

# Beyond SNP Heritability: Polygenicity and Discoverability of Phenotypes Estimated with a Univariate Gaussian Mixture Model

Dominic Holland<sup>a,b,\*</sup>, Oleksandr Frei<sup>d</sup>, Rahul Desikan<sup>c</sup>, Chun-Chieh Fan<sup>a,e,f</sup>, Alexey A. Shadrin<sup>d</sup>, Olav B. Smeland<sup>a,d,g</sup>, V. S. Sundar<sup>a,f</sup>, Paul Thompson<sup>h</sup>, Ole A. Andreassen<sup>d,g</sup>, Anders M. Dale<sup>a,b,f,i</sup>

<sup>a</sup>Center for Multimodal Imaging and Genetics, University of California at San Diego, La Jolla, CA 92037, USA,

<sup>b</sup>Department of Neurosciences, University of California, San Diego, La Jolla, CA 92093, USA,

<sup>c</sup>Department of Radiology, University of California, San Francisco, San Francisco, CA 94158, USA,

<sup>d</sup>NORMENT, KG Jebsen Centre for Psychosis Research, Institute of Clinical Medicine, University of Oslo 0424 Oslo, Norway,

<sup>e</sup>Department of Cognitive Sciences, University of California at San Diego, La Jolla, CA 92093, USA,

<sup>f</sup>Department of Radiology, University of California, San Diego, La Jolla, CA 92093, USA,

<sup>g</sup>Division of Mental Health and Addiction, Oslo University Hospital, 0407 Oslo, Norway,

<sup>h</sup>Keck School of Medicine, University of Southern California, Los Angeles, CA 90033, USA,

<sup>i</sup>Department of Psychiatry, University of California, San Diego, La Jolla, CA 92093, USA,

## Abstract

Of signal interest in the genetics of human traits is estimating their polygenicity (the proportion of causally associated single nucleotide polymorphisms (SNPs)) and the discoverability (or effect size variance) of the causal SNPs. Narrow-sense heritability is proportional to the product of these quantities. We present a basic model, using detailed linkage disequilibrium structure from an extensive reference panel, to estimate these quantities from genome-wide association studies (GWAS) summary statistics for SNPs with minor allele frequency  $>1\%$ . We apply the model to diverse phenotypes and validate the implementation with simulations. We find model polygenicities ranging from  $\simeq 2 \times 10^{-5}$  to  $\simeq 4 \times 10^{-3}$ , with discoverabilities similarly ranging over two orders of magnitude. A power analysis allows us to estimate the proportions of phenotypic variance explained additively by causal SNPs at current sample sizes, and map out sample sizes required to explain larger portions of additive SNP heritability. The model also allows for estimating residual inflation.

**Keywords:** GWAS, Polygenicity, Discoverability, Heritability, Causal SNPs, Effect size, Linkage Disequilibrium

## INTRODUCTION

The genetic components of complex human traits and diseases arise from hundreds to likely many thousands of single nucleotide polymorphisms (SNPs) (Visscher et al., 2012), most of which have weak effects. As sample sizes increase, more of the associated SNPs are identifiable (they reach genome-wide significance), though power for discovery varies widely across phenotypes. Of particular interest are estimating the proportion of SNPs (polygenicity) involved in any particular phenotype; their effective strength of association (discoverability); the proportion of variation in susceptibility, or phenotypic variation, captured additively by all common causal SNPs (approximately, the narrow sense heritability), and the fraction of that captured by genome-wide significant SNPs – all of which are active areas of research (Stahl et al., 2012; Yang et al., 2015; So et al., 2011; Speed et al., 2012; Lee et al., 2011; Yang et al., 2011a; Kumar et al., 2016; Palla and Dudbridge, 2015). However, the effects of population structure (Price et al., 2010), combined with high polygenicity and linkage

disequilibrium (LD), leading to spurious degrees of SNP association, or inflation, considerably complicate matters, and are also areas of much focus (Yang et al., 2011c; Bulik-Sullivan et al., 2015; Kang et al., 2010). Yet, despite recent significant advances, it has been difficult to develop a mathematical model of polygenic architecture based on GWAS that can be used for power estimated across human phenotypes.

Here, in a unified approach explicitly taking into account LD, we present a model relying on genome-wide association studies (GWAS) summary statistics (z-scores for SNP associations with a phenotype (Pasaniuc and Price, 2016)) to estimate polygenicity ( $\pi_1$ ) and discoverability ( $\sigma_\beta^2$ ), as well as elevation of z-scores due to any residual inflation of the z-scores arising from variance distortion induced by cryptic relatedness ( $\sigma_0^2$ ), which remains a concern in large-scale studies (Price et al., 2010). We estimate  $\pi_1$ ,  $\sigma_\beta^2$ , and  $\sigma_0^2$ , by postulating a z-score probability distribution function (pdf) that explicitly depends on them, and fitting it to the actual distribution of GWAS z-scores.

Estimates of polygenicity and discoverability allow one to estimate compound quantities, like narrow-sense heritability captured by the SNPs (Witte et al., 2014); to predict the power of larger-scale GWAS to discover genome-wide significant loci; and to understand why some pheno-

\*Corresponding author:

email: dominic.holland@gmail.com

Phone: 858-822-1776

Fax: 858-534-1078

types have higher power for SNP discovery and proportion of heritability explained than other phenotypes.

In previous work (Holland et al., 2016) we presented a related model that treated the overall effects of LD on z-scores in an approximate way. Here we take the details of LD explicitly into consideration, resulting in a conceptually more basic model to predict the distribution of z-scores. We apply the model to multiple phenotypes, in each case estimating the three model parameters and auxiliary quantities, including the overall inflation factor  $\lambda$ , (traditionally referred to as genomic control (Devlin and Roeder, 1999)) for SNP sets, and narrow sense heritability,  $h^2$ . We also perform extensive simulations on genotypes with realistic LD structure in order to validate the interpretation of the model parameters.

## METHODS

### Overview

Our basic model is a simple postulate for the distribution of causal effects (denoted  $\beta$  below). Our model assumes that only a fraction of all SNPs are in some sense causally related to any given phenotype. We work with a reference panel of approximately 11 million SNPs, and assume that all common causal SNPs ( $\text{MAF} > 0.01$ ) are contained in it. Any given GWAS will have z-scores for a subset of these reference SNPs. When a z-score partially involves a latent causal component (i.e., not pure noise), we assume that it arises through LD with neighboring causal SNPs, or that it itself is causal. Stating the model is straightforward, but solving it is more complicated. We construct a pdf for z-scores that directly follows from the underlying distribution of effects. For any given tag SNP's z-score, it is dependent on the other SNPs the focal SNP is in LD with, taking into account their LD with the focal SNP and their heterozygosity (i.e., it depends not just on the focal tag SNP's total LD and heterozygosity, but also on the distribution of neighboring reference SNPs in LD with it and their heterozygosities). We present two ways of constructing the model pdf for z-scores, using multinomial expansion, and using convolution. The former is more intuitive, but the latter is more numerically tractable and is used here to obtain all reported results. The problem then is finding the three model parameters that give a maximum likelihood best fit for the model's prediction of the distribution of z-scores to the actual distribution of z-scores. Because we are fitting three parameters using  $\gtrsim 10^6$  data points, it is appropriate to incorporate some data reduction to facilitate the computations. To that end, we bin the data (z-scores) into a  $10 \times 10$  grid of heterozygosity-by-total LD (having tested different grid sizes to ensure convergence of results). Also, when building the LD and heterozygosity structures of reference SNPs, we fine-grained the LD range ( $0 \leq r^2 \leq 1$ ), again ensuring that bins were small enough that results were well converged. To fit the model to the data we bin the z-scores (within each heterozygosity/total

LD window) and calculate the multinomial probability for having the actual distribution of z-scores (numbers of z-scores in the z-score bins) given the model pdf for the distribution of z-scores, and adjusting the model parameters using a multidimensional unconstrained nonlinear minimization (Nelder-Mead).

A visual summary of the predicted and actual distribution of z-scores is obtained by making quantile-quantile plots showing, for a wide range of significance thresholds going well beyond genome-wide significance, the proportion (x-axis) of tag SNPs exceeding any given threshold (y-axis) in the range.

With the pdf in hand, various quantities can be calculated: the number of causal SNPs; the expected genetic effect (denoted  $\delta$  below) at the current sample size for a tag SNP given the SNP's z-score and its full LD and heterozygosity structure; the SNP heritability; and the sample size required to explain any percentage of that with genome-wide significant SNPs. The model can easily be extended using a more complex distribution for the underlying  $\beta$ 's, with multiple-component mixtures and incorporating negative selection through both heterozygosity and linkage disequilibrium.

### The Model: Probability Distribution for Z-Scores

To establish notation, we note briefly that in GWAS, ignoring covariates, one traditionally performs simple linear regression in relating genotype to phenotype (Holland et al., 2016; Thompson et al., 2015). That is, assume a linear vector equation (no summation over repeated indices)

$$y = g_i \beta_i + e_i \quad (1)$$

for phenotype vector  $y$  over  $N$  samples (mean-centered and normalized to unit variance), mean-centered genotype vector  $g_i$  for the  $i^{\text{th}}$  of  $n$  SNPs, true effect (regression coefficient)  $\beta_i$ , and residual vector  $e_i$  containing the effects of all the other causal SNPs, the independent environmental component, and random error. For SNP  $i$ , the estimated simple linear regression coefficient is

$$\hat{\beta}_i = g_i^T y / (g_i^T g_i) = \text{cov}(g_i, y) / \text{var}(g_i), \quad (2)$$

where  $T$  denotes transpose and  $g_i^T g_i / N = \text{var}(g_i) = H_i$  is the SNP's heterozygosity (frequency of the heterozygous genotype):  $H_i = 2p_i(1 - p_i)$  where  $p_i$  is the frequency of either of the SNP's alleles.

Consistent with the work of others (Yang et al., 2011c; Zeng et al., 2018), we assume the causal SNPs are distributed randomly throughout the genome (an assumption that can be relaxed when explicitly considering different SNP categories, but that in the main is consistent with the additive variation explained by a given part of the genome being proportional to the length of DNA (Yang et al., 2011b)), and that their  $\beta$  coefficients in the GWAS framework (and in the absence of inflation due to population structure) are distributed normally with variance

given by a constant,  $\sigma_{\beta}^2$ :

$$\beta \sim \mathcal{N}(0, \sigma_{\beta}^2) \quad (3)$$

Taking into account all SNPs (the remaining ones are all null by definition), this is equivalent to the two-component Gaussian mixture model

$$\beta \sim \pi_1 \mathcal{N}(0, \sigma_{\beta}^2) + (1 - \pi_1) \mathcal{N}(0, 0) \quad (4)$$

where  $\mathcal{N}(0, 0)$  is the Dirac delta function, so that considering all SNPs, the net variance is  $\text{var}(\beta) = \pi_1 \sigma_{\beta}^2$ . If there is no LD, the association z-scores for SNPs can be decomposed into an effect  $\delta$  and a residual environment and error term,  $\epsilon$ , which is assumed to be independent of  $\delta$ , and which in the absence of inflation is  $\epsilon \sim \mathcal{N}(0, 1)$  (Holland et al., 2016):

$$z = \delta + \epsilon \quad (5)$$

with

$$\delta = \sqrt{NH} \beta \quad (6)$$

so that

$$\begin{aligned} \text{var}(z) &= \text{var}(\delta) + \text{var}(\epsilon) \\ &\equiv \sigma^2 + 1 \end{aligned} \quad (7)$$

where

$$\sigma^2 = \sigma_{\beta}^2 NH. \quad (8)$$

By construction, when there is no genetic effect,  $\delta = 0$ , so that  $\text{var}(\epsilon) = 1$  unless the z-scores are elevated by variance distortion due to cryptic relatedness in the sample.

If there is no cryptic relatedness in the sample, but the sample is composed of two or more subpopulations with different allele frequencies for a subset of markers, the marginal distribution of an individual's genotype at any of those markers will be inflated. This situation describes pure population stratification in the sample (Laird and Lange, 2010). The squared z-score for such a marker will then follow a noncentral chi-square distribution; the noncentrality parameter will contain the causal genetic effect, if any, but biased up or down (confounding or loss of power, depending on the relative sign of the genetic effect and the bias term), and will be inversely proportional to the inflated genotype variance. Thus, the effect of stratification is nontrivial; we assume in this study that this effect has largely been accounted for using standard methods (Wu et al., 2011).

Pure cryptic relatedness in the sample (drawn from a population mixture with at least one subpopulation with identical-by-descent marker alleles, but no population stratification), leads to variance distortion in the distribution of z-scores (Devlin and Roeder, 1999). If  $z_u$  denotes the uninflated z-scores, then the inflated z-scores are  $z = \sigma_0 z_u$ , where  $\sigma_0 \geq 1$  characterizes the inflation. Thus, from Eq. 7, in the presence of inflation due to cryptic relatedness

$$\begin{aligned} \text{var}(z) &= \sigma_0^2 (\sigma^2 + 1) \\ &\equiv \tilde{\sigma}^2 + \sigma_0^2 \\ &\equiv \tilde{\sigma}_{\beta}^2 NH + \sigma_0^2 \end{aligned} \quad (9)$$

where  $\tilde{\sigma}_{\beta}^2 \equiv \sigma_0^2 \sigma_{\beta}^2$ , so that with inflation  $\text{var}(\delta) = \tilde{\sigma}^2$  and  $\epsilon \sim \mathcal{N}(0, \sigma_0^2)$ . In the presence of inflation one is dealing with inflated random variables  $\tilde{\beta} \sim \mathcal{N}(0, \tilde{\sigma}_{\beta}^2)$ , but we will drop the tilde on the  $\beta$ 's in what follows.

Implicit in Eq. 8 is approximating the denominator,  $1 - q^2$ , of the  $\chi^2$  statistic noncentrality parameter to be 1, where  $q^2$  is the proportion of phenotypic variance explained by the causal variant, i.e.,  $q \equiv \sqrt{H} \beta$ . So a more correct  $\delta$  is

$$\delta = \sqrt{N} q / \sqrt{1 - q^2}. \quad (10)$$

Taylor expanding in  $q$  and then taking the variance gives

$$\text{var}(\delta) = \sigma_{\beta}^2 NH [1 + (15/4) \sigma_{\beta}^4 H^2 + O(\sigma_{\beta}^8 H^4)]. \quad (11)$$

The additional terms will be vanishingly small and so do not contribute in a distributional sense; (quasi-) Mendelian or outlier genetic effects represent an extreme scenario where the model is not expected to be accurate, but SNPs for such traits are by definition easily detectable. So Eq. 8 remains valid for the polygenicity of complex traits.

Now consider the effects of LD on z-scores. The simple linear regression coefficient estimate for tag SNP  $i$ ,  $\hat{\beta}_i$ , and hence the GWAS z-score, implicitly incorporates contributions due to LD with neighboring causal SNPs. (A tag SNP is a SNP with a z-score, imputed or otherwise; generally these will compose a smaller set than that available in reference panels like 1000 Genomes used here for calculating the LD structure of tag SNPs.) In Eq. 1,  $e_i = \sum_{j \neq i} g_j \beta_j + \epsilon$ , where  $g_j$  is the genotype vector for SNP  $j$ ,  $\beta_j$  is its true regression coefficient, and  $\epsilon$  is the independent true environmental and error residual vector (over the  $N$  samples). Thus, explicitly including all causal true  $\beta$ 's, Eq. 2 becomes

$$\begin{aligned} \hat{\beta}_i &= \frac{\sum_j g_i^T g_j \beta_j}{NH_i} + \frac{g_i^T \epsilon}{g_i^T g_i} \\ &\equiv \beta'_i + \epsilon'_i \end{aligned} \quad (12)$$

(the sum over  $j$  now includes SNP  $i$  itself). This is the simple linear regression expansion of the estimated regression coefficient for SNP  $i$  in terms of the independent latent (true) causal effects and the latent environmental (plus error) component;  $\beta'_i$  is the effective simple linear regression expression for the true genetic effect of SNP  $i$ , with contributions from neighboring causal SNPs mediated by LD. Note that  $g_i^T g_j / N$  is simply  $\text{cov}(g_i, g_j)$ , the covariance between genotypes for SNPs  $i$  and  $j$ . Since correlation is covariance normalized by the variances,  $\beta'_i$  in Eq. 12 can be written as

$$\beta'_i = \sum_j \sqrt{\frac{H_j}{H_i}} r_{ij} \beta_j. \quad (13)$$

where  $r_{ij}$  is the correlation between genotypes at SNP  $j$  and tag SNP  $i$ . Then, from Eq. 5, the z-score for the tag

SNP's association with the phenotype is given by:

$$\begin{aligned} z_i &= \sqrt{NH_i}\beta'_i + \epsilon_i \\ &= \sqrt{N} \sum_j \sqrt{H_j} r_{ij} \beta_j + \epsilon_i. \end{aligned} \quad (14)$$

We noted that in the absence of LD, the distribution of the residual in Eq. 5 is assumed to be univariate normal. But in the presence of LD (Eq. 14) there are induced correlations, so the appropriate extension would be multivariate normal for  $\epsilon_i$  (Zhu and Stephens, 2017). A limitation of the present work is that we do not consider this complexity. This may account for the relatively minor misfit in the simulation results for cases of high polygenicity – see below.

Thus, for example, if the SNP itself is not causal but is in LD with  $k$  causal SNPs that all have heterozygosity  $H$ , and where its LD with each of these is the same, given by some value  $r^2$  ( $0 < r^2 \leq 1$ ), then  $\tilde{\sigma}^2$  in Eq. 9 will be given by

$$\tilde{\sigma}^2 = kr^2\tilde{\sigma}_\beta^2NH. \quad (15)$$

For this idealized case, the marginal distribution, or pdf, of z-scores for a set of such associated SNPs is

$$f_1(z; N, \mathcal{H}, \sigma_\beta, \sigma_0) = \phi(z; 0, kr^2\tilde{\sigma}_\beta^2NH + \sigma_0^2) \quad (16)$$

where  $\phi(\cdot; \mu, \sigma^2)$  is the normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and  $\mathcal{H}$  is shorthand for the LD and heterozygosity structure of such SNPs (in this case, denoting exactly  $k$  causals with LD given by  $r^2$  and heterozygosity given by  $H$ ). If a proportion  $\alpha$  of all tag SNPs are similarly associated with the phenotype while the remaining proportion are all null (not causal and not in LD with causal SNPs), then the marginal distribution for all SNP z-scores is the Gaussian mixture

$$f(z) = (1 - \alpha)\phi(z; 0, \sigma_0^2) + \alpha f_1(z), \quad (17)$$

dropping the parameters for convenience.

For real genotypes, however, the LD and heterozygosity structure is far more complicated, and of course the causal SNPs are generally numerous and unknown. Thus, more generally, for each tag SNP  $\mathcal{H}$  will be a two-dimensional histogram over LD ( $r^2$ ) and heterozygosity ( $H$ ), each grid element giving the number of SNPs falling within the edges of that  $(r^2, H)$  bin. Alternatively, for each tag SNP it can be built as two one-dimensional histograms, one giving the LD structure (counts of neighboring SNPs in each LD  $r^2$  bin), and the other giving, for each  $r^2$  bin, the mean heterozygosity for those neighboring SNPs, which should be accurate for sufficiently fine binning. We use the latter in what follows. We present two consistent ways of expressing the *a posteriori* pdf for z-scores, based on multinomial expansion and on convolution, that provide complementary views. The multinomial approach perhaps gives a more intuitive feel for the problem, but the convolution approach is considerably more tractable numerically and is used here to obtain all reporter results.

## Model PDF: Multinomial Expansion

As in our previous work, we incorporate the model parameter  $\pi_1$  for the fraction of all SNPs that are causal (Holland et al., 2016). Additionally, we calculate the actual LD and heterozygosity structure for each SNP. That is, for each SNP we build a histogram of the numbers of other SNPs in LD with it for  $w$  equally-spaced  $r^2$ -windows between  $r_{min}^2$  and 1 where  $r_{min}^2 = 0.05$  (approximately the noise floor for correlation when LD is calculated from the 503 samples in 1000 Genomes), and record the mean heterozygosity for each bin; as noted above, we use  $\mathcal{H}$  as shorthand to represent all this. We find that  $w \simeq 20$  is sufficient for converged results. For any given SNP, the set of SNPs thus determined to be in LD with it constitute its LD block, with their number given by  $n$  (LD with self is always 1, so  $n$  is at least 1). The pdf for z-scores, given  $N, \mathcal{H}$ , and the three model parameters  $\pi_1, \sigma_\beta, \sigma_0$ , will then be given by the sum of Gaussians that are generalizations of Eq. 16 for different combinations of numbers of causals among the  $w$  LD windows, each Gaussian scaled by the probability of the corresponding combination of causals among the LD windows, i.e., by the appropriate multinomial distribution term.

For  $w$   $r^2$ -windows, we must consider the possibilities where the tag SNP is in LD with all possible numbers of causal SNPs in each of these windows, or any combination thereof. There are thus  $w + 1$  categories of SNPs: null SNPs (which  $r^2$ -windows they are in is irrelevant), and causal SNPs, where it does matter which  $r^2$ -windows they reside in. If window  $i$  has  $n_i$  SNPs ( $\sum_{i=1}^w n_i = n$ ) and mean heterozygosity  $H_i$ , and the overall fraction of SNPs that are causal is  $\pi_1$ , then the probability of having simultaneously  $k_0$  null SNPs,  $k_1$  causal SNPs in window 1, and so on through  $k_w$  causal SNPs in window  $w$ , for a nominal total of  $K$  causals ( $\sum_{i=1}^w k_i = K$  and  $k_0 = n - K$ ), is given by the multinomial distribution, which we denote  $M(k_0, \dots, k_w; n_0, \dots, n_w; \pi_1)$ . For an LD block of  $n$  SNPs, the prior probability,  $p_i$ , for a SNP to be causal and in window  $i$  is the product of the independent prior probabilities of a SNP being causal and being in window  $i$ :  $p_i = \pi_1 n_i / n$ . The prior probability of being null (regardless of  $r^2$ -window) is simply  $p_0 = (1 - \pi_1)$ . The probability of a given breakdown  $k_0, \dots, k_w$  of the neighboring SNPs into the  $w + 1$  categories is then given by

$$M(k_0, \dots, k_w; n_0, \dots, n_w; \pi_1) = \frac{n!}{k_0! \dots k_w!} p_0^{k_0} \dots p_w^{k_w} \quad (18)$$

and the corresponding Gaussian is

$$\phi(z; 0, (k_1 H_1 r_1^2 + \dots + k_w H_w r_w^2) \tilde{\sigma}_\beta^2 N + \sigma_0^2). \quad (19)$$

For a SNP with LD and heterozygosity structure  $\mathcal{H}$ , the pdf for its z-score, given  $N$  and the model parameters, is then given by summing over all possible numbers of total causals in LD with the SNP, and all possible distributions

of those causals among the  $w$   $r^2$ -windows:

$$\text{pdf}(z; N, \mathcal{H}, \pi_1, \sigma_\beta, \sigma_0) = \sum_{K=0}^{K_{max}} \sum_{k_1, \dots, k_w} \frac{n!}{k_0! \dots k_w!} p_0^{k_0} \dots p_w^{k_w} \times \phi(z; 0, (k_1 H_1 r_1^2 + \dots + k_w H_w r_w^2) \tilde{\sigma}_\beta^2 N + \sigma_0^2), \quad (20)$$

where  $K_{max}$  is bounded above by  $n$ . Note again that  $\mathcal{H}$  is shorthand for the heterozygosity and linkage-disequilibrium structure of the SNP, giving the set  $\{n_i\}$  (as well as  $\{H_i\}$ ), and hence, for a given  $\pi_1, p_i$ . Also there is the constraint  $\sum_{i=1}^w k_i = K$  on the second summation, and, for all  $i$ ,  $\max(k_i) = \max(K, n_i)$ , though generally  $K_{max} \ll n_i$ . The number of ways of dividing  $K$  causal SNPs amongst  $w$  LD windows is given by the binomial coefficient  $\binom{a}{b}$ , where  $a \equiv K + w - 1$  and  $b \equiv w - 1$ , so the number of terms in the second summation grows rapidly with  $K$  and  $w$ . However, because  $\pi_1$  is small (often  $\leq 10^{-3}$ ), the upper bound on the first summation over total number of potential causals  $K$  in the LD block for the SNP can be limited to  $K_{max} < \min(20, n)$ , even for large blocks with  $n \simeq 10^3$ . That is,

$$\sum_{K=0}^{K_{max}} \sum_{k_1, \dots, k_w} M(k_0, \dots, k_w; n_0, \dots, n_w; \pi_1) \simeq 1. \quad (21)$$

Still, the number of terms is large; e.g., for  $K = 10$  and  $w = 10$  there are 92,378 terms.

### Model PDF: Convolution

From Eq. 14, there is an alternative and efficient procedure that allows for accurate calculation of a z-score's *a posteriori* pdf (given the SNP's heterozygosity and LD structure, and the phenotype's model parameters). Any GWAS z-score is a sum of unobserved random variables (LD-mediated contributions from neighboring causal SNPs, and the additive environmental component), and the pdf for such a composite random variable is given by the convolution of the pdfs for the component random variables. Since convolution is associative, and the Fourier transform of the convolution of two functions is just the product of the individual Fourier transforms of the two functions, one can obtain the *a posteriori* pdf for z-scores as the inverse Fourier transform of the product of the Fourier transforms of the individual random variable components.

From Eq. 14  $z$  is a sum of correlation- and heterozygosity-weighted random variables  $\{\beta_j\}$  and the random variable  $\epsilon$ , where  $\{\beta_j\}$  denotes the set of true causal parameters for each of the SNPs in LD with the tag SNP whose z-score is under consideration. The Fourier transform  $F(k)$  of a Gaussian  $f(x) = c \times \exp(-ax^2)$  is  $F(k) = c\sqrt{\pi/a} \times \exp(-\pi^2 k^2/a)$ . From Eq. 4, for each SNP  $j$  in LD with the tag SNP ( $1 \leq j \leq b$ , where  $b$  is the tag SNP's block size),

$$\sqrt{N H_j} r_j \beta_j \sim \pi_1 \mathcal{N}(0, N H_j r_j^2 \tilde{\sigma}_\beta^2) + (1 - \pi_1) \mathcal{N}(0, 0). \quad (22)$$

The Fourier transform (with variable  $k$  – see below) of the first term on the right hand side is

$$F(k) = \pi_1 \exp(-2\pi^2 k^2 N H_j r_j^2 \tilde{\sigma}_\beta^2), \quad (23)$$

while that of the second term is simply  $(1 - \pi_1)$ . Additionally, the environmental term is  $\epsilon \sim \mathcal{N}(0, \sigma_0^2)$  (ignoring LD-induced correlation, as noted earlier), and its Fourier transform is  $\exp(-2\pi^2 \sigma_0^2 k^2)$ . For each tag SNP, one could construct the *a posteriori* pdf based on these Fourier transforms. However, it is more practical to use a coarse-grained representation of the data. Thus, in order to fit the model to a data set, we bin the tag SNPs whose z-scores comprise the data set into a two-dimensional heterozygosity / total LD grid (whose elements we denote “H-L” bins), and fit the model with respect to this coarse gridding instead of with respect to every individual tag SNP z-score; in the section “Parameter Estimation” below we describe using a  $10 \times 10$  grid. Additionally, for each H-L bin the LD  $r^2$  and heterozygosity histogram structure for each tag SNP is built, using  $w_{max}$  equally-spaced  $r^2$  bins for  $r_{min}^2 \leq r^2 \leq 1$ ;  $w_{max} = 20$  is large enough to allow for converged results;  $r_{min}^2 = 0.05$  is generally small enough to capture true causal associations in weak LD while large enough to exclude spurious contributions to the pdf arising from estimates of  $r^2$  that are non-zero due to noise. This points up a minor limitation of the model stemming from the small reference sample size ( $N_R = 503$  for 1000 Genomes) from which  $\mathcal{H}$  is built. Larger  $N_R$  would allow for more precision in handling very low LD ( $r^2 < 0.05$ ), but this is an issue only for situations with extremely large  $\sigma_\beta^2$  (high heritability with low polygenicity) that we do not encounter for the 12 phenotypes we analyze here. In any case, this can be calibrated for using simulations.

For any H-L bin with mean heterozygosity  $H$  and mean total LD  $L$  there will be an average LD and heterozygosity structure with a mean breakdown for the tag SNPs having  $n_w$  SNPs (not all of which necessarily are tag SNPs) with LD  $r^2$  in the  $w^{\text{th}}$   $r^2$  bin whose average heterozygosity is  $H_w$ . Thus, one can re-express z-scores for an H-L bin as

$$z = \sqrt{N} \sum_{w=1}^{w_{max}} \left( \sqrt{H_w} r_w \sum_{j=0}^{n_w} \beta_j \right) + \epsilon \quad (24)$$

where  $\beta_j$  and  $\epsilon$  are unobserved random variables.

In the spirit of the discrete Fourier transform (DFT), discretize the set of possible z-scores into the ordered set of  $n$  (equal to a power of 2) values  $z_1, \dots, z_n$  with equal spacing between neighbors given by  $\Delta z$  ( $z_n = -z_1 - \Delta z$ , and  $z_{n/2+1} = 0$ ). Taking  $z_1 = -38$  allows for the minimum p-values of  $5.8 \times 10^{-316}$  (near the numerical limit); with  $n = 2^{10}$ ,  $\Delta z = 0.0742$ . Given  $\Delta z$ , the Nyquist critical frequency is  $f_c = \frac{1}{2\Delta z}$ , so we consider the Fourier transform function for the z-score pdf at  $n$  discrete values  $k_1, \dots, k_n$ , with equal spacing between neighbors given by  $\Delta k$ , where  $k_1 = -f_c$  ( $k_n = -k_1 - \Delta k$ , and  $k_{n/2+1} = 0$ ; the DFT pair

$\Delta z$  and  $\Delta k$  are related by  $\Delta z \Delta k = 1/n$ ). Define

$$A_w \equiv -2\pi^2 N H_w r_w^2 \tilde{\sigma}_\beta^2. \quad (25)$$

(see Eq. 23). Then the product (over  $r^2$  bins) of Fourier transforms for the genetic contribution to z-scores, denoted  $G_j \equiv G(k_j)$ , is

$$G(k_j) = \prod_{w=1}^{w_{max}} (\pi_1 \exp(A_w k_j^2) + (1 - \pi_1))^{n_w}, \quad (26)$$

and the Fourier transform of the environmental contribution, denoted  $E_j \equiv E(k_j)$ , is

$$E(k_j) = \exp(-2\pi^2 \sigma_0^2 k_j^2). \quad (27)$$

Let  $\mathbf{F}_z = (G_1 E_1, \dots, G_n E_n)$  denote the vector of products of Fourier transform values, and let  $\mathcal{F}^{-1}$  denote the inverse Fourier transform operator. Then, the vector of pdf values for z-score bins (indexed by  $i$ ) in the H-L bin with mean LD and heterozygosity structure  $\mathcal{H}$ ,  $\mathbf{pdf}_z = (f_1, \dots, f_n)$  where  $f_i \equiv \text{pdf}(z_i|\mathcal{H})$ , is

$$\mathbf{pdf}_z = \mathcal{F}^{-1}[\mathbf{F}_z]. \quad (28)$$

## Data Preparation

For real phenotypes, we calculated SNP minor allele frequency (MAF) and LD between SNPs using the 1000 Genomes phase 3 data set for 503 subjects/samples of European ancestry (Consortium et al., 2015, 2012; Sveinbjornsson et al., 2016). For simulations, we used HAPGEN2 (Li and Stephens, 2003; Spencer et al., 2009; Su et al., 2011) to generate genotypes; we calculated SNP MAF and LD structure from 1000 simulated samples. We elected to use the same intersecting set of SNPs for real data and simulation. For HAPGEN2, we eliminated SNPs with  $\text{maf} < 0.002$ ; for 1000 Genomes, we eliminated SNPs for which the call rate (percentage of samples with useful data) was less than 90%. This left  $n_{\text{SNP}} = 11,015,833$  SNPs.

Sequentially moving through each chromosome in contiguous blocks of 5,000 SNPs, for each SNP in the block we calculated its Pearson  $r^2$  correlation coefficients with all SNPs in the central block itself and with all SNPs in the pair of flanking blocks of size up to 50,000 each. For each SNP we calculated its total LD (TLD), given by the sum of LD  $r^2$ 's thresholded such that if  $r^2 < r_{\text{min}}^2$  we set that  $r^2$  to zero. For each SNP we also built a histogram giving the numbers of SNPs in  $w_{\text{max}}$  equally-spaced  $r^2$ -windows covering the range  $r_{\text{min}}^2 \leq r^2 \leq 1$ . These steps were carried out independently for 1000 Genomes phase 3 and for HAPGEN2 (for the latter, we used 1000 simulated samples).

Employing a similar procedure, we also built binary (logical) LD matrices identifying all pairs of SNPs for which LD  $r^2 > 0.8$ , a liberal threshold for SNPs being ‘‘synonymous’’.

In applying the model to summary statistics, we calculated histograms of TLD and LD block size (using 100

bins in both cases) and ignoring SNPs whose TLD or block size was so large that their frequency was less than a hundredth of the respective histogram peak; typically this amounted to restricting to SNPs for which TLD  $\leq 600$  and LD block size  $\leq 1,500$ . We also ignored SNPs for which MAF  $\leq 0.01$ ,

We analyzed summary statistics for twelve phenotypes (in what follows, where sample sizes varied by SNP, we quote the median value): (1) bipolar disorder ( $N_{\text{cases}} = 20,352$ ,  $N_{\text{controls}} = 31,358$ ) (Stahl et al., 2018); (2) schizophrenia ( $N_{\text{cases}} = 35,476$ ,  $N_{\text{controls}} = 46,839$ ) (Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2014); (3) major depressive disorder ( $N_{\text{cases}} = 59,851$ ,  $N_{\text{controls}} = 113,154$ ) (Wray and Sullivan, 2017); (4) late onset Alzheimer’s disease (LOAD;  $N_{\text{cases}} = 17,008$ ,  $N_{\text{controls}} = 37,154$ ) (Lambert et al., 2013) (in the Supplementary Material we present results for a more recent GWAS with  $N_{\text{cases}} = 71,880$  and  $N_{\text{controls}} = 383,378$  (Jansen et al., 2018)); (5) Crohn’s disease ( $N_{\text{cases}} = 12,194$ ,  $N_{\text{controls}} = 34,915$ ); (6) ulcerative colitis ( $N_{\text{cases}} = 12,366$ ,  $N_{\text{controls}} = 34,915$ ) (de Lange et al., 2017); (7) number of years of formal education ( $N = 293,723$ ) (Okbay et al., 2016); (8) intelligence ( $N = 262,529$ ) (Sniekers et al., 2017; Savage et al., 2018); (9) low- ( $N = 89,873$ ) and (10) high-density lipoprotein ( $N = 94,295$ ) (Willer et al., 2013); (11) height ( $N = 707,868$ ) (Yengo et al., 2018; Wood et al., 2014); and (12) putamen volume (normalized by intracranial volume,  $N = 11,598$ ) (Hibar et al., 2015). Most participants were of European ancestry.

In Supplementary Material we show the results for our analysis of four additional phenotypes: amyotrophic lateral sclerosis (ALS, restricted to chromosome 9) (Van Rheenen et al., 2016), body mass index (BMI) (Locke et al., 2015), coronary artery disease (CAD) (Nikpay et al., 2015), and total cholesterol (TC) (Willer et al., 2013).

For schizophrenia, for example, there were 6,610,991 SNPs with finite z-scores out of the 11,015,833 SNPs from the 1000 Genomes reference panel that underlie the model; the genomic control factor for these SNPs was  $\lambda_{GC} = 1.466$ . Of these, 314,857 were filtered out due to low maf or very large LD block size. The genomic control factor for the remaining SNPs was  $\lambda_{GC} = 1.468$ ; for the pruned subsets, with  $\simeq 1.49 \times 10^6$  SNPs each, it was  $\lambda = 1.30$ . (Note that genomic control values for pruned data are always lower than for unpruned data.)

A limitation in the current work is that we have not taken account of imputation inaccuracy, where lower MAF SNPs are, through lower LD, less certain. Thus, the effects from lower MAF causal variants will noisier than for higher MAF variants.

## Simulations

We generated genotypes for  $10^5$  unrelated simulated samples using HAPGEN2 (Su et al., 2011). For narrow-sense heritability  $h^2$  equal to 0.1, 0.4, and 0.7, we considered polygenicity  $\pi_1$  equal to  $10^{-5}$ ,  $10^{-4}$ ,  $10^{-3}$ , and  $10^{-2}$ . For each of these 12 combinations, we randomly selected  $n_{\text{causal}}$

$= \pi_1 \times n_{snp}$  “causal” SNPs and assigned them  $\beta$ -values drawn from the standard normal distribution (i.e., independent of  $H$ ), with all other SNPs having  $\beta = 0$ . We repeated this ten times, giving ten independent instantiations of random vectors of  $\beta$ 's. Defining  $Y_G = G\beta$ , where  $G$  is the genotype matrix and  $\beta$  here is the vector of true coefficients over all SNPs, the total phenotype vector is constructed as  $Y = Y_G + \varepsilon$ , where the residual random vector  $\varepsilon$  for each instantiation is drawn from a normal distribution such that  $h^2 = \text{var}(Y_G)/\text{var}(Y)$ . For each of the instantiations this implicitly defines the “true” value  $\sigma_\beta^2$ .

The sample simple linear regression slope,  $\hat{\beta}$ , and the Pearson correlation coefficient,  $\hat{r}$ , are assumed to be t-distributed. These quantities have the same t-value:  $t = \hat{\beta}/\text{se}(\hat{\beta}) = \hat{r}/\text{se}(\hat{r}) = \hat{r}\sqrt{N-2}/\sqrt{1-\hat{r}^2}$ , with corresponding p-value from Student's  $t$  cumulative distribution function (cdf) with  $N-2$  degrees of freedom:  $p = 2 \times \text{tcdf}(-|t|, N-2)$  (see Supplementary Material). Since we are not here dealing with covariates, we calculated  $p$  from correlation, which is slightly faster than from estimating the regression coefficient. The t-value can be transformed to a z-value, giving the z-score for this  $p$ :  $z = -\Phi^{-1}(p/2) \times \text{sign}(\hat{r})$ , where  $\Phi$  is the normal cdf ( $z$  and  $t$  have the same p-value).

### Parameter Estimation

We randomly pruned SNPs using the threshold  $r^2 > 0.8$  to identify “synonymous” SNPs, performing ten such iterations. That is, for each of ten iterations, we randomly selected a SNP (not necessarily the one with largest z-score) to represent each subset of synonymous SNPs. For schizophrenia, for example, pruning resulted in approximately 1.3 million SNPs in each iteration.

The postulated pdf for a SNP's z-score depends on the SNP's LD and heterozygosity structure (histogram),  $\mathcal{H}$ . Given the data – the set of z-scores for available SNPs, as well as their LD and heterozygosity structure – and the  $\mathcal{H}$ -dependent pdf for z-scores, the objective is to find the model parameters that best predict the distribution of z-scores. We bin the SNPs with respect to a grid of heterozygosity and total LD; for any given H-L bin there will be a range of z-scores whose distribution the model it intended to predict. We find that a  $10 \times 10$ -grid of equally spaced bins is adequate for converged results. (Using equally-spaced bins might seem inefficient because of the resulting very uneven distribution of z-scores among grid elements – for example, orders of magnitude more SNPs in grid elements with low total LD compared with high total LD. However, the objective is to model the effects of H and L: using variable grid element sizes so as to maximize balance of SNP counts among grid elements means that the true H- and L-mediated effects of the SNPs in a narrow range of H and L get subsumed with the effects of many more SNPs in a much wider range of H and L – a misspecification of the pdf leading to some inaccuracy.) In lieu of or in addition to total LD (L) binning, one can bin SNPs with respect to their total LD block size (total number of

SNPs in LD, ranging from 1 to  $\sim 1,500$ ).

To find the model parameters that best fit the data, for a given H-L bin we binned the selected SNPs z-scores into equally-spaced bins of width  $dz=0.0742$  (between  $z_{min}=-38$  and  $z_{max}=38$ , allowing for p-values near the numerical limit of  $10^{-316}$ ), and from Eq. 28 calculated the probability for z-scores to be in each of those z-score bins (the prior probability for “success” in each z-score bin). Then, knowing the actual numbers of z-scores (numbers of “successes”) in each z-score bin, we calculated the multinomial probability,  $p_m$ , for this outcome. The optimal model parameter values will be those that maximize the accrual of this probability over all H-L bins. We constructed a cost function by calculating, for a given H-L bin,  $-\ln(p_m)$  and averaging over prunings, and then accumulating this over all H-L bins. Model parameters minimizing the cost were obtained from Nelder-Mead multidimensional unconstrained nonlinear minimization of the cost function, using the Matlab function `fminsearch()`.

### Posterior Effect Sizes

Model posterior effect sizes, given  $z$  (along with  $N$ ,  $\mathcal{H}$ , and the model parameters), were calculated using numerical integration over the random variable  $\delta$ :

$$\begin{aligned} \delta_{expected} &\equiv E(\delta|z) = \int P(\delta|z)\delta d\delta \\ &= \frac{1}{P(z)} \int P(z|\delta)P(\delta)d\delta. \end{aligned} \quad (29)$$

Here, since  $z|\delta \sim \mathcal{N}(\delta, \sigma_0^2)$ , the posterior probability of  $z$  given  $\delta$  is simply

$$P(z|\delta) = \phi(z; \delta, \sigma_0^2). \quad (30)$$

$P(z)$  is shorthand for pdf( $z|N, \mathcal{H}, \pi_1, \sigma_\beta, \sigma_0$ ), given by Eq. 28.  $P(\delta)$  is calculated by a similar procedure that lead to Eq. 28 but ignoring the environmental contributions  $\{E_j\}$ . Specifically, let  $\mathbf{F}_\delta = (G_1, \dots, G_n)$  denote the vector of products of Fourier transform values. Then, the vector of pdf values for genetic effect bins (indexed by  $i$ ; numerically, these will be the same as the z-score bins) in the H-L bin,  $\mathbf{pdf}_\delta = (f_1, \dots, f_n)$  where  $f_i \equiv \text{pdf}(\delta_i|\mathcal{H})$ , is

$$\mathbf{pdf}_\delta = \mathcal{F}^{-1}[\mathbf{F}_\delta]. \quad (31)$$

Similarly,

$$\begin{aligned} \delta_{expected}^2 &\equiv E(\delta^2|z) = \int P(\delta|z)\delta^2 d\delta \\ &= \frac{1}{P(z)} \int P(z|\delta)P(\delta)\delta^2 d\delta, \end{aligned} \quad (32)$$

which is used in power calculations.

### GWAS Power

Chip heritability,  $h_{SNP}^2$ , is the proportion of phenotypic variance that in principle can be captured additively by the  $n_{snp}$  SNPs under study (Witte et al., 2014). It is

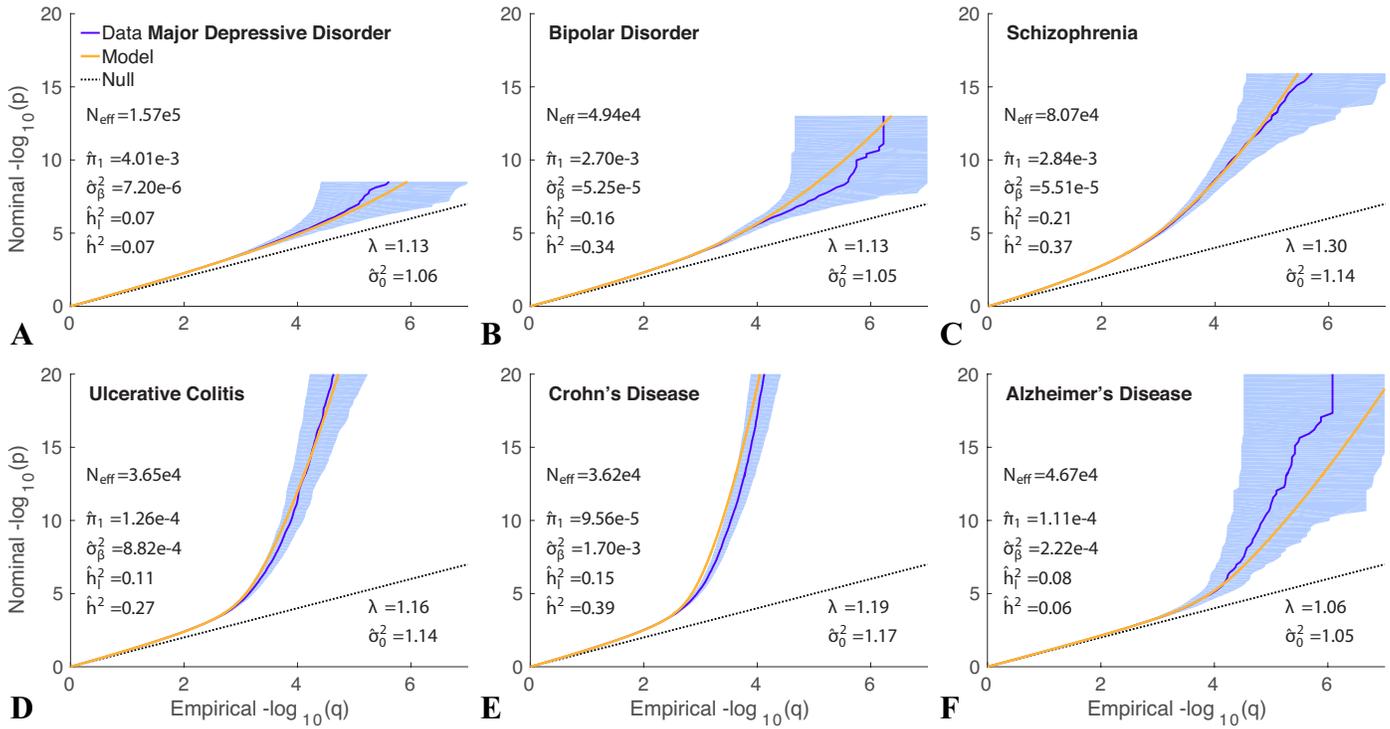


Figure 1: QQ plots of (pruned) z-scores for qualitative phenotypes (dark blue, 95% confidence interval in light blue) with model prediction (yellow): (A) major depressive disorder, (B) bipolar disorder, (C) schizophrenia, (D) ulcerative colitis, (E) Crohn’s disease, and (F) late onset Alzheimer’s disease (excluding APOE; see also Supplementary Material Fig. 10). The dashed line is the expected QQ plot under null (no SNPs associated with the phenotype).  $p$  is a nominal p-value for z-scores, and  $q$  is the proportion of z-scores with p-values exceeding that threshold.  $\lambda$  is the overall nominal genomic control factor for the pruned data (which is accurately predicted by the model in all cases). The three estimated model parameters are: polygenicity,  $\hat{\pi}_1$ ; discoverability,  $\hat{\sigma}_\beta^2$  (corrected for inflation); and SNP association  $\chi^2$ -statistic inflation factor,  $\hat{\sigma}_0^2$ .  $\hat{h}^2$  is the estimated narrow-sense chip heritability, re-expressed as  $h_l^2$  on the liability scale for these case-control conditions assuming a prevalence of: MDD 6.7% (NIMH, 2016), BIP 0.5% (Merikangas et al., 2011), SCZ 1% (Speed et al., 2017), UC 0.1% (Burisch et al., 2013), CD 0.1% (Burisch et al., 2013), AD 14% (for people aged 71 and older in the USA (Plassman et al., 2007; Alzheimer’s Association, 2018)). The estimated number of causal SNPs is given by  $\hat{n}_{causal} = \hat{\pi}_1 n_{snp}$  where  $n_{snp} = 11,015,833$  is the total number of SNPs, whose LD structure and MAF underlie the model; the GWAS z-scores are for subsets of these SNPs.  $N_{eff}$  is the effective case-control sample size – see text. Reading the plots: on the vertical axis, choose a p-value threshold (more extreme values are further from the origin), then the horizontal axis gives the proportion of SNPs exceeding that threshold (higher proportions are closer to the origin). Numerical values for the model parameters are also given in Table 1. See also Supplementary Material Figs. 13-19.

of interest to estimate the proportion of  $h_{SNP}^2$  that can be explained by SNPs reaching genome-wide significance,  $p \leq 5 \times 10^{-8}$  (i.e., for which  $|z| > z_t = 5.45$ ), at a given sample size (Pe’er et al., 2008; McCarthy et al., 2008). In Eq 1, for SNP  $i$  with genotype vector  $g_i$  over  $N$  samples, let  $y_{g_i} \equiv g_i \beta_i$ . If the SNP’s heterozygosity is  $H_i$ , then  $\text{var}(y_{g_i}) = \beta_i^2 H_i$ . If we knew the full set  $\{\beta_i\}$  of true  $\beta$ -values, then, for z-scores from a particular sample size  $N$ , the proportion of SNP heritability captured by genome-wide significant SNPs,  $A(N)$ , would be given by

$$A(N) = \frac{\sum_{i: |z_i| > z_t} \beta_i^2 H_i}{\sum_{\text{all } i} \beta_i^2 H_i}. \quad (33)$$

Now, from Eq. 14,  $\delta_i = \sqrt{N} \sum_j \sqrt{H_j} r_{ij} \beta_j$ . If SNP  $i$  is causal and sufficiently isolated so that it is not in LD with other causal SNPs, then  $\delta_i = \sqrt{N} \sqrt{H_i} \beta_i$ , and  $\text{var}(y_{g_i}) = \delta_i^2 / N$ . When all causal SNPs are similarly isolated, Eq.

33 becomes

$$A(N) = \frac{\sum_{i: |z_i| > z_t} \delta_i^2}{\sum_{\text{all } i} \delta_i^2}. \quad (34)$$

Of course, the true  $\beta_i$  are not known and some causal SNPs will likely be in LD with others. Furthermore, due to LD with causal SNPs, many SNPs will have a nonzero (latent or unobserved) effect size,  $\delta$ . Nevertheless, we can formulate an approximation to  $A(N)$  which, assuming the pdf for z-scores (Eq. 28) is reasonable, will be inaccurate to the degree that the average LD structure of genome-wide significant SNPs differs from the overall average LD structure. As before (see the subsection “Model PDF: Convolution”), consider a fixed set of  $n$  equally-spaced nominal z-scores covering a wide range of possible values (changing from the summations in Eq. 34 to the uniform summation spacing  $\Delta z$  now requires bringing the probability density into the summations). For each  $z$  from the fixed set (and, as before, employing data reduction by averaging so that H and L denote values for the  $10 \times 10$  grid), use

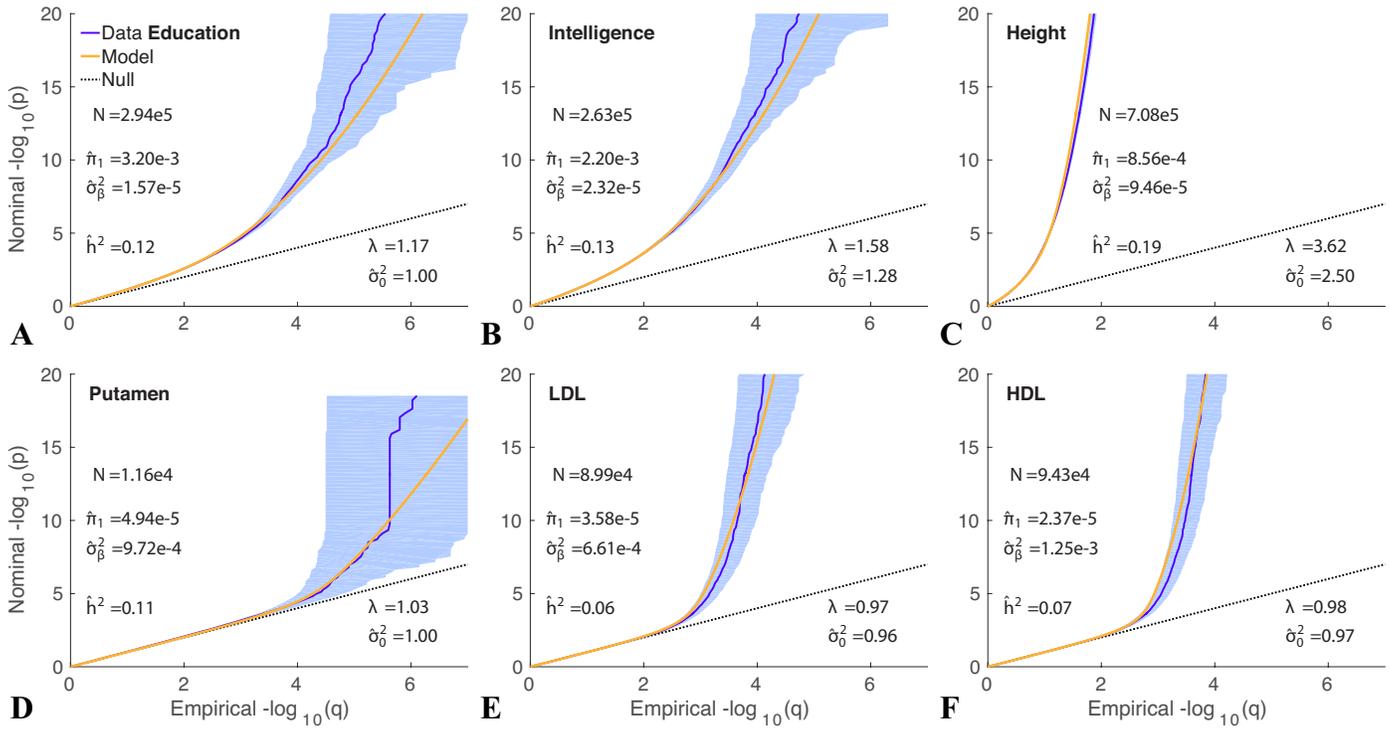


Figure 2: QQ plots of (pruned) z-scores for quantitative phenotypes: (A) educational attainment, (B) intelligence, (C) height, (D) putamen volume, (E) low-density lipoprotein (LDL), and (F) high-density lipoprotein (HDL).  $N$  is the sample size. See Fig. 1 for further description. Numerical values for the model parameters are also given in Table 1. See also Supplementary Material Figs. 20-25.

$E(\delta^2|z, N, H, L)$  given in Eq. 32 to define

$$C(z|N, H, L) \equiv E(\delta^2|z, N, H, L)P(z|N, H, L) \quad (35)$$

(emphasizing dependence on  $N$ ,  $H$ , and  $L$ ). Then, for any  $N$ ,  $A(N)$  can be estimated by

$$A(N) = \frac{\sum_{H,L} \sum_{z:|z|>z_t} C(z, N, H, L)}{\sum_{H,L} \sum_{all\ z} C(z, N, H, L)} \quad (36)$$

where  $\sum_{H,L}$  denotes sum over the H-L grid elements. The ratio in Eq. 36 should be accurate if the average effects of LD in the numerator and denominator cancel – which will always be true as the ratio approaches 1 for large  $N$ . Plotting  $A(N)$  gives an indication of the power of future GWAS to capture chip heritability.

### Quantile-Quantile Plots and Genomic Control

One of the advantages of quantile-quantile (QQ) plots is that on a logarithmic scale they emphasize behavior in the tails of a distribution, and provide a valuable visual aid in assessing the independent effects of polygenicity, strength of association, and cryptic relatedness – the roles played by the three model parameters – as well as showing how well a model fits data. QQ plots for the model were constructed using Eq. 28, replacing the normal pdf with the normal cdf, and replacing  $z$  with an equally-spaced vector  $\vec{z}_{nom}$  of length 10,000 covering a wide range of nominal  $|z|$  values (0 through 38). SNPs were divided into a  $10 \times 10$  grid of  $H \times L$  bins, and the cdf vector (with elements corresponding to

the z-values in  $\vec{z}_{nom}$ ) accumulated for each such bin (using mean values of  $H$  and  $L$  for SNPs in a given bin).

For a given set of samples and SNPs, the genomic control factor,  $\lambda$ , for the z-scores is defined as the median  $z^2$  divided by the median for the null distribution, 0.455 (Devlin and Roeder, 1999). This can also be calculated from the QQ plot. In the plots we present here, the abscissa gives the  $-\log_{10}$  of the proportion,  $q$ , of SNPs whose z-scores exceeded the two-tailed significance threshold  $p$ , transformed in the ordinate as  $-\log_{10}(p)$ . The median is at  $q_{med} = 0.5$ , or  $-\log_{10}(q_{med}) \simeq 0.3$ ; the corresponding empirical and model p-value thresholds ( $p_{med}$ ) for the z-scores – and equivalently for the z-scores-squared – can be read off from the plots. The genomic inflation factor is then given by  $\lambda = [\Phi^{-1}(p_{med}/2)]^2/0.455$ . Note that the values of  $\lambda$  reported here are for pruned SNP sets; these values will be lower than for the total GWAS SNP sets.

Knowing the total number,  $n_{tot}$ , of p-values involved in a QQ plot (number of GWAS z-scores from pruned SNPs), any point  $(q, p)$  (log-transformed) on the plot gives the number,  $n_p = q n_{tot}$ , of p-values that are as extreme as or more extreme than the chosen p-value. This can be thought of as  $n_p$  “successes” out of  $n_{tot}$  independent trials (thus ignoring LD) from a binomial distribution with prior probability  $q$ . To approximate the effects of LD, we estimate the number of independent SNPs as  $n_{tot}/f$  where  $f \simeq 10$ . The 95% binomial confidence interval for  $q$  is calculated as the exact Clopper-Pearson 95% inter-

val (Clopper and Pearson, 1934), which is similar to the normal approximation interval,  $q \pm 1.96\sqrt{q(1-q)/n_{tot}/f}$ .

### Number of Causal SNPs

The estimated number of causal SNPs is given by the polygenicity,  $\pi_1$ , times the total number of SNPs,  $n_{snp}$ :  $n_{causal} = \pi_1 n_{snp}$ .  $n_{snp}$  is given by the total number of SNPs that went into building the heterozygosity/LD structure,  $\mathcal{H}$  in Eq. 28, i.e., the approximately 11 million SNPs selected from the 1000 Genomes Phase 3 reference panel, not the number of tag SNPs in the particular GWAS. The parameters estimated are to be seen in the context of the reference panel, which we assume contains all common causal variants. Stable quantities (i.e., fairly independent of the reference panel size. e.g., using the full panel or ignoring every second SNP), are the estimated effect size variance and number of causal variants – which we demonstrate below – and hence the heritability. Thus, the polygenicity will scale inversely with the reference panel size. A reference panel with a substantially larger number of samples would allow for inclusion of more SNPs (non-zero maf), and thus the actual polygenicity estimated would change slightly.

### Narrow-sense Chip Heritability

Since we are treating the  $\beta$  coefficients as fixed effects in the simple linear regression GWAS formalism, with the phenotype vector standardized with mean zero and unit variance, from Eq. 1 the proportion of phenotypic variance explained by a particular causal SNP whose reference panel genotype vector is  $g$ ,  $q^2 = \text{var}(y; g)$ , is given by  $q^2 = \beta^2 H$ . The proportion of phenotypic variance explained additively by all causal SNPs is, by definition, the narrow sense chip heritability,  $h^2$ . Since  $E(\beta^2) = \sigma_\beta^2$  and  $n_{causal} = \pi_1 n_{snp}$ , and taking the mean heterozygosity over causal SNPs to be approximately equal to the mean over all SNPs,  $\bar{H}$ , the chip heritability can be estimated as

$$h^2 = \pi_1 n_{snp} \bar{H} \sigma_\beta^2. \quad (37)$$

Mean heterozygosity from the  $\sim 11$  million SNPs is  $\bar{H} = 0.2165$ .

For all-or-none traits like disease status, the estimated  $h^2$  from Eq. 37 for an ascertained case-control study is on the observed scale and is a function of the prevalence in the adult population,  $K$ , and the proportion of cases in the study,  $P$ . The heritability on the underlying continuous liability scale (Falconer, 1965),  $h_l^2$ , is obtained by adjusting for ascertainment (multiplying by  $K(1-K)/(P(1-P))$ , the ratio of phenotypic variances in the population and in the study) and rescaling based on prevalence (Dempster and Lerner, 1950; Lee et al., 2011):

$$h_l^2 = h^2 \frac{K(1-K)}{P(1-P)} \times \frac{K(1-K)}{a^2}, \quad (38)$$

where  $a$  is the height of the standard normal pdf at the truncation point  $z_K$  defined such that the area under the curve in the region to the right of  $z_K$  is  $K$ .

Phenotype	$\pi_1$	$\sigma_\beta^2$	$\sigma_0^2$	$n_{causal}$	$h_{(l)}^2$
MDD	4.01E-3	7.20E-6	1.06	4.4E4	0.07
Bipolar Disorder	2.70E-3	5.25E-5	1.05	3.0E4	0.16
Schizophrenia	2.84E-3	5.51E-5	1.14	3.1E4	0.21
Ulcerative Colitis	1.26E-4	8.82E-4	1.14	1.4E3	0.11
Crohn's Disease	9.56E-5	1.70E-3	1.17	1.1E3	0.15
AD	1.11E-4	2.22E-4	1.05	1.2E3	0.08
Education	3.20E-3	1.57E-5	1.00	3.5E4	0.12
Intelligence	2.20E-3	2.32E-5	1.28	2.4E4	0.13
Height	8.56E-4	9.46E-5	2.50	9.4E3	0.19
Putamen Volume	4.94E-5	9.72E-4	1.00	5.4E2	0.11
LDL	3.58E-5	6.61E-4	0.96	3.9E2	0.06
HDL	2.37E-5	1.25E-3	0.97	2.6E2	0.07

Table 1: Summary of model results for phenotypes shown in Figures 1 and 2. The subscript in  $h_{(l)}^2$  indicates that for the qualitative phenotypes (the first six) the reported SNP heritability is on the liability scale. MDD: Major Depressive Disorder; AD: Alzheimer's Disease (excluding APOE locus; for the full autosomal reference panel,  $h_l^2 = 0.15$  for AD – see Supporting Material Figure 10); LDL: low-density lipoproteins; HDL: high-density lipoproteins.

### GWAS Replication

Often a question arises whether z-scores for SNPs reaching genome-wide significance in a discovery-sample are compatible with the SNPs' z-scores in a replication-sample, particularly if any of those replication-sample z-scores are far from reaching genome-wide significance, or whether any apparent mismatch signifies some overlooked inconsistency. The model pdf allows one to make a principled statistical assessment in such cases. We present the details for this application, and results applied to studies of bipolar disorder, in the Supplementary Material.

## RESULTS

### Phenotypes

Figures 1 and 2 show QQ plots for the pruned z-scores for six qualitative and six quantitative phenotypes, along with model estimates (Supplementary Material Figs. 13-29 each show a  $4 \times 4$  grid breakdown by heterozygosity  $\times$  total-LD of QQ plots for all phenotypes studied here). In all cases, the model fit (yellow) closely tracks the data (dark blue). For the twelve phenotypes, estimates for the model polygenicity parameter range over two orders of magnitude, from  $\pi_1 \simeq 2 \times 10^{-5}$  to  $\pi_1 \simeq 4 \times 10^{-3}$ . The estimated SNP discoverability parameter (mean strength of SNP association with the phenotype, mean  $\beta^2$  for causals) also ranges over two orders of magnitude from  $\sigma_\beta^2 \simeq 7 \times 10^{-6}$  to  $\sigma_\beta^2 \simeq 2 \times 10^{-3}$  (in units where the variance of the phenotype is normalized to 1).

We find that schizophrenia and bipolar disorder appear to be similarly highly polygenic, with model polygenicities  $\simeq 2.84 \times 10^{-3}$  and  $\simeq 2.70 \times 10^{-3}$ , respectively. The model polygenicity of major depressive disorder, however, is 40%

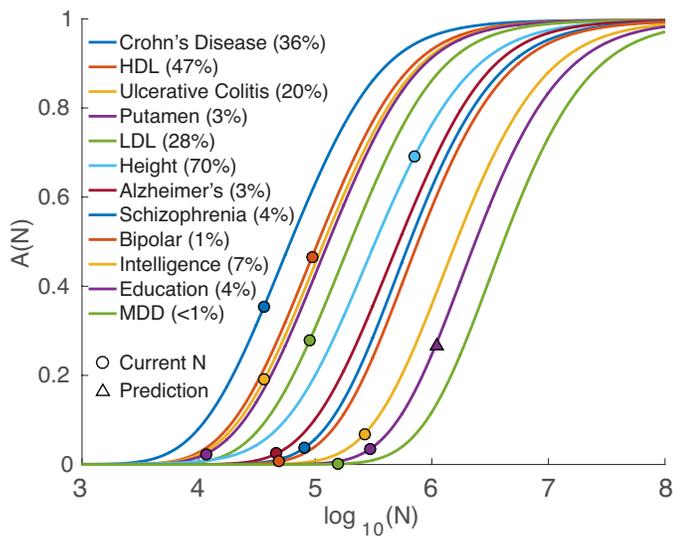


Figure 3: Proportion of narrow-sense chip heritability,  $A(N)$  (Eq. 36), captured by genome-wide significant SNPs as a function of sample size,  $N$ , for phenotypes shown in Figures 1 and Figure 2. Values for current sample sizes are shown in parentheses. Left-to-right curve order is determined by decreasing  $\sigma_\beta^2$ . The prediction for education at sample size  $N=1.1$  million is  $A(N) = 0.27$ , so that the proportion of phenotypic variance explained is predicted to be 3.5%, in good agreement with 3.2% reported in Lee et al. (2018). (The curve for AD excludes the APOE locus.)

higher,  $\pi_1 \simeq 4 \times 10^{-3}$  – the highest value among the twelve phenotypes. In contrast, the model polygenicities of late onset Alzheimer’s disease and Crohn’s disease are almost thirty times smaller than that of schizophrenia.

In Supplementary Material Figure 10 we show results for Alzheimer’s disease exclusively for chromosome 19 (which contains APOE), and for all autosomal chromosomes excluding chromosome 19. We also show results with the same chromosomal breakdown for a recent GWAS involving 455,258 samples that included 24,087 clinically diagnosed LOAD cases and 47,793 AD-by-proxy cases (individuals who were not clinically diagnosed with LOAD but for whom at least one parent had LOAD). These GWAS give consistent estimates of polygenicity:  $\pi_1 \sim 1 \times 10^{-4}$  excluding chromosome 19, and  $\pi_1 \sim 6 \times 10^{-5}$  for chromosome 19 exclusively.

Of the quantitative traits, educational attainment has the highest model polygenicity,  $\pi_1 = 3.2 \times 10^{-3}$ , similar to intelligence,  $\pi_1 = 2.2 \times 10^{-3}$ . Approximately two orders of magnitude lower in polygenicity are the endophenotypes putamen volume and low- and high-density lipoproteins; HDL had the lowest model polygenicity of the twelve phenotypes,  $\pi_1 = 2.37 \times 10^{-5}$ .

The estimated number of causal SNPs, in the reference template of 11 million SNPs, for schizophrenia is  $\hat{n}_{causal} \sim 3.1 \times 10^4$ , but major depressive disorder involves significantly more causal SNPs:  $\hat{n}_{causal} \simeq 4.4 \times 10^4$ . Alzheimer’s disease, however, involves only a relatively small number of causal SNPs,  $\hat{n}_{causal} \simeq 1,200$  (with  $\sim 10$  on chromosome 19 – see Supporting Material Figure 10),

while for height the model estimate is  $\hat{n}_{causal} \simeq 9,400$ . There is no underlying assumption in the model that the causal SNPs are in linkage equilibrium, and indeed for high-polygenicity phenotypes they can not all be in linkage equilibrium.

The model effective SNP discoverability for schizophrenia is  $\hat{\sigma}_\beta^2 = 5.51 \times 10^{-5}$ , similar to that for bipolar disorder. Major depressive disorder, which has the highest polygenicity, has the lowest SNP discoverability, approximately one-eighth that of schizophrenia; it is this low value, combined with high polygenicity that leads to the weak signal in Figure 1 (A) even though the sample size is relatively large. In contrast, SNP discoverability for Alzheimer’s disease is almost four times that of schizophrenia. The inflammatory bowel diseases, however, have much higher SNP discoverabilities, 16 and 31 times that of schizophrenia respectively for ulcerative colitis and Crohn’s disease – the latter having the highest value of the twelve phenotypes:  $\hat{\sigma}_\beta^2 = 1.7 \times 10^{-3}$ .

Additionally, for Alzheimer’s disease we show in Supplementary Material Figure 10 that the discoverability is two orders of magnitude greater for chromosome 19 than for the remainder of the autosome. Note that since two-thirds of the 2018 “cases” are AD-by-proxy, the discoverabilities for the 2018 data are, as expected, reduced relative to the values for the 2013 data (approximately 3.5 times smaller).

The effective SNP discoverability for education is  $\hat{\sigma}_\beta^2 = 1.57 \times 10^{-5}$ , about one-sixth that for height. However, high-density lipoprotein has a model discoverability an order of magnitude larger than that of height.

Note that for logistic linear regression coefficient  $\beta$ , the odds ratio for disease is  $OR = e^\beta$ ; for a rare disease, this is approximately equal to the genotypic relative risk:  $GRR \simeq OR$ . Since  $\mathbb{E}[\beta^2] = \sigma_\beta^2$ , the mean relative risk  $\mathbb{E}[GRR] \simeq 1 + \sigma_\beta^2/2$ . Thus, for schizophrenia for example, the mean relative risk is  $\simeq 1.00003$ .

The narrow sense SNP heritability from the ascertained case-control schizophrenia GWAS is estimated as  $h^2=0.37$ . Taking adult population prevalence of schizophrenia to be  $K=0.01$  (Purcell et al., 2009; Whiteford et al., 2013) (but see also Kinney et al. (2009), for  $K=0.005$ ), and given that there are 51,900 cases and 71,675 controls in the study, so that the proportion of cases in the study is  $P=0.42$ , the heritability on the liability scale for schizophrenia from Eq. 38 is  $\hat{h}_l^2=0.21$ . For bipolar disorder, with  $K=0.005$  (Merikangas et al., 2011), 20,352 cases and 31,358 controls,  $\hat{h}_l^2=0.16$ . Major depressive disorder appears to have a much lower model-estimated SNP heritability than schizophrenia:  $\hat{h}_l^2=0.07$ . The model estimate of SNP heritability for height is 19%, in contrast with reported values of 50%. The large discrepancy likely is related to the inflation which is generally simply assumed to be insignificant but which we estimate to be  $\sigma_0^2 = 2.5$ , which is in reasonable agreement with the intercepts obtained when directly plotting the data:  $E(z^2|H)$

$h^2$	$\hat{h}^2$	$\pi_1$	$\hat{\pi}_1$	$\sigma_\beta^2$	$\hat{\sigma}_\beta^2$	$\hat{\sigma}_0^2$	$n_{causal}$	$\hat{n}_{causal}$
0.1	0.12 (0.01)	1E-5	1.4E-5 (2E-6)	4.3E-3 (7E-4)	3.6E-3 (5E-4)	1.01 (0.002)	110	151 (20)
0.1	0.10 (0.01)	1E-4	1.0E-4 (2E-5)	4.2E-4 (2E-5)	4.1E-4 (5E-5)	1.01 (0.003)	1101	1130 (206)
0.1	0.09 (0.01)	1E-3	0.9E-3 (1E-4)	4.2E-5 (5E-7)	4.1E-5 (4E-6)	1.02 (0.003)	11015	10340 (1484)
0.1	0.09 (0.01)	1E-2	0.8E-2 (2E-3)	4.2E-6 (4E-8)	5.6E-6 (2E-6)	1.02 (0.002)	110158	83411 (25448)
0.4	0.52 (0.05)	1E-5	2.3E-5 (2E-6)	1.7E-2 (3E-3)	9.1E-2 (1E-3)	1.02 (0.002)	110	259 (20)
0.4	0.45 (0.02)	1E-4	1.2E-4 (8E-6)	1.7E-3 (7E-5)	1.5E-3 (9E-5)	1.04 (0.002)	1101	1310 (92)
0.4	0.39 (0.01)	1E-3	1.0E-3 (5E-5)	1.7E-4 (2E-6)	1.6E-4 (8E-6)	1.05 (0.003)	11015	10607 (578)
0.4	0.37 (0.01)	1E-2	0.9E-2 (1E-3)	1.7E-5 (2E-7)	1.7E-5 (2E-6)	1.06 (0.003)	110158	95135 (10851)
0.7	0.91 (0.09)	1E-5	2.9E-5 (2E-6)	3.0E-2 (5E-3)	1.3E-2 (2E-3)	1.02 (0.003)	110	324 (24)
0.7	0.82 (0.02)	1E-4	1.4E-4 (7E-6)	2.9E-3 (1E-4)	2.4E-3 (1E-4)	1.05 (0.002)	1101	1493 (79)
0.7	0.70 (0.01)	1E-3	1.0E-3 (4E-5)	2.9E-4 (4E-6)	2.8E-4 (1E-5)	1.08 (0.003)	11015	10866 (406)
0.7	0.66 (0.01)	1E-2	0.9E-2 (7E-4)	2.9E-5 (3E-7)	2.9E-5 (2E-6)	1.09 (0.003)	110158	95067 (8191)

Table 2: Simulation results: comparison of mean (std) true and estimated ( $\hat{\cdot}$ ) model parameters and derived quantities. Results for each line, for specified heritability  $h^2$  and fraction  $\pi_1$  of causal SNPs, are from 10 independent instantiations with random selection of the  $n_{causal}$  causal SNPs that are assigned a  $\beta$ -value from the standard normal distribution. Defining  $Y_g = G\beta$ , where  $G$  is the genotype matrix, the total phenotype vector is constructed as  $Y = Y_g + \varepsilon$ , where the residual random vector  $\varepsilon$  for each instantiation is drawn from a normal distribution such that  $\text{var}(Y) = \text{var}(Y_g)/h^2$  for predefined  $h^2$ . For each of the instantiations,  $i$ , this implicitly defines the true value  $\sigma_{\beta_i}^2$ , and  $\sigma_\beta^2$  is their mean. An example QQ plot for each line entry is shown in Supplementary Material, Figure 6.

and  $E(z^2|TLD)$  shown in Fig. 9.

Figure 3 shows the sample size required so that a given proportion of chip heritability is captured by genome-wide significant SNPs for the phenotypes (assuming equal numbers of cases and controls for the qualitative phenotypes:  $N_{eff} = 4/(1/N_{cases} + 1/N_{controls})$ , so that when  $N_{cases} = N_{controls}$ ,  $N_{eff} = N_{cases} + N_{controls} = N$ , the total sample size, allowing for a straightforward comparison with quantitative traits). At current sample sizes, only 4% of narrow-sense chip heritability is captured for schizophrenia and only 1% for bipolar disorder; using current methodologies, a sample size of  $N_{eff} \sim 1$  million would be required to capture the preponderance of SNP heritability for these phenotypes. Our model estimate for proportion of chip heritability captured for height by SNPs reaching  $p = 5 \times 10^{-8}$  at the current sample size is considerably higher, 70% (68% for the more stringent  $p = 10^{-8}$ ), and for HDL it is 47%. Major depressive disorder GWAS currently is greatly underpowered, as shown in Figure 3(A). For education, we predict that 3.5% of phenotypic variance would be explained at  $N = 1.1$  million, in good agreement with the value found from direct computation of 3.2% (Lee et al., 2018). For other phenotypes, the proportions of total SNP heritability captured at the available sample sizes are given in Figure 3. A comparison of our results with those of Zhang et al. (2018) is in Supporting Material Table 4.

We analyzed four additional phenotypes: amyotrophic lateral sclerosis (ALS, restricted to chromosome 9), body mass index (BMI), coronary artery disease (CAD), and total cholesterol (TC) – see Supplementary Material Figs. 11 and 12. The sample size for ALS was quite low, and we restricted the analysis to chromosome 9, which had most of the genome-wide significant tag SNPs; we estimate

that there are  $\sim 7$  causal SNPs with high discoverability on chromosome 9. For the three remaining phenotypes, the model estimates for the z-score distributions, Fig. 12, suggest that the single Gaussian with constant variance is less accurate for these phenotypes than for the others examined. Nevertheless, the results suggest that BMI is far more polygenic than TC, with the latter characterized by larger effects.

### Simulations

Table 2 shows the simulation results, comparing true and estimated values for the model parameters, heritability, and the number of causal SNPs, for twelve scenarios where  $\pi_1$  and  $\sigma_\beta^2$  both range over three orders of magnitude, encompassing the range of values for the phenotypes; in Supporting Material, Figure 6 shows QQ plots for a randomly chosen (out of 10)  $\beta$ -vector and phenotype instantiation for each of the twelve  $(\pi_1, h^2)$  scenarios. Most of the  $\hat{\pi}_1$  estimates are in very good agreement with the true values, though for the extreme scenario of high heritability and low polygenicity it is overestimated by factors of two-to-three. The numbers of estimated causal SNPs (out of  $\sim 11$  million) are in correspondingly good agreement with the true values, ranging in increasing powers of 10 from 110 through 110,158. The estimated discoverabilities ( $\hat{\sigma}_\beta^2$ ) are also in good agreement with the true values. In most cases,  $\hat{\sigma}_0^2$  is close to 1, indicating little or no global inflation, though it is elevated for high heritability with high polygenicity, suggesting it is capturing some ubiquitous effects. To test the modeling of inflation due to cryptic relatedness, we scaled the simulation z-scores as described earlier ( $z = \sigma_0 z_u$  with  $\sigma_0 > 1$ , where  $z_u$  are the original z-scores, i.e., not artificially inflated) and reran the model. E.g., for  $\sigma_0 = 1.2$ ,  $\hat{\pi}_1$  and  $\hat{\sigma}_\beta$  were found to be the same as

for the uninflated z-scores, and  $\hat{\sigma}_0 \simeq 1.2$  times the original (slightly inflated) estimates shown in Table 2.

In Supplementary Material, we examine the issue of model misspecification. Specifically, we assign causal effects  $\beta$  drawn from a Gaussian whose variance is not simply a constant but depends on heterozygosity, such that rarer causal SNPs will tend to have larger effects. The results – see Supplementary Material Table 3 – show that the model still makes reasonable estimates of the underlying genetic architecture.

### Dependence on Reference Panel

Our working assumption is that the reference panel of 11 million SNPs captures essentially all of the common SNPs that contribute in some causal fashion to phenotypes. What happens if, say, only half of the reference panel is used, for example by taking only every other SNP? We explicitly tested this, with a full rebuild of the culled reference panel SNP LD and heterozygosity structures, and rerunning the model on real phenotypes. The result is that all estimated parameters are as before except that  $\hat{\pi}_1$  doubles, leaving the estimated number of causal SNPs and heritability as before.

### DISCUSSION

Here we present a unified method based on GWAS summary statistics, incorporating detailed LD structure from a reference panel, for estimating phenotypic polygenicity,  $\pi_1$ , SNP discoverability or strength of association (the variance of the underlying causal effects),  $\sigma_\beta^2$ , and narrow-sense SNP heritability. In addition the model can be used to estimate residual inflation of the association statistics due to variance distortion induced by cryptic relatedness,  $\sigma_0^2$ .

We apply the model to twelve diverse phenotypes, six qualitative and six quantitative. From the estimated model parameters we also estimate the number of causal SNPs,  $n_{causal}$ , and the SNP heritability,  $h^2$  (for qualitative phenotypes, we re-express this as the proportion of population variance in disease liability,  $h_l^2$ , under a liability threshold model, adjusted for ascertainment). In addition, we estimate the proportion of SNP heritability captured by genome-wide significant SNPs at current sample sizes, and predict future sample sizes needed to explain the preponderance of SNP heritability.

We find that schizophrenia is highly polygenic, with  $\pi_1 = 2.8 \times 10^{-3}$ . This leads to an estimate of  $n_{causal} \simeq 31,000$ , which is in reasonable agreement with a recent estimate that the number of causals is  $>20,000$  (Loh et al., 2015). The SNP associations, however, are characterized by a narrow distribution,  $\sigma_\beta^2 = 6.27 \times 10^{-5}$ , indicating that most associations are of weak effect, i.e., have low discoverability. Bipolar disorder has similar parameters. The smaller sample size for bipolar disorder has led to fewer SNP discoveries compared with schizophrenia. How-

ever, from Figure 3, sample sizes for bipolar disorder are approaching a range where rapid increase in discoveries becomes possible. For educational attainment (Rietveld et al., 2013; Okbay et al., 2016; Cesarini and Visscher, 2017), the polygenicity is somewhat greater,  $\pi_1 = 3.2 \times 10^{-3}$ , leading to an estimate of  $n_{causal} \simeq 35,000$ , half a recent estimate,  $\simeq 70,000$ , for the number of loci contributing to heritability (Rietveld et al., 2013). The variance of the distribution for causal effect sizes is a quarter that of schizophrenia, indicating lower discoverability. Intelligence, a related phenotype (Sniekers et al., 2017; Plomin and von Stumm, 2018), has a larger discoverability than education while having lower polygenicity ( $\sim 10,000$  fewer causal SNPs).

In marked contrast are the lipoproteins and putamen volume which have very low polygenicity:  $\pi_1 < 5 \times 10^{-5}$ , so that only 250 to 550 SNPs (out of  $\sim 11$  million) are estimated to be causal. However, causal SNPs for putamen volume and HDL are characterized by relatively high discoverability, respectively 17-times and 23-times larger than for schizophrenia.

The QQ plots (which are sample size dependent) reflect these differences in genetic architecture. For example, the early departure of the schizophrenia QQ plot from the null line indicates its high polygenicity, while the steep rise for putamen volume after its departure corresponds to its high SNP discoverability.

Twin studies estimate the heritability of Alzheimer’s disease to be in the range 60-80% (Gatz et al., 2006). A recent raw genotype-based analysis (GCTA), including genes that contain rare variants that affect risk for AD, reported SNP heritability of 53% (Ridge et al., 2016; Yang et al., 2011a); an earlier related study that did not include rare variants and had only a quarter of the common variants estimated SNP heritability of 33% (Ridge et al., 2013). It should be noted that GCTA calculations of heritability are within the domain of the so-called infinitesimal model where all markers are assumed to be causal. Our model suggests, however, that phenotypes are characterized by polygenicities less than  $5 \times 10^{-3}$ ; for AD the polygenicity is  $\simeq 10^{-4}$ . Thus, though the GCTA approach yields a heritability estimate closer to the twin-based (broad sense) value, it is not clear how to interpret GCTA estimates of heritability when the underlying polygenicity in fact is  $\ll 1$ . For the 2013 data analyzed here (Lambert et al., 2013), a summary-statistics-based method applied to a subset of 54,162 of the 74,046 samples gave SNP heritability of almost 7% (Zheng et al., 2017; Bulik-Sullivan et al., 2015). Our estimate for the full 2013 dataset is 15%, half from APOE.

Onset and clinical progression of sporadic Alzheimer’s disease is strongly age-related (Holland et al., 2012; Desikan et al., 2017), with prevalence in differential age groups increasing at least up through the early 90s (Plassman et al., 2007). Thus, it would be more accurate to assess heritability (and its components, polygenicity and discoverability) with respect to, say, five-year age groups be-

ginning with age 65 years, and using a consistent control group of nonagenarians and centenarians. By the same token, comparisons among current and past AD GWAS are complicated because of potential differences in the age distributions of the respective case and the control cohorts. Additionally, choice of prevalence parameter in calculating liability-scale heritability, as well as the degree to which rare variants are included, will affect heritability estimates. The summary-statistic-based estimates of polygenicity that we report here are, however, likely to be robust for common SNPs:  $\pi_1 \simeq 1.1 \times 10^{-4}$ , with only a few causal SNPs on chromosome 19. The consistent deviations between the model fits and the data shown in Figure 1 (F) and Supporting Material Figure 10 (A) and (D) suggest that the model does not fully capture the genetic architecture of AD (ignoring APOE/chromosome 19). Future work will examine annotation-specific SNP categories, and the possibility of higher locality among causal SNPs – where one is, another might more likely be in the neighborhood than indicated by the random placement underlying the current version of the model. Regarding chromosome 19, the very low polygenicity coupled with the very high discoverability is a region where the model is less reliable in the precision of the parameter estimates; apart from proposed improvements just described, accuracy in this region of parameter space will be enhanced by more precise estimates of LD and heterozygosity structure of SNPs (enhanced signal to noise in the low LD ( $r^2 < 0.05$ ) and/or low heritability ( $H < 0.05$ ) range), so that  $r_{min}^2$  used in the model pdf (Eq. 28) can reliably be taken lower than 0.05. This can be achieved by using reference panels with larger sample sizes.

Our point estimate for the liability-scale SNP heritability of schizophrenia is  $h_l^2 = 0.21$  (assuming a population risk of 0.01), and that 4% of this (i.e., 1% of overall disease liability) is explainable based on common SNPs reaching genome-wide significance at the current sample size. This  $h_l^2$  estimate is in reasonable agreement with a recent result,  $h_l^2 = 0.27$  (Loh et al., 2015; Golan et al., 2014), also calculated from the PGC2 data set but using raw genotype data for 472,178 markers for a subset of 22,177 schizophrenia cases and 27,629 controls of European ancestry; and with an earlier result of  $h_l^2 = 0.23$  from PGC1 raw genotype data for 915,354 markers for 9,087 schizophrenia cases and 12,171 controls (Lee et al., 2012; Yang et al., 2011a). Our estimate of 1% of overall variation on the liability scale for schizophrenia explainable by genome-wide significant loci compares reasonably with the proportion of variance on the liability scale explained by Risk Profile Scores (RPS) reported as 1.1% using the MGS sample as target (the median for all 40 leave-one-out target samples analyzed is 1.19% – see Extended Data Figure 5 and Supplementary Tables 5 and 6 in Schizophrenia Working Group of the Psychiatric Genomics Consortium (2014); this was incorrectly reported as 3.4% in the main paper). These results show that current sample sizes need to increase substantially in order for RPSs to have predictive utility, as the vast major-

ity of associated SNPs remain undiscovered. Our power estimates indicate that  $\sim 500,000$  cases and an equal number of controls would be needed to identify these SNPs (note that there is a total of approximately 3 million cases in the US alone). The identified SNPs then need to be mapped to genes and their modality (e.g., regulatory or functional effects) determined, so that targeted therapeutics can be developed (Schubert et al., 2015). Greater power for discovery is achievable by using prior information involving SNP functional categories (Schork et al., 2013; Andreassen et al., 2013; Sveinbjornsson et al., 2016). However, it is not yet clear how significant a role genomics can play in psychiatric precision medicine (Breen et al., 2016). It is noteworthy that estimates of broad-sense heritability of schizophrenia from twin and family studies are in the range 0.6-0.8 (Sullivan et al., 2003; Lichtenstein et al., 2009), considerably higher than the narrow-sense chip heritability estimates from GWAS. Additionally, schizophrenia is considered a spectrum disorder with multiple phenotypic dimensions and diverse clinical presentation (MacDonald and Schulz, 2009; Peralta and Cuesta, 2001); GWAS might therefore benefit from considering continuous phenotypes rather than dichotomous variables in such situations (Edwards et al., 2016). More specifically in the context of the present model, if a nominally categorical phenotype can be decomposed into more than one subcategory, there is potential for enhanced power for discovery. The heritability estimated in a binary case-control design would be an average over heritabilities for the case subcategories. If those heritabilities are similar, then, since the union of the subcategory polygenicities gives the total polygenicity over all cases, the  $\sigma_\beta^2$  for any subcategory will be larger by a factor equal to the ratio of overall polygenicity to the subcategory polygenicity, and the corresponding power curve (as in Figure 3) will shift to the left.

For educational attainment, we estimate SNP heritability  $h^2 = 0.12$ , in good agreement with the estimate of 11.5% given in Okbay et al. (2016). As with schizophrenia, this is substantially less than the estimate of heritability from twin and family studies of  $\simeq 40\%$  of the variance in educational attainment explained by genetic factors (Branigan et al., 2013; Rietveld et al., 2013).

For putamen volume, we estimate the SNP heritability  $h^2 = 0.11$ , in reasonable agreement with an earlier estimate of 0.1 for the same overall data set (Hibar et al., 2015; So et al., 2011).

For height we find that its polygenicity is  $\pi_1 = 8.56 \times 10^{-4}$ , two-fifths that of intelligence (and very far from “omnigenic” (Boyle et al., 2017)), while its discoverability is four times that of intelligence, leading to a SNP heritability of 19%. This is in considerable disagreement with the GCTA SNP heritability estimate of 50% obtained from a subset of the raw genotype data (Wood et al., 2014). Power analysis for the model shows that for height 68% of the narrow-sense heritability arising from common SNPs is explained by more stringent genome-wide significant SNPs ( $p \leq 10^{-8}$ ) at the current sample size, i.e., 13% of total

phenotypic variance, which is substantially less than the 24.6% direct estimate from significant SNPs (Yengo et al., 2018). It is not clear why these large discrepancies exist. One relevant factor, however, is the very large estimated inflation that we report,  $\sigma_0^2 = 2.5$ , suggesting that z-scores are too large by a factor of 1.6. This inflation estimate comports with our intercept estimates when regressing  $z^2$  on total LD and heterozygosity – see Supporting Material Fig. 9. Also, our model is in a sense a basic mixture model, incorporating just one gaussian with constant variance for the association  $\beta$ -coefficients for the subset of causal SNPs. We will report on extensions to this model, incorporating heterozygosity in the variance (Zeng et al., 2018), additional gaussians (Zhang et al., 2018), and total LD dependence, in forthcoming work.

Our liability-scale heritability estimate for schizophrenia is substantially lower than the estimate from LD Score regression,  $h_l^2 = 0.555$  (Bulik-Sullivan et al., 2015). When the per-allele effect size,  $\beta$ , is independent of the MAF, LD Score regression will lead to a biased rotation of the regression line, increasing the slope, which is proportional to the estimated heritability, and decreasing the intercept, which gives the residual inflation. Thus, the discrepancy with LD Score regression in principle might arise due to the assumption in LD Score regression that effect sizes for *all* variants (the infinitesimal model) are drawn independently from normal distributions with variance inversely proportional to heterozygosity (so that the variance explained per SNP is uncorrelated with TLD, since TLD is positively correlated with MAF – Supporting Material Figure 7(B)). However, we show that the expected  $z^2$  increases linearly with H (as well as TLD, Supporting Material Figure 8), suggesting that  $E(\beta^2) = \sigma_\beta^2$  is, to a first approximation, independent of MAF (this is only an approximate argument because, strictly speaking,  $E(z^2) = \text{var}(z)$  depends on the heterozygosities of the causal reference SNPs the tag SNP is in LD with, not the heterozygosity of the tag SNP itself – see Eqs. 14 and 20). So the problem might arise due to assuming  $\beta \sim \mathcal{N}(0, H^S \sigma_\beta^2)$  with  $S = -1$  for *all* SNPs (polygenicity  $\pi_1 = 1$ ), whereas (1) only a small fraction of SNPs are causal; (2) of those only a subset might have the selection effect; and (3) the selection effect might well be weaker (closer to 0) than implied by  $S = -1$ . In forthcoming work we will report on extensions to our basic model, including heterozygosity-dependence of the genetic effect variance parameter, multiple Gaussian, and total LD dependence of prior probabilities.

Twin-based studies indicate that ALS heritability is around 65%. A recent GWAS, using a custom reference panel that considerably increased imputation accuracy (relative to 1000 Genomes) for low MAF SNPs, identified three new and one known associated loci. The SNP heritability was reported as 8.5% from GCTA, and 8.2% from LD score regression. Although these estimates implicitly are based on the infinitesimal model (polygenicity  $\pi_1 = 1$ ), which is not likely to reflect the actual genetic architecture, the dominant contribution to heritability was found

to come from low-frequency variants ( $\text{MAF} < 0.1$ ), where our data are scant. This points to the importance of the underlying reference panel, particularly with regard to rare variants, when assessing results from a causal effects mixture model.

## CONCLUSION

The SNP-level causal effects model we have presented is based on GWAS summary statistics and detailed LD structure of an underlying reference panel, and assumes a Gaussian distribution of effect sizes at a fraction of SNPs randomly distributed across the autosomal genome. We have shown that it captures the broad genetic architecture of diverse complex traits, where polygenicities and the variance of the effect sizes range over orders of magnitude. In addition, the model provides a roadmap for discovery in future GWAS. Future extensions and refinements include modeling specific polygenicities and effect size variances for different SNP functional annotation categories (Schork et al., 2013; Andreassen et al., 2013; Sveinbjornsson et al., 2016), possible modified pdf for non-Gaussian distribution of effects at the tails of the z-score distributions, examining individual chromosomes and possible allele frequency dependencies in different phenotypes, and extension to pleiotropic analyses. Higher accuracy in characterizing causal alleles in turn will enable greater power for SNP discovery. The current model (essentially Eq. 4) and its implementation (essentially Eq. 28) are basic elements for building a more refined model of SNP effects using summary statistics.

## Acknowledgments

We thank the consortia for making available their GWAS summary statistics, and the many people who provided DNA samples.

## Funding

Research Council of Norway (262656, 248984, 248778, 223273) and KG Jebsen Stiftelsen; ABCD-USA Consortium (5U24DA041123).

## References

- Alzheimer’s Association, 2018. 2018 alzheimer’s disease facts and figures. *Alzheimer’s & Dementia* 14 (3), 367–429.
- Andreassen, O. A., Djurovic, S., Thompson, W. K., Schork, A. J., Kendler, K. S., O’Donovan, M. C., Rujescu, D., Werge, T., van de Bunt, M., Morris, A. P., et al., 2013. 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.
- Boyle, E. A., Li, Y. I., Pritchard, J. K., 2017. An expanded view of complex traits: From polygenic to omnigenic. *Cell* 169 (7), 1177–1186.

- Branigan, A. R., McCallum, K. J., Freese, J., 2013. Variation in the heritability of educational attainment: An international meta-analysis. *Social Forces*, 109–140.
- Breen, G., Li, Q., Roth, B. L., O'Donnell, P., Didriksen, M., Dolmetsch, R., O'Reilly, P. F., Gaspar, H. A., Manji, H., Huebel, C., et al., 2016. Translating genome-wide association findings into new therapeutics for psychiatry. *Nature Neuroscience* 19 (11), 1392–1396.
- Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Patterson, N., Daly, M. J., Price, A. L., Neale, B. M., of the Psychiatric Genomics Consortium, S. W. G., et al., 2015. Ld score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature genetics* 47 (3), 291–295.
- Burisch, J., Jess, T., Martinato, M., Lakatos, P. L., ECCO-EpiCom, 2013. The burden of inflammatory bowel disease in europe. *Journal of Crohn's and Colitis* 7 (4), 322–337.
- Cesarini, D., Visscher, P. M., 2017. Genetics and educational attainment. *npj Science of Learning* 2 (1), 4.
- Clopper, C. J., Pearson, E. S., 1934. The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika* 26 (4), 404–413.
- Consortium, . G. P., et al., 2012. An integrated map of genetic variation from 1,092 human genomes. *Nature* 491 (7422), 56–65.
- Consortium, . G. P., et al., 2015. A global reference for human genetic variation. *Nature* 526 (7571), 68–74.
- de Lange, K. M., Moutsianas, L., Lee, J. C., Lamb, C. A., Luo, Y., Kennedy, N. A., Jostins, L., Rice, D. L., Gutierrez-Achury, J., Ji, S.-G., et al., 2017. Genome-wide association study implicates immune activation of multiple integrin genes in inflammatory bowel disease. *Nature genetics* 49 (2), 256.
- Dempster, E. R., Lerner, I. M., 1950. Heritability of threshold characters. *Genetics* 35 (2), 212.
- Desikan, R. S., Fan, C. C., Wang, Y., Schork, A. J., Cabral, H. J., Cupples, L. A., Thompson, W. K., Besser, L., Kukull, W. A., Holland, D., et al., 2017. Genetic assessment of age-associated alzheimer disease risk: Development and validation of a polygenic hazard score. *PLoS medicine* 14 (3), e1002258.
- Devlin, B., Roeder, K., Dec 1999. Genomic control for association studies. *Biometrics* 55 (4), 997–1004.
- Edwards, A. C., Bigdeli, T. B., Docherty, A. R., Bacanu, S., Lee, D., De Candia, T. R., Moscati, A., Thiselton, D. L., Maher, B. S., Wormley, B. K., et al., 2016. Meta-analysis of positive and negative symptoms reveals schizophrenia modifier genes. *Schizophrenia bulletin* 42 (2), 279–287.
- Falconer, D. S., 1965. The inheritance of liability to certain diseases, estimated from the incidence among relatives. *Annals of human genetics* 29 (1), 51–76.
- Gatz, M., Reynolds, C. A., Fratiglioni, L., Johansson, B., Mortimer, J. A., Berg, S., Fiske, A., Pedersen, N. L., 2006. Role of genes and environments for explaining alzheimer disease. *Archives of general psychiatry* 63 (2), 168–174.
- Golan, D., Lander, E. S., Rosset, S., 2014. Measuring missing heritability: inferring the contribution of common variants. *Proceedings of the National Academy of Sciences* 111 (49), E5272–E5281.
- Hibar, D. P., Stein, J. L., Renteria, M. E., Arias-Vasquez, A., Desrivieres, S., Jahanshad, N., Toro, R., Wittfeld, K., Abramovic, L., Andersson, M., et al., 2015. Common genetic variants influence human subcortical brain structures. *Nature*.
- Holland, D., Desikan, R. S., Dale, A. M., McEvoy, L. K., Initiative, A. D. N., et al., 2012. Rates of decline in alzheimer disease decrease with age. *PloS one* 7 (8), e42325.
- Holland, D., Wang, Y., Thompson, W. K., Schork, A., Chen, C. H., Lo, M. T., Witoelar, A., Werge, T., O'Donovan, M., Andreassen, O. A., Dale, A. M., 2016. Estimating Effect Sizes and Expected Replication Probabilities from GWAS Summary Statistics. *Front Genet* 7, 15.
- Jansen, I., Savage, J., Watanabe, K., Bryois, J., Williams, D., Steinberg, S., Sealock, J., Karlsson, I., Hagg, S., Athanasiu, L., et al., 2018. Genetic meta-analysis identifies 10 novel loci and functional pathways for alzheimer's disease risk. *bioRxiv*, 258533.
- Kang, H. M., Sul, J. H., Zaitlen, N. A., Kong, S.-y., Freimer, N. B., Sabatti, C., Eskin, E., et al., 2010. Variance component model to account for sample structure in genome-wide association studies. *Nature genetics* 42 (4), 348–354.
- Kinney, D. K., Teixeira, P., Hsu, D., Napoleon, S. C., Crowley, D. J., Miller, A., Hyman, W., Huang, E., 2009. Relation of schizophrenia prevalence to latitude, climate, fish consumption, infant mortality, and skin color: a role for prenatal vitamin d deficiency and infections? *Schizophrenia bulletin*, sbp023.
- Kumar, S. K., Feldman, M. W., Rehkopf, D. H., Tuljapurkar, S., 2016. Limitations of gcta as a solution to the missing heritability problem. *Proceedings of the National Academy of Sciences* 113 (1), E61–E70.
- Laird, N. M., Lange, C., 2010. The fundamentals of modern statistical genetics. Springer Science & Business Media.
- Lambert, J.-C., Ibrahim-Verbaas, C. A., Harold, D., Naj, A. C., Sims, R., Bellenguez, C., Jun, G., DeStefano, A. L., Bis, J. C., Beecham, G. W., et al., 2013. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for alzheimer's disease. *Nature genetics* 45 (12), 1452–1458.
- Lee, J. J., Wedow, R., Okbay, A., Kong, E., Maghziyan, O., Zacher, M., Nguyen-Viet, T. A., Bowers, P., Sidorenko, J., Linnér, R. K., et al., 2018. Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nature genetics* 50 (8), 1112.
- Lee, S. H., DeCandia, T. R., Ripke, S., Yang, J., Sullivan, P. F., Goddard, M. E., Keller, M. C., Visscher, P. M., Wray, N. R., Consortium, S. P. G.-W. A. S., et al., 2012. Estimating the proportion of variation in susceptibility to schizophrenia captured by common snps. *Nature genetics* 44 (3), 247–250.
- Lee, S. H., Wray, N. R., Goddard, M. E., Visscher, P. M., 2011. Estimating missing heritability for disease from genome-wide association studies. *The American Journal of Human Genetics* 88 (3), 294–305.
- Li, N., Stephens, M., 2003. Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics* 165 (4), 2213–2233.
- Lichtenstein, P., Yip, B. H., Björk, C., Pawitan, Y., Cannon, T. D., Sullivan, P. F., Hultman, C. M., 2009. Common genetic determinants of schizophrenia and bipolar disorder in swedish families: a population-based study. *The Lancet* 373 (9659), 234–239.
- Locke, A. E., Kahali, B., Berndt, S. I., Justice, A. E., Pers, T. H., Day, F. R., Powell, C., Vedantam, S., Buchkovich, M. L., Yang, J., et al., 2015. Genetic studies of body mass index yield new insights for obesity biology. *Nature* 518 (7538), 197.
- Loh, P.-R., Bhatia, G., Gusev, A., Finucane, H. K., Bulik-Sullivan, B. K., Pollack, S. J., de Candia, T. R., Lee, S. H., Wray, N. R., Kendler, K. S., et al., 2015. Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nature genetics*.
- MacDonald, A. W., Schulz, S. C., 2009. What we know: findings that every theory of schizophrenia should explain. *Schizophrenia Bulletin* 35 (3), 493–508.
- McCarthy, M. I., Abecasis, G. R., Cardon, L. R., Goldstein, D. B., Little, J., Ioannidis, J. P., Hirschhorn, J. N., 2008. Genome-wide association studies for complex traits: consensus, uncertainty and challenges. *Nature reviews genetics* 9 (5), 356–369.
- Mehta, P., Kaye, W., Raymond, J. e. a., 2018. Prevalence of amyotrophic lateral sclerosis 2014 united states. *MMWR Morb Mortal Wkly Rep* 67, 216–218.
- Merikangas, K. R., Jin, R., He, J.-P., Kessler, R. C., Lee, S., Sampson, N. A., Viana, M. C., Andrade, L. H., Hu, C., Karam, E. G., et al., 2011. Prevalence and correlates of bipolar spectrum disorder in the world mental health survey initiative. *Archives of general psychiatry* 68 (3), 241–251.
- Nikpay, M., Goel, A., Won, H.-H., Hall, L. M., Willenborg, C., Kanoni, S., Saleheen, D., Kyriakou, T., Nelson, C. P., Hopewell, J. C., et al., 2015. A comprehensive 1000 genomes-based genome-wide association meta-analysis of coronary artery disease. *Nature genetics* 47 (10), 1121.
- NIMH, 2016. Prevalence of Major Depressive Episode Among Adults. (accessed December 27, 2018).

- URL <https://www.nimh.nih.gov/health/statistics/major-depression.shtml>
- Okbay, A., Beauchamp, J. P., Fontana, M. A., Lee, J. J., Pers, T. H., Rietveld, C. A., Turley, P., Chen, G.-B., Emilsson, V., Meddens, S. F. W., et al., 2016. Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* 533 (7604), 539–542.
- Palla, L., Dudbridge, F., 2015. A fast method that uses polygenic scores to estimate the variance explained by genome-wide marker panels and the proportion of variants affecting a trait. *The American Journal of Human Genetics* 97 (2), 250–259.
- Pasaniuc, B., Price, A. L., 2016. Dissecting the genetics of complex traits using summary association statistics. *Nature Reviews Genetics*.
- Pe'er, I., Yelensky, R., Altshuler, D., Daly, M. J., 2008. Estimation of the multiple testing burden for genomewide association studies of nearly all common variants. *Genetic epidemiology* 32 (4), 381–385.
- Peralta, V., Cuesta, M. J., 2001. How many and which are the psychopathological dimensions in schizophrenia? issues influencing their ascertainment. *Schizophrenia research* 49 (3), 269–285.
- Plassman, B. L., Langa, K. M., Fisher, G. G., Heeringa, S. G., Weir, D. R., Ofstedal, M. B., Burke, J. R., Hurd, M. D., Potter, G. G., Rodgers, W. L., et al., 2007. Prevalence of dementia in the united states: the aging, demographics, and memory study. *Neuroepidemiology* 29 (1-2), 125–132.
- Plomin, R., von Stumm, S., 2018. The new genetics of intelligence. *Nature Reviews Genetics*.
- Price, A. L., Zaitlen, N. A., Reich, D., Patterson, N., 2010. New approaches to population stratification in genome-wide association studies. *Nature Reviews Genetics* 11 (7), 459–463.
- Purcell, S. M., Wray, N. R., Stone, J. L., Visscher, P. M., O'Donovan, M. C., Sullivan, P. F., Sklar, P., Purcell, S. M., Stone, J. L., Sullivan, P. F., et al., 2009. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature* 460 (7256), 748–752.
- Ridge, P. G., Hoyt, K. B., Boehme, K., Mukherjee, S., Crane, P. K., Haines, J. L., Mayeux, R., Farrer, L. A., Pericak-Vance, M. A., Schellenberg, G. D., et al., 2016. Assessment of the genetic variance of late-onset alzheimer's disease. *Neurobiology of aging* 41, 200–e13.
- Ridge, P. G., Mukherjee, S., Crane, P. K., Kauwe, J. S., et al., 2013. Alzheimer2019s disease: analyzing the missing heritability. *PLoS one* 8 (11), e79771.
- Rietveld, C. A., Medland, S. E., Derringer, J., Yang, J., Esko, T., Martin, N. W., Westra, H.-J., Shakhbazov, K., Abdellaoui, A., Agrawal, A., et al., 2013. Gwas of 126,559 individuals identifies genetic variants associated with educational attainment. *science* 340 (6139), 1467–1471.
- Sanchis-Gomar, F., Perez-Quilis, C., Leischik, R., Lucia, A., 2016. Epidemiology of coronary heart disease and acute coronary syndrome. *Annals of translational medicine* 4 (13).
- Savage, J., Jansen, P., Stringer, S., Watanabe, K., Bryois, J., de Leeuw, C., Nagel, M., Awasthi, S., Barr, P., Coleman, J., Grasby, K., Hammerschlag, A., Kaminski, J., Karlsson, R., et al., 2018. Genome-wide association meta-analysis (n=269,867) identifies new genetic and functional links to intelligence. *Nature genetics* forthcoming.  
URL [https://ctg.cncr.nl/software/summary\\_statistics](https://ctg.cncr.nl/software/summary_statistics)
- Schizophrenia Working Group of the Psychiatric Genomics Consortium, Jul 2014. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511 (7510), 421–427.
- Schork, A. J., Thompson, W. K., Pham, P., Torkamani, A., Roddey, J. C., Sullivan, P. F., Kelsoe, J. R., O'Donovan, M. C., Furburg, H., Schork, N. J., et al., 2013. All snps are not created equal: genome-wide association studies reveal a consistent pattern of enrichment among functionally annotated snps. *PLoS genetics* 9 (4), e1003449.
- Schubert, C. R., O'Donnell, P., Quan, J., Wendland, J. R., Xi, H. S., Winslow, A. R., Domenici, E., Essioux, L., Kam-Thong, T., Airey, D. C., et al., 2015. Brainseq: neurogenomics to drive novel target discovery for neuropsychiatric disorders. *Neuron* 88 (6), 1078.
- Sniekers, S., Stringer, S., Watanabe, K., Jansen, P. R., Coleman, J. R., Krapohl, E., Taskesen, E., Hammerschlag, A. R., Okbay, A., Zabaneh, D., et al., 2017. Genome-wide association meta-analysis of 78,308 individuals identifies new loci and genes influencing human intelligence. *Nature genetics* 49 (7), 1107.
- So, H.-C., Li, M., Sham, P. C., 2011. Uncovering the total heritability explained by all true susceptibility variants in a genome-wide association study. *Genetic epidemiology* 35 (6), 447–456.
- Speed, D., Cai, N., Johnson, M. R., Nejentsev, S., Balding, D. J., Consortium, U., et al., 2017. Reevaluation of snp heritability in complex human traits. *Nature genetics* 49 (7), 986.
- Speed, D., Hemani, G., Johnson, M. R., Balding, D. J., 2012. Improved heritability estimation from genome-wide snps. *The American Journal of Human Genetics* 91 (6), 1011–1021.
- Spencer, C. C., Su, Z., Donnelly, P., Marchini, J., 2009. Designing genome-wide association studies: sample size, power, imputation, and the choice of genotyping chip. *PLoS Genet* 5 (5), e1000477.
- Stahl, E., Breen, G., Forstner, A., McQuillin, A., Ripke, S., Cichon, S., Scott, L., Ophoff, R., Andreassen, O. A., Kelsoe, J., Sklar, P., 2018. Genomewide association study identifies 30 loci associated with bipolar disorder. bioRxiv.  
URL <https://www.biorxiv.org/content/early/2018/01/24/173062>
- Stahl, E. A., Wegmann, D., Trynka, G., Gutierrez-Achury, J., Do, R., Voight, B. F., Kraft, P., Chen, R., Kallberg, H. J., Kurzeeman, F. A., et al., 2012. Bayesian inference analyses of the polygenic architecture of rheumatoid arthritis. *Nature genetics* 44 (5), 483–489.
- Su, Z., Marchini, J., Donnelly, P., 2011. Hapgen2: simulation of multiple disease snps. *Bioinformatics* 27 (16), 2304–2305.
- Sullivan, P. F., Kendler, K. S., Neale, M. C., 2003. Schizophrenia as a complex trait: evidence from a meta-analysis of twin studies. *Archives of general psychiatry* 60 (12), 1187–1192.
- Sveinbjornsson, G., Albrechtsen, A., Zink, F., Gudjonsson, S. A., Oddson, A., Måsson, G., Holm, H., Kong, A., Thorsteinsdottir, U., Sulem, P., et al., 2016. Weighting sequence variants based on their annotation increases power of whole-genome association studies. *Nature genetics*.
- Thompson, W. K., Wang, Y., Schork, A., Zuber, V., Andreassen, O. A., Dale, A. M., Holland, D., Shujing, X., 2015. An empirical bayes method for estimating the distribution of effects in genome-wide association studies. *PLoS Genetics* [in press].
- Van Rheenen, W., Shatunov, A., Dekker, A. M., McLaughlin, R. L., Diekstra, F. P., Pulit, S. L., Van Der Spek, R. A., Vösa, U., De Jong, S., Robinson, M. R., et al., 2016. Genome-wide association analyses identify new risk variants and the genetic architecture of amyotrophic lateral sclerosis. *Nature genetics* 48 (9), 1043.
- Visscher, P. M., Brown, M. A., McCarthy, M. I., Yang, J., 2012. Five years of gwas discovery. *The American Journal of Human Genetics* 90 (1), 7–24.
- Whiteford, H. A., Degenhardt, L., Rehm, J., Baxter, A. J., Ferrari, A. J., Erskine, H. E., Charlson, F. J., Norman, R. E., Flaxman, A. D., Johns, N., et al., 2013. Global burden of disease attributable to mental and substance use disorders: findings from the global burden of disease study 2010. *The Lancet* 382 (9904), 1575–1586.
- Willer, C. J., Schmidt, E. M., Sengupta, S., Peloso, G. M., Gustafsson, S., Kanoni, S., Ganna, A., Chen, J., Buchkovich, M. L., Mora, S., et al., 2013. Discovery and refinement of loci associated with lipid levels. *Nature genetics* 45 (11), 1274.
- Witte, J. S., Visscher, P. M., Wray, N. R., 2014. The contribution of genetic variants to disease depends on the ruler. *Nature Reviews Genetics* 15 (11), 765–776.
- Wood, A. R., Esko, T., Yang, J., Vedantam, S., Pers, T. H., Gustafsson, S., Chu, A. Y., Estrada, K., Luan, J., Kutalik, Z., et al., 2014. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature genetics* 46 (11), 1173–1186.
- Wray, N. R., Sullivan, P. F., 2017. Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. bioRxiv.

URL <https://www.biorxiv.org/content/early/2017/07/24/167577>

- Wu, C., DeWan, A., Hoh, J., Wang, Z., 2011. A comparison of association methods correcting for population stratification in case-control studies. *Annals of human genetics* 75 (3), 418–427.
- Yang, J., Bakshi, A., Zhu, Z., Hemani, G., Vinkhuyzen, A. A., Lee, S. H., Robinson, M. R., Perry, J. R., Nolte, I. M., van Vliet-Ostaptchouk, J. V., et al., 2015. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nature genetics*.
- Yang, J., Lee, S. H., Goddard, M. E., Visscher, P. M., 2011a. Gcta: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics* 88 (1), 76–82.
- Yang, J., Manolio, T. A., Pasquale, L. R., Boerwinkle, E., Caporaso, N., Cunningham, J. M., De Andrade, M., Feenstra, B., Feingold, E., Hayes, M. G., et al., 2011b. Genome partitioning of genetic variation for complex traits using common snps. *Nature genetics* 43 (6), 519–525.
- Yang, J., Weedon, M. N., Purcell, S., Lettre, G., Estrada, K., Willer, C. J., Smith, A. V., Ingelsson, E., O'connell, J. R., Mangino, M., et al., 2011c. Genomic inflation factors under polygenic inheritance. *European Journal of Human Genetics* 19 (7), 807–812.
- Yengo, L., Sidorenko, J., Kemper, K. E., Zheng, Z., Wood, A. R., Weedon, M. N., Frayling, T. M., Hirschhorn, J., Yang, J., Visscher, P. M., et al., 2018. Meta-analysis of genome-wide association studies for height and body mass index in ~ 700,000 individuals of european ancestry. *bioRxiv*, 274654.
- Zeng, J., Vlaming, R., Wu, Y., Robinson, M. R., Lloyd-Jones, L. R., Yengo, L., Yap, C. X., Xue, A., Sidorenko, J., McRae, A. F., et al., 2018. Signatures of negative selection in the genetic architecture of human complex traits. *Nature genetics* 50 (5), 746.
- Zhang, Y., Qi, G., Park, J.-H., Chatterjee, N., 2018. Estimation of complex effect-size distributions using summary-level statistics from genome-wide association studies across 32 complex traits. *Nature genetics* 50 (9), 1318.
- Zheng, J., Erzurumluoglu, A. M., Elsworth, B. L., Kemp, J. P., Howe, L., Haycock, P. C., Hemani, G., Tansey, K., Laurin, C., Pourcain, B. S., et al., 2017. Ld hub: a centralized database and web interface to perform ld score regression that maximizes the potential of summary level gwas data for snp heritability and genetic correlation analysis. *Bioinformatics* 33 (2), 272–279.
- Zhu, X., Stephens, M., 2017. Bayesian large-scale multiple regression with summary statistics from genome-wide association studies. *The annals of applied statistics* 11 (3), 1561.