in the disease model. We allowed complex inter-relationships among baseline disease risk (α_i), PGS mean (μ_i) and exposure means (γ_{i1} and γ_{i2}) across families, in manners that are known to affect the estimation of gene-environment interactions using unrelated individuals.¹⁴

We designed the second simulation setting to investigate the robustness of different methods for the estimation of simulated effects of educational attainment (EA)-PGS, a score which is known to be highly confounded with environmental factors due to population substructure^{15,16} in the UK Biobank (UKB) study (www.ukbiobank.ac.uk). We built EA-PGS using independent SNPs ($r^2 < 0.01$ within 1000kb) and weights reported in previous work (PGS Catalog ID: PGS002012)¹⁷ for the UKB participants. To simulate geographical population structure, we matched unrelated UKB males and females within assessment centers (UKB Field ID: 54), and by birth locations defined by north and east coordinates (UKB Field ID: 129 and 130). Additionally, we matched pairs based on both geographical proximity and similarity in educational attainment (UKB Field ID: 6138) to simulate additional effect of assortative mating. For each matched pair of UKB participants, we simulated children by generating their genotype values based on Mendel's law of transmission and their disease status based on Equation (5), where a "hidden" environmental variable, defined as the average of BMI values (UKB Field ID: 21001) of the two parents, were introduced to influence disease risk. Additional details on

both simulation settings, including the matching algorithm used, can be found in Supplemental Notes.

Data Analysis for Application of PGS-TRI to Autism Spectrum Disorders (ASD)

We illustrate the versatile capabilities of PGS-TRI by using it to analyze genotype and epidemiologic data available on trios from the Simons Foundation Powering Autism Research for Knowledge (SPARK) study¹⁸ (<https://sparkforautism.org/>) (Table S1). The key objectives included: (1) obtaining transmission-based estimates of the DE of the most advanced ASD-PGS to date built from prior studies, and comparing such effect size from that reported from prior case-control studies, (2) assessing the portability of European-derived PGS to non-European ancestry families, (3) evaluating the association of PGS for multiple neurocognitive traits with ASD risk using a transmission-based method, (4) characterizing the nature of interaction of ASD-PGS and prenatal exposure on ASD risk, and (5) discovering potential novel association of ASD risk with DE and δ -IDE of PGS for gene expression and obesity-related metabolite traits.

Specifically, we examined ASD risk associated with direct and indirect effects of pre-constructed PGSs for ASD¹⁹ and 7 other cognitive and mental health traits, including educational attainment,¹⁷ schizophrenia,²⁰ strictly defined lifetime major depressive disorder,²¹ bipolar disorder,²² neuroticism,¹⁷ insomnia,¹⁷ and attention-deficit/hyperactivity disorder (ADHD)²³ across different ancestry groups. As a negative control, we used body mass index (BMI).¹⁷ The PGS for ASD itself was defined based on the largest GWAS to date conducted by the iPSYCH consortium and involved a total of 28,017 SNPs. We standardized all of the PGSs by ancestry-specific standard errors obtained from the 1000 Genomes (1000G) Phase 3 Project.²⁴ This allows interpreting underlying risk parameters, i.e., relative risks, in a standard unit scale. We further examined gene-environment interactions of the ASD-PGS with several prenatal environmental factors, information collected in SPARK participants using questionnaires designed by the Genes and Environment Autism Research Study (GEARS). Finally, we examined the association of ASD risk with genetically-predicted levels of gene expression and metabolites using thousands of genetic scores generated by the OMICSPRED project.²⁵ Here the underlying hypothesis is that genetically-predicted biomarker levels, as captured by the

underlying PGS, could affect ASD risk through direct or indirect effects. We do note a caveat that because the PGSs for biomarkers have been derived based on adult samples, direct effect of PGSs on children's outcome is only possible if the same PGS also predicts biomarker levels in fetal state or/and early childhood. Anticipating limited power for this analysis, we only included those biomolecular traits that can be predicted with accuracy $R^2 \geq 0.1$ by the underlying genetic scores according to the internal validation in OMICSPRED and those which include at least 5 SNPs in the underlying model. This criterion resulted in the evaluation of a total of 4,991 gene expression levels and 27 highly correlated obesity-related metabolites. Additional details of data pre-processing and covariate coding can be found in Supplemental Notes.

Data Analysis for Application of PGS-TRI to Orofacial Clefts (OFCs)

We also used PGS-TRI to investigate the risk of non-syndromic OFCs using case-parent trio data from the Gene Environment Association Studies (GENEVA).^{26,27} This analysis included a total of 778 European and 1,126 East-Asian ancestry trios (see Tables S2-3 for distribution of trios by subtypes and exposure). We used the trio data to examine the effects of a pre-defined OFC-PGS on the risk of OFCs across different subtypes and ancestry groups, and its interaction with prenatal exposure to maternal smoking, maternal alcohol consumption, use of multivitamins during pregnancy, and prenatal environmental tobacco smoke exposure.²⁶⁻²⁸ PGS for cleft lip with or without cleft palate (CL/P) were constructed using 24 SNPs and their respective weights sourced from the PGScatalog.^{29,30} Similar to ASD analysis, we standardized the PGS using ancestry-specific standard errors estimated from the 1000G project. Finally, we also examined the risk of OFCs (CL/P) associated with the DE and δ -IDE of genetically-predicted gene expressions and metabolite levels using genetic scores generated from the OMICSPRED project.²⁵ Additional details of data pre-processing and covariate coding can be found in Supplemental Notes.

Results

Simulation Studies

In simulation studies, where we directly generate data under the assumed model, we observe that PGS-TRI produces unbiased effect-size estimates, well-controlled type-I error rates, and calibrated confidence intervals for all of the different types of parameters across a realistic range of population stratification scenarios (Fig.2, Extended Fig.1). The pTDT, being a transmission-based method, also produces unbiased tests for direct genetic effects (DE) and has identical power to PGS-TRI across different scenarios. Standard logistic regression analysis of unrelated case-control participants produces biased inference for DE in the presence of correlations between PGS mean and disease risks across families. Logistic regression and case-only analysis also produce significant bias for inference on gene-environment interaction parameters in the presence of complex population substructures across which disease risk, exposure distribution, and PGS distribution co-vary. If parental genotype data were available for unrelated case-control participants, then logistic regression model could also be used to estimate the magnitude of differential indirect effects (δ -IDE). In the absence of population stratification, both methods are valid for the estimation of DE and δ -IDE, but logistic regression is more powerful for detecting DE while PGS-TRI is more powerful for detecting δ -IDE (Extended Fig.2). Further, logistic regression can produce biased inference for DE not only in the presence of population stratification but also in the presence of δ -IDE when parental data are not available to account for such effects (Extended Fig.3).

We further confirm the robustness of PGS-TRI to realistic patterns of population substructure and assortative mating by considering analyses of educational attainment (EA)-PGS in the UK Biobank study. In this setting, we matched unrelated male and female participants to form “parents” and children’s genotype data were simulated under Mendel’s law of inheritance. We observe that when participants are matched by geographical proximity, there is significant across-family correlation ($P < 5 \times 10^{-6}$) between EA-PGS and parental-BMI, which is treated as a hidden “environmental” confounder for the simulations of disease risk in children (Extended Fig.4a-4d). In this setting, logistic regression analysis of unrelated cases and controls, even after adjustment for genetic PCs and geographical coordinates, can produce significantly biased inference for DE (Fig.3). When we further added EA to the matching criterion, i.e., allowing for assortative mating in addition to geographical population structure, there was increased bias in

logistic regression. PGS-TRI produces both unbiased tests and effect size estimates in all scenarios.

Polygenic Risk for Autism Spectrum Disorders (ASD)

We first examined the association of ASD-PGS derived from the European-ancestry (EUR) iPSYCH study¹⁹ with ASD risk across different ancestry groups (Fig.4a; Table S4). We observe that for the EUR population, PGS-TRI produced a transmission-based estimate of DE (RR = 1.32, 95% CI = [1.22,1.44]) closely matching the reported estimate (OR = 1.33, 95% CI = [1.30,1.36]) from the iPSYCH study,¹⁹ which predominantly used unrelated case-control samples. Thus, our estimate of effect size suggests no evidence of significant bias in prior population-based GWAS due to unadjusted population stratification. We also detected evidence of a significant direct effect of the PGS on ASD risk within Asian- and African-ancestry families with the corresponding effect size estimates being of similar magnitude as those derived from EUR families. However, we did not find any evidence of direct effects of PGS on ASD risk in Americas ancestry families.

We further observed that PGS-TRI produced estimates of DE of PGS for EA and several psychiatric traits on ASD risk in the cross-ancestry population (combined ancestry groups of EUR, Asian, African, and Americas) (Fig.4b) which have similar patterns as has been reported in prior studies.¹³ Most of these associations, however, could not be detected specifically for the non-European populations due to low power and perhaps lack of transportability of the PGS to non-EUR populations (Table S5). We did not detect any evidence of the differential indirect effect of the parental PGS values, including that of ASD itself, on the offspring risk with the exception of schizophrenia-PGS, for which borderline evidence ($P = 0.048$) is seen. We did not detect evidence ($FDR < 0.05$) of any DE or δ -IDE of OMICSPRED generated gene-expression-PGSs on ASD risk in the cross-ancestry or EUR-only analysis (Table S7-8). The corresponding q-q plots indicate PGS-TRI controls type-I error rates well in genome-wide analysis (Extended Fig.5). However, in the analysis of a group of 27 obesity-related metabolites, we found pervasive evidence of maternal mediated indirect effects of underlying genetic scores on the risk of ASD in offsprings (Fig.4c; Table S9) where the strongest signal was observed for Apolipoprotein

B ($P = 1.6 \times 10^{-3}$, FDR = 0.02). However, because the lipoprotein groups are highly correlated with each other (Extended Fig.6), it is difficult to pinpoint the specific metabolites causing the underlying effects. We did not find evidence of DE of any of the metabolite-PGSs on ASD risk (Table S10).

Finally, we explored gene-environment interactions of the ASD-PGS and several specific pre- and peri-natal environmental factors on ASD risk (Table S6). We generally did not observe any strong evidence of interactions, indicating that polygenic risk and environmental factors generally act multiplicatively on the risk of ASD. We observed nominal evidence that maternal alcohol consumption, before and during pregnancy, modified the DE of ASD PGS on offspring risk in EUR families ($P = 0.04$ ever, before pregnancy; $P = 0.05$, during pregnancy; compared to never drinkers).

Polygenic Risk of Orofacial Clefts (OFCs)

We applied PGS-TRI to analyze the risk of orofacial clefts using EUR- and Asian- ancestry trios available from the GENEVA study (Fig.5; Table S11-12). We first examined the risk of various OFCs subtypes with a PGS incorporating 24 SNPs defined by an earlier study (PGS Catalog ID: PGS002266).^{29,30} We found highly significant and consistent levels of DE of this PGS on the risk of OFCs across different subtypes, including CL-alone, CL&P, and CL/P (CL-alone and CL&P combined), and across both populations. The strengths of the associations are of similar magnitude as those reported in previous studies using both population- and family-based samples.²⁹ In our analysis, we did not find any evidence of association of this PGS with the CP-alone subtype, an anatomically and embryologically distinct subtype, which was not surprising, considering the original PGS was developed based on studies of CL/P subtypes only. In the cross-ancestry population (combined EUR and Asian population) analysis, we did not observe evidence of differential indirect effects of the parental PGS on the risk OFCs in offspring. In the EUR-only analysis, however, we observed nominal-level evidence of maternally mediated δ -IDE of the PGS on the risk of CL-alone ($P = 0.03$).

We further used PGS-TRI to explore the interaction of the DE of OFC-PGS with several known prenatal environmental risk factors and offspring's sex on the OFCs subtypes (Fig.5a). In the EUR population, we detected evidence of interaction of the PGS

with maternal smoking during pregnancy on the risk of CL-alone and combined CL/P; and to a lesser extent with maternal environmental smoking exposure in CL&P. These interactions are not found to be significant in Asians, but the power for detecting interaction in this population was very low as only a small proportion (3%) of the affected probands were exposed to maternal smoking (Table S3). We further observed evidence of PGS by sex interactions in EUR on CL&P risk, but an opposite direction for the same interaction effect in the Asian population, which cancelled each other in cross-ancestry population analysis.

Finally, we examined the risk of OFCs (CL/P, and combined OFCs) associated with DE and δ -IDE of transcriptomic-PGSs generated by the OMICSPRED study.²⁵ We detected strong evidence of DE of genetically predicted expression of *TRAF3IP3* on the risk of CL/P (cross-ancestry $P = 2.0 \times 10^{-12}$) with the strength of association appearing to be stronger in the Asian population compared with the EUR population (heterogeneity test $P = 0.01$, Fig.5c; Table S13). The genetic score of *TRAF3IP3* expression involved 44 SNPs within the cis-region and had a prediction $R^2 = 0.139$ in the EUR population. An intronic SNP rs2235370 of *TRAF3IP3* has been previously reported as a sentinel variant associated with the risk of CL/P in prior GWAS cross-ancestry meta-analysis.³¹ The DE of *TRAF3IP3*-PGS on CL/P risk became insignificant when we removed 8 SNPs in linkage disequilibrium with rs2235370 ($r^2 > 0.3$), but the DE of genetic score based only on the 8 SNPs remained highly significant (cross-ancestry $P = 6.5 \times 10^{-13}$). Application of genotypic TDT³² further showed a strong association of each of the individual 8 SNPs with the risk of CL/P (Table S14). These results combined suggested the presence of a haplotype in this region which has protective effect on OFC risk, likely mediated by the expression level of the gene *TRAF3IP3*. We did not detect any evidence of δ -IDE of the transcriptomic-PGS on OFCs risk (Extended Fig. 7; Table S15). We also did not find evidence of DE or δ -IDE of metabolomic-PGSs on OFCs risk (Table S16-17).

Discussion

We have developed a new analytic framework, PGS-TRI, for the analysis of polygenic scores in case-parent trio studies. It enables transmission-based estimation of the risk of

the outcome in offspring, accounting for the direct effects of inherited PGS and its interaction with environmental factors. Further, the method leverages observed asymmetry in PGS value between parents within families to estimate maternally or paternally mediated indirect effects on the offspring outcomes. We conducted simulation studies across various realistic scenarios, including one involving EA-PGS in the UK Biobank study, to demonstrate the robustness of the proposed method against population structure and assortative mating. Our applications of PGS-TRI to 2 early developmental health conditions, ASD and OFCs, provide the first transmission-based estimates of effect sizes for established PGSs and address concerns over biases in previous studies caused by uncorrected population structure, assortative mating, or indirect genetic effects. We also applied PGS-TRI to novel analyses exploring polygenic gene-environment interactions on these two conditions. Finally, we used PGSs for gene expression and metabolite traits to examine any evidence of their direct and indirect effects on the two conditions.

Case-parent trio studies and related analytic methods have a significant history and long-standing utility in the field of genetic epidemiology, especially for developmental health conditions in children. Originally, the transmission disequilibrium test (TDT)³³ was proposed as an allelic test for linkage and association, and was believed that it would form the basis of future GWAS.³⁴ Several other key studies noted that transmission-based testing and risk estimation can be conducted based on marker genotype data without having to assume multiplicative effects of underlying alleles.^{35,36} Subsequently, a series of methods were proposed for the transmission-based analysis of multi-allelic markers,³⁷ complex pedigrees,³⁸ gene-environment interactions,³⁹⁻⁴¹ indirect (also referred to in prior studies as “nurturing”),^{42,43} and imprinting/parent-of-origin- effects,^{41,43} and a more general class of distribution-free methods for family-based association testing also emerged.⁴⁴⁻⁴⁶ As GWAS required very large sample sizes for the detection of small polygenic effects, studies based on unrelated cases and controls became widely popular due to the ease of recruiting. In the post-GWAS era, however, there has been now renewed interest in the use of case-parent trios and other family-based studies for more robust characterizations of risks associated with GWAS-identified genetic effects.^{7,11,12} Additionally, there are

practical considerations that make it more feasible to collect data from parents in families with young affected children than collecting large samples of unrelated cases and controls. Thus, a family trio design is particularly well suited for modern and growing developmental outcomes such as ASD and OFCs, among others.

The first transmission-based method for PGS analysis in case-parent trios was proposed based on the deviation of observed PGS values of children from its expected value under Mendel's law of transmission, which is represented by the mid-parental PGS values.¹² While the method provided a valid test, the underlying statistics do not provide valid effect estimates. Here, we show that unbiased risk estimation requires scaling of transmission-disequilibrium statistics by estimates of within-family variance of the PGS, which we now incorporate in the PGS-TRI method. Further, our method, which is motivated by the likelihood of the trio-data, allows for the modeling of PGS-environment interactions within the transmission-based framework and can incorporate indirect effects of parental PGS on children's outcomes. In particular, we show that under our framework, an estimate of asymmetric maternal and paternal indirect effects of PGS can be obtained by the average difference in PGS values between the two parents across families, scaled by average within-family variance parameters.

We used data from the SPARK consortium to obtain transmission-based estimates of ASD risk associated with pre-defined PGS of ASD and several other neurocognitive traits defined by prior association studies. Our results validate previously reported risks in EUR ancestry populations. There is also a critical need to assess the portability of ASD-PGS in non-EUR ancestries and our analyses confirm transmission-based estimates of effect sizes for PGS are of similar magnitude in Asian- and African-ancestry groups. Further, the lack of evidence of any DE of PGS in Americas ancestry population highlights SNP discovery efforts are needed to better identify and capture polygenic risks in this group. Our results further show that PGS and the subset of pre- and peri-natal environmental risk factors examined in this study generally have multiplicative effects on the risk of ASD. Finally, we explored whether genetic scores of molecular traits, including genome-wide gene expressions and cardiometabolic metabolites, were associated with ASD risk through direct or indirect effects. We found statistically significant evidence of

indirect maternal genetic effects of PGS for a set of obesity-related metabolites on the ASD risk in offspring. While this finding requires further validation in larger studies, it is particularly noteworthy given the consistent associations between maternal obesity (both pre-pregnancy and during pregnancy) and ASD risk in offspring, as reported in epidemiology studies.^{47,48} Interestingly, we did not detect evidence of an indirect maternal effect of BMI-PGS itself on ASD risk in offspring, which may be due to lack of power or due to lack of underlying causal effect. Overall, this analysis showcases how case-parent trio studies can be leveraged to test for potential causal effects of parental exposures on children's outcomes, with parental PGS serving as "instruments" within the Mendelian Randomization framework.

We also applied our new method to analyze the polygenic risk of OFCs using data from the GENEVA study. This analysis established transmission-based risk estimates for a pre-defined PGS across different OFCs subtypes in both EUR and Asian-ancestry groups. We identified modest evidence of non-multiplicative interaction of the DE of OFC-PGS and maternal smoking, as well as environmental maternal tobacco smoke exposure during pregnancy. Prior genome-wide SNP-environment interaction studies,²⁶ including but not limited to GENEVA samples, have not revealed genome-wide significant findings. Joint tests for genetic associations and gene-environment interactions had indicated modest evidence of gene-by-maternal smoking interactions for rs7541797 near PAX7, but not much evidence was found for gene-by-environmental tobacco smoking.²⁶ In our PGS-based analysis, on the other hand, we find some evidence of aggregated SNPs by smoking interaction effects across different known loci. In the future, large and diverse studies including experimental model systems are needed to further characterize gene-environment interactions in the etiology of OFCs.

While our study has many major strengths as described above, our study has several limitations as well. Our framework allows estimation of differences in the indirect genetic effects that two parents may have on children, but it cannot identify indirect effects of the father and the mother separately. In our setting, we allow family-specific disease risk and PGS distributional parameters to have arbitrary distribution across families. It is possible that by incorporating stronger parametric assumptions one can extend our

framework for the estimation of the indirect genetic effects for each parent separately. We have modeled gene-environment interactions only with respect to DE of PGS, and further research is merited to extend the model to allow the possibility of interactions of environmental factors with indirect effect of parental PGS on children's outcome. Our assumption of gene-environment independence within "homogeneous" families is highly robust to the presence of population structure and assortative mating. However, there could be direct correlation between PGS and environmental exposures due to pleiotropic effects. Finally, the OMICSPRED study primarily trained PGS for molecular traits using adult samples, meaning any DE on developmental outcomes would depend on the PGS being predictive of early-life molecular traits. While molecular QTLs often have robust effects across life stages,⁴⁹ our DE analysis may lack power if early-life molecular levels are not as predictable by PGSs generated by OMICSPRED.

The proposed framework for conducting PGS analysis in case-parent trios opens numerous avenues for further research, including applying analogous analyses in other types of ascertained families such as mother-child dyads, case-parent trios with unaffected siblings, and more complex pedigrees which may be ascertained through multiple probands. Additionally, there is potential to extend the model for joint analysis of multiple possibly correlated PGSs, which could be useful for applications such as multivariable Mendelian randomization analysis.⁵⁰ Case-parent trio designs are widely used for many developmental and early childhood conditions due to the feasibility of collecting parental data and biosamples. Moreover, for many late-onset diseases, studies have collected genetic data on extended families ascertained with specific conditions like breast cancer for understanding risks associated with rare high-penetrant mutations. Our proposed method, along with its future extensions, can facilitate PGS analysis in ascertained families, enabling robust characterization of associated risks, both by a PGS itself and in conjunction with rare mutations and non-genetic risk factors.

Figure Legends

Figure 1. Figure illustrating our PGS-TRI model framework for case-parent trio family study designs. $\beta_G, \beta_{GE}, \beta_E$: population-level PGS direct genetic effect (β_G), direct PGS-E interactions (β_{GE}), and direct environmental effects (β_E), associated with offspring's outcome. β_M, β_F : mother and father indirect genetic effects associated with the offspring's outcome risk. $PGS_{iC}, PGS_{iM}, PGS_{iF}$: PGS values of child, mother, father in family $i = 1, \dots, N$. D_{iC} : outcome status of the offspring in family $i = 1, \dots, N$.

Figure 2. Performance of PGS-TRI and alternative methods for estimating parameters of PGS direct effect (DE), differential parental indirect genetic effect (δ -IDE), and PGSxE interactions in simulation studies. Results are shown for (a) Type I error of PGS DE and δ -IDE, when underlying true effects are 0; (b) Type I error of PGS-E interaction, when underlying true main effect DE is 0.4; (c) bias in parameter estimates and magnitude of standard deviation (SD) of estimates for DE and δ -IDE; (d) bias in parameter estimates and magnitude of SD of estimate for PGS-E interaction. pTDT is implemented as an alternative method for testing DE. Logistic regression is implemented for testing and estimation of DE assuming unrelated controls are available of the same size as the number of cases. Logistic regression is also implemented for testing and estimation of δ -IDE further assuming parental genotypes are available for the unrelated cases and controls. Further, a case-only method is also implemented for testing PGSxE interaction. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population.

Figure 3. Performance of PGS-TRI and alternative methods for the estimation of genetic effects associated with educational attainment PGS in the UK Biobank-based simulation study. Results are shown for (a) Type I error associated with direct-effect (DE) of PGS (b) Type I error associated with differential parental indirect genetic effect (δ -IDE) of PGS (c) bias and SD associated with estimation of DE and (d) bias and SD associated with estimation of δ -IDE. Among alternative methods, pTDT is implemented for testing of DE. For testing of DE, logistic regression is implemented for the analysis of unrelated cases and controls without any adjustment, or adjustment for top 10 genetic principal components (PCs) constructed from the parental data, or top 10 PCs plus the assessment centres, and north and east birth coordinates of the parents. For the testing and estimation of δ -IDE using logistic regression, we assume parental genotype data are available on unrelated cases and controls. Data on children for matched pairs of UKB participants are repeatedly simulated, and then a set of 2000 case-parent trios, or a set of 2000 unrelated cases and 2000 unrelated controls, are further sampled for subsequent analysis. Parents were either matched by only geographic proximity to simulate effect of population stratification or geographical proximity and educational attainment level to simulate effect of both population stratification and assortative mating. **PS**: population stratification bias; **AM**: assortative mating; **PC10**: logistic regression with the top 10 genetic PCs as covariates; **PC10 Geo**: logistic regression with the top 10 genetic PCs, north and east birth co-ordinates, and assessment centres as covariates in the model.

Figure 4. SPARK study results for autism spectrum disorder (ASD). (a) Relative risk estimates for direct (DE) and maternally mediated indirect effects (δ -IDE) of ASD-PGS (PGS ID: PGS000327) on ASD risk across multiple ancestry groups; (b) Relative risk estimates of DE and δ -IDE of PGS for multiple neurocognitive traits on ASD risk (only shown for combined population analysis due to sample size in ancestral subpopulations); (c) Heatmap of estimates of maternally

mediated indirect effects (δ -IDE) of PGS for 27 obesity-related metabolite PGS on the risk of ASD in offsprings. In (a) and (b), PGS are standardized within each ancestry group by population mean and SD calculated using 1000G reference data of independent individuals so that relative risk (RR) corresponds to an increase in risk per SD unit increase in PGS value. In (c), PGSs are standardized using population mean and SD by pooling all individuals within each ancestry group in the current study. Asterisks indicate significant metabolites after Bonferonni correction for the effective number of independent tests ($P < 0.025$). **ADHD**: attention-deficit/hyperactivity disorder; **EAS/SAS**: combined East Asian and South Asian populations; **CI**: confidence interval; **BMI**: body mass index, used as a negative control; **IDL**: intermediate-density lipoproteins; **LDL**: low-density lipoprotein; **VLDL**: very small low-density lipoprotein.

Figure 5. Results from the analysis in the GENEVA study of orofacial clefts (OFCs) trios.

(a) Estimates of relative risk of different OFC subtypes associated with the DE and δ -IDE of OFC-PGS (PGS ID: PGS002266) and the interaction of OFC-PGS with several maternally medicated risk factors (b) Miami plot showing results from the transcriptome-wide association study using PGS-TRI of the risk of the combined OFC CL/P associated with the DE of PGS for gene-expression traits available from the OMNISPRED study. In (a), PGS were standardized within each ancestry group by population mean and SD calculated using 1000G reference data of independent individuals so that relative risk (RR) corresponded to an increase in risk per SD unit increase in PGS value. In (b), PGSs are standardized using population mean and SD by pooling all individuals within each ancestry group in the current study.

Extended Figure 1. Coverage level of 95% confidence intervals of PGS-TRI and alternative methods for estimation of (a) DE, (b) δ -IDE, and (c-d) PGSxE interactions in simulation studies.

For each type of parameter, results are shown in scenarios in the absence and the presence of population stratification bias. Logistic regression is implemented for the estimation of DE and PGS by E interaction effects assuming unrelated controls are available of the same size as the number of cases. Logistic regression is also implemented for the estimation of δ -IDE further assuming that parental genotypes are available for the unrelated cases and controls. Additionally, a case-only method is implemented for testing PGSxE interaction terms. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population.

Extended Figure 2. Power of PGS-TRI compared to alternative methods for the detection of (a) DE, (b) δ -IDE, (c)(d) PGS by E interaction terms.

pTDT is implemented as an alternative method for testing DE. The case-only method is implemented as an alternative method for testing PGSxE. Logistic regression is also implemented for testing of DE and PGSxE assuming unrelated controls are available of the same size as the number of cases. Logistic regression is further implemented for the testing of δ -IDE assuming parental genotypes are available for the unrelated cases and controls. For fair comparisons, power results are only presented in the absence of population stratification when logistic regression and the case-only method have no bias. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population. All tests were two-sided and were conducted at a significance level of 0.05.

Extended Figure 3. Biases and SD of estimates of DE of PGS using PGS-TRI and logistic regression when δ -IDE of parental PGS is not incorporated into modeling.

bottom left panels, data were simulated assuming no δ -IDE. In the top right and bottom right panels, data were simulated in the presence of δ -IDE. Both PGS-TRI and logistic regression were fitted without the parental indirect effect parameters in the model.

Extended Figure 4. Geographical distributions of (a) Educational Attainment (EA) PGS, (b) Educational Attainment, (c) Body Mass Index (BMI), (d) Townsend Index, (e)-(h) Top 4 Principal Components of unrelated UK Biobank participants (considered as “parents” in our simulation study). Individuals were grouped into 100 clusters based on their east and north co-ordinates of birthplaces using the K-means clustering. Colors represent the mean values of grouped individuals in each cluster on the map. The between-cluster correlation between BMI and EA-PGS is -0.48 ($P < 5 \times 10^{-6}$), the within-cluster correlation is -0.018 ($P > 0.05$). Parental BMI was treated as a hidden environmental confounding variable while simulating outcome status in children.

Extended Figure 5. Quantile-quantile plot (QQ plot) of p-values generated by PGS-TRI for the transcriptome-wide association study of autism risk in the SPARK study. The DE and δ -IDE of PGS associated with a total of 4,989 transcriptomic traits are tested. The diagonal lines correspond to expected p-values percentiles under the null hypothesis and the shaded regions represent 95% confidence bands. **DE:** PGS direct effect; **δ -IDE:** differential parental indirect genetic effect; λ_{GC} : genomic inflation factor, a value that equals 1 represents no inflation.

Extended Figure 6. Principal component (PC) analysis of 27 obesity-related metabolite PGS using reference genotype data from the 1000 Genomes European population. Results indicate that more than 98% of variations are explained by a single PGS indicating the number of underlying effective tests is close to 1.

Extended Figure 7. Quantile-quantile plot (QQ plot) of p-values generated by PGS-TRI for the transcriptome-wide association study of the risk of OFC CL/P subtype in the GENEVA study. The DE and δ -IDE of PGS associated with a total of 4,991 transcriptomic traits are tested. The diagonal lines correspond to expected p-values percentiles under the null hypothesis and the shaded regions represent 95% confidence bands. **DE:** PGS direct effect; **δ -IDE:** differential parental indirect genetic effect; λ_{GC} : genomic inflation factor, a value that equals 1 represents no inflation.

Data and Code Availability

GENEVA datasets are available in dbGaP at <https://www.ncbi.nlm.nih.gov> through dbGaP accession number phs000094.v1.p1.

The software PGS-TRI and the polygenic transmission disequilibrium test are publicly available at <https://github.com/ziquaow/PGS.TRI>

Acknowledgments

Drs. Wang and Chatterjee were supported by the National Institutes of Health (NIH) grant 1R01HG010480-01. Dr. Wang was also supported by 1K99HG013674-01. The research of Dr. Chatterjee was also supported by U01CA249866, U01HG011719, and 1U24OD023382-01. The GEARS project and Drs. Volk and Ladd-Acosta were supported by funding from R01ES034554 for this work. Dr. Ray was supported by R35GM150836. Drs. Beaty and Ruczinski were supported by NIDCR R01DE031855.

This research has been conducted using the UK Biobank Resource under Application Number 17731.

We are grateful to all of the families in SPARK, the SPARK clinical sites and SPARK staff. We appreciate obtaining access to phenotypic and genetic data on SFARI Base. Approved researchers can obtain the SPARK population dataset described in this study by applying at <https://base.sfari.org>. We appreciate obtaining access to recruit participants through SPARK research match on SFARI Base.

Funding support for the study entitled “International Consortium to Identify Genes and Interactions Controlling Oral Clefts” was provided by several previous grants from the National Institute of Dental and Craniofacial Research (NIDCR). Data and samples were drawn from several studies awarded to members of this consortium. Funding to support original data collection, previous genotyping and analysis came from several sources to individual investigators. Funding for individual investigators include: R21-DE-013707 and R01-DE-014581 (Beaty); R37-DE-08559 and P50-DE-016215 (Murray, Marazita) and the Iowa Comprehensive Program to Investigate Craniofacial and Dental Anomalies (Murray); R01-DE-09886 (Marazita), R01-DE-012472 (Marazita), R01-DE-014677 (Marazita), R01-DE-016148 (Marazita), R21-DE-016930 (Marazita); R01-DE-013939 (Scott). Parts of this research were supported in part by the Intramural Research Program of the NIH, National Institute of Environmental Health Sciences (Wilcox, Lie). Additional recruitment was supported by the Smile Train Foundation for recruitment in China (Jabs, Beaty, Shi) and a grant from the Korean government (Jee).

The genome-wide association study, also known the Cleft Consortium, is part of the Gene Environment Association Studies (GENEVA) program of the trans-NIH Genes, Environment and Health Initiative [GEI] supported by U01-DE-018993. Genotyping services were provided by the Center for Inherited Disease Research (CIDR). CIDR is

fully funded through a federal contract from the National Institutes of Health (NIH) to The Johns Hopkins University, contract number HHSN268200782096C. Funds for genotyping were provided by the NIDCR through CIDR's NIH contract. Assistance with genotype cleaning, as well as with general study coordination, was provided by the GENEVA Coordinating Center (U01-HG-004446) and by the National Center for Biotechnology Information (NCBI). We sincerely thank all of the patients and families at each recruitment site for participating in this study, and we gratefully acknowledge the invaluable assistance of clinical, field and laboratory staff who contributed to this effort over the years.

Trans-Omics in Precision Medicine (TOPMed) program imputation panel (version Freeze5) was supported by the National Heart, Lung and Blood Institute (NHLBI); see www.nhlbiwgs.org. TOPMed study investigators contributed data to the reference panel, which can be accessed through the Michigan Imputation Server; see <https://imputationserver.sph.umich.edu>. The panel was constructed and implemented by the TOPMed Informatics Research Center at the University of Michigan (3R01HL-117626-02S1; contract HHSN268201800002I). The TOPMed Data Coordinating Center (3R01HL-120393-02S1; contract HHSN268201800001I) provided additional data management, sample identity checks, and overall program coordination and support. We gratefully acknowledge the studies and participants who provided biological samples and data for TOPMed.

References

1. Kerminen S, Martin AR, Koskela J, et al. Geographic variation and bias in the polygenic scores of complex diseases and traits in Finland. *Am J Hum Genet.* 2019;104(6):1169–1181. doi: 10.1016/j.ajhg.2019.05.001.
2. Berg JJ, Harpak A, Sinnott-Armstrong N, et al. Reduced signal for polygenic adaptation of height in UK Biobank. *Elife.* 2019;8:10.7554/eLife.39725. doi: 10.7554/eLife.39725.
3. Sohail M, Maier RM, Ganna A, et al. Polygenic adaptation on height is overestimated due to uncorrected stratification in genome-wide association studies. *Elife.* 2019;8:10.7554/eLife.39702. doi: 10.7554/eLife.39702.
4. Morris TT, Davies NM, Hemani G, Smith GD. Population phenomena inflate genetic associations of complex social traits. *Sci Adv.* 2020;6(16):eaay0328. doi: 10.1126/sciadv.aay0328.
5. Abdellaoui A, Dolan CV, Verweij KJH, Nivard MG. Gene-environment correlations across geographic regions affect genome-wide association studies. *Nat Genet.* 2022;54(9):1345–1354. doi: 10.1038/s41588-022-01158-0.
6. Davies NM, Howe LJ, Brumpton B, Havdahl A, Evans DM, Davey Smith G. Within family Mendelian randomization studies. *Hum Mol Genet.* 2019;28(R2):R170–R179. doi: 10.1093/hmg/ddz204.

7. Howe LJ, Nivard MG, Morris TT, et al. Within-sibship genome-wide association analyses decrease bias in estimates of direct genetic effects. *Nat Genet.* 2022;54(5):581–592. doi: 10.1038/s41588-022-01062-7.
8. Trejo S, Domingue BW. Genetic nature or genetic nurture? introducing social genetic parameters to quantify bias in polygenic score analyses. *Biodemography Soc Biol.* 2018;64(3-4):187–215. doi: 10.1080/19485565.2019.1681257.
9. Warrington NM, Beaumont RN, Horikoshi M, et al. Maternal and fetal genetic effects on birth weight and their relevance to cardio-metabolic risk factors. *Nat Genet.* 2019;51(5):804–814. doi: 10.1038/s41588-019-0403-1.
10. Wu Y, Zhong X, Lin Y, et al. Estimating genetic nurture with summary statistics of multigenerational genome-wide association studies. *Proc Natl Acad Sci U S A.* 2021;118(25):e2023184118. doi: 10.1073/pnas.2023184118. doi: 10.1073/pnas.2023184118.
11. Kong A, Thorleifsson G, Frigge ML, et al. The nature of nurture: Effects of parental genotypes. *Science.* 2018;359(6374):424–428. doi: 10.1126/science.aan6877.
12. Weiner DJ, Wigdor EM, Ripke S, et al. Polygenic transmission disequilibrium confirms that common and rare variation act additively to create risk for autism spectrum disorders. *Nat Genet.* 2017;49(7):978–985. doi: 10.1038/ng.3863.

13. Bulik-Sullivan B, Finucane HK, Anttila V, et al. An atlas of genetic correlations across human diseases and traits. *Nat Genet.* 2015;47(11):1236–1241. doi: 10.1038/ng.3406.
14. An J, Won S, Lutz SM, Hecker J, Lange C. Effect of population stratification on SNP-by-environment interaction. *Genet Epidemiol.* 2019;43(8):1046–1055. <https://doi.org/10.1002/gepi.22250>. doi: 10.1002/gepi.22250.
15. Abdellaoui A, Hugh-Jones D, Yengo L, et al. Genetic correlates of social stratification in great britain. *Nature Human Behaviour.* 2019;3(12):1332–1342. <https://doi.org/10.1038/s41562-019-0757-5>. doi: 10.1038/s41562-019-0757-5.
16. Veller C, Coop GM. Interpreting population- and family-based genome-wide association studies in the presence of confounding. *PLOS Biology.* 2024;22(4):e3002511. <https://doi.org/10.1371/journal.pbio.3002511>.
17. Prive F, Aschard H, Carmi S, et al. Portability of 245 polygenic scores when derived from the UK biobank and applied to 9 ancestry groups from the same cohort. *Am J Hum Genet.* 2022;109(1):12–23. doi: 10.1016/j.ajhg.2021.11.008.
18. SPARK Consortium. SPARK: A US cohort of 50,000 families to accelerate autism research. *Neuron.* 2018;97(3):488–493. doi: 10.1016/j.neuron.2018.01.015.
19. Grove J, Ripke S, Als TD, et al. Identification of common genetic risk variants for autism spectrum disorder. *Nat Genet.* 2019;51(3):431–444. doi: 10.1038/s41588-019-0344-8.

20. Zheutlin AB, Dennis J, Karlsson Linner R, et al. Penetrance and pleiotropy of polygenic risk scores for schizophrenia in 106,160 patients across four health care systems. *Am J Psychiatry*. 2019;176(10):846–855. doi: 10.1176/appi.ajp.2019.18091085.
21. Cai N, Revez JA, Adams MJ, et al. Minimal phenotyping yields genome-wide association signals of low specificity for major depression. *Nat Genet*. 2020;52(4):437–447. doi: 10.1038/s41588-020-0594-5.
22. Gui Y, Zhou X, Wang Z, et al. Sex-specific genetic association between psychiatric disorders and cognition, behavior and brain imaging in children and adults. *Transl Psychiatry*. 2022;12(1):347–6. doi: 10.1038/s41398-022-02041-6.
23. Lahey BB, Tong L, Pierce B, et al. Associations of polygenic risk for attention-deficit/hyperactivity disorder with general and specific dimensions of childhood psychological problems and facets of impulsivity. *J Psychiatr Res*. 2022;152:187–193. doi: 10.1016/j.jpsychires.2022.06.019.
24. Auton A, Abecasis GR, Altshuler DM, et al. A global reference for human genetic variation. *Nature*. 2015;526(7571):68–74. <https://doi.org/10.1038/nature15393>. doi: 10.1038/nature15393.
25. Xu Y, Ritchie SC, Liang Y, et al. An atlas of genetic scores to predict multi-omic traits. *Nature*. 2023;616(7955):123–131. doi: 10.1038/s41586-023-05844-9.

26. Zhang W, Venkataraghavan S, Hetmanski JB, et al. Detecting gene-environment interaction for maternal exposures using case-parent trios ascertained through a case with non-syndromic orofacial cleft. *Front Cell Dev Biol.* 2021;9:621018. doi: 10.3389/fcell.2021.621018.
27. Beaty TH, Murray JC, Marazita ML, et al. A genome-wide association study of cleft lip with and without cleft palate identifies risk variants near MAFB and ABCA4. *Nat Genet.* 2010;42(6):525–529. doi: 10.1038/ng.580.
28. Beaty TH, Ruczinski I, Murray JC, et al. Evidence for gene-environment interaction in a genome wide study of nonsyndromic cleft palate. *Genet Epidemiol.* 2011;35(6):469–478. doi: 10.1002/gepi.20595.
29. Yu Y, Alvarado R, Petty LE, et al. Polygenic risk impacts PDGFRA mutation penetrance in non-syndromic cleft lip and palate. *Hum Mol Genet.* 2022;31(14):2348–2357. doi: 10.1093/hmg/ddac037.
30. Lambert SA, Gil L, Jupp S, et al. The polygenic score catalog as an open database for reproducibility and systematic evaluation. *Nat Genet.* 2021;53(4):420–425. doi: 10.1038/s41588-021-00783-5.
31. Yang Y, Suzuki A, Iwata J, Jun G. Secondary genome-wide association study using novel analytical strategies disentangle genetic components of cleft lip and/or cleft palate in 1q32.2. *Genes (Basel).* 2020;11(11):1280. doi: 10.3390/genes11111280. doi: 10.3390/genes11111280.

32. Schwender H, Taub MA, Beaty TH, Marazita ML, Ruczinski I. Rapid testing of SNPs and Gene–Environment interactions in Case–Parent trio data based on exact analytic parameter estimation. *Biometrics*. 2012;68(3):766–773. <https://doi.org/10.1111/j.1541-0420.2011.01713.x>. doi: 10.1111/j.1541-0420.2011.01713.x.
33. Spielman RS, McGinnis RE, Ewens WJ. Transmission test for linkage disequilibrium: The insulin gene region and insulin-dependent diabetes mellitus (IDDM). *Am J Hum Genet*. 1993;52(3):506–516.
34. Risch R, Merikangas M. The future of genetic studies of complex human diseases. *Science*. 1996;273(5281):1516–1517. <https://doi.org/10.1126/science.273.5281.1516>. doi: 10.1126/science.273.5281.1516.
35. Ayres KL, Curnow RN. Detecting non-multiplicative genotype relative risks from transmissions of parental alleles to affected children. *J Hum Genet*. 2005;50(1):46–48. <https://doi.org/10.1007/s10038-004-0217-5>. doi: 10.1007/s10038-004-0217-5.
36. Clayton D, Jones H. Transmission/disequilibrium tests for extended marker haplotypes. *The American Journal of Human Genetics*. 1999;65(4):1161–1169. <https://www.sciencedirect.com/science/article/pii/S0002929707626196>. doi: 10.1086/302566.
37. SHAM PC, CURTIS D. An extended transmission/disequilibrium test (TDT) for multi-allele marker loci. *Ann Hum Genet*. 1995;59(3):323–336. <https://doi.org/10.1111/j.1469-1809.1995.tb00751.x>. doi: 10.1111/j.1469-1809.1995.tb00751.x.

38. Abecasis GR, Cookson WOC, Cardon LR. Pedigree tests of transmission disequilibrium. *European Journal of Human Genetics*. 2000;8(7):545–551. <https://doi.org/10.1038/sj.ejhg.5200494>. doi: 10.1038/sj.ejhg.5200494.
39. Schaid DJ. Case-parents design for gene-environment interaction. *Genet Epidemiol*. 1999;16(3):261–273. doi: 10.1002/(SICI)1098-2272(1999)16:33.0.CO;2-M.
40. Umbach DM, Weinberg CR. The use of case-parent triads to study joint effects of genotype and exposure. *Am J Hum Genet*. 2000;66(1):251–261. doi: 10.1086/302707.
41. Cordell HJ, Barratt BJ, Clayton DG. Case/pseudocontrol analysis in genetic association studies: A unified framework for detection of genotype and haplotype associations, gene-gene and gene-environment interactions, and parent-of-origin effects. *Genet Epidemiol*. 2004;26(3):167–185. <https://doi.org/10.1002/gepi.10307>. doi: 10.1002/gepi.10307.
42. Mitchell LE. Differentiating between fetal and maternal genotypic effects, using the transmission test for linkage disequilibrium. *Am J Hum Genet*. 1997;60(4):1006–1007.
43. Weinberg CR, Wilcox AJ, Lie RT. A log-linear approach to case-parent-triad data: Assessing effects of disease genes that act either directly or through maternal effects and that may be subject to parental imprinting. *Am J Hum Genet*. 1998;62(4):969–978. doi: 10.1086/301802.
44. Horvath S, Xu X, Laird NM. The family based association test method: Strategies for studying general genotype–phenotype associations. *European Journal of Human*

Genetics. 2001;9(4):301–306. <https://doi.org/10.1038/sj.ejhg.5200625>. doi: 10.1038/sj.ejhg.5200625.

45. Lange C, Laird NM. On a general class of conditional tests for family-based association studies in genetics: The asymptotic distribution, the conditional power, and optimality considerations. *Genet Epidemiol*. 2002;23(2):165–180.

<https://doi.org/10.1002/gepi.209>. doi: 10.1002/gepi.209.

46. Lange C, DeMeo D, Silverman EK, Weiss ST, Laird NM. PBAT: Tools for family-based association studies. *The American Journal of Human Genetics*. 2004;74(2):367–369. <https://www.sciencedirect.com/science/article/pii/S0002929707618473>. doi: 10.1086/381563.

47. Zhang S, Lin T, Zhang Y, Liu X, Huang H. Effects of parental overweight and obesity on offspring's mental health: A meta-analysis of observational studies. *PLOS ONE*. 2022;17(12):e0276469. <https://doi.org/10.1371/journal.pone.0276469>.

48. Li Y, Ou J, Liu L, Zhang D, Zhao J, Tang S. Association between maternal obesity and autism spectrum disorder in offspring: A meta-analysis. *J Autism Dev Disord*. 2016;46(1):95–102. <https://doi.org/10.1007/s10803-015-2549-8>. doi: 10.1007/s10803-015-2549-8.

49. Gaunt TR, Shihab HA, Hemani G, et al. Systematic identification of genetic influences on methylation across the human life course. *Genome Biol*. 2016;17:61–z. doi: 10.1186/s13059-016-0926-z.

50. Burgess S, Thompson SG. Multivariable mendelian randomization: The use of pleiotropic genetic variants to estimate causal effects. *Am J Epidemiol*.

2015;181(4):251–260. doi: 10.1093/aje/kwu283.

Main Figures and Extended Figures

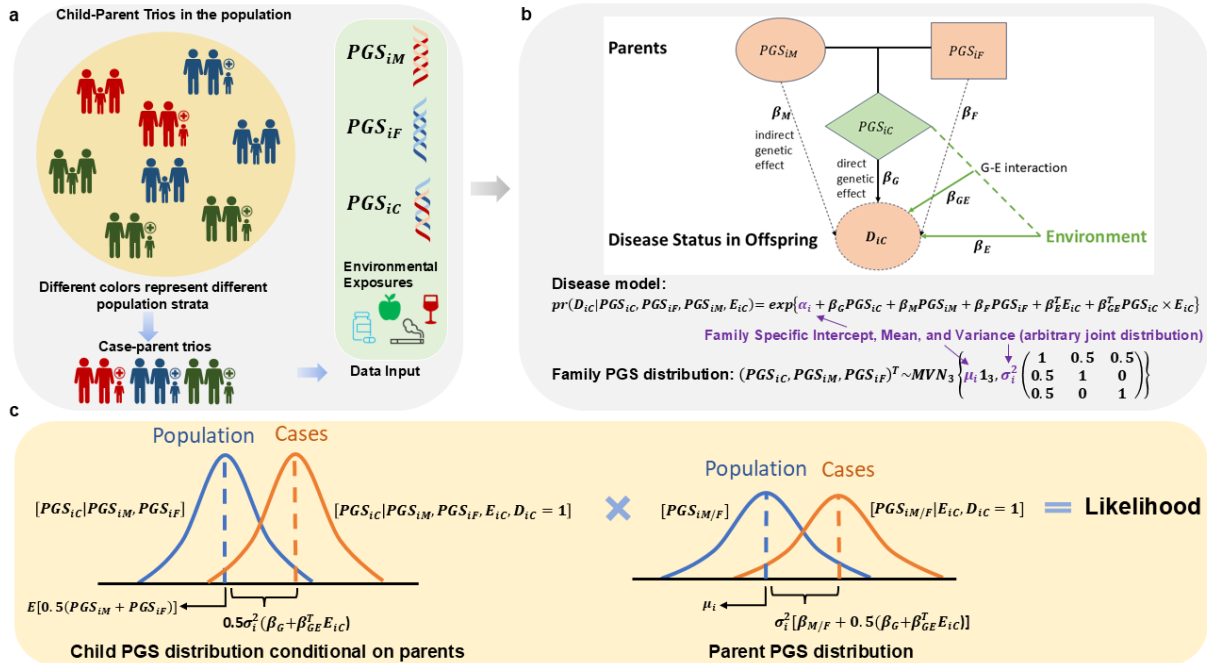


Figure 1. Figure illustrating our PGS-TRI model framework for case-parent trio family study designs. $\beta_G, \beta_{GE}, \beta_E$: population-level PGS direct genetic effect (β_G), direct PGS-E interactions (β_{GE}), and direct environmental effects (β_E), associated with offspring's outcome. β_M, β_F : mother and father indirect genetic effects associated with the offspring's outcome risk. $PGS_{iC}, PGS_{iM}, PGS_{iF}$: PGS values of child, mother, father in family $i = 1, \dots, N$. D_{iC} : outcome status of the offspring in family $i = 1, \dots, N$.

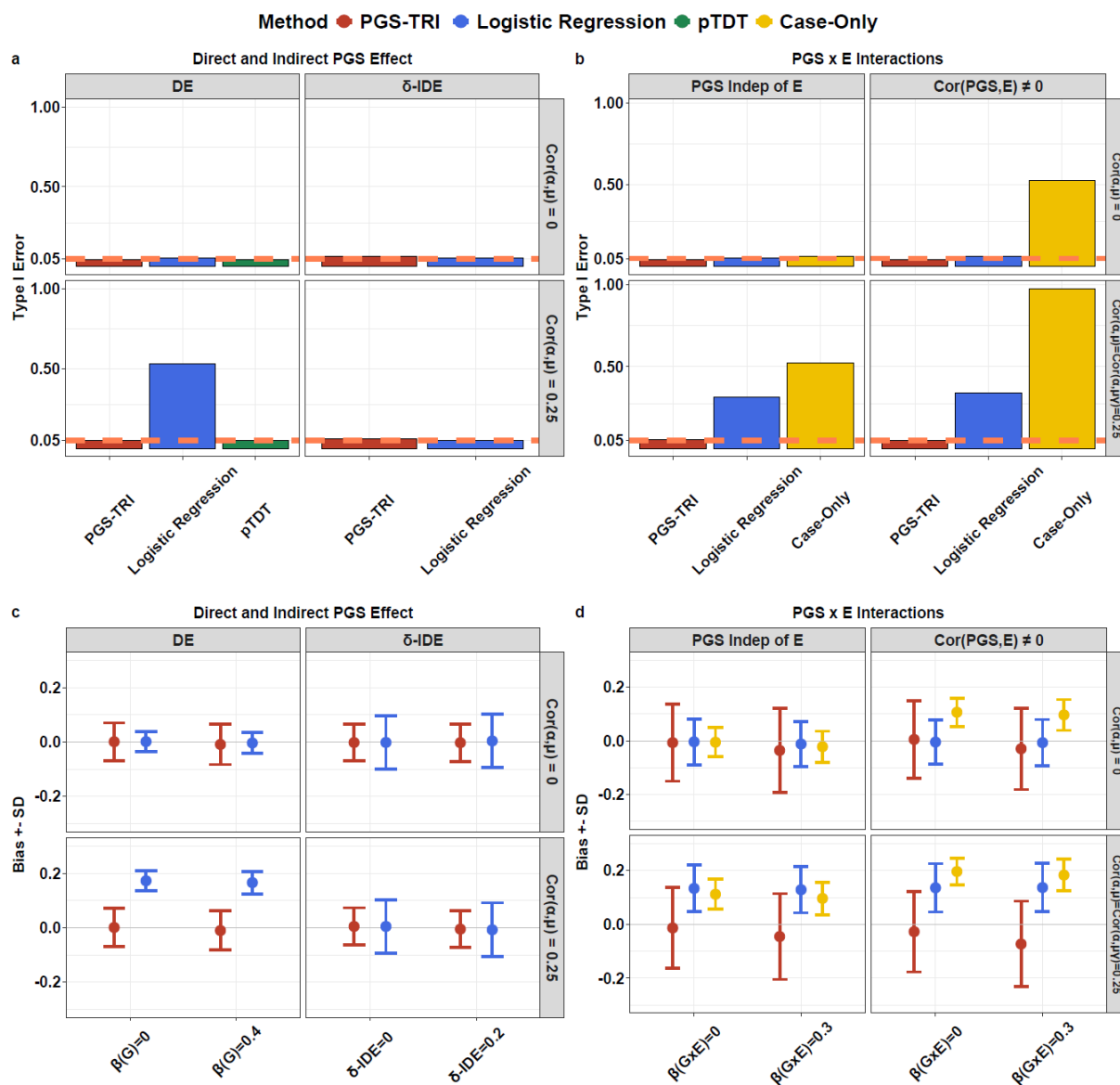


Figure 2. Performance of PGS-TRI and alternative methods for estimating parameters of PGS direct effect (DE), differential parental indirect genetic effect (δ -IDE), and PGSxE interactions in simulation studies. Results are shown for (a) Type I error of PGS DE and δ -IDE, when underlying true effects are 0; (b) Type I error of PGS-E interaction, when underlying true main effect DE is 0.4; (c) bias in parameter estimates and magnitude of standard deviation (SD) of estimates for DE and δ -IDE; (d) bias in parameter estimates and magnitude of SD of estimate for PGS-E interaction. pTDT is implemented as an alternative method for testing DE. Logistic regression is implemented for testing and estimation of DE assuming unrelated controls are available of the same size as the number of cases. Logistic regression is also implemented for testing and estimation of δ -IDE further assuming parental genotypes are available for the unrelated cases and controls. Further, a case-only method is also implemented for testing

PGSxE interaction. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population.

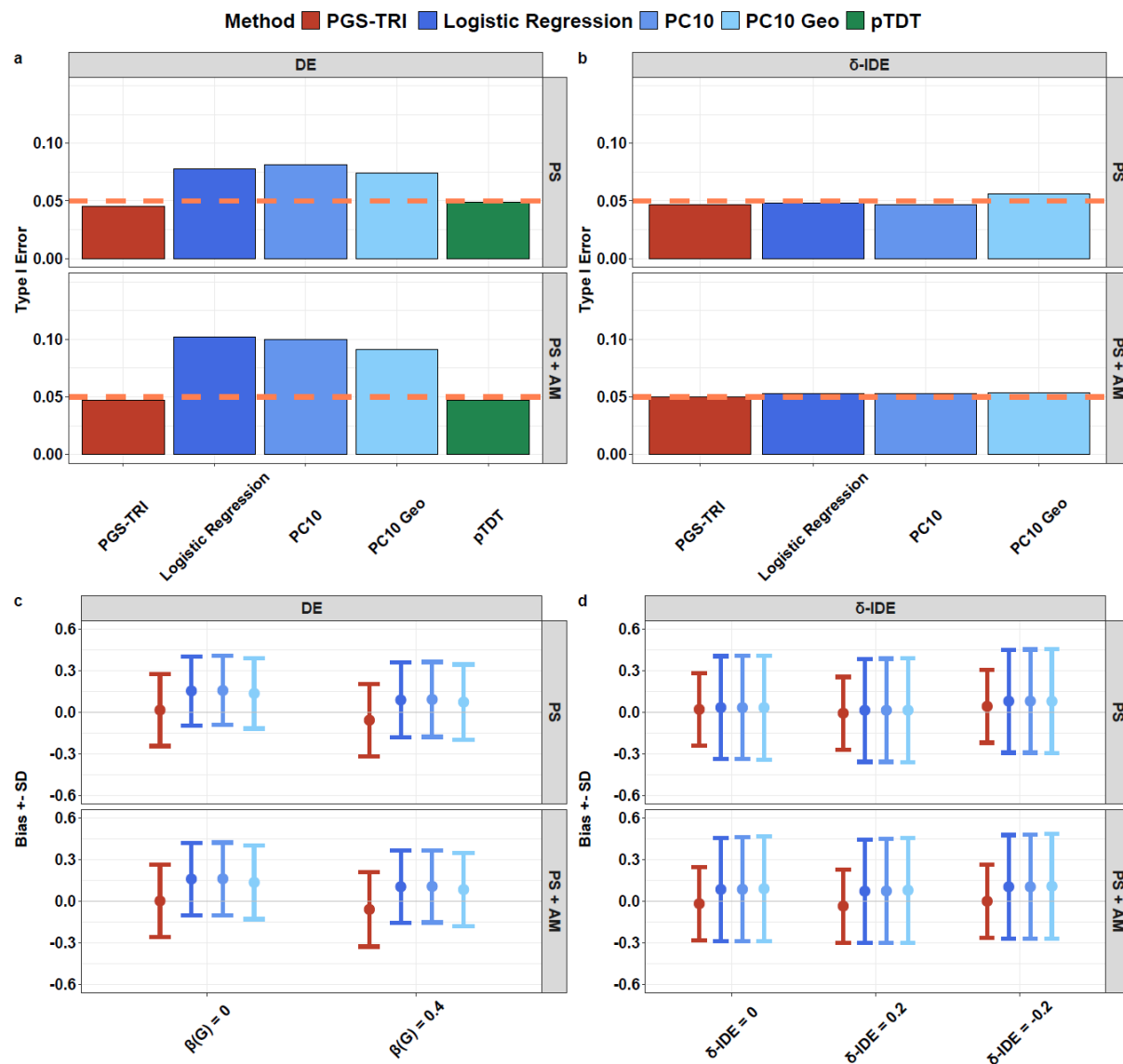


Figure 3. Performance of PGS-TRI and alternative methods for the estimation of genetic effects associated with educational attainment PGS in the UK Biobank-based simulation study. Results are shown for (a) Type I error associated with direct-effect (DE) of PGS (b) Type I error associated with differential parental indirect genetic effect (δ -IDE) of PGS (c) bias and SD associated with estimation of DE and (d) bias and SD associated with estimation of δ -IDE. Among alternative methods, pTDT is implemented for testing of DE. For testing of DE, logistic regression is implemented for the analysis of unrelated cases and controls without any adjustment, or adjustment for top 10 genetic principal components (PCs) constructed from the

parental data, or top 10 PCs plus the assessment centres, and north and east birth coordinates of the parents. For the testing and estimation of δ -IDE using logistic regression, we assume parental genotype data are available on unrelated cases and controls. Data on children for matched pairs of UKB participants are repeatedly simulated, and then a set of 2000 case-parent trios, or a set of 2000 unrelated cases and 2000 unrelated controls, are further sampled for subsequent analysis. Parents were either matched by only geographic proximity to simulate effect of population stratification or geographical proximity and educational attainment level to simulate effect of both population stratification and assortative mating. **PS**: population stratification bias; **AM**: assortative mating; **PC10**: logistic regression with the top 10 genetic PCs as covariates; **PC10 Geo**: logistic regression with the top 10 genetic PCs, north and east birth co-ordinates, and assessment centres as covariates in the model.

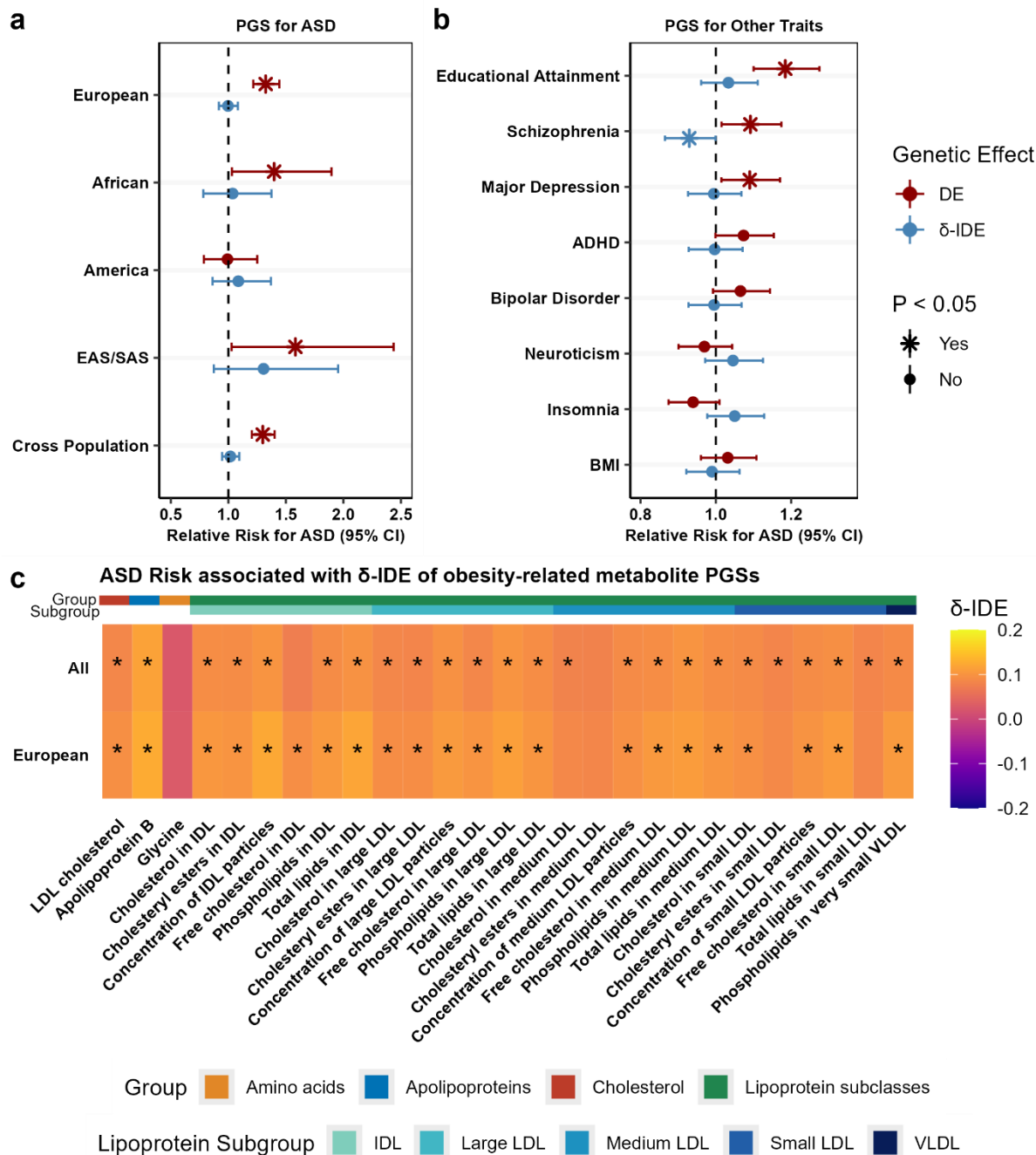
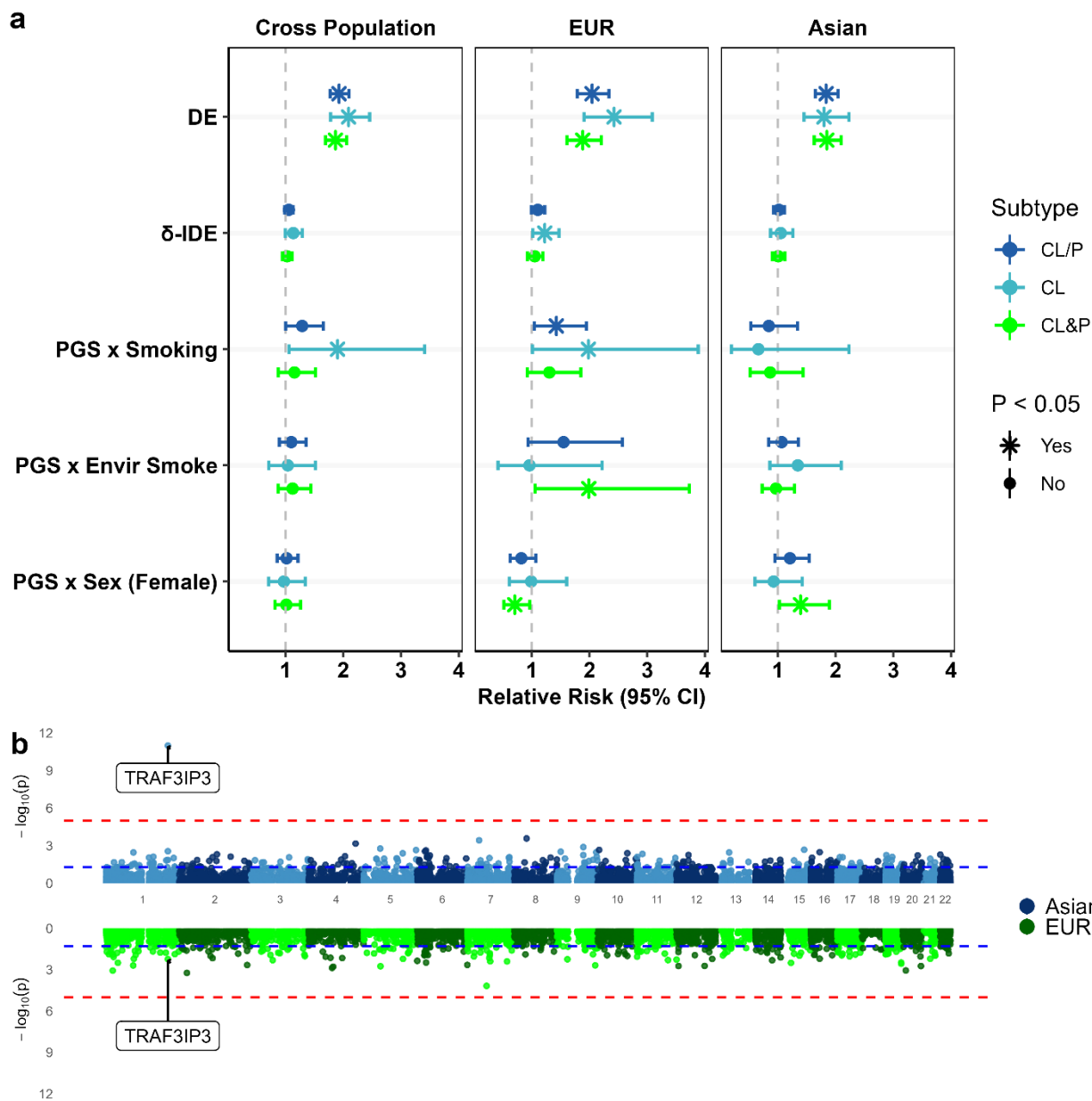
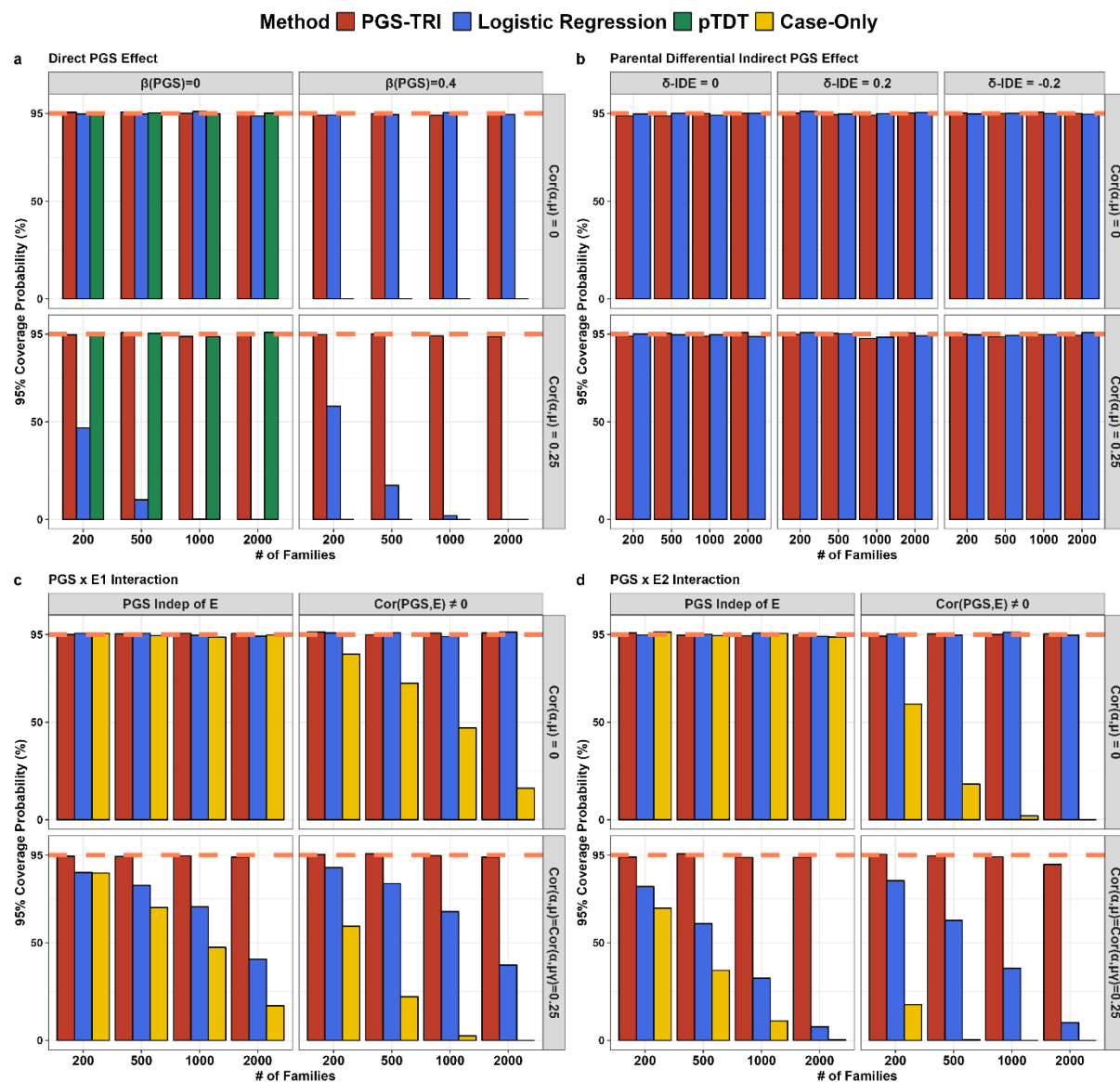


Figure 4. SPARK study results for autism spectrum disorder (ASD). (a) Relative risk estimates for direct (DE) and maternally mediated indirect effects (δ -IDE) of ASD-PGS (PGS ID: PGS000327) on ASD risk across multiple ancestry groups; (b) Relative risk estimates of DE and δ -IDE of PGS for multiple neurocognitive traits on ASD risk (only shown for combined population analysis due to sample size in ancestral subpopulations); (c) Heatmap of estimates of maternally mediated indirect effects (δ -IDE) of PGS for 27 obesity-related metabolite PGS on the risk of ASD in offsprings. In (a) and (b), PGS are standardized within each ancestry group by population mean and SD calculated using 1000G reference data of independent individuals so that relative risk (RR) corresponds to an increase in risk per SD unit increase in PGS value. In (c), PGSs are standardized using population mean and SD by pooling all individuals within each ancestry group in the current study. Asterisks indicate significant

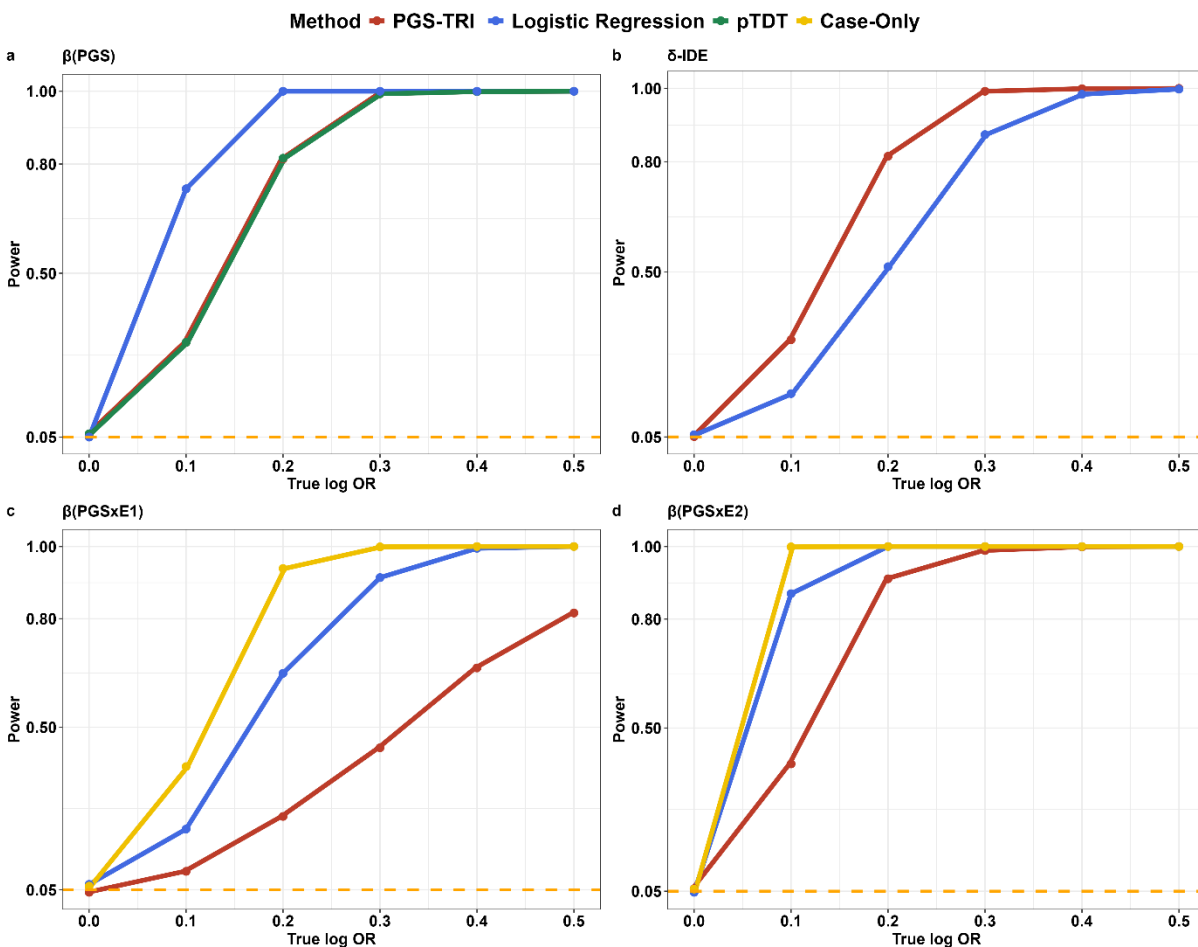
metabolites after Bonferonni correction for the effective number of independent tests ($P < 0.025$). **ADHD**: attention-deficit/hyperactivity disorder; **EAS/SAS**: combined East Asian and South Asian populations; **CI**: confidence interval; **BMI**: body mass index, used as a negative control; **IDL**: intermediate-density lipoproteins; **LDL**: low-density lipoprotein; **VLDL**: very small low-density lipoprotein.



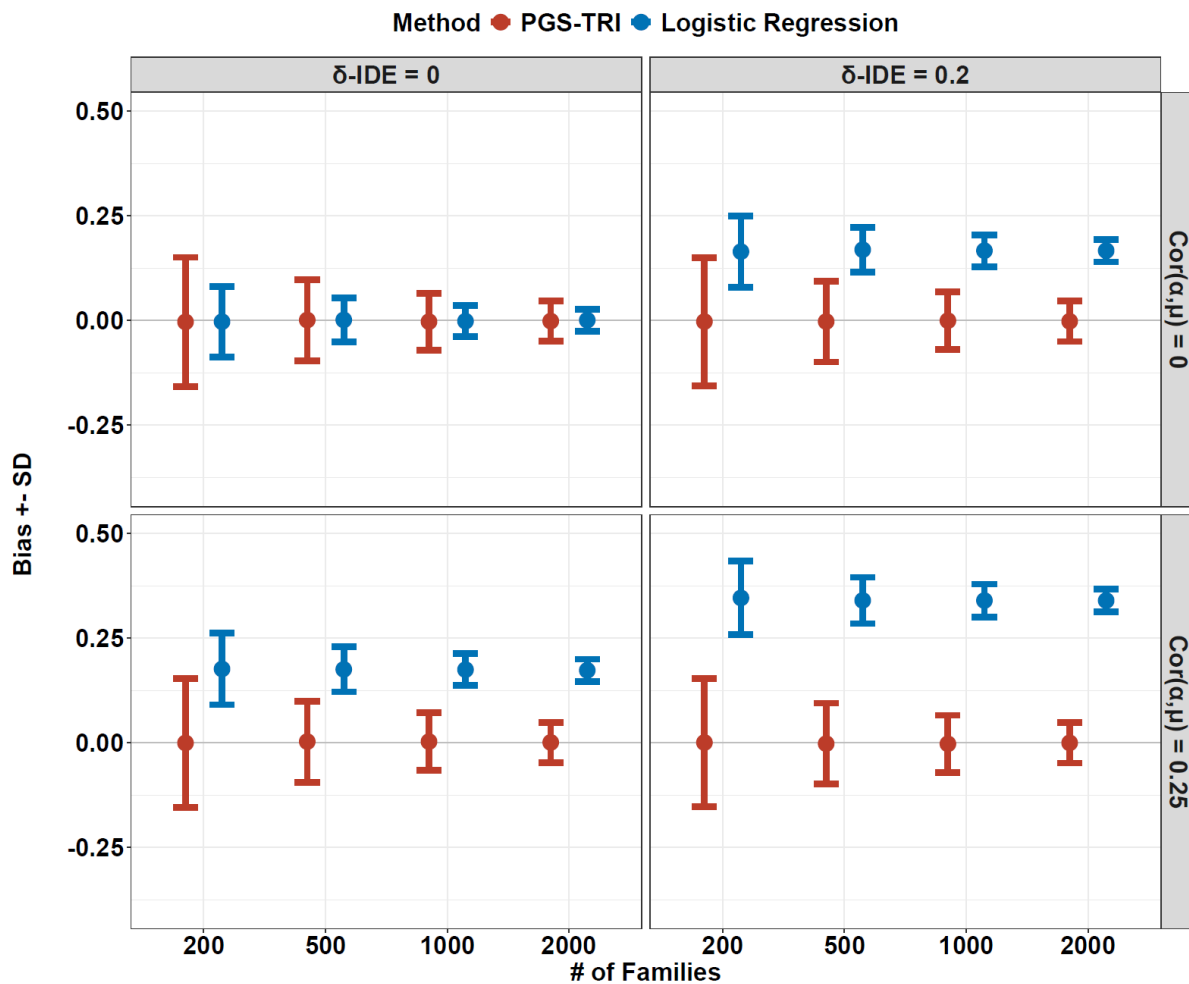
Extended Figures



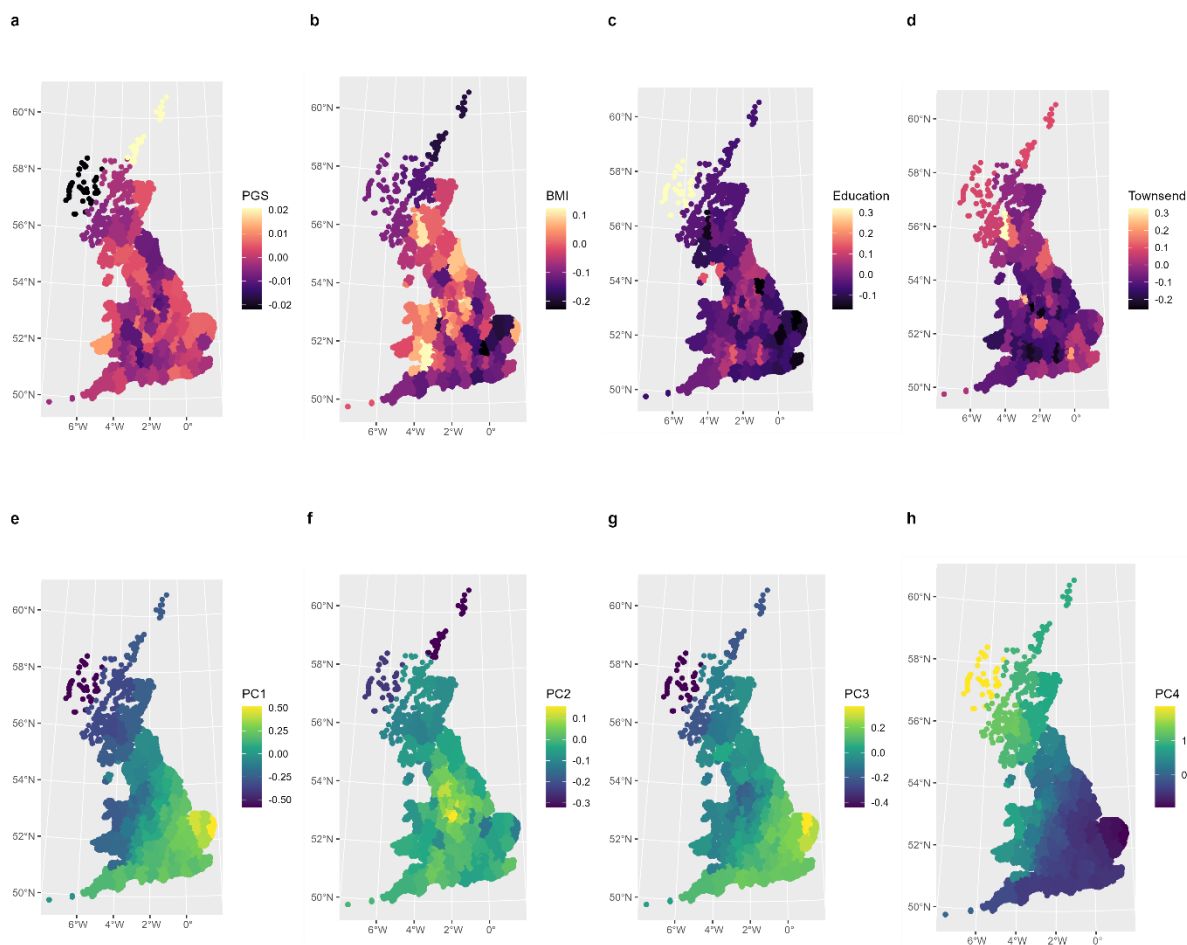
Extended Figure 1. Coverage level of 95% confidence intervals of PGS-TRI and alternative methods for estimation of (a) DE, (b) δ -IDE, and (c-d) PGSxE interactions in simulation studies. For each type of parameter, results are shown in scenarios in the absence and the presence of population stratification bias. Logistic regression is implemented for the estimation of DE and PGS by E interaction effects assuming unrelated controls are available of the same size as the number of cases. Logistic regression is also implemented for the estimation of δ -IDE further assuming that parental genotypes are available for the unrelated cases and controls. Additionally, a case-only method is implemented for testing PGSxE interaction terms. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population.



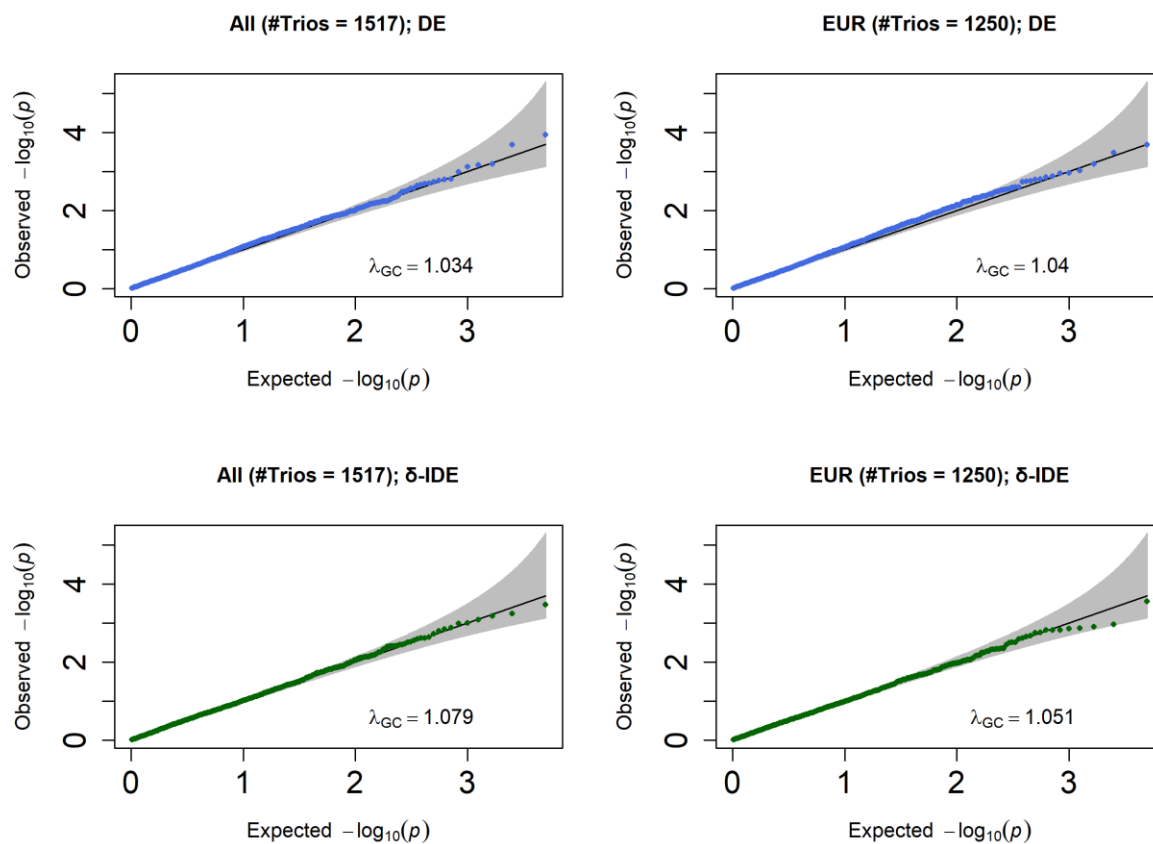
Extended Figure 2. Power of PGS-TRI compared to alternative methods for the detection of (a) DE, (b) δ -IDE, (c)(d) PGS by E interaction terms. pTDT is implemented as an alternative method for testing DE. The case-only method is implemented as an alternative method for testing PGSxE. Logistic regression is also implemented for testing of DE and PGSxE assuming unrelated controls are available of the same size as the number of cases. Logistic regression is further implemented for the testing of δ -IDE assuming parental genotypes are available for the unrelated cases and controls. For fair comparisons, power results are only presented in the absence of population stratification when logistic regression and the case-only method have no bias. Data are repeatedly simulated for 1000 trios, or 1000 unrelated cases and 1000 unrelated controls from the underlying population. All tests were two-sided and were conducted at a significance level of 0.05.



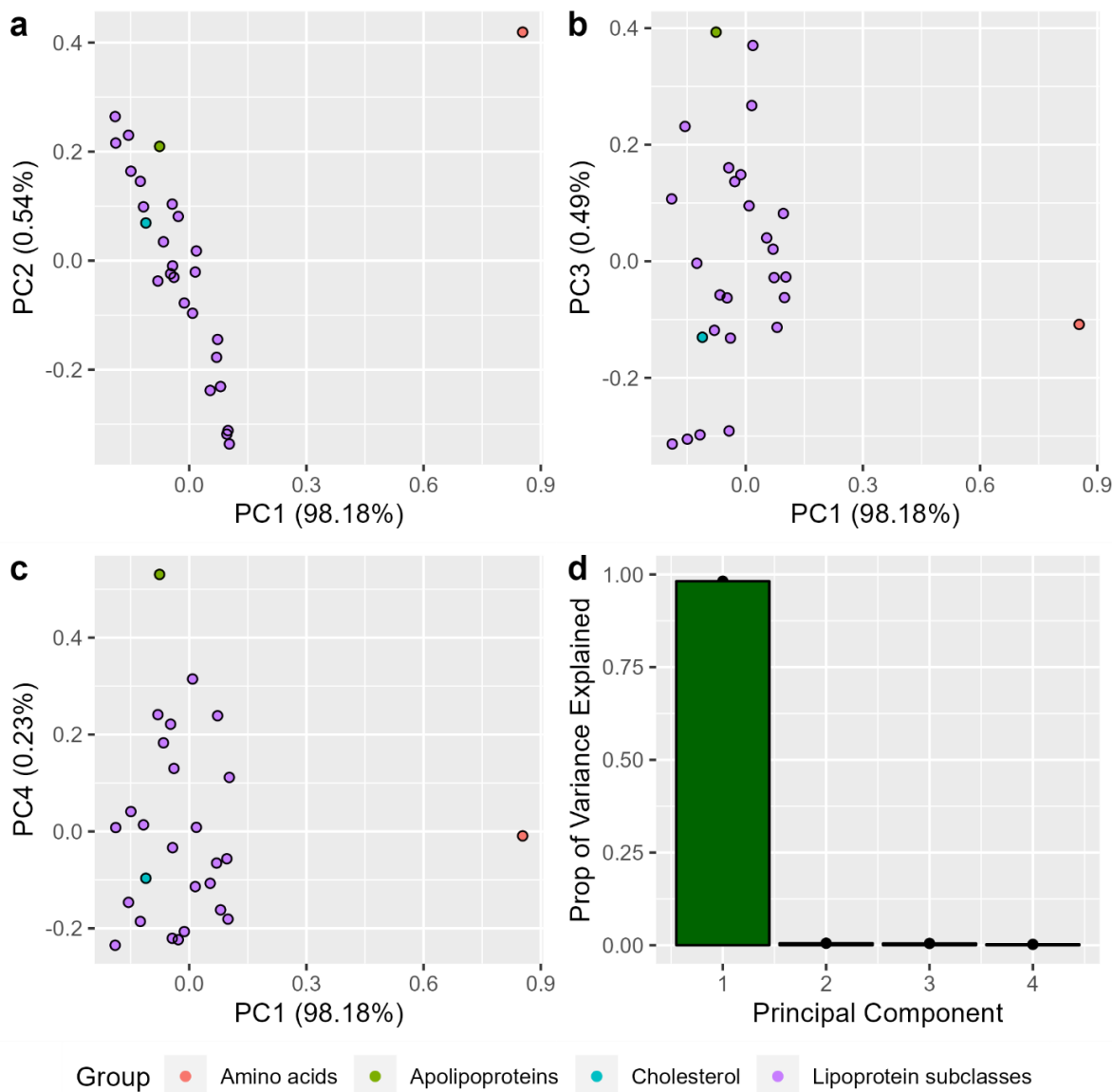
Extended Figure 3. Biases and SD of estimates of DE of PGS using PGS-TRI and logistic regression when δ -IDE of parental PGS is not incorporated into modeling. In the top left and bottom left panels, data were simulated assuming no δ -IDE. In the top right and bottom right panels, data were simulated in the presence of δ -IDE. Both PGS-TRI and logistic regression were fitted without the parental indirect effect parameters in the model.



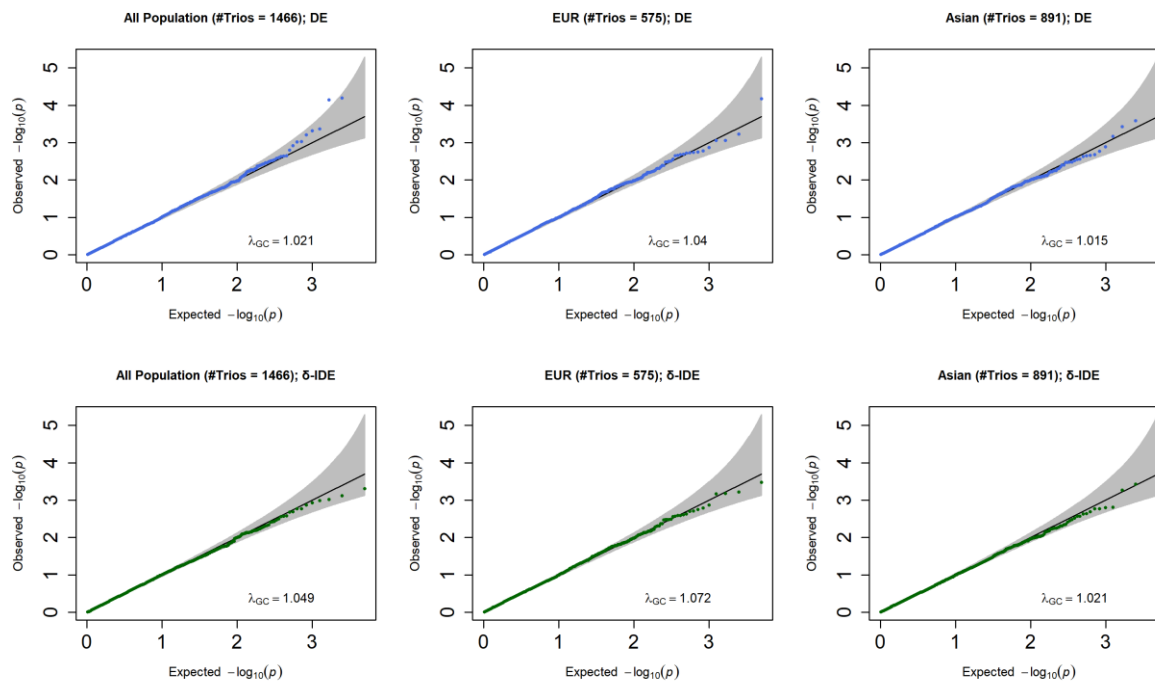
Extended Figure 4. Geographical distributions of (a) Educational Attainment (EA) PGS, (b) Educational Attainment, (c) Body Mass Index (BMI), (d) Townsend Index, (e)-(h) Top 4 Principal Components of unrelated UK Biobank participants (considered as “parents” in our simulation study). Individuals were grouped into 100 clusters based on their east and north co-ordinates of birthplaces using the K-means clustering. Colors represent the mean values of grouped individuals in each cluster on the map. The between-cluster correlation between BMI and EA-PGS is -0.48 ($P < 5 \times 10^{-6}$), the within-cluster correlation is -0.018 ($P > 0.05$). Parental BMI was treated as a hidden environmental confounding variable while simulating outcome status in children.



Extended Figure 5. Quantile-quantile plot (QQ plot) of p-values generated by PGS-TRI for the transcriptome-wide association study of autism risk in the SPARK study. The DE and δ -IDE of PGS associated with a total of 4,989 transcriptomic traits are tested. The diagonal lines correspond to expected p-values percentiles under the null hypothesis and the shaded regions represent 95% confidence bands. **DE**: PGS direct effect; **δ -IDE**: differential parental indirect genetic effect; λ_{GC} : genomic inflation factor, a value that equals 1 represents no inflation.



Extended Figure 6. Principal component (PC) analysis of 27 obesity-related metabolite PGS using reference genotype data from the 1000 Genomes European population. Results indicate that more than 98% of variations are explained by a single PGS indicating the number of underlying effective tests is close to 1.



Extended Figure 7. Quantile-quantile plot (QQ plot) of p-values generated by PGS-TRI for the transcriptome-wide association study of the risk of OFC CL/P subtype in the GENEVA study. The DE and δ -IDE of PGS associated with a total of 4,991 transcriptomic traits are tested. The diagonal lines correspond to expected p-values percentiles under the null hypothesis and the shaded regions represent 95% confidence bands. **DE**: PGS direct effect; **δ -IDE**: differential parental indirect genetic effect; λ_{GC} : genomic inflation factor, a value that equals 1 represents no inflation.