

# An Atlas of Genetic Correlations across Human Diseases and Traits

Brendan Bulik-Sullivan<sup>†\*,1,2,3</sup>, Hilary K Finucane<sup>\*,4</sup>, Verner Anttila<sup>1,2,3</sup>, Alexander Gusev<sup>5,6</sup>, Felix R. Day<sup>7</sup>, ReproGen Consortium<sup>8</sup>, Psychiatric Genomics Consortium<sup>8</sup>, Genetic Consortium for Anorexia Nervosa of the Wellcome Trust Case Control Consortium<sup>3,8</sup>, John R.B. Perry<sup>7</sup>, Nick Patterson<sup>1</sup>, Elise Robinson<sup>1,2,3</sup>, Mark J Daly<sup>1,2,3</sup>, Alkes L Price<sup>\*\*,1,5,6</sup>, and Benjamin M Neale<sup>†\*\*,1,2,3</sup>

<sup>1</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>2</sup>Stanley Center for Psychiatric Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>3</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts, USA.

<sup>4</sup>Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA, USA.

<sup>5</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA.

<sup>6</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA.

<sup>7</sup>MRC Epidemiology Unit, University of Cambridge School of Clinical Medicine, Institute of Metabolic Science, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK

<sup>8</sup>A list of members and affiliations appears in the Supplementary Note.

## Abstract

The major challenges preventing estimation of genetic correlation from GWAS data are the lack of availability of individual genotype data and widespread sample overlap. We circumvent these difficulties by introducing a technique for estimating genetic correlation that requires only GWAS summary statistics and is not biased by sample overlap. We use our method to estimate 300 genetic correlations among 25 traits, totaling more than 1.5 million unique phenotype measurements. Our results include a positive genetic correlation between anorexia nervosa and schizophrenia and a negative genetic correlation between anorexia nervosa and body mass index, as well as a large number of replications and positive controls. These results highlight the power of a polygenic modeling framework, since there currently are no genome-wide significant SNPs for anorexia nervosa.

---

\*Co-first authors

\*\*Co-last authors

<sup>†</sup>Address correspondence to BBS ([bulik@broadinstitute.org](mailto:bulik@broadinstitute.org)) or BMN ([bneale@broadinstitute.org](mailto:bneale@broadinstitute.org)).

# Introduction

Discovering correlations between phenotypes is a fundamental goal of epidemiology, with applications to classification and treatment of disease, as well as development of pharmaceutical drugs. One classical strategy in epidemiology is to search for correlations between phenotypes via observational studies; however, these studies can be biased by confounding and reverse causation [1, 2]. For heritable traits, an alternative strategy that is more robust to confounding is to scan for pairs of phenotypes with shared genetic etiology.

The first methods for testing for genetic overlap were family studies [3–7]. The disadvantage of these methods is the requirement to measure all traits on the same individuals, which scales poorly to studies of multiple traits, especially traits that are difficult or costly to measure (*e.g.*, low-prevalence diseases). Genome-wide association studies (GWAS) produce effect-size estimates for specific genetic variants, so it is possible to test for shared genetics by looking for correlations in effect-sizes across traits, which does not require measuring multiple traits per individual.

A widely-used technique for testing for relationships between phenotypes using GWAS data is Mendelian randomization (MR) [1, 2], which is the specialization to genetics of instrumental variables [8]. MR is effective for traits that are influenced by large-effect common variants [9, 10]; however, for many complex traits, heritability is distributed over thousands of variants with small effects [11]. For these traits, MR suffers from low power and weak instrument bias [8, 12].

In this paper, our goal is to estimate genetic correlation, a quantity that includes the effects of all SNPs, including those that do not reach genome-wide significance (Methods). In cases where two phenotypes have a cause-effect relationship, genetic correlation can be interpreted as a re-scaling of the MR estimate from a hypothetical instrument built using all SNPs (Methods). Genetic correlation is also meaningful for pairs of diseases, in which case it can be interpreted as the genetic analogue of comorbidity. The two main existing techniques for estimating genetic correlation from GWAS data are restricted maximum likelihood (REML) [13–18] and polygenic scores [19, 20]. These methods have only been applied to a few traits, because they require individual genotype data, which are difficult to obtain due to informed consent limitations.

We introduce a computationally fast technique based on LD Score regression [21] for estimating genetic correlation using only GWAS summary statistics. We apply this method to data from 25 GWAS and report genetic correlations for 300 pairs of phenotypes.

## Results

### Overview of Methods

Our method for estimating genetic correlation from summary statistics relies on the fact that the GWAS effect-size estimate for a given SNP incorporates the effects of all SNPs in linkage disequilibrium (LD) with that SNP [21, 22]. For a polygenic trait, SNPs with high LD will have higher  $\chi^2$ -statistics on average than SNPs with low LD [21]. A similar relationship holds if we replace  $\chi^2$ -statistics for a single study with the product of  $z$ -scores from two studies. Precisely, under a polygenic model [13, 15], the expected value of  $z_{1j}z_{2j}$  is

$$\mathbb{E}[z_{1j}z_{2j}] = \frac{\sqrt{N_1N_2}\rho_g}{M}\ell_j + \frac{\rho N_s}{\sqrt{N_1N_2}}, \quad (1)$$

where  $N_i$  is the sample size for study  $i$ ,  $\rho_g$  is genetic covariance (defined in Methods),  $\ell_j$  is LD Score [21],  $N_s$  is the number of individuals included in both studies, and  $\rho$  is the phenotypic correlation among the  $N_s$  overlapping samples. We derive this equation in the Supplementary Note. If study 1 and study 2 are the same study, then equation 1 reduces to the single-trait result from [21], because genetic covariance between a trait and itself is heritability, and  $\chi^2 = z^2$ . A similar result holds if one or both studies is a case/control study, in which case genetic covariance is on the observed scale. As a consequence of equation 1, we can estimate genetic covariance using the slope from the regression of  $z_{1j}z_{2j}$  on LD Score, which is computationally very fast (Methods). Normalizing genetic covariance to lie in  $[-1, 1]$  yields genetic correlation:  $r_g := \rho_g / \sqrt{h_1^2 h_2^2}$ , where  $h_i^2$  denotes the SNP-heritability [13] from study  $i$ . Theory and simulation confirm that our method produces robust estimates of genetic correlation, even when summary statistics are affected by population stratification [21] or sample overlap.

## Simulations

We performed a series of simulations to evaluate the robustness of the model to potential confounders such as sample overlap and model misspecification, and to verify the accuracy of the standard error estimates. Details of our simulation setup are provided in the methods. Table 1 shows LD Score regression estimates and standard errors from 1000 simulations of quantitative traits with complete sample overlap. The LD Score regression standard errors (Methods) matched the standard deviation across simulations, and the LD Score regression estimates were not biased by sample overlap. For comparison, if we (incorrectly) estimate genetic correlation using LD Score regression with intercept constrained to zero, we obtain out-of-bounds estimates around  $\hat{r}_g \approx 2$ . LD Score regression with incorrectly contained intercept is a representative example of an estimator that is not robust to sample overlap, similar to polygenic scores [20]. Table S1 shows results from simulations with one quantitative trait and one ascertained binary trait. These simulations confirm the derivations in the Supplementary Note and demonstrate that LD Score regression is not biased by oversampling of cases.

Parameter	Truth	Estimate	SD	SE
$h^2$	0.58	0.58	0.07	0.07
$\rho_g$	0.29	0.29	0.06	0.06
$r_g$	0.50	0.49	0.08	0.07

Table 1: *Simulations with complete sample overlap. Truth shows the true parameter values. Estimate shows the average LD Score regression estimate across 100 simulations. SD shows the standard deviation of the estimates across 100 simulations, and SE shows the mean LD Score regression SE across 100 simulations. Further details of the simulation setup are given in the Methods.*

Estimates of heritability and genetic covariance can be biased if the underlying model of genetic architecture is misspecified, *e.g.*, if variance explained is correlated with LD Score or MAF [21, 23]. Because genetic correlation is estimated as a ratio, it is more robust: biases that affect the numerator and the denominator in the same direction tend to cancel. We obtain approximately correct estimates of genetic correlation even in simulations with models of genetic architecture where our estimates of heritability and genetic covariance are biased (Table S2).

## Replication of Psychiatric Cross-Disorder Results

As technical validation, we replicated the estimates of genetic correlations among psychiatric disorders obtained with individual genotypes and REML in [16], by applying LD Score regression to summary statistics from the same data [24]. These summary statistics were generated from non-overlapping samples, so we applied LD Score regression using both unconstrained and constrained intercepts (Methods). Results from these analyses are shown in Figure 1. As expected, the results from LD Score regression were similar to the results from REML. LD Score regression with constrained intercept gave standard errors that were only slightly larger than those from REML, while the standard errors from LD Score regression with intercept were substantially larger, especially for traits with small sample sizes (*e.g.*, ADD, ASD).

## Application to Summary Statistics From 25 Phenotypes

We used cross-trait LD Score regression to estimate genetic correlations among 25 phenotypes (URLs, Methods). Genetic correlation estimates for all 300 pairwise combinations of the 25 traits are shown in Figure 2. For clarity of presentation, the 25 phenotypes were restricted to contain only one phenotype from each cluster of closely related phenotypes (Methods). Genetic correlations among the educational, anthropometric, smoking, and insulin-related phenotypes that were excluded from Figure 2 are shown in Table S4 and Figures S1, S2 and S3, respectively. References and sample sizes are shown in Table S3.

For the majority of pairs of traits in Figure 2, no GWAS-based genetic correlation estimate has been reported; however, many associations have been described informally based on the observation of overlap among genome-wide significant loci. Examples of genetic correlations that are consistent with overlap among top loci include the correlations between plasma lipids and cardiovascular disease [10]; age at onset of menarche and obesity [25]; type 2 diabetes, obesity, fasting glucose, plasma lipids and cardiovascular disease [26]; birth weight, adult height and type 2 diabetes [27,28]; birth length, adult height and infant head circumference [29,30]; and childhood obesity and adult obesity [29]. For many of these pairs of traits, we can reject the null hypothesis of zero genetic correlation with overwhelming statistical significance (*e.g.*,  $p < 10^{-20}$  for age at onset of menarche and obesity).

The first section of table 2 lists genetic correlation results that are consistent with epidemiological associations, but, as far as we are aware, have not previously been reported using genetic data. Our estimates of the genetic correlation between age at onset of menarche and adult height [31], cardiovascular disease [32] and type 2 diabetes [32,33] are consistent with the epidemiological associations. Our estimate of a negative genetic correlation between anorexia nervosa and obesity (and a similar genetic correlation with BMI) suggests that the same genetic factors influence normal variation in BMI as well as dysregulated BMI in psychiatric illness. This result is consistent with our observation that BMI GWAS findings implicate neuronal, rather than metabolic, cell-types and epigenetic marks [34]. The negative genetic correlation between adult height and coronary artery disease agrees with a replicated epidemiological association [35–37]. We observe several significant associations with the educational attainment phenotypes from Rietveld *et al.* [38]: we estimate a statistically significant negative genetic correlation between college and Alzheimer’s disease, which agrees with epidemiological results [39,40]. The positive genetic correlation between college and bipolar disorder is consistent with the psychiatric literature [41,42]. Our estimate of a negative genetic correlation between smoking and college is consistent with the fact that smoking is more

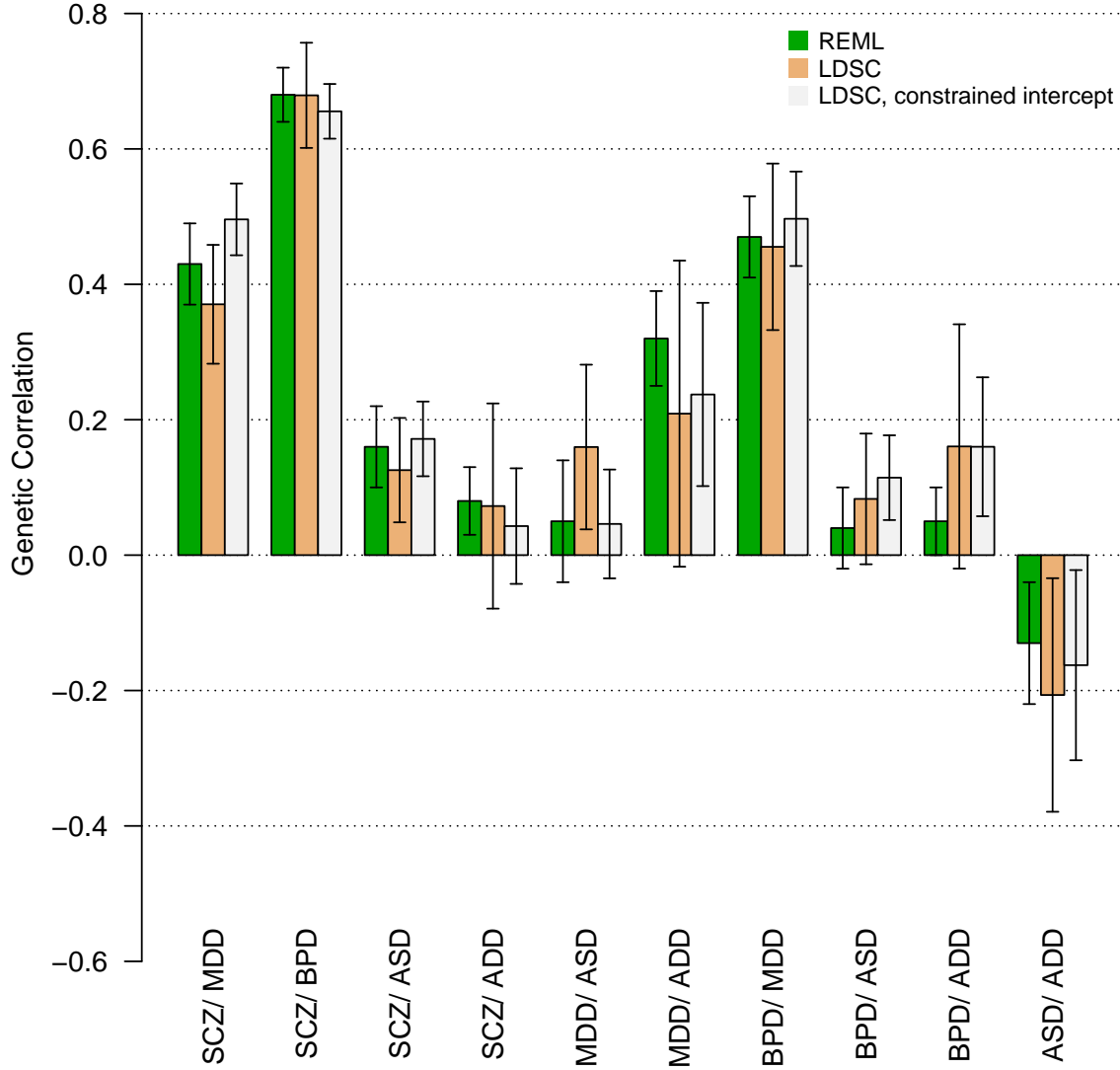


Figure 1: *Replication of Psychiatric Cross-Disorder Results.* This plot compares LD Score regression estimates of genetic correlation using the summary statistics from [24] to estimates obtained from REML with the same data [16]. The horizontal axis indicates pairs of phenotypes, and the vertical axis indicates genetic correlation. Error bars are standard errors. Green is REML; orange is LD Score with intercept and white is LD Score with constrained intercept. The estimates of genetic correlation among psychiatric phenotypes in figure 2 use larger sample sizes; this analysis is intended as a technical validation. Abbreviations: ADD = attention deficit disorder; ASD = autism spectrum disorder; BPD = bipolar disorder; MDD = major depressive disorder; SCZ = schizophrenia.

prevalent among less-educated groups [43].

The second section of table 2 lists three results that are, to the best of our knowledge, new both

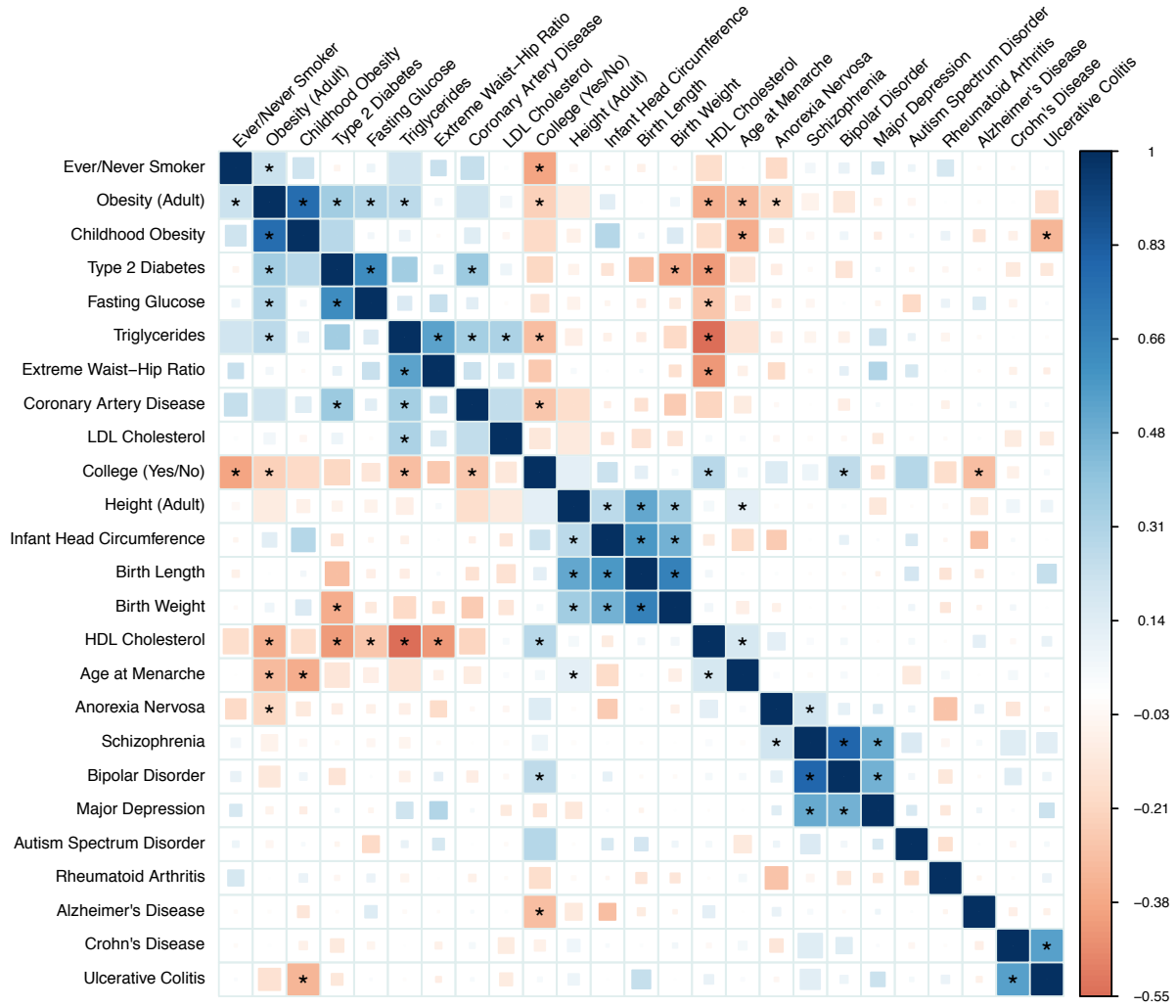


Figure 2: *Genetic Correlations among 25 GWAS. Blue represents positive genetic correlations; red represents negative. Larger squares correspond to more significant p-values. Genetic correlations that are different from zero at 1% FDR are shown as full-sized squares. Genetic correlations that are significantly different from zero after Bonferroni correction for the 300 tests in this figure have an asterisk. We show results that do not pass multiple testing correction as smaller squares in order to avoid whitening out positive controls where the estimate points in the expected direction, but does not achieve statistical significance due to small sample size. This multiple testing correction is conservative, since the tests are not independent.*

to genetics and epidemiology. One, we find a positive genetic correlation between anorexia nervosa and schizophrenia. Comorbidity between eating and psychotic disorders has not been thoroughly investigated in the psychiatric literature [44,45], and this result raises the possibility of similarity

	Phenotype 1	Phenotype 2	$r_g$ (se)	$p$ -value
Epidemiological	Age at Menarche	Height (Adult)	0.11 (0.03)	$6 \times 10^{-5}$ **
	Age at Menarche	Type 2 Diabetes	-0.13 (0.04)	$3 \times 10^{-3}$
	Age at Menarche	Triglycerides	-0.15 (0.04)	$1 \times 10^{-3}$ *
	Coronary Artery Disease	Age at Menarche	-0.11 (0.05)	$4 \times 10^{-2}$
	Coronary Artery Disease	College (Yes/No)	-0.278 (0.07)	$1 \times 10^{-4}$ **
	Coronary Artery Disease	Height (Adult)	-0.17 (0.05)	$2 \times 10^{-4}$ *
	Alzheimer's	College (Yes/No)	-0.30 (0.08)	$1 \times 10^{-4}$ **
	Bipolar Disorder	College (Yes/No)	0.026 (0.064)	$6 \times 10^{-5}$ **
	Obesity (Adult)	College (Yes/No)	-0.23 (0.04)	$2 \times 10^{-8}$ **
	Triglycerides	College (Yes/No)	-0.30 (0.04)	$5 \times 10^{-12}$ **
	Anorexia Nervosa	Obesity (Adult)	-0.20 (0.04)	$4 \times 10^{-6}$ **
	Ever/Never Smoker	College (Yes/No)	-0.39 (0.07)	$1 \times 10^{-9}$ **
	Ever/Never Smoker	Obesity (Adult)	0.22 (0.05)	$7 \times 10^{-5}$ **
New/Nonzero	Autism Spectrum Disorder	College (Yes/No)	0.28 (0.08)	$5 \times 10^{-4}$ *
	Ulcerative Colitis	Childhood Obesity	-0.33 (0.08)	$3.9 \times 10^{-5}$ **
	Anorexia Nervosa	Schizophrenia	0.19 (0.04)	$1.5 \times 10^{-5}$ **
New/Zero	Schizophrenia	Alzheimer's	0.05 (0.05)	0.58
	Schizophrenia	Ever/Never Smoker	0.03 (0.06)	0.26
	Schizophrenia	Triglycerides	-0.05 (0.04)	0.21
	Schizophrenia	LDL Cholesterol	-0.02 (0.03)	0.64
	Schizophrenia	HDL Cholesterol	0.03 (0.04)	0.50
	Schizophrenia	Rheumatoid Arthritis	-0.05 (0.05)	0.38
	Crohn's Disease	Rheumatoid Arthritis	-0.02 (0.09)	0.83
	Ulcerative Colitis	Rheumatoid Arthritis	-0.09 (0.09)	0.33

Table 2: Genetic correlation estimates, standard errors and  $p$ -values for selected pairs of traits. Results are grouped into genetic correlations that are new genetic results, but are consistent with established epidemiological associations (“Epidemiological”), genetic correlations that are new both to genetics and epidemiology (“New/Nonzero”) and interesting null results (“New/Zero”). The  $p$ -values are uncorrected  $p$ -values. Results that pass multiple testing correction for the 300 tests in Figure 2 at 1% FDR have a single asterisk; results that pass Bonferroni correction have two asterisks. We present some genetic correlations that agree with epidemiological associations but that do not pass multiple testing correction in our data.

between these classes of disease. Two, we estimate a negative genetic correlation between ulcerative colitis (UC) and childhood obesity. The relationship between premorbid BMI and ulcerative colitis is not well-understood; exploring this relationship may be a fruitful direction for further investigation. Three, we estimate a positive genetic correlation between autism spectrum disorder (ASD) and educational attainment, which itself has very high genetic correlation with IQ [38, 46, 47]. The ASD summary statistics were generated using a case-pseudocontrol study design, so this result cannot be explained by the tendency for the parents of children who receive a diagnosis of ASD to be better educated than the general population [48]. The distribution of IQ among individuals with ASD has lower mean than the general population, but with heavy tails [49] (*i.e.*, an excess of individuals with low and high IQ). There is evidence that the genetic architectures of high IQ and low IQ ASD are dissimilar [50]. We are unable to offer an explanation for this result, but propose that further exploration of this genetic correlation may be an important direction for future research.

The third section of table 2 lists interesting examples where the genetic correlation is close to zero with small standard error. The lack of genetic correlation between schizophrenia and rheumatoid arthritis is interesting because schizophrenia has been observed to be protective for rheumatoid arthritis [51]. The absence of genetic correlation between schizophrenia and smoking is notable because of the high prevalence of smoking among individuals with schizophrenia [52]. The absence of genetic correlation between schizophrenia and plasma lipid levels contrasts with a previous report of extensive pleiotropy between schizophrenia and triglycerides [53]. However, this observation from Andreassen, *et al.* [53] could be explained the sensitivity of the method used to the properties of a few regions with strong LD, rather than trait biology (Table S5). We estimate near-zero genetic correlation between Alzheimer’s disease and schizophrenia. The genetic correlations between Alzheimers disease and the other psychiatric traits (anorexia nervosa, bipolar, major depression, ASD) are also close to zero, but with larger standard errors, due to smaller sample sizes. This suggests that the genetic basis of Alzheimer’s disease is distinct from psychiatric conditions. Last, we estimate near zero genetic correlation between rheumatoid arthritis (RA) and both Crohn’s disease (CD) and UC. Although these diseases share many associated loci [54,55], there appears to be no directional trend. Some RA risk alleles are also risk alleles for UC and CD, but many RA risk alleles are protective for UC and CD [54], yielding near-zero genetic correlation. This is an example of pleiotropy without genetic correlation (Methods).

Finally, our estimates of genetic correlations among metabolic traits are consistent with the estimates obtained using REML in Vattikuti *et al.* [17] (Supplementary Table S4), and are directionally consistent with the recent Mendelian randomization results from Wuertz *et al.* [56]. Our estimate of 0.57 (0.074) for the genetic correlation between CD and UC is consistent with the estimate of 0.62 (0.042) from Chen *et al.* [18].

## Discussion

We have described a new method for estimating genetic correlation from GWAS summary statistics, which we applied to a dataset of GWAS summary statistics consisting of 25 traits and more than 1.5 million unique phenotype measurements. We reported several new findings that would have been difficult or impossible to obtain with existing methods, including a positive genetic correlation between educational attainment and autism spectrum disorder and a positive genetic correlation between anorexia nervosa and schizophrenia. Our method replicated many previously-reported GWAS-based genetic correlations, and confirmed observations of overlap among genome-wide significant SNPs, MR results and epidemiological associations.

This method is an advance for several reasons: it does not require individual genotypes, genome-wide significant SNPs or LD-pruning (which loses information if causal SNPs are in LD). Our method is not biased by sample overlap and is computationally fast. Furthermore, our approach does not require measuring multiple traits on the same individuals, so it scales easily to studies of thousands of pairs of traits. These advantages allow us to estimate genetic correlation for many more pairs of phenotypes than was possible with existing methods.

The challenges in interpreting genetic correlation are similar to the challenges in MR. We highlight two difficulties. First, genetic correlation is immune to environmental confounding, but is subject to genetic confounding, analogous to confounding by pleiotropy in MR. For example, the genetic correlation between HDL and CAD in Figure 2 could result from a causal effect  $HDL \rightarrow CAD$ , but could also be mediated by triglycerides (TG) [10,57], represented graphically [58]



as  $\text{HDL} \leftarrow \text{G} \rightarrow \text{TG} \rightarrow \text{CAD}$ , where  $\text{G}$  is the set of genetic variants with effects on both HDL and TG. Extending genetic correlation to multiple genetically correlated phenotypes is an important direction for future work. Second, although genetic correlation estimates are not biased by over-sampling of cases, they are affected by other forms of selection bias, such as misclassification [16].

We note several limitations of LD Score regression as an estimator of genetic correlation. First, LD Score regression requires larger sample sizes than methods that use individual genotypes in order to achieve acceptable standard error. Second, LD Score regression is not currently applicable to samples from recently-admixed populations. Third, methods built from polygenic models, such as LD Score regression and REML, are most effective when applied to traits with polygenic genetic architectures. For traits where significant SNPs account for a sizable proportion of heritability, analyzing only these SNPs can be more powerful. Developing methods that make optimal use of both large-effect SNPs and diffuse polygenic signal is a direction for future research.

Despite these limitations, we believe that the LD Score regression estimator of genetic correlation will be a useful addition to the epidemiological toolbox, since it allows for rapid screening for correlations among a diverse set of traits, without the need for measuring multiple traits on the same individuals or genome-wide significant SNPs.

## Methods

### Definition of Genetic Covariance and Correlation

All definitions refer to narrow-sense heritabilities and genetic covariances. Let  $S$  denote a set of  $M$  SNPs, let  $X$  denote a vector of additively (0-1-2) coded genotypes for the SNPs in  $S$ , and let  $y_1$  and  $y_2$  denote phenotypes. Define  $\beta := \operatorname{argmax}_{\alpha \in \mathbb{R}^M} \operatorname{Cor}[y_1, X\alpha]^2$ , where the maximization is performed in the population (*i.e.*, in the infinite data limit). Let  $\gamma$  denote the corresponding vector for  $y_2$ . This is a projection, so  $\beta$  is unique modulo SNPs in perfect LD. Define  $h_S^2$ , the heritability explained by SNPs in  $S$ , as  $h_S^2(y_1) := \sum_j \beta_j^2$  and  $\rho_S(y_1, y_2)$ , the genetic covariance among SNPs in  $S$ , as  $\rho_S(y_1, y_2) := \sum_{j \in S} \beta_j \gamma_j$ . The genetic correlation among SNPs in  $S$  is  $r_S(y_1, y_2) := \rho_S(y_1, y_2) / \sqrt{h_S^2(y_1) h_S^2(y_2)}$ , which lies in  $[-1, 1]$ . Following [13], we use subscript  $g$  (as in  $h_g^2, \rho_g, r_g$ ) when the set of SNPs is genotyped and imputed SNPs in GWAS.

SNP genetic correlation ( $r_g$ ) is different from family study genetic correlation. In a family study, the relationship matrix captures information about all genetic variation, not just common SNPs. As a result, family studies estimate the total genetic correlation ( $S$  equals all variants). Unlike the relationship between SNP-heritability [13] and total heritability, for which  $h_g^2 \leq h^2$ , no similar relationship holds between SNP genetic correlation and total genetic correlation. If  $\beta$  and  $\gamma$  are more strongly correlated among common variants than rare variants, then the total genetic correlation will be less than the SNP genetic correlation.

Genetic correlation is (asymptotically) proportional to Mendelian randomization estimates. If we use a genetic instrument  $g_i := \sum_{j \in S} X_{ij} \beta_j$  to estimate the effect  $b_{12}$  of  $y_1$  on  $y_2$ , the 2SLS estimate is  $\hat{b}_{2SLS} := g^\top y_2 / g^\top y_1$  [8]. The expectations of the numerator and denominator are  $\mathbb{E}[g^\top y_2] = \rho_S(y_1, y_2)$  and  $\mathbb{E}[g^\top y_1] = h_S^2(y_1)$ . Thus,  $\operatorname{plim}_{N \rightarrow \infty} \hat{b}_{2SLS} = r_S(y_2, y_1) \sqrt{h_S^2(y_1) / h_S^2(y_2)}$ . If we use the same set  $S$  of SNPs to estimate  $b_{12}$  and  $b_{21}$  (*e.g.*, if  $S$  is the set of all common SNPs, as in the genetic correlation analyses in this paper), then this procedure is symmetric in  $y_1$  and  $y_2$ .

Genetic correlation is different from pleiotropy. Two traits have a pleiotropic relationship if many variants affect both. Genetic correlation is a stronger condition than pleiotropy: to exhibit genetic correlation, the directions of effect must also be consistently aligned.

### Cross-Trait LD Score Regression

We estimate genetic covariance by regressing  $z_{1j} z_{2j}$  against  $\ell_j \sqrt{N_{1j} N_{2j}}$ , (where  $N_{ij}$  is the sample size for SNP  $j$  in study  $i$ ) then multiplying the resulting slope by  $M$ , the number of SNPs in the reference panel with MAF between 5% and 50% (technically, this is an estimate of  $\rho_{5-50\%}$ , see the Supplementary Note).

If we know the amount of sample overlap ahead of time, we can reduce the standard error by constraining the intercept with the `--constrain-intercept` flag in `ldsc`. This works even if there is nonzero sample overlap, in which case the intercept should be constrained to  $N_s \rho / \sqrt{N_1 N_2}$ .

### Regression Weights

For heritability estimation, we use the regression weights from [21]. If effect sizes for both phenotypes are drawn from a bivariate normal distribution, then the optimal regression weights for genetic

covariance estimation are

$$\text{Var}[z_{1j}z_{2j} | \ell_j] = \left( \frac{N_1 h_1^2 \ell_j}{M} + 1 \right) \left( \frac{N_2 h_2^2 \ell_j}{M} + 1 \right) + 2 \left( \frac{\sqrt{N_1 N_2} \rho_g}{M} \ell_j + \frac{\rho N_s}{\sqrt{N_1 N_2}} \right)^2 \quad (2)$$

(Supplementary Note). This quantity depends on several parameters ( $h_1^2, h_2^2, \rho_g, \rho, N_s$ ) which are not known a priori, so it is necessary to estimate them from the data. We compute the weights in two steps:

1. The first regression is weighted using heritabilities from the single-trait LD Score regressions,  $\rho N_s = 0$ , and  $\rho_g$  estimated as  $\hat{\rho}_g := (\bar{\ell} \sqrt{N_1 N_2})^{-1} \sum_j z_{1j} z_{2j}$ .
2. The second regression is weighted using the estimates of  $\rho N_s$  and  $\rho_g$  from step 1. The genetic covariance estimate that we report is the estimate from the second regression.

Linear regression with weights estimated from the data is called feasible generalized least squares (FGLS). FGLS has the same limiting distribution as WLS with optimal weights, so WLS  $p$ -values are valid for FGLS [8]. We multiply the heteroskedasticity weights by  $1/\ell_j$  (where  $\ell_j$  is LD Score with sum over regression SNPs) in order to downweight SNPs that are overcounted. This is a heuristic: the optimal approach is to rotate the data so that it is de-correlated, but this rotation matrix is difficult to compute.

## Assessment of Statistical Significance via Block Jackknife

Summary statistics for SNPs in LD are correlated, so the OLS standard error will be biased downwards. We estimate a heteroskedasticity-and-correlation-robust standard error with a block jackknife over blocks of adjacent SNPs. This is the same procedure used in [21], and gives accurate standard errors in simulations (Table 1). We obtain a standard error for the genetic correlation by using a ratio block jackknife over SNPs. The default setting in `ldsc` is 200 blocks per genome, which can be adjusted with the `--num-blocks` flag.

## Computational Complexity

Let  $N$  denote sample size and  $M$  the number of SNPs. The computational complexity of the steps involved in LD Score regression are as follows:

1. Computing summary statistics takes  $\mathcal{O}(MN)$  time.
2. Computing LD Scores takes  $\mathcal{O}(MN)$  time, though the  $N$  for computing LD Scores need not be large. We use the  $N = 378$  Europeans from 1000 Genomes.
3. LD Score regression takes  $\mathcal{O}(M)$  time and space.

For a user who has already computed summary statistics and downloads LD Scores from our website (URLs), the computational cost of LD Score regression is  $\mathcal{O}(M)$  time and space. For comparison, REML takes time  $\mathcal{O}(MN^2)$  for computing the GRM and  $\mathcal{O}(N^3)$  time for maximizing the likelihood.

Practically, estimating LD Scores takes roughly an hour parallelized over chromosomes, and LD Score regression takes about 15 seconds per pair of phenotypes on a 2014 MacBook Air with 1.7 GhZ Intel Core i7 processor.

## Simulations

We simulated quantitative traits under an infinitesimal model in 2062 controls from a Swedish study. To simulate the standard scenario where many causal SNPs are not genotyped, we simulated phenotypes by drawing casual SNPs from 622,146 best-guess imputed 1000 Genomes SNPs on chromosome 2, then retained only the 90,980 HM3 SNPs with MAF above 5% for LD Score regression. We used in-sample LD Scores for LD Score regression.

## Summary Statistic Datasets

We selected traits for inclusion in the main text via the following procedure:

1. Begin with all publicly available non-sex-stratified European-only summary statistics.
2. Remove studies that do not provide signed summary statistics.
3. Remove studies not imputed to at least HapMap 2.
4. Remove studies that include heritable covariates.
5. Remove all traits with heritability  $z$ -score below 4. Genetic correlation estimates for traits with heritability  $z$ -score below 4 are generally too noisy to interpret.
6. Prune clusters of correlated phenotypes (*e.g.*, obesity classes 1-3) by picking the trait from each cluster with the highest heritability  $z$ -score.

We then applied the following filters (implemented in the script `sumstats_to_chisq.py` included with `ldsc`):

1. For studies that provide a measure of imputation quality, filter to INFO above 0.9.
2. For studies that provide sample MAF, filter to sample MAF above 1%.
3. In order to restrict to well-imputed SNPs in studies that do not provide a measure of imputation quality, filter to HapMap3 [60] SNPs with 1000 Genomes EUR MAF above 5%, which tend to be well-imputed in most studies. This step should be skipped if INFO scores are available for all studies.
4. If sample size varies from SNP to SNP, remove SNPs with effective sample size less than 0.67 times the 90th percentile of sample size.
5. Remove indels and structural variants.
6. Remove strand-ambiguous SNPs.
7. Remove SNPs whose alleles do not match the alleles in 1000 Genomes.
8. Because the presence of outliers can increase the regression standard error, we also removed SNPs with extremely large effect sizes ( $\chi^2 > 80$ , as in [21]).

Genomic control (GC) correction at any stage biases the heritability and genetic covariance estimates downwards (see the Supplementary Note of [21]). The biases in the numerator and denominator of genetic correlation cancel exactly, so genetic correlation is not biased by GC correction. A majority of the studies analyzed in this paper used GC correction, so we do not report genetic covariance and heritability.

Data on Alzheimer’s disease were obtained from the following source:

*International Genomics of Alzheimer's Project (IGAP) is a large two-stage study based upon genome-wide association studies (GWAS) on individuals of European ancestry. In stage 1, IGAP used genotyped and imputed data on 7,055,881 single nucleotide polymorphisms (SNPs) to meta-analyze four previously-published GWAS datasets consisting of 17,008 Alzheimer's disease cases and 37,154 controls (The European Alzheimer's Disease Initiative, EADI; the Alzheimer Disease Genetics Consortium, ADGC; The Cohorts for Heart and Aging Research in Genomic Epidemiology consortium, CHARGE; The Genetic and Environmental Risk in AD consortium, GERAD). In stage 2, 11,632 SNPs were genotyped and tested for association in an independent set of 8,572 Alzheimer's disease cases and 11,312 controls. Finally, a meta-analysis was performed combining results from stages 1 and 2.*

We only used stage 1 data for LD Score regression.

## URLs

1. ldsc software:  
[github.com/bulik/ldsc](https://github.com/bulik/ldsc)
2. This paper:  
[github.com/bulik/gencor\\_tex](https://github.com/bulik/gencor_tex)
3. PGC (psychiatric) summary statistics:  
[www.med.unc.edu/pgc/downloads](http://www.med.unc.edu/pgc/downloads)
4. GIANT (anthropometric) summary statistics:  
[www.broadinstitute.org/collaboration/giant/index.php/GIANT\\_consortium\\_data\\_files](http://www.broadinstitute.org/collaboration/giant/index.php/GIANT_consortium_data_files)
5. EGG (Early Growth Genetics) summary statistics:  
[www.egg-consortium.org/](http://www.egg-consortium.org/)
6. MAGIC (insulin, glucose) summary statistics:  
[www.magicinvestigators.org/downloads/](http://www.magicinvestigators.org/downloads/)
7. CARDIoGRAM (coronary artery disease) summary statistics:  
[www.cardiogramplusc4d.org](http://www.cardiogramplusc4d.org)
8. DIAGRAM (T2D) summary statistics:  
[www.diagram-consortium.org](http://www.diagram-consortium.org)
9. Rheumatoid arthritis summary statistics:  
[www.broadinstitute.org/ftp/pub/rheumatoid\\_arthritis/Stahl.etal\\_2010NG/](http://www.broadinstitute.org/ftp/pub/rheumatoid_arthritis/Stahl.etal_2010NG/)
10. IGAP (Alzheimers) summary statistics:  
[www.pasteur-lille.fr/en/recherche/u744/igap/igap\\_download.php](http://www.pasteur-lille.fr/en/recherche/u744/igap/igap_download.php)
11. IIBDGC (inflammatory bowel disease) summary statistics:  
[www.ibdgenetics.org/downloads.html](http://www.ibdgenetics.org/downloads.html)  
We used a newer version of these data with 1000 Genomes imputation.
12. Plasma lipid summary statistics:  
[www.broadinstitute.org/mpg/pubs/lipids2010/](http://www.broadinstitute.org/mpg/pubs/lipids2010/)
13. SSGAC (educational attainment) summary statistics:  
[www.ssgac.org/](http://www.ssgac.org/)
14. Beans:  
[www.barismo.com](http://www.barismo.com)  
[www.bluebottlecoffee.com](http://www.bluebottlecoffee.com)

## Acknowledgements

We would like to thank P. Sullivan, C. Bulik and S. Caldwell for helpful discussion. This work was supported by NIH grants R01 MH101244 (ALP), R03 CA173785 (HKF) and by the Fannie and John Hertz Foundation (HKF). The coffee that Brendan drank while writing this paper was roasted by Barismo in Arlington, MA and Blue Bottle Coffee in Oakland, CA.

Data on anorexia nervosa were obtained by funding from the WTCCC3 WT088827/Z/09 titled “A genome-wide association study of anorexia nervosa”.

Data on glycaemic traits have been contributed by MAGIC investigators and have been downloaded from [www.magicinvestigators.org](http://www.magicinvestigators.org).

Data on coronary artery disease / myocardial infarction have been contributed by CARDIOGRAMplusC4D investigators and have been downloaded from [www.CARDIOGRAMPLUSC4D.ORG](http://www.CARDIOGRAMPLUSC4D.ORG)

We thank the International Genomics of Alzheimer’s Project (IGAP) for providing summary results data for these analyses. The investigators within IGAP contributed to the design and implementation of IGAP and/or provided data but did not participate in analysis or writing of this report. IGAP was made possible by the generous participation of the control subjects, the patients, and their families. The i-Select chips was funded by the French National Foundation on Alzheimer’s disease and related disorders. EADI was supported by the LABEX (laboratory of excellence program investment for the future) DISTALZ grant, Inserm, Institut Pasteur de Lille, Universit de Lille 2 and the Lille University Hospital. GERAD was supported by the Medical Research Council (Grant 503480), Alzheimer’s Research UK (Grant 503176), the Wellcome Trust (Grant 082604/2/07/Z) and German Federal Ministry of Education and Research (BMBF): Competence Network Dementia (CND) grant 01GI0102, 01GI0711, 01GI0420. CHARGE was partly supported by the NIH/NIA grant R01 AG033193 and the NIA AG081220 and AGES contract N01-AG-12100, the NHLBI grant R01 HL105756, the Icelandic Heart Association, and the Erasmus Medical Center and Erasmus University. ADGC was supported by the NIH/NIA grants: U01 AG032984, U24 AG021886, U01 AG016976, and the Alzheimer’s Association grant ADGC-10-196728.

## Author Contributions

MJD provided reagents. BMN and ALP provided reagents. CL, ER, VA, JP and FD aided in the interpretation of results. JP and FD provided data on age at onset of menarche. The caffeine molecule is responsible for all that is good about this manuscript. BBS and HKF are responsible for the rest. All authors revised and approved the final manuscript.

## Competing Financial Interests

Unfortunately, we have no financial conflicts of interest to declare.

# 1 References

- [1] George Davey Smith and Shah Ebrahim. Mendelian randomization: can genetic epidemiology contribute to understanding environmental determinants of disease? *International journal of epidemiology*, 32(1):1–22, 2003.
- [2] George Davey Smith and Gibran Hemani. Mendelian randomization: genetic anchors for causal inference in epidemiological studies. *Human molecular genetics*, 23(R1):R89–R98, 2014.
- [3] SG Vandenberg. Multivariate analysis of twin differences. *Methods and goals in human behavior genetics*, pages 29–43, 1965.
- [4] Oscar Kempthorne and Richard H Osborne. The interpretation of twin data. *American journal of human genetics*, 13(3):320, 1961.
- [5] John C Loehlin and Steven Gerritjan Vandenberg. *Genetic and environmental components in the covariation of cognitive abilities: An additive model*. Louisville Twin Study, University of Louisville, 1966.
- [6] Michael Neale and Lon Cardon. *Methodology for genetic studies of twins and families*. Number 67. Springer, 1992.
- [7] Paul Lichtenstein, Benjamin H Yip, Camilla Björk, Yudi Pawitan, Tyrone D Cannon, Patrick F Sullivan, and Christina M Hultman. Common genetic determinants of schizophrenia and bipolar disorder in swedish families: a population-based study. *The Lancet*, 373(9659):234–239, 2009.
- [8] Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics: An empiricist’s companion*. Princeton university press, 2008.
- [9] Benjamin F Voight, Gina M Peloso, Marju Orho-Melander, Ruth Frikke-Schmidt, Maja Barbalic, Majken K Jensen, George Hindy, Hilma Hólm, Eric L Ding, Toby Johnson, et al. Plasma hdl cholesterol and risk of myocardial infarction: a mendelian randomisation study. *The Lancet*, 380(9841):572–580, 2012.
- [10] Ron Do, Cristen J Willer, Ellen M Schmidt, Sebanti Sengupta, Chi Gao, Gina M Peloso, Stefan Gustafsson, Stavroula Kanoni, Andrea Ganna, Jin Chen, et al. Common variants associated with plasma triglycerides and risk for coronary artery disease. *Nature genetics*, 45(11):1345–1352, 2013.
- [11] Peter M Visscher, Matthew A Brown, Mark I McCarthy, and Jian Yang. Five years of gwas discovery. *The American Journal of Human Genetics*, 90(1):7–24, 2012.
- [12] Stephen Burgess, Simon G Thompson, et al. Avoiding bias from weak instruments in mendelian randomization studies. *International journal of epidemiology*, 40(3):755–764, 2011.
- [13] Jian Yang, Beben Benyamin, Brian P McEvoy, Scott Gordon, Anjali K Henders, Dale R Nyholt, Pamela A Madden, Andrew C Heath, Nicholas G Martin, Grant W Montgomery, et al. Common snps explain a large proportion of the heritability for human height. *Nature Genetics*, 42(7):565–569, 2010.



- [14] Jian Yang, S Hong Lee, Michael E Goddard, and Peter M Visscher. Gcta: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics*, 88(1):76–82, 2011.
- [15] Sang Hong Lee, Jian Yang, Michael E Goddard, Peter M Visscher, and Naomi R Wray. Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics*, 28(19):2540–2542, 2012.
- [16] Cross-Disorder Group of the Psychiatric Genomics Consortium et al. Genetic relationship between five psychiatric disorders estimated from genome-wide snps. *Nature Genetics*, 2013.
- [17] Shashaank Vattikuti, Juen Guo, and Carson C Chow. Heritability and genetic correlations explained by common snps for metabolic syndrome traits. *PLoS genetics*, 8(3):e1002637, 2012.
- [18] Guo-Bo Chen, Sang Hong Lee, Marie-Jo A Brion, Grant W Montgomery, Naomi R Wray, Graham L Radford-Smith, Peter M Visscher, et al. Estimation and partitioning of (co) heritability of inflammatory bowel disease from gwas and immunochip data. *Human molecular genetics*, page ddu174, 2014.
- [19] Shaun M Purcell, Naomi R Wray, Jennifer L Stone, Peter M Visscher, Michael C O’Donovan, Patrick F Sullivan, Pamela Sklar, Shaun M Purcell, Jennifer L Stone, Patrick F Sullivan, et al. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*, 460(7256):748–752, 2009.
- [20] Frank Dudbridge. Power and predictive accuracy of polygenic risk scores. *PLoS genetics*, 9(3):e1003348, 2013.
- [21] Brendan Bulik-Sullivan, Po-Ru Loh, Hilary Finucane, Stephan Ripke, Jian Yang, Nick Patterson, Mark J Daly, Alkes L Price, and Benjamin M Neale. Ld score regression distinguishes confounding from polygenicity in genome-wide association studies. *bioRxiv*, 2014.
- [22] Jian Yang, Michael N Weedon, Shaun Purcell, Guillaume Lettre, Karol Estrada, Cristen J Willer, Albert V Smith, Erik Ingelsson, Jeffrey R O’Connell, Massimo Mangino, et al. Genomic inflation factors under polygenic inheritance. *European Journal of Human Genetics*, 19(7):807–812, 2011.
- [23] Doug Speed, Gibran Hemani, Michael R Johnson, and David J Balding. Improved heritability estimation from genome-wide snps. *The American Journal of Human Genetics*, 91(6):1011–1021, 2012.
- [24] Cross-Disorder Group of the Psychiatric Genomics Consortium et al. Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet*, 381(9875):1371, 2013.
- [25] John RB Perry, Felix Day, Cathy E Elks, Patrick Sulem, Deborah J Thompson, Teresa Ferreira, Chunyan He, Daniel I Chasman, Tõnu Esko, Gudmar Thorleifsson, et al. Parent-of-origin-specific allelic associations among 106 genomic loci for age at menarche. *Nature*, 514(7520):92–97, 2014.

- [26] Andrew P Morris, Benjamin F Voight, Tanya M Teslovich, Teresa Ferreira, Ayellet V Segre, Valgerdur Steinthorsdottir, Rona J Strawbridge, Hassan Khan, Harald Grallert, Anubha Mahajan, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nature genetics*, 44(9):981, 2012.
- [27] Momoko Horikoshi, Hanieh Yaghootkar, Dennis O Mook-Kanamori, Ulla Sovio, H Rob Taal, Branwen J Hennig, Jonathan P Bradfield, Beate St Pourcain, David M Evans, Pimphen Charoen, et al. New loci associated with birth weight identify genetic links between intrauterine growth and adult height and metabolism. *Nature genetics*, 45(1):76–82, 2013.
- [28] Rachel M Freathy, Amanda J Bennett, Susan M Ring, Beverley Shields, Christopher J Groves, Nicholas J Timpson, Michael N Weedon, Eleftheria Zeggini, Cecilia M Lindgren, Hana Lango, et al. Type 2 diabetes risk alleles are associated with reduced size at birth. *Diabetes*, 58(6):1428–1433, 2009.
- [29] Early Growth Genetics (EGG) Consortium et al. A genome-wide association meta-analysis identifies new childhood obesity loci. *Nature genetics*, 44(5):526–531, 2012.
- [30] H Rob Taal, Beate St Pourcain, Elisabeth Thiering, Shikta Das, Dennis O Mook-Kanamori, Nicole M Warrington, Marika Kaakinen, Eskil Kreiner-Møller, Jonathan P Bradfield, Rachel M Freathy, et al. Common variants at 12q15 and 12q24 are associated with infant head circumference. *Nature genetics*, 44(5):532–538, 2012.
- [31] NC Onland-Moret, PHM Peeters, CH Van Gils, F Clavel-Chapelon, T Key, A Tjønneland, A Trichopoulou, R Kaaks, Jonas Manjer, S Panico, et al. Age at menarche in relation to adult height the epic study. *American journal of epidemiology*, 162(7):623–632, 2005.
- [32] Felix Day et al. Puberty timing associated with diabetes, cardiovascular disease and also diverse health outcomes in men and women: the uk biobank study. *Submitted*, 2014.
- [33] Cathy E Elks, Ken K Ong, Robert A Scott, Yvonne T van der Schouw, Judith S Brand, Petra A Wark, Pilar Amiano, Beverley Balkau, Aurelio Barricarte, Heiner Boeing, et al. Age at menarche and type 2 diabetes risk the epic-interact study. *Diabetes care*, 36(11):3526–3534, 2013.
- [34] Hilary K. Finucane, Brendan Bulik-Sullivan, Alexander Gusev, Gosia Trynka, Yakir Reshef, Po-Ru Loh, Verner Anttila, Han Xu, Chongzhi Zang, Kyle Farh, Stephan Ripke, Felix R. Day, The ReproGen Consortium, Schizophrenia Working Group of the Psychiatric Genomics Consortium, Shaun Purcell, Eli Stahl, Sara Lindstrom, John R. B. Perry, Yukinori Okada, Brad Bernstein, Soumya Raychaudhuri, Mark Daly, Nick Patterson, Benjamin M. Neale, and Alkes L. Price. Polygenic effects of cell-type-specific functional elements in 17 traits and 1.3 million phenotyped samples. *In preparation*, 2014.
- [35] Na Wang, Xianglan Zhang, Yong-Bing Xiang, Gong Yang, Hong-Lan Li, Jing Gao, Hui Cai, Yu-Tang Gao, Wei Zheng, and Xiao-Ou Shu. Associations of adult height and its components with mortality: a report from cohort studies of 135 000 chinese women and men. *International journal of epidemiology*, 40(6):1715–1726, 2011.

- [36] Patricia R Hebert, Janet W Rich-Edwards, JE Manson, Paul M Ridker, Nancy R Cook, Gerald T O'Connor, Julie E Buring, and Charles H Hennekens. Height and incidence of cardiovascular disease in male physicians. *Circulation*, 88(4):1437–1443, 1993.
- [37] Janet W Rich-Edwards, JoAnn E Manson, Meir J Stampfer, Graham A Colditz, Walter C Willett, Bernard Rosner, Frank E Speizer, and Charles H Hennekens. Height and the risk of cardiovascular disease in women. *American journal of epidemiology*, 142(9):909–917, 1995.
- [38] Cornelius A Rietveld, Sarah E Medland, Jaime Derringer, Jian Yang, Tõnu Esko, Nicolas W Martin, Harm-Jan Westra, Konstantin Shakhbazov, Abdel Abdellaoui, Arpana Agrawal, et al. Gwas of 126,559 individuals identifies genetic variants associated with educational attainment. *Science*, 340(6139):1467–1471, 2013.
- [39] Deborah E Barnes and Kristine Yaffe. The projected effect of risk factor reduction on alzheimer’s disease prevalence. *The Lancet Neurology*, 10(9):819–828, 2011.
- [40] Sam Norton, Fiona E Matthews, Deborah E Barnes, Kristine Yaffe, and Carol Brayne. Potential for primary prevention of alzheimer’s disease: an analysis of population-based data. *The Lancet Neurology*, 13(8):788–794, 2014.
- [41] James H MacCabe, Mats P Lambe, Sven Cnattingius, Pak C Sham, Anthony S David, Abraham Reichenberg, Robin M Murray, and Christina M Hultman. Excellent school performance at age 16 and risk of adult bipolar disorder: national cohort study. *The British Journal of Psychiatry*, 196(2):109–115, 2010.
- [42] Jari Tiihonen, Jari Haukka, Markus Henriksson, Mary Cannon, Tuula Kieseppä, Ilmo Laaksonen, Juhani Sinivuori, and Jouko Lönnqvist. Premorbid intellectual functioning in bipolar disorder and schizophrenia: results from a cohort study of male conscripts. *American Journal of Psychiatry*, 162(10):1904–1910, 2005.
- [43] John P Pierce, Michael C Fiore, Thomas E Novotny, Evridiki J Hatziandreu, and Ronald M Davis. Trends in cigarette smoking in the united states: educational differences are increasing. *Jama*, 261(1):56–60, 1989.
- [44] Ruth H Striegel-Moore, Vicki Garvin, Faith-Anne Dohm, and Robert A Rosenheck. Psychiatric comorbidity of eating disorders in men: a national study of hospitalized veterans. *International Journal of Eating Disorders*, 25(4):399–404, 1999.
- [45] Barton J Blinder, Edward J Cumella, and Visant A Sanathara. Psychiatric comorbidities of female inpatients with eating disorders. *Psychosomatic Medicine*, 68(3):454–462, 2006.
- [46] Ian J Deary, Steve Strand, Pauline Smith, and Cres Fernandes. Intelligence and educational achievement. *Intelligence*, 35(1):13–21, 2007.
- [47] Catherine M Calvin, Cres Fernandes, Pauline Smith, Peter M Visscher, and Ian J Deary. Sex, intelligence and educational achievement in a national cohort of over 175,000 11-year-old schoolchildren in england. *Intelligence*, 38(4):424–432, 2010.

- [48] Maureen S Durkin, Matthew J Maenner, F John Meaney, Susan E Levy, Carolyn DiGuseppi, Joyce S Nicholas, Russell S Kirby, Jennifer A Pinto-Martin, and Laura A Schieve. Socioeconomic inequality in the prevalence of autism spectrum disorder: evidence from a us cross-sectional study. *PLoS One*, 5(7):e11551, 2010.
- [49] Elise B Robinson, Kaitlin E Samocha, Jack A Kosmicki, Lauren McGrath, Benjamin M Neale, Roy H Perlis, and Mark J Daly. Autism spectrum disorder severity reflects the average contribution of de novo and familial influences. *Proceedings of the National Academy of Sciences*, 111(42):15161–15165, 2014.
- [50] Kaitlin E Samocha, Elise B Robinson, Stephan J Sanders, Christine Stevens, Aniko Sabo, Lauren M McGrath, Jack A Kosmicki, Karola Rehnström, Swapan Mallick, Andrew Kirby, et al. A framework for the interpretation of de novo mutation in human disease. *Nature genetics*, 46(9):944–950, 2014.
- [51] Alan J Silman and Jacqueline E Pearson. Epidemiology and genetics of rheumatoid arthritis. *Arthritis Res*, 4(Suppl 3):S265–S272, 2002.
- [52] Jose de Leon and Francisco J Diaz. A meta-analysis of worldwide studies demonstrates an association between schizophrenia and tobacco smoking behaviors. *Schizophrenia research*, 76(2):135–157, 2005.
- [53] Ole A Andreassen, Srdjan Djurovic, Wesley K Thompson, Andrew J Schork, Kenneth S Kendler, Michael C O’Donovan, Dan Rujescu, Thomas Werge, Martijn van de Bunt, Andrew P Morris, et al. Improved detection of common variants associated with schizophrenia by leveraging pleiotropy with cardiovascular-disease risk factors. *The American Journal of Human Genetics*, 92(2):197–209, 2013.
- [54] Chris Cotsapas, Benjamin F Voight, Elizabeth Rossin, Kasper Lage, Benjamin M Neale, Chris Wallace, Gonçalo R Abecasis, Jeffrey C Barrett, Timothy Behrens, Judy Cho, et al. Pervasive sharing of genetic effects in autoimmune disease. *PLoS genetics*, 7(8):e1002254, 2011.
- [55] Kyle Kai-How Farh, Alexander Marson, Jiang Zhu, Markus Klei, William J Housley, Samantha Beik, Noam Shores, Holly Whitton, Russell JH Ryan, Alexander A Shishkin, et al. Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature*, 2014.
- [56] Peter Wurtz et al. Metabolic signatures of adiposity in young adults: Mendelian randomization analysis and effects of weight change. *PLoS Medicine*, 2014.
- [57] Stephen Burgess, Daniel F Freitag, Hassan Khan, Donal N Gorman, and Simon G Thompson. Using multivariable mendelian randomization to disentangle the causal effects of lipid fractions. *PloS one*, 9(10):e108891, 2014.
- [58] Sander Greenland, Judea Pearl, and James M Robins. Causal diagrams for epidemiologic research. *Epidemiology*, pages 37–48, 1999.
- [59] Christopher C Chang, Carson C Chow, Laurent CAM Telier, Shashaank Vattikuti, Shaun M Purcell, and James J Lee. Second-generation plink: rising to the challenge of larger and richer datasets. *arXiv preprint arXiv:1410.4803*, 2014.

- [60] International HapMap 3 Consortium et al. Integrating common and rare genetic variation in diverse human populations. *Nature*, 467(7311):52–58, 2010.
- [61] Karl Pearson and Alice Lee. On the inheritance of characters not capable of exact quantitative measurement. *Philosophical Transactions of the Royal Society of London, A (195)*, pages 79–150, 1901.
- [62] Schizophrenia Working Group of the Psychiatric Genomics Consortium et al. Biological insights from 108 schizophrenia-associated genetic loci. *Nature*, 511(7510):421–427, 2014.
- [63] Pamela Sklar, Stephan Ripke, Laura J Scott, Ole A Andreassen, Sven Cichon, Nick Craddock, Howard J Edenberg, John I Nurnberger, Marcella Rietschel, Douglas Blackwood, et al. Large-scale genome-wide association analysis of bipolar disorder identifies a new susceptibility locus near *od4*. *Nature genetics*, 43(10):977, 2011.
- [64] Stephan Ripke, Naomi R Wray, Cathryn M Lewis, Steven P Hamilton, Myrna M Weissman, Gerome Breen, Enda M Byrne, Douglas HR Blackwood, Dorret I Boomsma, Sven Cichon, et al. A mega-analysis of genome-wide association studies for major depressive disorder. *Molecular psychiatry*, 18(4):497–511, 2012.
- [65] Vesna Boraska, Christopher S Franklin, James AB Floyd, Laura M Thornton, Laura M Huckins, Lorraine Southam, N William Rayner, Ioanna Tachmazidou, Kelly L Klump, Janet Treasure, et al. A genome-wide association study of anorexia nervosa. *Molecular psychiatry*, 2014.
- [66] Tobacco, Genetics Consortium, et al. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nature genetics*, 42(5):441–447, 2010.
- [67] Jean-Charles Lambert, Carla A Ibrahim-Verbaas, Denise Harold, Adam C Naj, Rebecca Sims, Céline Bellenguez, Gyungah Jun, Anita L DeStefano, Joshua C Bis, Gary W Beecham, et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for alzheimer’s disease. *Nature genetics*, 2013.
- [68] Hana Lango Allen, Karol Estrada, Guillaume Lettre, Sonja I Berndt, Michael N Weedon, Fernando Rivadeneira, Cristen J Willer, Anne U Jackson, Sailaja Vedantam, Soumya Raychaudhuri, et al. Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature*, 467(7317):832–838, 2010.
- [69] Sonja I Berndt, Stefan Gustafsson, Reedik Mägi, Andrea Ganna, Eleanor Wheeler, Mary F Feitosa, Anne E Justice, Keri L Monda, Damien C Croteau-Chonka, Felix R Day, et al. Genome-wide meta-analysis identifies 11 new loci for anthropometric traits and provides insights into genetic architecture. *Nature genetics*, 45(5):501–512, 2013.
- [70] Heribert Schunkert, Inke R König, Sekar Kathiresan, Muredach P Reilly, Themistocles L Assimes, Hilma Holm, Michael Preuss, Alexandre FR Stewart, Maja Barbalic, Christian Gieger, et al. Large-scale association analysis identifies 13 new susceptibility loci for coronary artery disease. *Nature genetics*, 43(4):333–338, 2011.
- [71] Tanya M Teslovich, Kiran Musunuru, Albert V Smith, Andrew C Edmondson, Ioannis M Stylianou, Masahiro Koseki, James P Pirruccello, Samuli Ripatti, Daniel I Chasman, Cristen J

- Willer, et al. Biological, clinical and population relevance of 95 loci for blood lipids. *Nature*, 466(7307):707–713, 2010.
- [72] Alisa K Manning, Marie-France Hivert, Robert A Scott, Jonna L Grimsby, Nabila Bouatia-Naji, Han Chen, Denis Rybin, Ching-Ti Liu, Lawrence F Bielak, Inga Prokopenko, et al. A genome-wide approach accounting for body mass index identifies genetic variants influencing fasting glyceimic traits and insulin resistance. *Nature genetics*, 44(6):659–669, 2012.
  - [73] Ralf JP van der Valk, Eskil Kreiner-Møller, Marjolein N Kooijman, Mònica Guxens, Evangelia Stergiakouli, Annika Sääf, Jonathan P Bradfield, Frank Geller, M Geoffrey Hayes, Diana L Cousminer, et al. A novel common variant in *dcst2* is associated with length in early life and height in adulthood. *Human molecular genetics*, page ddu510, 2014.
  - [74] Luke Jostins, Stephan Ripke, Rinse K Weersma, Richard H Duerr, Dermot P McGovern, Ken Y Hui, James C Lee, L Philip Schumm, Yashoda Sharma, Carl A Anderson, et al. Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature*, 491(7422):119–124, 2012.
  - [75] Eli A Stahl, Soumya Raychaudhuri, Elaine F Remmers, Gang Xie, Stephen Eyre, Brian P Thomson, Yonghong Li, Fina AS Kurreeman, Alexandra Zhernakova, Anne Hinks, et al. Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. *Nature genetics*, 42(6):508–514, 2010.
  - [76] Ole A Andreassen, Wesley K Thompson, Andrew J Schork, Stephan Ripke, Morten Mattingsdal, John R Kelsoe, Kenneth S Kendler, Michael C O’Donovan, Dan Rujescu, Thomas Werge, et al. Improved detection of common variants associated with schizophrenia and bipolar disorder using pleiotropy-informed conditional false discovery rate. *PLoS Genetics*, 9(4):e1003455, 2013.

# Supplementary Note

## Quantitative Traits

Suppose we sample two cohorts for two phenotypes,  $y_1$  and  $y_2$ , with sample sizes  $N_1$  and  $N_2$ . We model phenotypes as  $y_1 = Y\beta + \delta$ , and  $y_2 = Z\gamma + \epsilon$ , where  $Y$  and  $Z$  are matrices of genotypes normalized to mean zero and variance one<sup>1</sup>, with dimensions  $N_1 \times M$  and  $N_2 \times M$ , respectively;  $\beta$  and  $\gamma$  are vectors of per-normalized genotype effect sizes, and  $\delta$  and  $\epsilon$  are vectors of residuals, representing environmental effects and non-additive genetic effects. In this model,  $Y$  and  $Z$  are unobserved matrices of all SNPs, including SNPs that are not genotyped.

We treat all of  $Y, Z, \beta, \gamma, \delta$  and  $\epsilon$  as random. Suppose that  $(\beta, \gamma)$  has mean zero and covariance matrix<sup>2</sup>

$$\text{Var}[(\beta, \gamma)] = \frac{1}{M} \begin{pmatrix} h_1^2 I & \rho_g I \\ \rho_g I & h_2^2 I \end{pmatrix},$$

and  $(\delta, \epsilon)$  has mean zero and covariance matrix

$$\text{Var}[(\delta, \epsilon)] = \begin{pmatrix} (1 - h_1^2) I & \rho_e I \\ \rho_e I & (1 - h_2^2) I \end{pmatrix}.$$

Let  $\rho := \rho_g + \rho_e$ . Genotypes are *i.i.d.* draw from a distribution with covariance matrix  $R$  (*i.e.*,  $r$  is an LD matrix with  $r_{jk} = \mathbb{E}[Y_{ij}Y_{ik}]$ ). There are  $N_s$  individuals included in both studies.

**Lemma 1.** *Under this model, the expected genetic covariance between phenotypes is  $\rho_g$ , justifying our use of the notation  $\rho_g$ .*

*Proof.* Let  $X$  denote an  $1 \times M$  vector of normalized, centered genotypes for an arbitrary individual. Under the model, the additive genetic component of  $y_1$  for this individual is  $\sum_j X_j \beta_j$ , and the additive genetic component of  $y_2$  for this individual is  $\sum_j X_j \gamma_j$ . Thus, the genetic covariance between  $y_1$  and  $y_2$  is

$$\begin{aligned} \text{Cov} \left[ \sum_j X_j \beta_j, \sum_j X_j \gamma_j \right] &= \mathbb{E} \left[ \left( \sum_j X_j \beta_j \right) \left( \sum_j X_j \gamma_j \right) \right] \\ &= \sum_j \sum_k \mathbb{E}[X_j X_k \beta_j \gamma_k] \\ &= \sum_j \mathbb{E}[X_j^2 \beta_j \gamma_j] \\ &= \sum_j \mathbb{E}[X_j^2] \mathbb{E}[\beta_j \gamma_j] \\ &= \rho_g. \end{aligned}$$

□

<sup>1</sup>We ignore the distinction between normalizing and centering in the population and in the sample, since this introduces only  $\mathcal{O}(1/N)$  error.

<sup>2</sup>The assumption that all  $\beta$  is drawn with equal variance for all SNPs hides an implicit assumption that rare SNPs have larger per-allele effect sizes than common SNPs. As discussed in the simulations section of the main text and in our earlier work [21], LD Score regression is robust to moderate violations of this assumption, though it may break down in extreme cases, *e.g.*, if all causal variants are rare. In situations where a different model for  $\text{Var}[\beta]$  is more appropriate, all proofs in this note go through with LD Score replaced by weighted LD Scores,  $\ell_j = \sum_k \text{Var}[\beta_j] r_{jk}^2$ .

We compute linear regression  $z$ -scores  $z_{1j} := Y_j^\top y_1 / \sqrt{N_1}$  and  $z_{2j} := Y_j^\top y_2 / \sqrt{N_2}$  for genotyped SNPs  $j$ .

**Proposition 1.** *Let  $j$  denote a genotyped SNP. Under the model described above,*

$$\mathbb{E}[z_{1j}z_{2j}] = \frac{\sqrt{N_1N_2}\rho_g}{M}\ell_j + \frac{N_s\rho}{\sqrt{N_1N_2}}. \quad (3)$$

*Proof.* By the law of total expectation,

$$\mathbb{E}[z_{1j}z_{2j}] = \mathbb{E}[\mathbb{E}[z_{1j}z_{2j} \mid Y, Z]] \quad (4)$$

First we compute the inner expectation from equation 4, with  $Z$  and  $Y$  fixed.

$$\begin{aligned} \mathbb{E}[z_{1j}z_{2j} \mid Y, Z] &= \frac{1}{\sqrt{N_1N_2}} \mathbb{E}[Y_j^\top y_1 y_2^\top Z_j] \\ &= \frac{1}{\sqrt{N_1N_2}} Y_j^\top \mathbb{E}[(Y\beta + \delta)(Z\gamma + \epsilon)^\top] Z_j \\ &= \frac{1}{\sqrt{N_1N_2}} Y_j^\top \left( Y \mathbb{E}[\beta^\top \gamma] Z + \mathbb{E}[\delta^\top Z \gamma] + \mathbb{E}[\beta^\top Y^\top \epsilon] + \mathbb{E}[\delta^\top \epsilon] \right) Z_j \\ &= \frac{1}{\sqrt{N_1N_2}} Y_j^\top \left( Y \mathbb{E}[\beta^\top \gamma] Z + \mathbb{E}[\delta^\top \epsilon] \right) Z_j \\ &= \frac{1}{\sqrt{N_1N_2}} \left( \frac{\rho_g}{M} Y_j^\top Y Z_j^\top Z + \rho_e Y_j^\top Z_j \right). \end{aligned} \quad (5)$$

Next, we remove the conditioning on  $Y$  and  $Z$ .

$$\frac{1}{\sqrt{N_1N_2}} \mathbb{E}[Y_j^\top Z_j] = \frac{N_s}{\sqrt{N_1N_2}}, \quad (6)$$

and

$$\frac{1}{\sqrt{N_1N_2}} \mathbb{E}[Y_j^\top Y Z_j^\top Z] = \ell_j + \frac{MN_s}{\sqrt{N_1N_2}}. \quad (7)$$

Substituting equations 6 and 7 into equation 5,

$$\begin{aligned} \mathbb{E}[z_{1j}z_{2j}] &= \frac{\sqrt{N_1N_2}\rho_g}{M}\ell_j + \frac{N_s(\rho_g + \rho_e)}{\sqrt{N_1N_2}} \\ &= \frac{\sqrt{N_1N_2}\rho_g}{M}\ell_j + \frac{N_s\rho}{\sqrt{N_1N_2}}. \end{aligned} \quad (8)$$

If study 1 and study 2 are the same study, then  $N_1 = N_2 = N_s$ ,  $\rho_g = h_g^2$  and  $\rho = 1$ , so equation 8 reduces to the LD Score regression equation for a single trait from [21].  $\square$

## Regression Weights

We can improve the efficiency of LD Score regression weighting by the reciprocal of the conditional variance function (CVF),  $\text{Var}[z_{1j}z_{2j} \mid \ell_j]$ . The CVF is not uniquely determined by the assumptions



about the first and second moments of  $\beta$  and  $\gamma$  used to derive proposition 1. Therefore we derive the CVF for the case where  $z_{1j}$  and  $z_{2j}$  are jointly distributed as bivariate normal<sup>3</sup>.

From a standard formula for double second moments of the bivariate normal, the CVF is

$$\begin{aligned}\text{Var}[z_{1j}z_{2j} | \ell_j] &= \mathbb{E}[z_{1j}^2 z_{2j}^2] \\ &= \text{Var}[z_{1j}]\text{Var}[z_{2j}] + 2\mathbb{E}[z_{1j}z_{2j}]^2 \\ &= \left(\frac{N_1 h_1^2 \ell_j}{M} + 1\right) \left(\frac{N_2 h_2^2 \ell_j}{M} + 1\right) + 2 \left(\frac{\sqrt{N_1 N_2} \rho_g}{M} \ell_j + \frac{\rho N_s}{\sqrt{N_1 N_2}}\right)^2\end{aligned}\quad (9)$$

In cases where the normality assumption does not hold, LD Score regression will remain unbiased, but may be inefficient, because the regression weights will be suboptimal.

## Liability Threshold Model

In the liability threshold (probit) model [61], binary traits are determined by an unobserved continuous liability  $\psi$ . The observed trait is  $y := \mathbf{1}[\psi > \tau]$ , where  $\tau$  is the liability threshold. If  $\psi$  is normally distributed, then setting  $\tau := \Phi^{-1}(1 - K)$  (where  $\Phi$  is the standard normal cdf) yields a population prevalence of  $K$ .

The liability threshold construction is surprisingly general. If  $y$  is an arbitrary binary phenotype, then we can define a normally distributed liability for  $y$  by setting  $\psi_i$  to be a random draw from a normal distribution truncated at the left by  $\tau$  if  $y_i = 1$  and truncated at the right by  $\tau$  if  $y_i = 0$ . We prefer to think of liability scale heritability as a useful way to compare the heritabilities of binary phenotypes with different prevalences, rather than a quantity that is only meaningful if one is willing to make strong assumptions about the data generating process.

For phenotypes generated according to the liability threshold model, we can estimate not only the heritability and genetic covariance of the observed phenotype, but also the heritability and genetic covariance of the unobserved liability.

In the next lemma, we derive population case and control allele frequencies in terms of the heritability of liability when liability is generated by a polygenic model. Since we are only modeling additive effects and are willing to assume Hardy-Weinberg equilibrium, we lose no generality and simplify notation considerably by stating the proofs in terms of haploid genotypes.

We state this lemma in terms of marginal per-allele effect sizes, instead of the per-normalized-genotype effect sizes considered in the earlier results on quantitative traits. Here marginal means that these are the effect sizes obtained by univariate regression of phenotype against genotype in the infinite data limit. Haploid normalized genotypes are defined  $X_j := (G_j - p_j)/\sqrt{p_j(1 - p_j)}$ . If  $\beta_j$  is the marginal per-normalized-genotype effect and  $\zeta_j$  is the marginal per-allele effect, we have  $X_j \beta_j = G_j \zeta_j$ . Thus, setting  $G_{ij} = 1$  yields  $\zeta_j = \beta_j \sqrt{(1 - p_j)/p_j}$ .

**Lemma 2.** *Suppose unobserved liabilities  $\psi, \varphi$  for traits  $y_1, y_2$  with thresholds  $\tau_1, \tau_2$  corresponding to prevalences  $K_1, K_2$  are generated according to the usual polygenic model for quantitative traits, i.e.,  $\psi_i = \sum_j X_{ij} \beta_j + \delta$ ,  $\varphi_i = \sum_j X_{ij} \gamma_j + \epsilon$ , with*

$$\text{Var}[(\beta, \gamma)] = \frac{1}{M} \begin{pmatrix} h_1^2 I & \rho_g I \\ \rho_g I & h_2^2 I \end{pmatrix},$$

<sup>3</sup>For instance, it is sufficient but not necessary to assume that  $\beta, \gamma, \delta$  and  $\epsilon$  are multivariate normal. More generally, the z-scores will be approximately normal if  $\beta$  and  $\gamma$  are reasonably polygenic. If the distribution of effect sizes is heavy-tailed, e.g., if there are few casual SNPs, then the CVF may be larger.

and

$$\text{Var}[(\delta, \epsilon)] = \begin{pmatrix} (1 - h_1^2)I & \rho_e I \\ \rho_e I & (1 - h_2^2)I \end{pmatrix}.$$

Let  $\zeta_j$  and  $\xi_j$  denote the marginal per-allele effect sizes of  $j$  on  $\psi$  and  $\varphi$ , and let  $p_{cas,kj} := \mathbb{P}[G_{ij} = 1 \mid y_{ik} = 1]$  for  $k = 1, 2$ . Then

$$\begin{aligned} \mathbb{E}[p_{cas,1j} - p_{con,1j}] &= 0, \\ \mathbb{E}[p_{cas,2j} - p_{con,2j}] &= 0, \\ \text{Var}[p_{cas,1j} - p_{con,1j}] &= \frac{p_j(1 - p_j)\phi(\tau_1)^2 h_1^2}{MK_1^2(1 - K_1)^2} \ell_j, \\ \text{Var}[p_{cas,2j} - p_{con,2j}] &= \frac{p_j(1 - p_j)\phi(\tau_2)^2 h_2^2}{MK_2^2(1 - K_2)^2} \ell_j, \\ \text{Cov}[p_{cas,1j} - p_{con,1j}, p_{cas,2j} - p_{con,2j}] &= \frac{p_j(1 - p_j)\phi(\tau_1)\phi(\tau_2)\rho_g}{MK_1(1 - K_1)K_2(1 - K_2)} \ell_j, \end{aligned}$$

where  $\phi$  is the standard normal density. These results apply to population allele frequencies, not allele frequencies in a finite sample. We deal with ascertained finite samples in the next section.

*Proof.* This proof is accomplished in two steps. First, we compute allele frequencies conditional on the marginal effects on liability. To do this, we reverse the conditional probability using Bayes' theorem, which reduces the problem to a series of [Taylor approximations to] Gaussian integrals. Second, we remove the conditioning on the marginal effects on liability in order to express the allele frequencies in terms of  $h_1^2, h_2^2, \rho_g$  and  $\ell_j$ . Since liability is just a quantitative trait, we need only apply the LD Score regression equation for quantitative traits.

By Bayes' rule,

$$\begin{aligned} \mathbb{P}[G_{ij} = 1 \mid y_{i1} = 1, \zeta_j] &= \frac{\mathbb{P}[y_{i1} = 1 \mid G_{ij} = 1, \zeta_j]\mathbb{P}[G_{ij} = 1]}{\mathbb{P}[y_{i1} = 1]} \\ &= \frac{p_j}{K_1} \mathbb{P}[y_{i1} = 1 \mid G_{ij} = 1, \zeta_j] \\ &= \frac{p_j}{K_1} \mathbb{P}[\psi_i > \tau_1 \mid G_{ij} = 1, \zeta_j]. \end{aligned} \tag{10}$$

The distribution of  $\psi$  is  $\psi \mid G_{ij} = 1, \zeta_j = N(\zeta_j, 1)$  (where the approximation that the variance equals one holds when the marginal heritability explained by  $j$  is small, which is the typical case in GWAS). Thus  $\mathbb{P}[\psi_i > \tau_1 \mid G_{ij} = 1]$  is simply a Gaussian integral. We approximate this probability with a first-order Taylor expansion around  $\tau_1$ .

$$\begin{aligned} \mathbb{P}[\psi_i > \tau_1 \mid G_{ij} = 1, \zeta_j] &= 1 - \Phi(\tau_1 - \zeta_j) \\ &\approx K_1 + \phi(\tau_1)\zeta_j, \end{aligned} \tag{11}$$

Substituting equation 11 into equation 10,

$$\mathbb{P}[G_{ij} = 1 \mid y_{i1} = 1, \zeta_j] = \frac{p_j}{K} (K + \phi(\tau_1)\zeta_j). \tag{12}$$

A similar argument shows that

$$\mathbb{P}[G_{ij} = 1 \mid y_{i1} = 0, \zeta_j] = \frac{p_j}{1 - K_1} (1 - K - \phi(\tau_1)\zeta_j). \tag{13}$$

Subtracting equation 12 from equation 13,

$$\mathbb{P}[G_{ij} = 1 \mid y_{i1} = 1, \zeta_j] - \mathbb{P}[G_{ij} = 1 \mid y_{i1} = 0, \zeta_j] = p_j \frac{\phi(\tau_1)\zeta_j}{K_1(1 - K_1)}. \quad (14)$$

Similar results hold for trait 2, replacing  $\zeta$  with  $\xi$  and subscript 1 with subscript 2.

We have written the probabilities in question in terms of constants and marginal effects on liability. Since liability is simply a quantitative trait, the means, variances, and covariances of the marginal effects on liability are described by the LD Score regression equation for quantitative traits. Precisely,  $\mathbb{E}[\xi_j] = \mathbb{E}[\zeta_j] = 0$ ,  $\text{Var}[\xi_j] = (1 - p_j)h_1^2\ell_j/p_kM$ ,  $\text{Var}[\zeta_j] = (1 - p_j)h_2^2\ell_j/p_jM$  and  $\text{Cov}[\zeta_j, \xi_j] = (1 - p_j)\rho_g\ell_j/p_jM$ . If we combine these results with equation 14, we find that  $\mathbb{E}[p_{cas,1j} - p_{con,1j}] = 0$ ;

$$\begin{aligned} \text{Var}[p_{cas,1j} - p_{con,1j}] &= \text{Var} \left[ \frac{p_j\phi(\tau_1)\zeta_j}{K_1(1 - K_1)} \right] \\ &= \frac{p_j(1 - p_j)\phi(\tau_1)^2h_1^2}{MK_1^2(1 - K_1)^2}\ell_j, \end{aligned} \quad (15)$$

(similarly for trait two) and

$$\begin{aligned} \text{Cov}[p_{cas,1j} - p_{con,1j}, p_{cas,2j} - p_{con,2j}] &= \text{Cov} \left[ \frac{p_j\phi(\tau_1)\zeta_j}{K_1(1 - K_1)}, \frac{p_j\phi(\tau_2)\xi_j}{K_2(1 - K_2)} \right] \\ &= \frac{p_j(1 - p_j)\phi(\tau_1)\phi(\tau_2)\rho_g}{MK_1(1 - K_1)K_2(1 - K_2)}\ell_j. \end{aligned} \quad (16)$$

□

## Ascertained Studies of Liability Threshold Traits

In the next proposition, we derive an LD Score regression equation for ascertained case/control studies.

Let  $P_i$  denote the sample prevalence of  $y_i$  in study  $i$  for  $i = 1, 2$ . We compute  $z$ -scores

$$z_j := \frac{\sqrt{NP(1 - P)}(\hat{p}_{cas} - \hat{p}_{con})}{\sqrt{\hat{p}_j(1 - \hat{p}_j)}},$$

where  $\hat{p}_j$  denotes allele frequency in the entire sample<sup>4</sup>,  $\hat{p}_{cas}$  denotes sample case allele frequency and  $\hat{p}_{con}$  denotes sample control allele frequency.

We emphasize one subtlety before stating the main proposition. The results in this section allow for study  $i$  to select samples based on  $y_j$  only if  $i = j$ . If study 1 ascertains on  $y_2$  – for example, if all cases in study 1 have  $y_1 = y_2 = 1$  – then  $\hat{p}_{cas,1j}$  will not be an unbiased estimate of  $p_{cas,1j}$ . Indeed, in this example,  $\mathbb{E}[\hat{p}_{cas,1j}] = \mathbb{P}[G_{ij} = 1 \mid y_1 = y_2 = 1]$ , which will not equal  $p_{cas,1j} = \mathbb{P}[G_{ij} = 1 \mid y_1 = 1]$  unless  $\rho = 1$  or  $\rho = 0$ . This follows from the fact that the conditionals and marginals of a bivariate normal are equal iff  $\rho = 0$  or  $\rho = 1$ . We do not derive formulae describing the bias, except to note that the most common scenario, the “healthy controls” model – cases are sampled independently but all controls are controls for both traits – is probably nothing to worry about. If cases are rare,  $\mathbb{P}[G_{ij} = 1 \mid y_1 = 0] \approx \mathbb{P}[G_{ij} = 1 \mid y_1 = y_2 = 0]$ . Conditioning on  $y_2 = 0$  hardly changes the distribution, because  $y_2 = 0$  most of the time, anyway.

<sup>4</sup>The expected value of  $\hat{p}_j$  is not equal to  $p_j$  unless  $P = K$  or  $j$  has zero marginal effect.

**Proposition 2.** Under the polygenic liability threshold model from lemma 2,

$$\mathbb{E}[z_{1j}z_{2j}] \approx \frac{\sqrt{N_1N_2}\rho_{g,obs}}{M}\ell_j + \sqrt{N_1N_2P_1(1-P_1)P_2(1-P_2)} \left( \sum_{a,b \in \{cas,con\}} \frac{N_{a,b}(-1)^{1+\mathbf{1}[a=b]}}{N_{a,1}N_{b,2}} \right) \quad (17)$$

where

$$\rho_{g,obs} := \rho_g \left( \frac{\phi(\tau_1)\phi(\tau_2)\sqrt{P_1(1-P_1)P_2(1-P_2)}}{K_1(1-K_1)K_2(1-K_2)} \right)$$

denotes observed scale genetic covariance,  $N_{a,b}$  denotes the number of individuals with phenotype  $a$  in study 1 and  $b$  in study two for  $a, b \in \{cas, con\}$  (e.g.,  $N_{cas,con}$  is the number of individuals who are a case in study 1 but a control in study 2),  $N_i$  denotes total sample size in study  $i$  and  $N_{a,i}$  for  $a \in \{cas, con\}$  and  $i = 1, 2$  denotes the number of individuals with phenotype  $a$  in study  $i$ .

*Proof.* The full form of  $z_{1j}z_{2j}$  is

$$z_{1j}z_{2j} = \frac{\sqrt{cN_1N_2}(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j})}{\sqrt{\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j})}},$$

where  $c := P_1(1-P_1)P_2(1-P_2)$ . Our strategy for obtaining the expectation is

$$\mathbb{E}[z_{1j}z_{2j}] \approx \sqrt{cN_1N_2} \frac{\mathbb{E}[(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j})]}{\mathbb{E}[\sqrt{\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j})}]} \quad (18)$$

$$\approx \sqrt{cN_1N_2} \frac{\mathbb{E}[(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j})]}{\sqrt{\mathbb{E}[\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j})]}} \quad (19)$$

$$= \sqrt{cN_1N_2} \frac{\mathbb{E}[\mathbb{E}[(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j}) | \zeta_j, \xi_j]]}{\sqrt{\mathbb{E}[\mathbb{E}[\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j}) | \zeta_j, \xi_j]]}}, \quad (20)$$

where  $\zeta_j$  and  $\xi_j$  denote the marginal per-allele effects of  $j$ . Approximation 18 hides  $\mathcal{O}(1/N)$  error from moving from the expectation of a ratio to a ratio of expectations. Approximation 19 hides  $\mathcal{O}(1/N)$  error from moving from the expectation of a square root to a square root of expectations. Equality 20 follows from applying the law of total expectation to the numerator and denominator.

First, we compute the numerator. By linearity of expectation,

$$\begin{aligned} \mathbb{E}[(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j}) | \zeta_j, \xi_j] &= \mathbb{E}[\hat{p}_{cas,1j}\hat{p}_{cas,2j}] - \mathbb{E}[\hat{p}_{cas,1j}\hat{p}_{con,2j}] \\ &\quad - \mathbb{E}[\hat{p}_{con,1j}\hat{p}_{cas,2j}] + \mathbb{E}[\hat{p}_{con,1j}\hat{p}_{con,2j}] \end{aligned} \quad (21)$$

After conditioning on the marginal effects  $\zeta_j$  and  $\xi_j$ , the only source of variance in the sample allele frequencies  $\hat{p}_{cas,1}, \hat{p}_{con,1}, \hat{p}_{cas,2}, \hat{p}_{con,2}$  is sampling error. Write  $\hat{p}_{cas,1j}\hat{p}_{cas,2j} = (p_{cas,1j} + \eta)(p_{cas,2j} + \nu)$ , where  $\eta$  and  $\nu$  denote sampling error. If study 1 and study 2 share samples,  $\nu$  and  $\eta$  will be correlated:

$$\begin{aligned} \mathbb{E}[\hat{p}_{cas,1j}\hat{p}_{cas,2j} | \zeta_j, \xi_j] &= p_{cas,1j}p_{cas,2j} + \mathbb{E}[\eta\nu] \\ &\approx p_{cas,1j}p_{cas,2j} + \frac{N_{cas,cas}\sqrt{p_{cas,1j}(1-p_{cas,1j})p_{cas,2j}(1-p_{cas,2j})}}{N_{cas,1}N_{cas,2}} \end{aligned} \quad (22)$$

$$\approx p_{cas,1j}p_{cas,2j} \left( 1 + \frac{N_{cas,cas}}{N_{cas,1}N_{cas,2}} \right), \quad (23)$$

where approximation 22 is the (bivariate) central limit theorem, and approximation 23 comes from ignoring the difference between  $\sqrt{p_{cas,1j}(1-p_{cas,1j})p_{cas,2j}(1-p_{cas,2j})}$  and  $p_j(1-p_j)$ . This step is justified in the derivation of the denominator. Similar relationships hold for the other terms in equation 21.

We can remove the conditioning on  $\zeta_j$  and  $\xi_j$  using equation 16.

$$\mathbb{E}[(p_{cas,1j} - p_{con,1j})(p_{cas,2j} - p_{con,2j})] = \frac{p_j(1-p_j)\phi(\tau_1)\phi(\tau_2)\rho_g}{MK_1(1-K_1)K_2(1-K_2)}\ell_j. \quad (24)$$

If we combine equations 23 and 24, we obtain

$$\mathbb{E}[(\hat{p}_{cas,1j} - \hat{p}_{con,1j})(\hat{p}_{cas,2j} - \hat{p}_{con,2j})] \approx p_j(1-p_j) \left( \frac{\phi(\tau_1)\phi(\tau_2)\rho_g}{c'M}\ell_j + \sum_{a,b \in \{cas, con\}} \frac{N_{a,b}(-1)^{1+\mathbf{1}[a=b]}}{N_{a,1}N_{b,2}} \right), \quad (25)$$

where  $c' := K_1(1-K_1)K_2(1-K_2)$ .

Next, we derive the expectation of the denominator. Conditional on  $\zeta_j$  and  $\xi_j$ ,  $\hat{p}_{1j}(1-\hat{p}_{1j})$  is  $P_1p_{cas,1j} + (1-P_1)p_{con,1j}$  plus  $\mathcal{O}(1/N)$  sampling variance. If studies 1 and 2 share samples, the  $\mathcal{O}(1/N)$  sampling variance in  $\hat{p}_{1j}(1-\hat{p}_{1j})$  and  $\hat{p}_{2j}(1-\hat{p}_{2j})$  will be correlated, but this still only amounts to  $\mathcal{O}(N_s/N_1N_2)$  error. If we remove the conditioning on  $\zeta_j$  and  $\xi_j$ , then  $P_1p_{cas,1j} + (1-P_1)p_{con,1j}$  is equal to  $p_j(1-p_j)$  plus  $\mathcal{O}(h_{1,obs}^2\ell_j/M)$  error from uncertainty in  $\zeta_j$ . The covariance between uncertainty in  $\zeta_j$  and uncertainty in  $\xi_j$  is driven by  $\rho_{g,obs}$ , so the expectation of the denominator is  $\mathbb{E}[\sqrt{\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j})}] = p_j(1-p_j)(1 + \mathcal{O}(N_s/N_1N_2) + \mathcal{O}(\rho_{g,obs}\ell_j/M))$ . We make the approximation<sup>5</sup> that

$$\mathbb{E} \left[ \sqrt{\hat{p}_{1j}(1-\hat{p}_{1j})\hat{p}_{2j}(1-\hat{p}_{2j})} \right] \approx p_j(1-p_j). \quad (26)$$

We obtain the desired result by dividing  $\sqrt{cN_1N_2}$  times equation 25 by equation 26.  $\square$

**Corollary 1.** *If study 1 is an ascertained study of a binary trait, and study 2 is a non-ascertained quantitative study, then proposition 2 holds, except with genetic covariance on the half-observed scale*

$$\rho_{g,obs} := \rho_g \left( \frac{\phi(\tau_1)\sqrt{P_1(1-P_1)}}{K_1(1-K_1)} \right).$$

**Corollary 2.** *For a single binary trait,*

$$\mathbb{E}[\chi_j^2] = \frac{Nh_{obs}^2}{M}\ell_j + 1, \quad (27)$$

where  $h_{obs}^2 = h^2\phi(\tau)^2P(1-P)/K^2(1-K)^2$ .

<sup>5</sup>For  $\ell_j = 100$  (roughly the median 1kG LD Score),  $M = 10^7$  and  $\rho_{g,obs} = 1$ , we get  $\rho_{g,obs}\ell_j/M = 10^{-5}$ . A worst-case value for  $N_s/N_1N_2$  might be  $N_s = N_1 = N_2 = 10^3$ , in which case  $N_s/N_1N_2 = 10^{-3}$ . Thus,  $\rho_{g,obs}\ell_j/M$  and  $N_s/N_1N_2$  will generally be at least 3 orders of magnitude smaller than 1.

*Proof.* This follows from proposition 2 if we set study 1 equal to study 2 and note that the observed scale genetic covariance between a trait and itself is observed scale heritability. To show that the intercept is one, observe that if study 1 and study 2 are the same, then

$$\begin{aligned} \sqrt{cN_1N_2} \left( \sum_{a,b \in \{cas, con\}} \frac{N_{a,b}(-1)^{1+\mathbf{1}[a=b]}}{N_{a,1}N_{b,2}} \right) &= NP(1-P) \left( \frac{1}{N_{cas}} + \frac{1}{N_{con}} \right) \\ &= \frac{NP(1-P)(N_{cas} + N_{con})}{N_{cas}N_{con}} \\ &= \frac{N^2P(1-P)}{N_{cas}N_{con}}. \end{aligned} \quad (28)$$

But  $NP = N_{cas}$  and  $N(1-P) = N_{con}$ , so equation 28 simplifies to 1.  $\square$

Finally, observe that  $\rho_{g,obs}/\sqrt{h_{1,obs}^2 h_{2,obs}^2} = \rho_g/\sqrt{h_1^2 h_2^2} = r_g$ . Put another way, the natural definition for “observed scale genetic correlation” turns out to be the same as regular genetic correlation, because the scale transformation factors in the numerator and denominator cancel. This is convenient: we can compute genetic correlations for binary traits on a sensible scale without having to worry about sample and population prevalences.

## Flavors of Heritability and Genetic Correlation

The heritability parameter estimated by **ldsc** is subtly different than the heritability parameter  $h_g^2$  estimated by **GCTA**. If  $g$  denotes the set of all genotyped SNPs in some GWAS, define  $\beta_{GCTA} := \operatorname{argmax}_{\alpha \in \mathbb{R}^{|g|}} \operatorname{Cor}[y_1, X_g \alpha]^2$ , where  $X_g$  is a random vector of genotypes for SNPs in  $g$ . Then the heritability parameter estimated by **GCTA** is defined

$$h_g^2 := \sum_{j \in g} \beta_{GCTA,j}^2.$$

Let  $S$  denote the set of SNPs used to compute LD Scores (*i.e.*,  $\ell_j = \sum_{k \in S} r_{jk}^2$ ), and let  $\beta_S := \operatorname{argmax}_{\alpha \in \mathbb{R}^{|S|}} \operatorname{Cor}[y_1, X_S \alpha]^2$ . Generally  $\beta_{S,j} \neq \beta_{GCTA,j}$  unless all SNPs in  $S \setminus g$  are not in LD with SNPs in  $g$ . Define

$$h_S^2 := \sum_{j \in S} \beta_{S,j}^2.$$

Let  $S'$  denote the set of SNPs in  $S$  with MAF above 5%. Define

$$h_{5-50\%}^2 := \sum_{j \in S'} \beta_{S,j}^2. \quad (29)$$

The default setting in **ldsc** is to report  $h_{5-50\%}^2$ , estimated as the slope from LD Score regression times  $M_{5-50\%}$ , the number of SNPs with MAF above 5%.

The reason for this is the following: suppose that  $h^2$  per SNP is not constant as a function of MAF. Then the slope of LD Score regression will represent some sort of weighted average of the values of  $h^2$  per SNP, with more weight given to classes of SNPs that are well-represented among the regression SNPs. In a typical GWAS setting, the regression SNPs are mostly common SNPs, so multiplying the slope from LD Score regression by  $M$  (which includes rare SNPs) amounts to

extrapolating that  $h^2$  per SNP among common variants is the same as  $h^2$  per SNP among rare variants. This extrapolation is particularly risky, because there are many more rare SNPs than common SNPs.

It is probably reasonable to treat  $h^2$  per SNP as a constant function of MAF for SNPs with MAF above 5%, but we have very little information about  $h^2$  per SNP for SNPs with MAF below 5%. Therefore we report  $h_{5-50\%}^2$  instead of  $h_S^2$  to avoid excessive extrapolation error. This lower bound can be pushed lower with larger sample sizes and better rare variant coverage, either from sequencing or imputation.

There are two main distinctions between  $h_{5-50\%}^2$  and  $h_g^2$ . First,  $h_g^2$  does not include the effects of common SNPs that are not tagged by the set of genotyped SNPs  $g$ . Second, the effects of causal 4% SNPs are not counted towards  $h_{5-50\%}^2$ . In practice, neither of these distinctions makes a large difference, since most GWAS arrays focus on common variation and manage to assay or tag almost all common variants, which is why we do not emphasize this distinction in the main text.

The relationship between the genetic covariance parameter estimated by LD Score regression and the genetic covariance parameter estimated by **GCTA** is similar to the relationship between  $h_{5-50\%}^2$  and  $h_g^2$ . Choice of  $M$  is not important for genetic correlation, because the factors of  $M$  in the numerator and denominator cancel.

## Supplementary Tables

### Simulations with one Binary Trait and one Quantitative Trait

Prevalence	$\hat{h}^2$	$\hat{h}_{liab}^2$	$\hat{r}_g$
0.01	0.72 (0.1)	0.59 (0.04)	0.51 (0.4)
0.05	0.72 (0.12)	0.59 (0.07)	0.45 (0.17)
0.2	0.72 (0.11)	0.6 (0.08)	0.46 (0.14)
0.5	0.73 (0.11)	0.59 (0.08)	0.42 (0.17)

Table S1: *Simulations with one binary trait and one quantitative trait. The prevalence column describes the population prevalence of the binary trait. We ran 100 simulations for each prevalence. The  $\hat{h}^2$  column shows the mean heritability estimate for the quantitative trait. The  $\hat{h}_{liab}^2$  column shows the mean liability-scale heritability estimate for the binary trait. The  $\hat{r}_g$  column shows the mean genetic correlation estimate. Standard deviations across 100 simulations in parentheses. The true parameter values were  $r_g = 0.46$ ,  $h^2 = 0.7$  for the quantitative trait and  $h_{liab}^2 = 0.6$  for the binary trait. For all simulations, the quantitative trait sample size was 1000, the binary trait sample size was 1000 cases and 1000 controls, and there were 500 overlapping samples. There were 1000 effective independent SNPs. The environmental covariance was 0.2. We simulated case/control ascertainment using simulated LD block genotypes and a rejection sampling model of ascertainment. This is the same strategy used to simulate case/control ascertainment in [21].*



## Simulations with MAF- and LD-Dependent Genetic Architecture

LD Score	$h^2(5-50\%)$	$\rho_g(5-50\%)$	$r_g(5-50\%)$
Truth	0.83	0.42	0.5
HM3	0.53 (0.08)	0.28 (0.07)	0.52 (0.1)
PNG	0.36 (0.08)	0.18 (0.06)	0.5 (0.13)
30 Bins	0.81 (0.12)	0.41 (0.08)	0.51 (0.09)
60 Bins	0.81 (0.12)	0.41 (0.09)	0.51 (0.09)

Table S2: *Simulations with MAF- and LD-dependent genetic architecture. Effect sizes were drawn from normal distributions such that the variance of per-allele effect sizes was uncorrelated with MAF, and variants with LD Score below 100 were fourfold enriched for heritability. Sample size was 2062 with complete overlap between studies; causal SNPs were about 600,000 best-guess imputed 1kG SNPs on chr 2, and the SNPs retained for the LD Score regression were the subset of about 100,000 of these SNPs that were included in HM3. True parameter values are shown in the top line of the table. Estimates are averages across 100 simulations. Standard deviations (in parentheses) are standard deviations across 100 simulations. LD Scores were estimated using in-sample LD and a 1cM window. HM3 means LD Score with sum taken over SNPs in HM3. PNG (per-normalized-genotype) means LD Score with the sum taken over all SNPs in 1kG as in [21]. 30 bins means per-allele LD Score binned on a MAF by LD Score grid with MAF breaks at 0.05, 0.1, 0.2, 0.3 and 0.4 and LD Score breaks at 35, 75, 150 and 400. 60 bins means per-allele LD Score binned on a MAF by LD Score grid with MAF breaks at 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4 and 0.45 and LD Score breaks at 30, 60, 120, 200 and 300. These simulations demonstrate that naive (HM3, PNG) LD Score regression gives correct genetic correlation estimates even when heritability and genetic covariance estimates are biased, so long as genetic correlation does not depend on LD.*

## Sample Sizes and References

Trait	Reference	Sample Size
Schizophrenia	PGC Schizophrenia Working Group, <i>Nature</i> , 2014 [62]	70,100
Bipolar disorder	PGC Bipolar Working Group, <i>Nat Genet</i> , 2011 [63]	16,731
Major depression	PGC MDD Working Group, <i>Mol Psych</i> , 2013 [64]	18,759
Anorexia Nervosa	Boraska, <i>et al.</i> , <i>Mol Psych</i> , 2014 [65]	17,767
Autism Spectrum Disorder	PGC Cross-Disorder Group, <i>Lancet</i> , 2013 [24]	10,263
Ever/Never Smoked	TAG Consortium, 2010 <i>Nat Genet</i> , [66]	74,035
Alzheimer's	Lambert, <i>et al.</i> , <i>Nat Genet</i> , 2013 [67]	54,162
College	Rietveld, <i>et al.</i> , <i>Science</i> , 2013 [38]	101,069
Height	Lango Allen, <i>et al.</i> , <i>Nature</i> 2010 [68]	133,858
Obesity Class 1	Berndt, <i>et al.</i> , <i>Nat Genet</i> , 2013 [69]	98,000
Extreme Waist-Hip Ratio	Berndt, <i>et al.</i> , <i>Nat Genet</i> , 2013 [69]	10,000
Coronary Artery Disease	Schunkert, <i>et al.</i> , <i>Nat Genet</i> , 2011 [70]	86,995
Triglycerides	Teslovich, <i>et al.</i> , <i>Nature</i> , 2010 [71]	96,598
LDL Cholesterol	Teslovich, <i>et al.</i> , <i>Nature</i> , 2010 [71]	95,454
HDL Cholesterol	Teslovich, <i>et al.</i> , <i>Nature</i> , 2010 [71]	99,900
Type-2 Diabetes	Morris, <i>et al.</i> , <i>Nat Genet</i> , 2012 [26]	69,033
Fasting Glucose	Manning, <i>et al.</i> , <i>Nat Genet</i> , 2012 [72]	46,186
Childhood Obesity	EGG Consortium, <i>Nat Genet</i> , 2012 [29]	13,848
Birth Length	van der Valk, <i>et al.</i> , <i>HMG</i> , 2014 [73]	22,263
Birth Weight	Horikoshi, <i>et al.</i> , <i>Nat Genet</i> , 2013 [27]	26,836
Infant Head Circumference	Taal, <i>et al.</i> , <i>Nat Genet</i> , 2012 [30]	10,767
Age at Menarche	Perry, <i>et al.</i> , <i>Nature</i> , 2014 [25]	132,989
Crohn's Disease	Jostins, <i>et al.</i> , <i>Nature</i> , 2012 [74]	20,883
Ulcerative Colitis	Jostins, <i>et al.</i> , <i>Nature</i> , 2012 [74]	27,432
Rheumatoid Arthritis	Stahl, <i>et al.</i> , <i>Nat Genet</i> , 2010 [75]	25,708

Table S3: *Sample sizes and references for traits analyzed in the main text.*

## Genetic Correlation between Educational Attainment Phenotypes

Phenotype 1	Phenotype 2	$r_g$	$se$
College (Yes/No)	Years of Education	1.00	0.014

Table S4: *Genetic correlation between the two educational attainment phenotypes from Rietveld, et al. [38].*

## Supplementary Figures

### Genetic Correlations among Anthropometric Traits

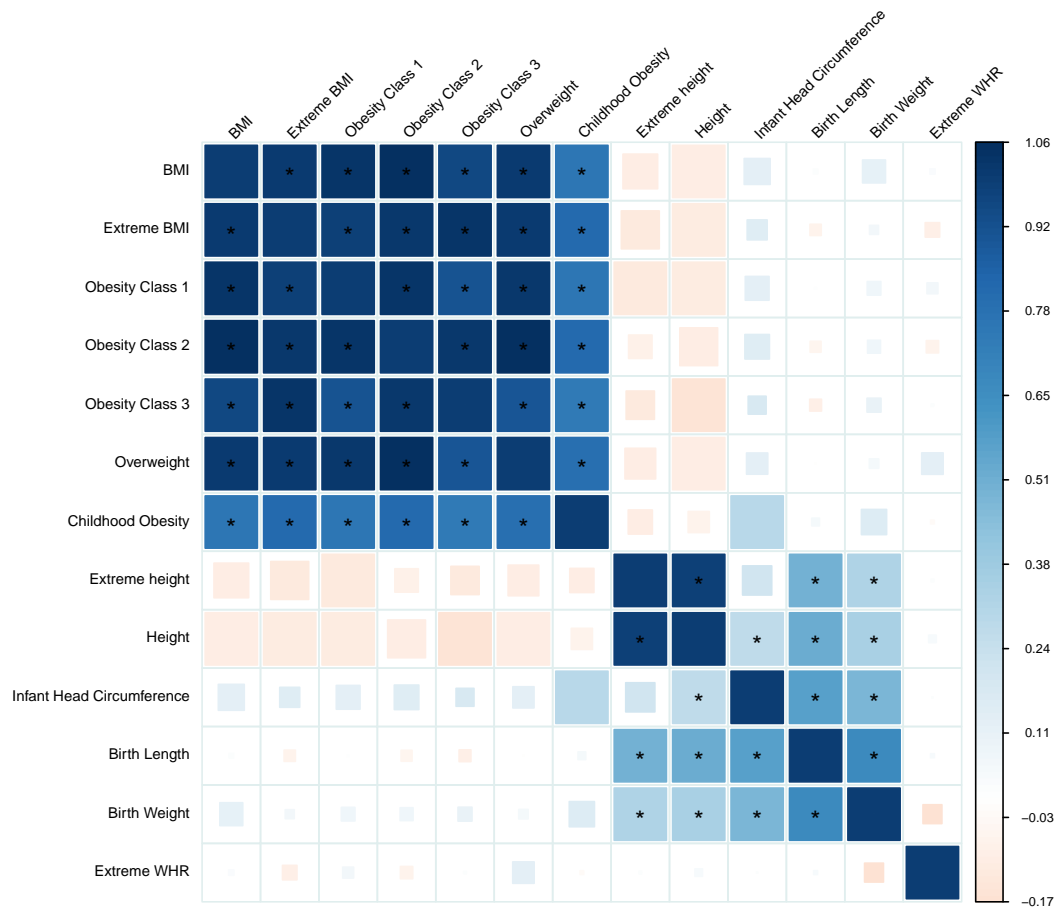


Figure S1: Genetic correlations among anthropometric traits from studies by the GIANT and EGG consortia. The structure of the figure is the same as Figure 2 in the main text: blue corresponds to positive genetic correlations; red corresponds to negative genetic correlation. Larger squares correspond to more significant p-values. Genetic correlations that are different from zero at 1% FDR are shown as full-sized squares. Genetic correlations that are significantly different from zero at significance level 0.05 after Bonferroni correction are given an asterisk.

## Genetic Correlations among Smoking Traits

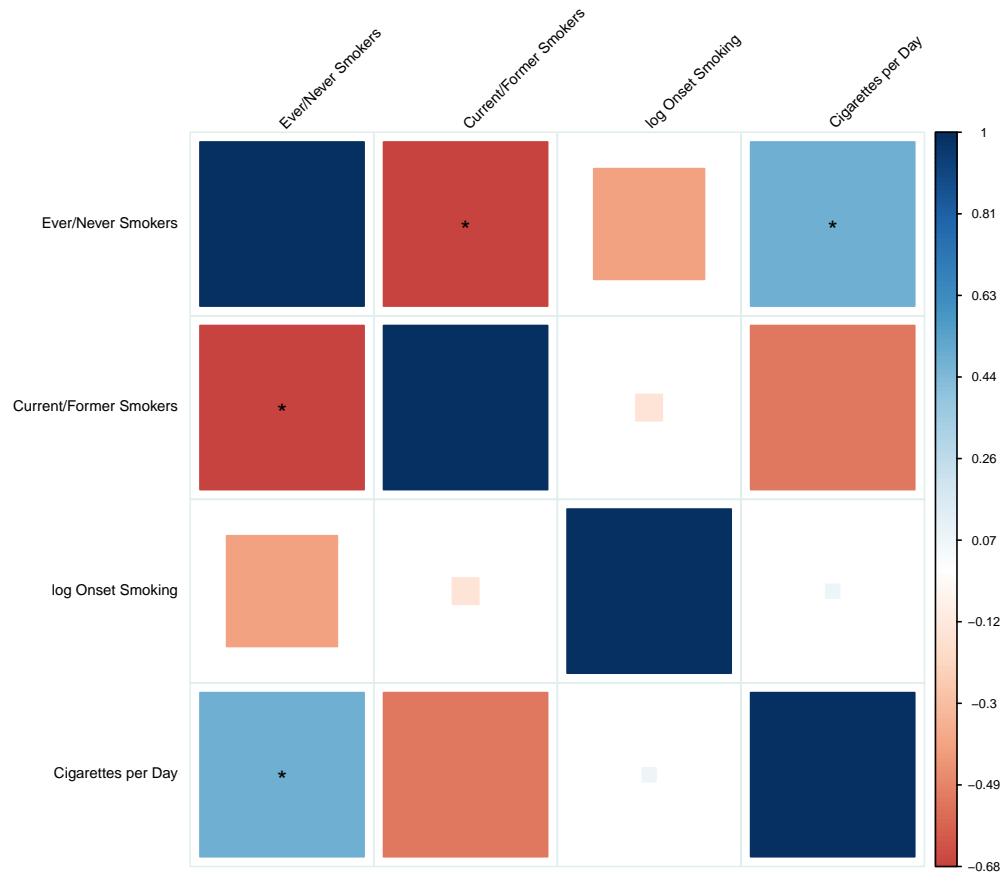


Figure S2: Genetic correlations among smoking traits from the Tobacco and Genetics (TAG) consortium. The structure of the figure is the same as Figure 2 in the main text: blue corresponds to positive genetic correlations; red corresponds to negative genetic correlation. Larger squares correspond to more significant p-values. Genetic correlations that are different from zero at 1% FDR are shown as full-sized squares. Genetic correlations that are significantly different from zero at significance level 0.05 after Bonferroni correction are given an asterisk.

## Genetic Correlations among Insulin-Related Traits

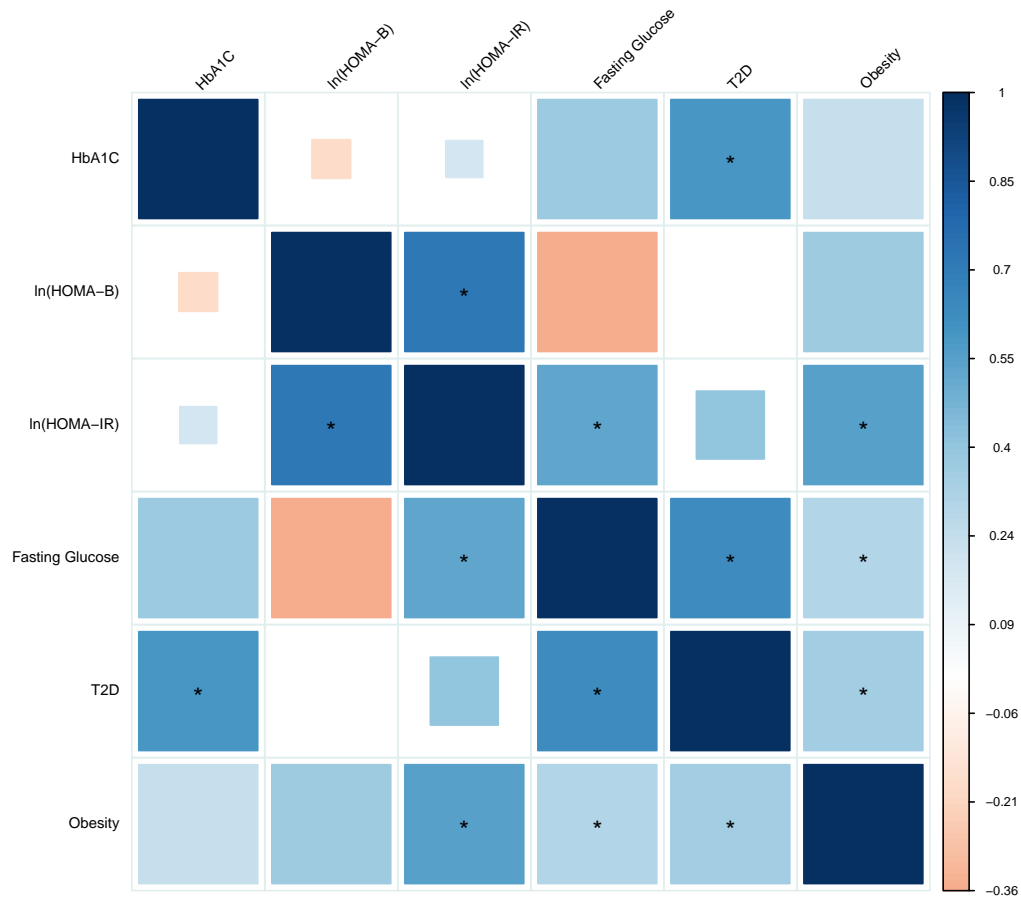


Figure S3: Genetic correlations among insulin-related traits from studies by the MAGIC consortium. The structure of the figure is the same as Figure 2 in the main text: blue corresponds to positive genetic correlations; red corresponds to negative genetic correlation. Larger squares correspond to more significant p-values. Genetic correlations that are different from zero at 1% FDR are shown as full-sized squares. Genetic correlations that are significantly different from zero at significance level 0.05 after Bonferroni correction are given an asterisk.

## Metabolic Genetic Correlations from Vattikuti, et al. and LD Score

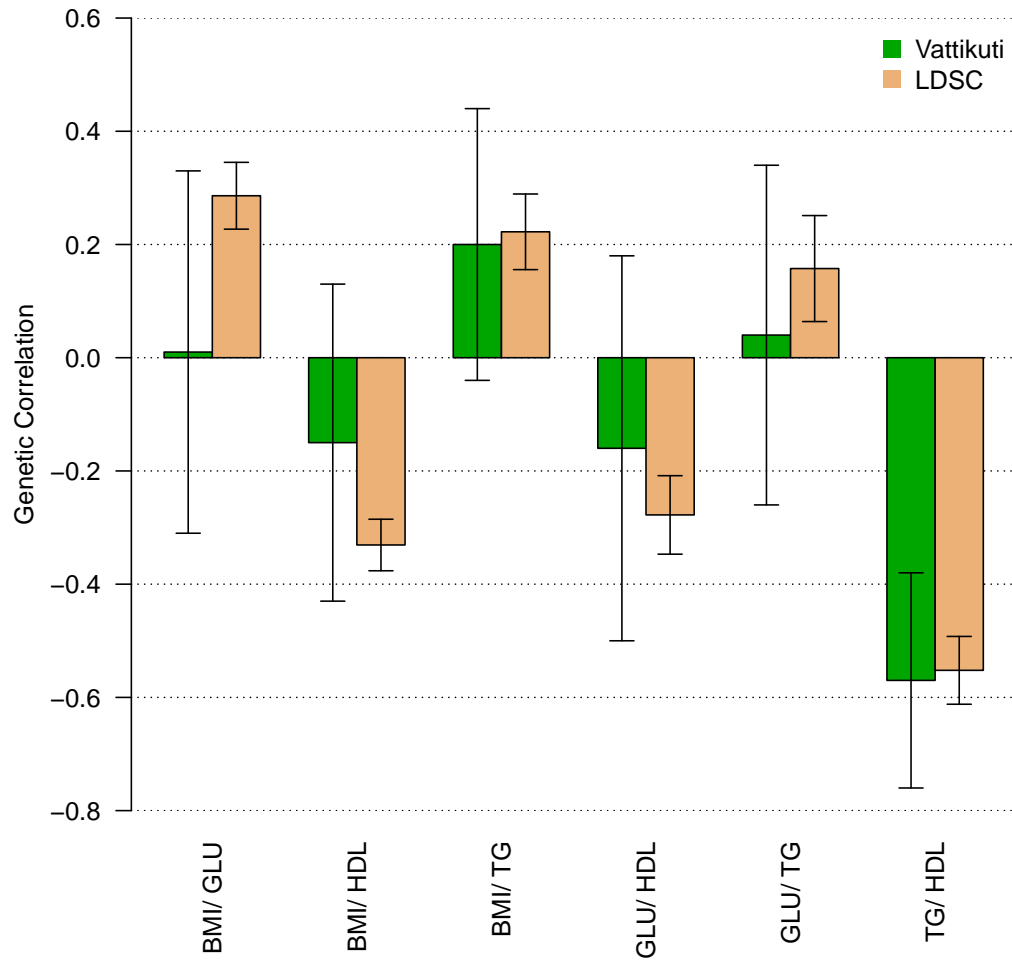


Figure S4: This figure compares estimates of genetic correlations among metabolic traits from table 3 of Vattikuti et al. [17] to the estimates from LD Score regression (and larger sample sizes). Error bars are standard errors.

## Schizophrenia / TG Conditional QQ Plot with and without the MHC

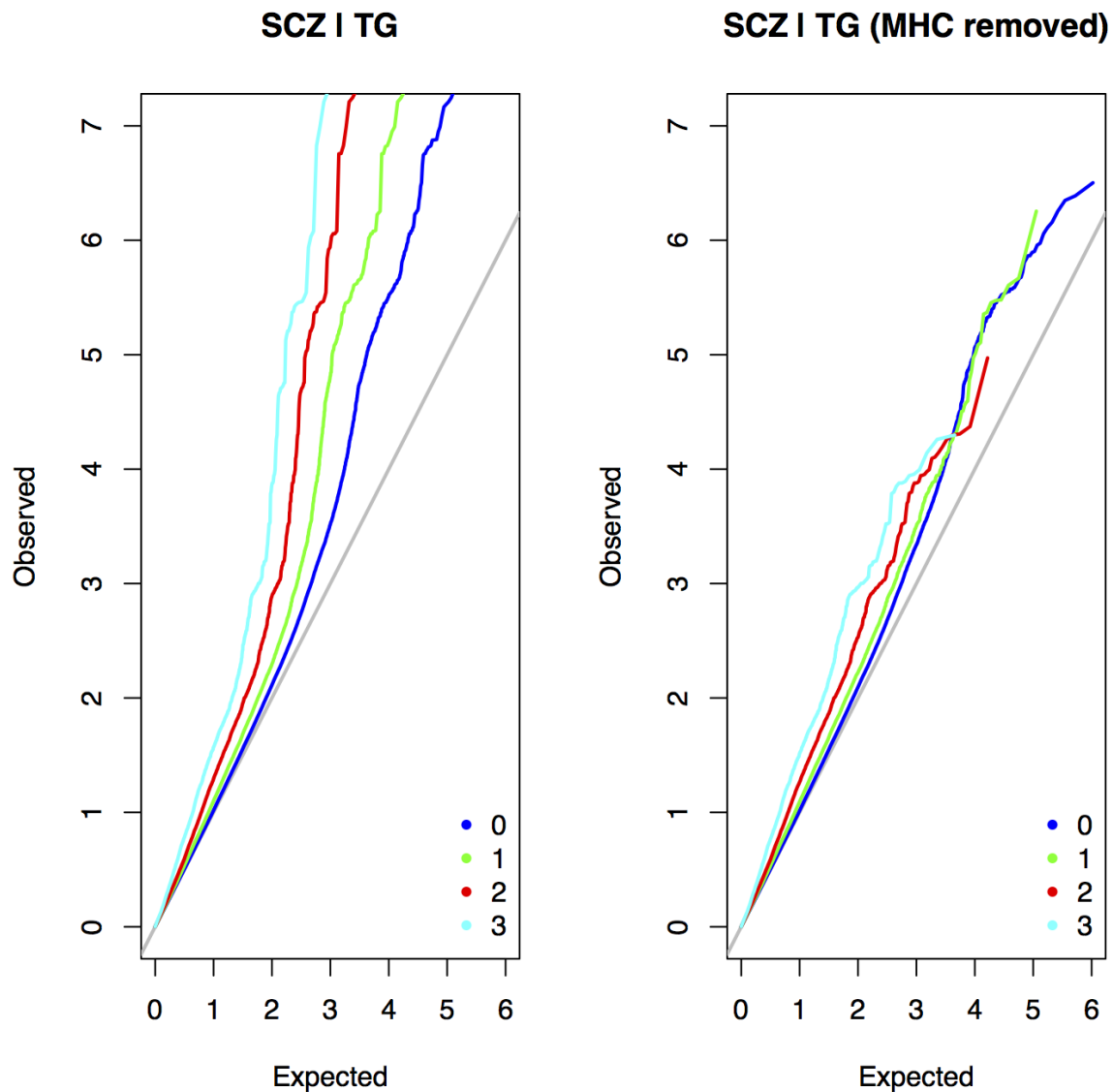


Figure S5: We reproduced the conditional QQ plot comparing schizophrenia (SCZ) and triglycerides (TG) from Andreassen et al. [76] (left). The major histocompatibility complex (MHC, chr6, 25-35 MB) is a genomic region containing SNPs with exceptionally high LD Scores and the strongest GWAS association for schizophrenia [62] as well as a GWAS association to TG [71]. Removing the MHC removes all signal of enrichment from the conditional QQ plot (right).



## Collaborators

The members of the Schizophrenia Working Group of the Psychiatric Genetics Consortium are Stephan Ripke, Benjamin M. Neale, Aiden Corvin, James T.R. Walters, Kai-How Farh, Peter A. Holmans, Phil Lee, Brendan Bulik-Sullivan, David A. Collier, Hailiang Huang, Tune H. Pers, Ingrid Agartz, Esben Agerbo, Margot Albus, Madeline Alexander, Farooq Amin, Silviu A. Bacanu, Martin Begemann, Richard A. Belliveau, Jr., Judit Bene, Sarah E. Bergen, Elizabeth Bevilacqua, Tim B. Bigdeli, Donald W. Black, Anders D. Brglum, Richard Bruggeman, Nancy G. Buccola, Randy L. Buckner, William Byerley, Wiepke Cahn, Guiqing Cai, Dominique Campion, Rita M. Cantor, Vaughan J. Carr, Noa Carrera, Stanley V. Catts, Kimberly D. Chambert, Raymond C.K. Chan, Ronald Y.L. Chen, Eric Y.H. Chen, Wei Cheng, Eric F.C. Cheung, Siow Ann Chong, C. Robert Cloninger, David Cohen, Nadine Cohen, Paul Cormican, Nick Craddock, James J. Crowley, David Curtis, Michael Davidson, Kenneth L. Davis, Franziska Degenhardt, Jurgen Del Favero, Lynn E. DeLisi, Ditte Demontis, Dimitris Dikeos, Timothy Dinan, Srdjan Djurovic, Gary Donohoe, Elodie Drapeau, Jubao Duan, Frank Dudbridge, Naser Durmishi, Peter Eichhammer, Johan Eriksson, Valentina Escott-Price, Laurent Essioux, Ayman H. Fanous, Marttila S. Farrell, Josef Frank, Lude Franke, Robert Freedman, Nelson B. Freimer, Marion Friedl, Joseph I. Friedman, Menachem Fromer, Giulio Genovese, Lyudmila Georgieva, Elliot S. Gershon, Ina Giegling, Paola Giusti-Rodriguez, Stephanie Godard, Jacqueline I. Goldstein, Vera Golimbet, Srihari Gopal, Jacob Gratten, Jakob Grove, Lieuwe de Haan, Christian Hammer, Marian L. Hamshere, Mark Hansen, Thomas Hansen, Vahram Haroutunian, Annette M. Hartmann, Frans A. Henskens, Stefan Herms, Joel N. Hirschhorn, Per Hoffmann, Andrea Hofman, Mads V. Hollegaard, David M. Hougaard, Masashi Ikeda, Inge Joa, Antonio Julia, Rene S. Kahn, Luba Kalaydjieva, Sena Karachanak-Yankova, Juha Karjalainen, David Kavanagh, Matthew C. Keller, Brian J. Kelly, James L. Kennedy, Andrey Khrunin, Yunjung Kim, Janis Klovins, James A. Knowles, Bettina Konte, Vaidutis Kucinskas, Zita Ausrele Kucinskiene, Hana Kuzelova-Ptackova, Anna K. Kahler, Claudine Laurent, Jimmy Lee Chee Keong, S. Hong Lee, Sophie E. Legge, Bernard Lerer, Miaoxin Li, Tao Li, Kung-Yee Liang, Jeffrey Lieberman, Svetlana Limborska, Carmel M. Loughland, Jan Lubinski, Jouko Lonnqvist, Milan Macek, Jr., Patrik K.E. Magnusson, Brion S. Maher, Wolfgang Maier, Jacques Mallet, Sara Marsal, Manuel Mattheisen, Morten Mattingsdal, Robert W. McCarley, Colm McDonald, Andrew M. McIntosh, Sandra Meier, Carin J. Meijer, Bela Melegh, Ingrid Melle, Raquelle I. Meshulam-Gately, Andres Metspalu, Patricia T. Michie, Lili Milani, Vihra Milanova, Younes Mokrab, Derek W. Morris, Ole Mors, Preben B. Mortensen, Kieran C. Murphy, Robin M. Murray, Inez Myin-Germeys, Bertram Mller-Myhsok, Mari Nelis, Igor Nenadic, Deborah A. Nertney, Gerald Nestadt, Kristin K. Nicodemus, Liene Nikitina-Zake, Laura Nisenbaum, Annelie Nordin, Eadbhard O'Callaghan, Colm O'Dushlaine, F. Anthony O'Neill, Sang-Yun Oh, Ann Olincy, Line Olsen, Jim Van Os, Psychosis Endophenotypes International Consortium, Christos Pantelis, George N. Papadimitriou, Sergi Papiol, Elena Parkhomenko, Michele T. Pato, Tiina Paunio, Milica Pejovic-Milovancevic, Diana O. Perkins, Olli Pietilinen, Jonathan Pimm, Andrew J. Pocklington, John Powell, Alkes Price, Ann E. Pulver, Shaun M. Purcell, Digby Quested, Henrik B. Rasmussen, Abraham Reichenberg, Mark A. Reimers, Alexander L. Richards, Joshua L. Roffman, Panos Roussos, Douglas M. Ruderfer, Veikko Salomaa, Alan R. Sanders, Ulrich Schall, Christian R. Schubert, Thomas G. Schulze, Sibylle G. Schwab, Edward M. Scolnick, Rodney J. Scott, Larry J. Seidman, Jianxin Shi, Engilbert Sigurdsson, Teimuraz Silagadze, Jeremy M. Silverman, Kang Sim, Petr Slominsky, Jordan W. Smoller, Hon-Cheong So, Chris C.A. Spencer, Eli A. Stahl, Hreinn Stefansson, Stacy Steinberg, Elisabeth Stogmann, Richard E. Straub, Eric Strengman, Jana Strohmaier, T. Scott Stroup, Mythily Subra-

maniam, Jaana Suvisaari, Dragan M. Svrakic, Jin P. Szatkiewicz, Erik Sderman, Srinivas Thirumalai, Draga Toncheva, Paul A. Tooney, Sarah Tosato, Juha Veijola, John Waddington, Dermot Walsh, Dai Wang, Qiang Wang, Bradley T. Webb, Mark Weiser, Dieter B. Wildenauer, Nigel M. Williams, Stephanie Williams, Stephanie H. Witt, Aaron R. Wolen, Emily H.M. Wong, Brandon K. Wormley, Jing Qin Wu, Hualin Simon Xi, Clement C. Zai, Xuebin Zheng, Fritz Zimprich, Naomi R. Wray, Kari Stefansson, Peter M. Visscher, Wellcome Trust Case Control Consortium, Rolf Adolfsson, Ole A. Andreassen, Douglas H.R. Blackwood, Elvira Bramon, Joseph D. Buxbaum, Anders D. Brglum, Sven Cichon, Ariel Darvasi, Enrico Domenici, Hannelore Ehrenreich, Tonu Esko, Pablo V. Gejman, Michael Gill, Hugh Gurling, Christina M. Hultman, Nakao Iwata, Assen V. Jablensky, Erik G. Jonsson, Kenneth S. Kendler, George Kirov, Jo Knight, Todd Lencz, Douglas F. Levinson, Qingqin S. Li, Jianjun Liu, Anil K. Malhotra, Steven A. McCarroll, Andrew McQuillin, Jennifer L. Moran, Preben B. Mortensen, Bryan J. Mowry, Markus M. Nthen, Roel A. Ophoff, Michael J. Owen, Aarno Palotie, Carlos N. Pato, Tracey L. Petryshen, Danielle Posthuma, Marcella Rietschel, Brien P. Riley, Dan Rujescu, Pak C. Sham, Pamela Sklar, David St. Clair, Daniel R. Weinberger, Jens R. Wendland, Thomas Werge, Mark J. Daly, Patrick F. Sullivan, and Michael C. O'Donovan.

The members of the Genetic Consortium for Anorexia Nervosa (GCAN) are V Boraska, C S Franklin, J A B Floyd, L M Thornton, L M Huckins, L Southam, N William Rayner, I Tachmazidou, K L Klump, J Treasure, C M Lewis, U Schmidt, F Tozzi, K Kiezebrink, J Hebebrand, P Gorwood, R A H Adan, M J H Kas, A Favaro, P Santonastaso, F Fernandez-Aranda, M Gratacos, F Rybakowski, M Dmitrzak-Weglarz, J Kaprio, A Keski-Rahkonen, A Raevuori, E F Van Furth, M C T Slof-Oot Landt, J I Hudson, T Reichborn-Kjennerud, G P S Knudsen, P Monteleone, A S Kaplan, A Karwautz, H Hakonarson, W H Berrettini, Y Guo, D Li, N J Schork, G Komaki, T Ando, H Inoko, T Esko, K Fischer, K Mnnik, A Metspalu, J H Baker, R D Cone, J Dackor, J E DeSocio, C E Hilliard, J K O'Toole, J Pantel, J P Szatkiewicz, C Taico, S Zerwas, S E Trace, O S P Davis, S Helder, K Bhren, R Burghardt, M de Zwaan, K Egberts, S Ehrlich, B Herpertz-Dahlmann, W Herzog, H Imgart, A Scherag, S Scherag, S Zipfel, C Boni, N Ramoz, A Versini, M K Brandys, U N Danner, C de Kove, J Hendriks, B P C Koeleman, R A Ophoff, E Strengman, A A van Elburg, A Bruson, M Clementi, D Degortes, M Forzan, E Tenconi, E Docampo, G Escarams, S Jimnez-Murcia, J Lissowska, A Rajewski, N Szeszenia-Dabrowska, A Slopian, J Hauser, L Karhunen, I Meulenbelt, P E Slagboom, A Tortorella, M Maj, G Dedoussis, D Dikeos, F Gonidakis, K Tziouvas, A Tsitsika, H Papezova, L Slachtova, D Martaskova, J L Kennedy, R D Levitan, Z Yilmaz, J Huemer, D Koubek, E Merl, G Wagner, P Lichtenstein, G Breen, S Cohen-Woods, A Farmer, P McGuffin, S Cichon, I Giegling, S Herms, D Rujescu, S Schreiber, H-E Wichmann, C Dina, R Sladek, G Gambaro, N Soranzo, A Julia, S Marsal, R a Rabionet, V Gaborieau, D M Dick, A Palotie, S Ripatti, E Widn, O A Andreassen, T Espeseth, A Lundervold, I Reinvang, V M Steen, S Le Hellard, M Mattingsda, I Ntalla, V Bencko, L Foretova, V Janout, M Navratilova, S Gallinger, D Pinto, S W Scherer, H Aschauer, L Carlberg, A Schosser, L Alfredsson, B Ding, L Klareskog, L Padyukov, C Finan, G Kalsi, M Roberts, D W Logan, L Peltonen, G R S Ritchie, P Courtet, S Guilleme, I Jaussent, J C Barrett, X Estivill, A Hinney, P F Sullivan, D A Collier, E Zeggini, and C M Bulik.

The members of the Wellcome Trust Case Control Consortium 3 (WTCCC3) are C A Anderson, J C Barrett, J A B Floyd, C S Franklin, R McGinnis, N Soranzo, E Zeggini, J Sambrook, J Stephens, W H Ouwehand, W L McArdle, S M Ring, D P Strachan, G Alexander, C M Bulik, D A Collier, P J Conlon, A Dominiczak, A Duncanson, A Hill, C Langford, G Lord, A P Maxwell, L Morgan, L Peltonen, R N Sandford, N Sheerin, N Soranzo, F O Vannberg, J C Barrett, D N A Genotyping, H Blackburn, W-M Chen, S Edkins, M Gillman, E Gray, S E Hunt, C Langford, a Nengut-Gumuscu,

S Potter, S S Rich, D Simpkin, and Pa Whittaker.

The members of the ReproGen consortium are John RB Perry, Felix Day, Cathy E Elks, Patrick Sulem, Deborah J Thompson, Teresa Ferreira, Chunyan He, Daniel I Chasman, Tnu Esko, Gudmar Thorleifsson, Eva Albrecht, Wei Q Ang, Tanguy Corre, Diana L Cousminer, Bjarke Feenstra, Nora Franceschini, Andrea Ganna, Andrew D Johnson, Sanela Kjellqvist, Kathryn L Lunetta, George McMahon, Ilja M Nolte, Lavinia Paternoster, Eleonora Porcu, Albert V Smith, Lisette Stolk, Alexander Teumer, Natalia Ternikova, Emmi Tikkanen, Sheila Ulivi, Erin K Wagner, Najaf Amin, Laura J Bierut, Enda M Byrne, JoukeJan Hottenga, Daniel L Koller, Massimo Mangino, Tune H Pers, Laura M YergesArmstrong, Jing Hua Zhao, Irene L Andrulis, Hoda AntonCulver, Femke Atsma, Stefania Bandinelli, Matthias W Beckmann, Javier Benitez, Carl Blomqvist, Stig E Bojesen, Manjeet K Bolla, Bernardo Bonanni, Hiltrud Brauch, Hermann Brenner, Julie E Buring, Jenny ChangClaude, Stephen Chanock, Jinhui Chen, Georgia ChenevixTrench, J. Margriet Colle, Fergus J Couch, David Couper, Andrea D Coveillo, Angela Cox, Kamila Czene, Adamo Pio D'adamo, George Davey Smith, Immaculata De Vivo, Ellen W Demerath, Joe Dennis, Peter Devilee, Aida K Dieffenbach, Alison M Dunning, Gudny Eiriksdottir, Johan G Eriksson, Peter A Fasching, Luigi Ferrucci, Dieter FleschJanys, Henrik Flyger, Tatiana Foroud, Lude Franke, Melissa E Garcia, Montserrat GarcaClosas, Frank Geller, Eco EJ de Geus, Graham G Giles, Daniel F Gudbjartsson, Vilmundur Gudnason, Pascal Gunel, Suqun Guo, Per Hall, Ute Hamann, Robin Haring, Catharina A Hartman, Andrew C Heath, Albert Hofman, Maartje J Hooning, John L Hopper, Frank B Hu, David J Hunter, David Karasik, Douglas P Kiel, Julia A Knight, VeliMatti Kosma, Zoltan Kutalik, Sandra Lai, Diether Lambrechts, Annika Lindblom, Reedik Mgi, Patrik K Magnusson, Arto Mannermaa, Nicholas G Martin, Gisli Masson, Patrick F McArdle, Wendy L McArdle, Mads Melbye Kyriaki Michailidou, Evelin Mihailov, Lili Milani, Roger L Milne, Heli Nevanlinna, Patrick Neven, Ellen A Nohr, Albertine J Oldehinkel, Ben A Oostra, Aarno Palotie,, Munro Peacock, Nancy L Pedersen, Paolo Peterlongo, Julian Peto, Paul DP Pharoah, Dirkje S Postma, Anneli Pouta, Katri Pylks, Paolo Radice, Susan Ring, Fernando Rivadeneira, Antonietta Robino, Lynda M Rose, Anja Rudolph, Veikko Salomaa, Serena Sanna, David Schlessinger, Marjanka K Schmidt, Mellissa C Southey, Ulla Sovio Meir J Stampfer, Doris Stckl Anna M Storniolo, Nicholas J Timpson Jonathan Tyrer, Jenny A Visser, Peter Vollenweider, Henry Vlzke, Gerard Waeber, Melanie Waldenberger, Henri Wallaschofski, Qin Wang, Gonneke Willemsen, Robert Winqvist, Bruce HR Wolffenbuttel, Margaret J Wright, Australian Ovarian Cancer Study The GENICA Network, kConFab, The LifeLines Cohort Study, The InterAct Consortium, Early Growth Genetics (EGG) Consortium, Dorret I Boomsma, Michael J Econs, KayTee Khaw, Ruth JF Loos, Mark I McCarthy, Grant W Montgomery, John P Rice, Elizabeth A Streeten, Unnur Thorsteinsdottir, Cornelia M van Duijn, Behrooz Z Alizadeh, Sven Bergmann, Eric Boerwinkle, Heather A Boyd, Laura Crisponi, Paolo Gasparini, Christian Gieger, Tamara B Harris, Erik Ingelsson, MarjoRiitta Jrvlin, Peter Kraft, Debbie Lawlor, Andres Metspalu, Craig E Pennell, Paul M Ridker, Harold Snieder, Thorkild IA Srensen, Tim D Spector, David P Strachan, Andr G Uitterlinden, Nicholas J Wareham, Elisabeth Widen, Marek Zygmunt, Anna Murray, Douglas F Easton, Kari Stefansson, Joanne M Murabito, Ken K Ong.