RESEARCH ARTICLE

# Estimation of non-null SNP effect size distributions enables the detection of enriched genes underlying complex traits

Wei Cheng[1,2], Sohini Ramachandran[1,2]*, Lorin Crawford[2,3,4]*

**1** Department of Ecology and Evolutionary Biology, Brown University, Providence, Rhode Island, United States of America, **2** Center for Computational Molecular Biology, Brown University, Providence, Rhode Island, United States of America, **3** Department of Biostatistics, Brown University, Providence, Rhode Island, United States of America, **4** Center for Statistical Sciences, Brown University, Providence, Rhode Island, United States of America

* sramachandran@brown.edu (SR); lorin_crawford@brown.edu (LC)

## Abstract

Traditional univariate genome-wide association studies generate false positives and negatives due to difficulties distinguishing associated variants from variants with spurious non-zero effects that do not directly influence the trait. Recent efforts have been directed at identifying genes or signaling pathways enriched for mutations in quantitative traits or case-control studies, but these can be computationally costly and hampered by strict model assumptions. Here, we present gene-ε, a new approach for identifying statistical associations between sets of variants and quantitative traits. Our key insight is that enrichment studies on the gene-level are improved when we reformulate the genome-wide SNP-level null hypothesis to identify spurious small-to-intermediate SNP effects and classify them as non-causal. gene-ε efficiently identifies enriched genes under a variety of simulated genetic architectures, achieving greater than a 90% true positive rate at 1% false positive rate for polygenic traits. Lastly, we apply gene-ε to summary statistics derived from six quantitative traits using European-ancestry individuals in the UK Biobank, and identify enriched genes that are in biologically relevant pathways.

## Author summary

Enrichment tests augment the standard univariate genome-wide association (GWA) framework by identifying groups of biologically interacting mutations that are enriched for associations with a trait of interest, beyond what is expected by chance. These analyses model local linkage disequilibrium (LD), allow many different mutations to be disease-causing across patients, and generate biologically interpretable hypotheses for disease mechanisms. However, existing enrichment analyses are hampered by high computational costs, and rely on GWA summary statistics despite the high false positive rate of the standard univariate GWA framework. Here, we present the gene-level association framework gene-ε (pronounced "genie"), an empirical Bayesian approach for identifying statistical associations between sets of mutations and quantitative traits. The central innovation

of gene-$\varepsilon$ is reformulating the GWA null model to distinguish between *(i)* mutations that are statistically associated with the disease but are unlikely to directly influence it, and *(ii)* mutations that are most strongly associated with a disease of interest. We find that, with our reformulated SNP-level null hypothesis, our gene-level enrichment model outperforms existing enrichment methods in simulation studies and scales well for application to emerging biobank datasets. We apply gene-$\varepsilon$ to six quantitative traits in the UK Biobank and recover novel and functionally validated gene-level associations.

## Introduction

Over the last decade, there has been an evolving debate about the types of insight genome-wide single-nucleotide polymorphism (SNP) genotype data offer into the genetic architecture of complex traits [1–5]. In the traditional genome-wide association (GWA) framework, individual SNPs are tested independently for association with a trait of interest. While this approach can have drawbacks [2, 3, 6], more recent approaches that combine SNPs within a region have gained power to detect biologically relevant genes and pathways enriched for correlations with complex traits [7–14]. Reconciling these two observations is crucial for biomedical genomics.

In the traditional GWA model, each SNP is assumed to either *(i)* directly influence (or perfectly tag a variant that directly influences) the trait of interest; or *(ii)* have no affect on the trait at all (see Fig 1A). Throughout this manuscript, for simplicity, we refer to SNPs under the former as "associated" and those under latter as "non-associated". These classifications are based on ordinary least squares (OLS) effect size estimates for each SNP in a regression framework, where the null hypothesis assumes that the true effects of non-associated SNPs are zero ($H_0$: $\beta_j$ = 0). The traditional GWA model is agnostic to trait architecture, and is underpowered with a high false-positive rate for "polygenic" traits or traits which are generated by many mutations of small effect [5, 15–17].

Suppose that in truth each SNP in a GWA dataset instead belongs to one of *three* categories depending on the underlying distribution of their effects on the trait of interest: *(i)* associated SNPs; *(ii)* non-associated SNPs that emit spurious nonzero statistical signals; and *(iii)* non-associated SNPs with zero-effects (Fig 1B) [18]. Associated SNPs may lie in enriched genes that directly influence the trait of interest. The phenomenon of a non-associated SNP emitting nonzero statistical signal can occur due to multiple reasons. For example, spurious nonzero SNP effects can be due to some varying degree of linkage disequilibrium (LD) with associated SNPs [19]; or alternatively, non-associated SNPs can have a trans-interaction effect with SNPs located within an enriched gene. In either setting, spurious SNPs can emit small-to-intermediate statistical noise (in some cases, even appearing indistinguishable from truly associated SNPs), thereby confounding traditional GWA tests (Fig 1B). Hereafter, we refer to this noise as "epsilon-genic effects" (denoted in shorthand as "$\varepsilon$-genic effects"). There is a need for a computational framework that has the ability to identify mutations associated with a wide range of traits, regardless of whether narrow-sense heritability is sparsely or uniformly distributed across the genome.

Here, we develop a new and scalable quantitative approach for testing aggregated sets of SNP-level GWA summary statistics for enrichment of associated mutations in a given quantitative trait. In practice, our approach can be applied to any user-specified set of genomic regions, such as regulatory elements, intergenic regions, or gene sets. In this study, for simplicity, we refer to our method as a gene-level test (i.e., an annotated collection of SNPs within the
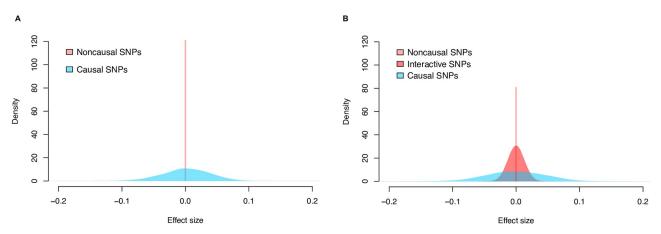
**Fig 1. Illustration of null hypothesis assumptions for the distribution of GWA SNP-level effect sizes according to different views on underlying genetic architectures.** The effect sizes of "non-associated" (pink), "spurious non-associated" (red), and "associated" (blue) SNPs were drawn from normal distributions with successively larger variances. **(A)** The traditional GWA model of complex traits simply assumes SNPs are associated or non-associated. Under the corresponding null hypothesis, associated SNPs are likely to emit nonzero effect sizes while non-associated SNPs will have effect sizes of zero. When there are many causal variants, we refer to the traits as polygenic. **(B)** Under our reformulated GWA model, there are three categories: associated SNPs, non-associated SNPs that emit spurious nonzero effect sizes, and non-associated SNPs with effect sizes of zero. We propose a multi-component framework (see also [18]), in which null SNPs can emit different levels of statistical signals based on (*i*) different degrees of connectedness (e.g., through linkage disequilibrium), or (*ii*) its regulated gene interacts with an enriched gene. While truly associated SNPs are still more likely to emit large effect sizes than SNPs in the other categories, null SNPs can have intermediate effect sizes. Here, our goal is to treat spurious SNPs with small-to-intermediate nonzero effects as being non-associated with the trait of interest.

boundary of a gene). The key contribution of our approach is that gene-level association tests should treat spurious SNPs with $\varepsilon$-genic effects as non-associated variants. Conceptually, this requires assessing whether SNPs explain more than some "epsilon" proportion of the phenotypic variance. In this generalized model, we reformulate the GWA null hypothesis to assume *approximately* no association for spurious non-associated SNPs where

$$H_0: \beta_j \approx 0, \qquad \beta_j \sim \mathcal{N}(0, \sigma_\varepsilon^2), \quad j = 1, \ldots, J \text{ SNPs.}$$

Here, $\sigma_\varepsilon^2$ denotes a "SNP-level null threshold" and represents the maximum proportion of phenotypic variance explained (PVE) that is contributed by spurious non-associated SNPs. This null hypothesis can be equivalently restated as $H_0: \mathbb{E}[\beta_j^2] \leq \sigma_\varepsilon^2$ (Fig 1B). Non-enriched genes are then defined as genes that only contain SNPs with $\varepsilon$-genic effects (i.e., $0 \leq \mathbb{E}[\beta_j^2] \leq \sigma_\varepsilon^2$ for every *j*-th SNP within that region). Enriched genes, on the other hand, are genes that contain at least one associated SNP (i.e., $\mathbb{E}[\beta_j^2] > \sigma_\varepsilon^2$ for at least one SNP *j* within that region). By accounting for the presence of spurious $\varepsilon$-genic effects (i.e., through different values of $\sigma_\varepsilon^2$ which the user can subjectively control), our approach flexibly constructs an appropriate GWA SNP-level null hypothesis for a wide range of traits with genetic architectures that land anywhere on the polygenic spectrum (see Materials and methods).

We refer to our gene-level association framework as "gene-$\varepsilon$" (pronounced "genie"). gene-$\varepsilon$ leverages our modified SNP-level null hypothesis to lower false positive rates and increases power for identifying gene-level enrichment within GWA studies. This happens via two key conceptual insights. First, gene-$\varepsilon$ regularizes observed (and inflated) GWA summary statistics so that SNP-level effect size estimates are positively correlated with the assumed generative model of complex traits. Second, it examines the distribution of regularized effect sizes to offer the user choices for an appropriate SNP-level null threshold $\sigma_\varepsilon^2$ to distinguish associated SNPs from spurious non-associated SNPs. This makes for an improved and refined hypothesis

testing strategy for identifying enriched genes underlying complex traits. With detailed simulations, we assess the power of gene-$\varepsilon$ to identify significant genes under a variety of genetic architectures, and compare its performance against multiple competing approaches [7, 10, 12, 14, 20]. We also apply gene-$\varepsilon$ to the SNP-level summary statistics of six quantitative traits assayed in individuals of European ancestry from the UK Biobank [21].

## Results

### Overview of gene-$\varepsilon$

The gene-$\varepsilon$ framework requires two inputs: GWA SNP-level effect size estimates, and an empirical linkage disequilibrium (LD, or variance-covariance) matrix. The LD matrix can be estimated directly from genotype data, or from an ancestry-matched set of samples if genotype data are not available to the user. We use these inputs to both estimate gene-level contributions to narrow-sense heritability $h^2$, and perform gene-level enrichment tests. After preparing the input data, there are three steps implemented in gene-$\varepsilon$, which are detailed below (Fig 2).

First, we shrink the observed GWA effect size estimates via regularized regression (Fig 2A and 2B; Eq (4) in Materials and methods). This shrinkage step reduces the inflation of OLS effect sizes for spurious SNPs [22], and increases their correlation with the assumed generative model for the trait of interest (particularly for traits with high heritability; S1 Fig). When assessing the performance of gene-$\varepsilon$ in simulations, we considered different types of regularization for the effect size estimates: the Least Absolute Shrinkage And Selection Operator (gene-$\varepsilon$-LASSO) [23], the Elastic Net solution (gene-$\varepsilon$-EN) [24], and Ridge Regression (gene-$\varepsilon$-RR) [25]. We also assessed our framework using the observed ordinary least squares (OLS) estimates without any shrinkage (gene-$\varepsilon$-OLS) to serve as motivation for having regularization as a step in the framework.

Second, we fit a $K$-mixture Gaussian model to all regularized effect sizes genome-wide with the goal of classifying SNPs as associated, non-associated with spurious statistical signal, or non-associated with zero-effects (Figs 1B and 2C; see also [18]). Each successive Gaussian mixture component has distinctly smaller variances ($\sigma_1^2 > \cdots > \sigma_K^2$) with the $K$-th component fixed at $\sigma_K^2 = 0$. Estimating these variance components helps determine an appropriate $k$-th category to serve as the cutoff for SNPs with null effects (i.e., choosing some variance component $\sigma_k^2$ to be the null threshold $\sigma_\varepsilon^2$). The gene-$\varepsilon$ software allows users to determine this cutoff subjectively. Intuitively, enriched genes are likely to contain important variants with relatively larger effects that are categorized in the early-to-middle mixture components. Since the biological interpretation of the middle components may not be consistent across trait architectures, we take a conservative approach in our selection of a cutoff when determining associated SNPs. Without loss of generality, we assume non-null SNPs appear in the first mixture component with the largest variance, while null SNPs appear in the latter components. By this definition, non-associated SNPs with spurious $\varepsilon$-genic or zero-effects then have PVEs that fall at or below the variance of the second component (i.e., $\sigma_\varepsilon^2 = \sigma_2^2$ and $H_0 : \mathbb{E}[\beta_j^2] \leq \sigma_2^2$ for the $j$-th SNP). gene-$\varepsilon$ allows for flexibility in the number of Gaussians that specify the range of null and non-null SNP effects. To achieve genome-wide scalability, we estimate parameters of the $K$-mixture model using an expectation-maximization (EM) algorithm.

Third, we group the regularized GWA summary statistics according to gene boundaries (or user-specified SNP-sets) and compute a gene-level enrichment statistic based on a commonly used quadratic form (Fig 2D) [7, 12, 20]. In expectation, these test statistics can be naturally interpreted as the contribution of each gene to the narrow-sense heritability. We use Imhof's method [26] to derive a $P$-value for assessing evidence in support of an association between a

**Fig 2. Schematic overview of gene-$\varepsilon$: Our new gene-level association approach accounting for spurious nonzero SNP-level effects.** **(A)** gene-$\varepsilon$ takes SNP-level GWA marginal effect sizes (OLS estimates $\hat{\boldsymbol{\beta}}$) and a linkage disequilibrium (LD) matrix ($\boldsymbol{\Sigma}$) as input. It is well-known that OLS effect size estimates are inflated due to LD (i.e., correlation structures) among genome-wide genotypes. **(B)** gene-$\varepsilon$ first uses its inputs to derive regularized effect size estimates ($\tilde{\boldsymbol{\beta}}$) through shrinkage methods (LASSO, Elastic Net and Ridge Regression; we explore performance of each solution under a variety of simulated trait architectures in Supporting Information). **(C)** A unique feature of gene-$\varepsilon$ is that it treats SNPs with spurious nonzero effects as non-associated. gene-$\varepsilon$ assumes a reformulated null distribution of SNP-level effects $\tilde{\beta}_j \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, where $\sigma_\varepsilon^2$ is the SNP-level null threshold and represents the maximum proportion of phenotypic variance explained (PVE) by a spurious or non-associated SNP. This leads to the reformulated SNP-level null hypothesis $H_0 : \mathbb{E}[\beta_j^2] \leq \sigma_\varepsilon^2$. To infer an appropriate $\sigma_\varepsilon^2$, gene-$\varepsilon$ fits a $K$-mixture of normal distributions over the regularized effect sizes with successively smaller variances ($\sigma_1^2 > \cdots > \sigma_K^2$; with $\sigma_K^2 = 0$). In this study (without loss of generality), we assume that associated SNPs will appear in the first set, while spurious and non-associated SNPs appear in the latter sets. By definition, the SNP-level null threshold is then $\sigma_\varepsilon^2 = \sigma_2^2$. **(D)** Lastly, gene-$\varepsilon$ computes gene-level association test statistics $\tilde{Q}_g$ using quadratic forms and corresponding $P$-values using Imhof's method. This assumes the common gene-level null $H_0 : Q_g = 0$, where the null distribution of $Q_g$ is dependent upon the SNP-level null threshold $\sigma_\varepsilon^2$. For more details, see Materials and methods.

https://doi.org/10.1371/journal.pgen.1008855.g002

given gene and the trait of interest. Details for each of these steps can be found in Materials and Methods, as well as in Supporting Information.

## Performance comparisons in simulation studies

To assess the performance of gene-$\varepsilon$, we simulated complex traits under multiple genetic architectures using real genotype data on chromosome 1 from individuals of European ancestry in the UK Biobank (Materials and methods). Following quality control procedures, our simulations included 36,518 SNPs (Supporting Information). Next, we used the NCBI's Reference Sequence (RefSeq) database in the UCSC Genome Browser [27] to annotate SNPs with

the appropriate genes. Simulations were conducted using two different SNP-to-gene assignments. In the first, we directly used the UCSC annotations which resulted in 1,408 genes to be used in the simulation study. In the second, we augmented the UCSC gene boundaries to include SNPs within ±50kb, which resulted in 1,916 genes in the simulation study. For both cases, we assumed a linear additive model for quantitative traits, while varying the following parameters: sample size ($N$ = 5,000 or 10,000); narrow-sense heritability ($h^2$ = 0.2 or 0.6); and the percentage of enriched genes (set to 1% or 10%). In each scenario, we considered traits being generated with and without additional population structure. In the latter setting, traits are simulated while also using the top ten principal components of the genotype matrix as covariates to create stratification. Regardless of the setting, GWA summary statistics were computed by fitting a single-SNP univariate linear model (via OLS) without any control for population structure. Comparisons were based on 100 different simulated runs for each parameter combination.

We compared the performance of gene-$\varepsilon$ against that of five competing gene-level association or enrichment methods: SKAT [20], VEGAS [7], MAGMA [10], PEGASUS [12], and RSS [14] (Supporting Information). As previously noted, we also explored the performance of gene-$\varepsilon$ while using various degrees of regularization on effect size estimates, with gene-$\varepsilon$-OLS being treated as a baseline. SKAT, VEGAS, and PEGASUS are frequentist approaches, in which SNP-level GWA $P$-values are drawn from a correlated chi-squared distribution with covariance estimated using an empirical LD matrix [28]. MAGMA is also a frequentist approach in which gene-level $P$-values are derived from distributions of SNP-level effect sizes using an $F$-test [10]. RSS is a Bayesian model-based enrichment method which places a likelihood on the observed SNP-level GWA effect sizes (using their standard errors and LD estimates), and assumes a spike-and-slab shrinkage prior on the true SNP effects [29]. Conceptually, SKAT, MAGMA, VEGAS, and PEGASUS assume null models under the traditional GWA framework, while RSS and gene-$\varepsilon$ allow for traits to have architectures with more complex SNP effect size distributions.

For all methods, we assess the power and false discovery rates (FDR) for identifying correct genes at a Bonferroni-corrected threshold ($P$ = 0.05/1408 genes = 3.55×10$^{-5}$ and $P$ = 0.05/1916 genes = 2.61×10$^{-5}$, depending on if the ±50kb buffer was used) or median probability model (posterior enrichment probability >0.5; see [30]) (S1–S16 Tables). We also compare their ability to rank true positives over false positives via receiver operating characteristic (ROC) and precision-recall curves (Fig 3 and S2–S16 Figs). While we find gene-$\varepsilon$ and RSS have the best tradeoff between true and false positive rates, RSS does not scale well for genome-wide analyses (Table 1). In many settings, gene-$\varepsilon$ has similar power to RSS (while maintaining a considerably lower FDR), and generally outperforms RSS in precision-versus-recall. gene-$\varepsilon$ also stands out as the best approach in scenarios where the observed OLS summary statistics were produced without first controlling for confounding stratification effects in more heritable traits (i.e., $h^2$ = 0.6). Computationally, gene-$\varepsilon$ gains speed by directly assessing evidence for rejecting the gene-level null hypothesis, whereas RSS must compute the posterior probability of being an enriched gene (which can suffer from convergence issues; Supporting Information). For context, an analysis of just 1,000 genes takes gene-$\varepsilon$ an average of 140 seconds to run on a personal laptop, while RSS takes around 9,400 seconds to complete.

When using GWA summary statistics to identify genotype-phenotype associations, modeling the appropriate trait architecture is crucial. As expected, all methods we compared in this study have relatively more power for traits with high $h^2$. However, our simulation studies confirm the expectation that the max utility for methods assuming the traditional GWA framework (i.e., SKAT, MAGMA, VEGAS, and PEGASUS) is limited to scenarios where heritability is low, phenotypic variance is dominated by just a few enriched genes with large effects, and

**Fig 3. Receiver operating characteristic (ROC) and precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations ($N$ = 10, 000; $h^2$ = 0.6).** We simulate complex traits under different genetic architectures and GWA study scenarios, varying the following parameters: narrow sense heritability, proportion of associated genes, and sample size (Supporting Information). Here, the sample size $N$ = 10, 000 and the narrow-sense heritability $h^2$ = 0.6. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of this step (labeled OLS; orange). We further compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% associated genes) and polygenic (10% associated genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% associated genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates.

https://doi.org/10.1371/journal.pgen.1008855.g003

summary statistics are not confounded by population structure (S2, S3, S9, and S10 Figs). RSS, gene-$\varepsilon$-EN, and gene-$\varepsilon$-LASSO robustly outperform these methods for the other trait architectures (Fig 3, S4–S8 and S11–S16 Figs). One major reason for this result is that shrinkage and penalized regression methods appropriately correct for inflation in GWA summary statistics (S1 Fig). For example, we find that the regularization used by gene-$\varepsilon$-EN and gene-$\varepsilon$-LASSO is able to recover effect size estimates that are almost perfectly correlated ($r^2 > 0.9$) with the true effect sizes used to simulate sparse architectures (e.g., simulations with 1% enriched genes). In S17–S24 Figs, we show a direct comparison between gene-$\varepsilon$ with and without regularization to show how inflated SNP-level summary statistics directly affect the ability to identify enriched genes across different trait architectures. Regularization also allows gene-$\varepsilon$ to preserve type 1

**Table 1. Computational time for running gene-$\varepsilon$ and other gene-level association approaches, as a function of the total number genes analyzed and the number of SNPs within each gene.** Methods compared include: gene-$\varepsilon$, PEGASUS [12], VEGAS [7], RSS [14], MAGMA [10], and SKAT [20]. Here, we simulated 10 datasets for each pair of parameter values (number of genes analyzed, and number of SNPs within each gene). Each table entry represents the average computation time (in seconds) it takes each approach to analyze a dataset of the size indicated. Run times were measured on a MacBook Pro (Processor: 3.1-gigahertz (GHz) Intel Core i5, Memory: 8GB 2133-megahertz (MHz) LPDDR3). Only a single core on the machine was used. PEGASUS, SKAT, and MAGMA are score-based methods and, thus, are expected to take the least amount of time to run. Both gene-$\varepsilon$ and RSS are regression-based methods, but gene-$\varepsilon$ is scalable in both the number of genes and the number of SNPs per gene. The increased computational burden of RSS results from its need to do Bayesian posterior inference; however, gene-$\varepsilon$ is able to scale because it leverages regularization and point estimation for hypothesis testing.

| # Total Genes | # SNPs per Gene | Average Time (sec) | | | | | |
|---|---|---|---|---|---|---|---|
| | | gene-$\varepsilon$ | PEGASUS | VEGAS | RSS | MAGMA | SKAT |
| 250 | 5 | 2.18 | 2.99 | 39.18 | 3.33 | <0.10 | 1.17 |
| | 10 | 4.34 | 1.55 | 57.22 | 13.81 | <0.10 | 1.90 |
| | 20 | 12.94 | 1.22 | 85.54 | 55.49 | <0.10 | 3.63 |
| 500 | 5 | 8.62 | 6.10 | 77.35 | 14.70 | <0.10 | 2.25 |
| | 10 | 16.00 | 3.37 | 106.05 | 56.38 | <0.10 | 4.08 |
| | 20 | 37.88 | 2.52 | 194.21 | 248.90 | <0.10 | 7.07 |
| 1000 | 5 | 25.89 | 11.81 | 152.12 | 60.11 | 0.28 | 4.87 |
| | 10 | 40.69 | 6.33 | 200.78 | 250.51 | 0.58 | 8.59 |
| | 20 | 136.96 | 6.87 | 284.97 | 9410.37 | 1.19 | 14.21 |

https://doi.org/10.1371/journal.pgen.1008855.t001

error when traits are generated under the null hypothesis of no gene enrichment. Importantly, our method is relatively conservative when GWA summary statistics are less precise and derived from studies with smaller sample sizes (e.g., $N = 5,000$; S17 Table).

## Characterizing genetic architecture of quantitative traits in the UK Biobank

We applied gene-$\varepsilon$ to 1,070,306 genome-wide SNPs and six quantitative traits—height, body mass index (BMI), mean red blood cell volume (MCV), mean platelet volume (MPV), platelet count (PLC), waist-hip ratio (WHR)—assayed in 349,414 European-ancestry individuals in the UK Biobank (Supporting Information) [21]. After quality control, we regressed the top ten principal components of the genotype data onto each trait to control for population structure, and then we derived OLS SNP-level effect sizes using the traditional GWA framework. For completeness, we then analyzed these GWA effect size estimates with the four different implementations of gene-$\varepsilon$. In the main text, we highlight results under the Elastic Net solution; detailed findings with the other gene-$\varepsilon$ approaches can be found in Supporting Information.

While estimating $\varepsilon$-genic effects, gene-$\varepsilon$ provides insight into to the genetic architecture of a trait (S18 Table). For example, past studies have shown human height to have a higher narrow-sense heritability (estimates ranging from 45-80%; [6, 31–39]). Using Elastic Net regularized effect sizes, gene-$\varepsilon$ estimated approximately 11% of SNPs in the UK Biobank to be statistically associated with height. This meant approximately 110,000 SNPs had marginal PVEs $\mathbb{E}[\beta_j^2] > 0$ (Materials and methods). This number is similar to the 93,000 and 100,000 height associated variants previously estimated by Goldstein [40] and Boyle et al. [4], respectively. Additionally, gene-$\varepsilon$ identified approximately 2% of SNPs to be "causal" (meaning they had PVEs greater than the SNP-level null threshold, $\mathbb{E}[\beta_j^2] > \sigma_2^2$); again similar to the Boyle et al. [4] estimate of 3.8% causal SNPs for height using data from the GIANT Consortium [32], and the Lello et al. [41] estimate of 3.1% causal SNPs for height using European-ancestry individuals in the UK Biobank.

Compared to body height, narrow-sense heritability estimates for BMI have been considered both high and low (estimates ranging from 25-60%; [31, 33, 34, 36, 37, 39, 42–45]). Such

inconsistency is likely due to difference in study design (e.g., twin, family, population-based studies), many of which have been known to produce different levels of bias [44]. Here, our results suggest BMI to have a lower narrow-sense heritability than height, with a slightly different distribution of null and non-null SNP effects. Specifically, we found BMI to have 13% associated SNPs and 6% causal SNPs.

In general, we found our genetic architecture characterizations in the UK Biobank to reflect the same general themes we saw in the simulation study. Less aggressive shrinkage approaches (e.g., OLS and Ridge) are subject to misclassifications of associated, spurious, and non-associated SNPs. As a result, these methods struggle to reproduce well-known narrow-sense heritability estimates from the literature, across all six traits. This once again highlights the need for computational frameworks that are able to appropriately correct for inflation in summary statistics.

## gene-$\varepsilon$ identifies refined list of genetic enrichments

Next, we applied gene-$\varepsilon$ to the summary statistics from the UK Biobank and generated genome-wide gene-level association $P$-values (Fig 4A and 4B, S25A–S29A and S25B–S29B Figs). As in the simulation study, we conducted two separate analyses using two different SNP-to-gene annotations: *(i)* we used the RefSeq database gene boundary definitions directly, or *(b)* we augmented the gene boundaries by adding SNPs within a ±50 kilobase (kb) buffer to account for possible regulatory elements. A total of 14,322 genes were analyzed when using the UCSC boundaries as defined, and a total of 17,680 genes were analyzed when including the 50kb buffer. The ultimate objective of gene-$\varepsilon$ is to identify enriched genes, which we define as



**Fig 4. Gene-level association results from applying gene-$\varepsilon$ to body height (panels A and C) and mean platelet volume (MPV; panels B and D), assayed in European-ancestry individuals in the UK Biobank.** Body height has been estimated to have a narrow-sense heritability $h^2$ in the range of 0.45 to 0.80 [6, 31–39]; while, MPV has been estimated to have $h^2$ between 0.50 and 0.70 [33, 34, 58]. Manhattan plots of gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes for **(A)** body height and **(B)** MPV. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 3.49 \times 10^{-6}$ correcting for 14,322 autosomal genes analyzed). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes overlapping with the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$. We highlight categories with $Q$-values (i.e., false discovery rates) less than 0.05 and annotate corresponding genes in the Manhattan plots in **(A)** and **(B)**, respectively. For height, the only significant dbGAP category is "Body Height", with nine of the genes identified by gene-$\varepsilon$ appearing in this category. For MPV, the two significant dbGAP categories are "Platelet Count" and "Face"—the first of which is directly connected to trait [57, 60, 61].

https://doi.org/10.1371/journal.pgen.1008855.g004

containing at least one associated SNP and achieving a gene-level association *P*-value below a Bonferroni-corrected significance threshold (in our two analyses, $P = 0.05/14322$ genes $= 3.49 \times 10^{-6}$ and $P = 0.05/17680$ genes $2.83 \times 10^{-6}$, respectively; S19–S24 Tables). As a validation step, we compared gene-$\varepsilon$ *P*-values to RSS posterior enrichment probabilities for each gene. We also used the gene set enrichment analysis tool Enrichr [46] to identify dbGaP categories with an overrepresentation of significant genes reported by gene-$\varepsilon$ (Fig 4C and 4D, S25C–S29C and S25D–S29D Figs). A comparison of gene-level associations and gene set enrichments between the different gene-$\varepsilon$ approaches are also listed (S25–S27 Tables).

Many of the candidate enriched genes we identified by applying gene-$\varepsilon$ were not previously annotated as having trait-specific associations in either dbGaP or the GWAS catalog (Fig 4); however, many of these same candidate genes have been identified by past publications as related to the phenotype of interest (Table 2). It is worth noting that multiple genes would not have been identified by standard GWA approaches since the top SNP in the annotated region had a marginal association below a genome-wide threshold (see Table 2 and highlighted rows in S19–S24 Tables). Additionally, 45% of the genes selected by gene-$\varepsilon$ were also selected by RSS. For example, gene-$\varepsilon$ reports *C1orf150* as having a significant gene-level association with MPV ($P = 1 \times 10^{-20}$ and RSS posterior enrichment probability of 1), which is known to be associated with germinal center signaling and the differentiation of mature B cells that mutually activate platelets [47–49]. Importantly, nearly all of the genes reported by gene-$\varepsilon$ had evidence of overrepresentation in gene set categories that were at least related to the trait of interest. As expected, the top categories with Enrichr *Q*-values smaller than 0.05 for height and MPV were "Body Height" and "Platelet Count", respectively. Even for the less heritable MCV, the top significant gene sets included hematological categories such as "Transferrin", "Erythrocyte Indices", "Hematocrit", "Narcolepsy", and "Iron"—all of which have verified and clinically relevant connections to trait [50–57].

Lastly, gene-$\varepsilon$ also identified genes with rare causal variants. For example, *ZNF628* (which is not mapped to height in the GWAS catalog) was detected by gene-$\varepsilon$ with a significant *P*-value of $1 \times 10^{-20}$ (and $P = 4.58 \times 10^{-8}$ when the gene annotation included a 50kb buffer). Previous studies have shown a rare variant *rs147110934* within this gene to significantly affect adult height [38]. Rare and low-frequency variants are generally harder to detect under the traditional GWA framework. However, rare variants have been shown to be important for explaining the variation of complex traits [28, 39, 80–83]. With regularization and testing for spurious $\varepsilon$-genic effects, gene-$\varepsilon$ is able to distinguish between rare variants that are causal and SNPs with larger effect sizes due various types of correlations. This only enhances the power of gene-$\varepsilon$ to identify potential novel enriched genes.

## Discussion

During the past decade, it has been repeatedly observed that the traditional GWA framework can struggle to accurately differentiate between associated and spurious SNPs (which we define as SNPs that covary with associated SNPs but do not directly influence the trait of interest). As a result, the traditional GWA approach is prone to generating false positives, and detects variant-level associations spread widely across the genome rather than aggregated sets in disease-relevant pathways [4]. While this observation has spurred to many interesting lines of inquiry—such as investigating the role of rare variants in generating complex traits [9, 28, 80, 81], comparing the efficacy of tagging causal variants in different ancestries [84, 85], and integrating GWA data with functional -omics data [86–88]—the focus of GWA studies and studies integrating GWA data with other -omics data is still largely based on the role of individual variants, acting independently.

**Table 2. Top three newly identified candidate genes reported by gene-$\varepsilon$ for the six quantitative traits studied in the UK Biobank (using imputed genotypes with gene boundaries defined by the NCBI's RefSeq database in the UCSC Genome Browser [27]).** We call these novel candidate genes because they are not listed as being associated with the trait of interest in either the GWAS catalog or dbGaP, and they have top posterior enrichment probabilities with the trait using RSS analysis. Each gene is annotated with past functional studies that link them to the trait of interest. We also report each gene's overall trait-specific significance rank (out of 14,322 autosomal genes analyzed for each trait), as well as their heritability estimates from gene-$\varepsilon$ using Elastic Net to regularize GWA SNP-level effect size estimates. The traits are: height; body mass index (BMI); mean corpuscular volume (MCV); mean platelet volume (MPV); platelet count (PLC); and waist-hip ratio (WHR). ♣: Enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P = 4.67 \times 10^{-8}$ correcting for 1,070,306 SNPs analyzed; see highlighted rows in S19–S24 Tables for complete list). *: Multiple genes were tied for this ranking.

| Trait | Gene | Chr | gene-$\varepsilon$ P-Value | Rank | $h_g^2$ | Post. Prob. | Biological Relevance to Trait | Ref(s) |
|-------|------|-----|----------------------------|------|---------|-------------|-------------------------------|--------|
| Height | EZH2 | 7 | $9.34 \times 10^{-8}$ | 61 | $7.23 \times 10^{-3}$ | 1.000 | Associated with diseases Adamantinoma of Long Bone and Weaver Syndrome (characterized by rapid growth). | [62] |
| Height | C17orf42 | 17 | $5.38 \times 10^{-9}$ | 52 | $4.54 \times 10^{-3}$ | 1.000 | Known as the transcription elongation factor of mitochondria (TEFM) which regulates transcription and can affect body height. | [63] |
| Height | KISS1R | 19 | $1 \times 10^{-20}$ | 1* | $5.27 \times 10^{-4}$ | 0.970 | Associated with disorders of puberty and final height. | [64] |
| BMI | ZC3H4 | 19 | $1.62 \times 10^{-14}$ | 20 | $7.84 \times 10^{-3}$ | 1.000 | BMI-inducer known to be associated with adiposity and obesity. | [65–68] |
| BMI | PTOV1 | 19 | $1 \times 10^{-20}$ | 1* | $2.26 \times 10^{-3}$ | 0.990 | Found to be overexpressed in prostate adenocarcinomas which can be induced by obesity. | [69] |
| BMI | FBXO45♣ | 3 | $6.52 \times 10^{-7}$ | 23 | $1.82 \times 10^{-3}$ | 0.029 | Reported to be involved in children syndromic obesity. | [70] |
| MCV | SLC24A1 | 15 | $1.74 \times 10^{-7}$ | 50 | $4.66 \times 10^{-3}$ | 0.140 | Encoded protein is involved in glucose transportation pathway and MCV is reported to be associated with glucose level. | [69] |
| MCV | PDX1♣ | 13 | $1 \times 10^{-20}$ | 1* | $2.31 \times 10^{-4}$ | 0.019 | Associated with Glycated hemoglobin which is affected by MCV | [71] |
| MCV | RHOD | 11 | $1 \times 10^{-20}$ | 1* | $3.35 \times 10^{-4}$ | 0.002 | Associated with Wiskott-Aldrich Syndrome which is characterized by abnormal immune system function (immune deficiency) and a reduced ability to form blood clots. | [69, 72] |
| MPV | C1orf150 | 1 | $1 \times 10^{-20}$ | 1* | $3.44 \times 10^{-2}$ | 1.000 | Known as GCSAML which is involved with germinal center signaling and differentiation of mature B cells that mutually activate platelets. | [47–49] |
| MPV | KIAA0922 | 4 | $3.20 \times 10^{-6}$ | 64 | $7.17 \times 10^{-3}$ | 1.000 | Known as TMEM131L which is associated with canonical Wnt signaling and can effect platelet formation. | [73, 74] |
| MPV | TPT1♣ | 13 | $1 \times 10^{-20}$ | 1* | $3.25 \times 10^{-4}$ | 0.051 | mRNA expression is identified in platelets. | [69] |
| PLC | C1orf150 | 1 | $1 \times 10^{-20}$ | 1* | $2.51 \times 10^{-2}$ | 1.000 | Known as GCSAML which is involved with germinal center signaling and differentiation of mature B cells that mutually activate platelets. | [47–49] |
| PLC | PSMD2 | 3 | $1.42 \times 10^{-9}$ | 29 | $7.40 \times 10^{-3}$ | 1.000 | Also known as the 26S proteasome which is found to be important for platelet production. | [69] |
| PLC | APOB48R | 16 | $1 \times 10^{-20}$ | 1* | $1.36 \times 10^{-3}$ | 0.003 | Involved in Lipoprotein metabolism pathway which can affect platelet. | [69] |
| WHR | TFAP2B | 6 | $3.92 \times 10^{-7}$ | 21 | $3.60 \times 10^{-3}$ | 1.000 | Dietary protein associated with weight maintenance. | [67, 75] |
| WHR | WDR68 | 17 | $1.05 \times 10^{-7}$ | 20 | $1.10 \times 10^{-3}$ | 0.990 | Also known as DCAF7 which has been shown to bind Huntingtin-associated protein 1 (HAP1) and affect weight. | [76] |
| WHR | MLL | 11 | $8.14 \times 10^{-8}$ | 19 | $2.43 \times 10^{-3}$ | 0.940 | Orthologous gene in mice that affects skeleton, body size, and growth. | [67, 77–79] |

Here, our objective is to identify biologically significant underpinnings of the genetic architecture of complex traits by modifying the traditional GWA null hypothesis from $H_0: \beta_j = 0$ (i.e., the $j$-th SNP has zero statistical association with the trait of interest) to $H_0: \beta_j \approx 0$. We accomplish this by testing for $\varepsilon$-genic effects: spurious small-to-intermediate effect sizes emitted by truly non-associated SNPs. We use an empirical Bayesian approach to learn the effect size distributions of null and non-null SNP effects, and then we aggregate (regularized) SNP-level association signals into a gene-level test statistic that represents the gene's contribution to the narrow-sense heritability of the trait of interest. Together, these two steps reduce false positives and increase power to identify the mutations, genes, and pathways that directly influence a trait's genetic architecture. By considering different thresholds for what constitutes a null SNP effect (i.e., different values of $\sigma_\varepsilon^2$ for spurious non-associated SNPs; Figs 1 and 2), gene-$\varepsilon$ offers the flexibility to construct an appropriate null hypothesis for a wide range of traits with

genetic architectures that land anywhere on the polygenic spectrum. It is important to stress that while we repeatedly point to our improved ability distinguish "causal" variants in enriched genes, gene-$\varepsilon$ is by no means a causal inference procedure. Instead, it is an association test which highlights genes in enriched pathways that are most likely to be associated with the trait of interest.

Through simulations, we showed the gene-$\varepsilon$ framework outperforms other widely used gene-level association methods (particularly for highly heritable traits), while also maintaining scalability for genome-wide analyses (Fig 3, S2–S24 Figs, Table 1, and S1–S17 Tables). Indeed, all the approaches we compared in this study showed improved performance when they used summary statistics derived from studies with larger sample sizes (i.e., simulations with $N = 10,000$). This is because the quality of summary statistics also improves in these settings (via the asymptotic properties of OLS estimates). Nonetheless, our results suggest that applying gene-$\varepsilon$ to summary statistics from previously published studies will increase the return made on investments in GWA studies over the last decade.

Like any aggregated SNP-set association method, gene-$\varepsilon$ has its limitations. Perhaps the most obvious limitation is that annotations can bias the interpretation of results and lead to erroneous scientific conclusions (i.e., might cause us to highlight the "wrong" gene [14, 89, 90]). We observed some instances of this during the UK Biobank analyses. For example, when studying MPV, *CAPN10* only appeared to be a significant gene after its UCSC annotated boundary was augmented by a ±50kb buffer window ($P = 1.85 \times 10^{-1}$ and $P = 1.17 \times 10^{-7}$ before and after the buffer was added, respectively; see S22 Table). After further investigation, this result occurred because the augmented definition of *CAPN10* included nearly all causal SNPs from the significant neighboring gene *RNPEPL1* ($P = 1 \times 10^{-20}$ and $P = 2.07 \times 10^{-9}$ before and after the buffer window was added, respectively). While this shows the need for careful biological interpretation of the results, it also highlights the power of gene-$\varepsilon$ to prioritize true genetic signal effectively.

Another limitation of gene-$\varepsilon$ is that it relies on the user to determine an appropriate SNP-level null threshold $\sigma_\varepsilon^2$ to serve as a cutoff between null and non-null SNP effects. In the current study, we use a *K*-mixture Gaussian model to classify SNPs into different categories and then (without loss of generality) we subjectively assume that associated SNPs only appear in the component with the largest variance (i.e., we choose $\sigma_\varepsilon^2 = \sigma_2^2$). Indeed, there can be many scenarios where this particular threshold choice is not optimal. For example, if there is one very strongly associated locus, the current implementation of the algorithm will assign it to its own mixture component and all other SNPs will be assumed to be not associated with the trait, regardless of the size of their corresponding variances. As previously mentioned, one practical guideline would be to select $\sigma_\varepsilon^2$ based on some *a priori* knowledge about a trait's architecture. However, a more robust approach would be to select the SNP-null hypothesis threshold based on the data at hand. One way to do this would be to take a fully Bayesian approach and allow posterior inference on $\sigma_\varepsilon^2$ to be dependent upon how much heritability is explained by SNPs placed in the top few largest components of the normal mixture. Recently, sparse Bayesian parametric [91] and nonparametric [92] Gaussian mixture models have been proposed for improved polygenic prediction with summary statistics. Combining these modeling strategies with our modified SNP-level null hypothesis could make for a more unified and data-driven implementation of the gene-$\varepsilon$ framework.

There are several other potential extensions for the gene-$\varepsilon$ framework. First, in the current study, we only focused on applying gene-$\varepsilon$ to quantitative traits (Fig 4, S25–S29 Figs, Table 2, and S18–S27 Tables). Future studies extending this approach to binary traits (e.g., case-control studies) should explore controlling for additional confounders that can occur within these

phenotypes, such as ascertainment [93–95]. Second, we only focus on data consisting of common variants; however, it would be interesting to extend gene-$\varepsilon$ for (*i*) rare variant association testing and (*ii*) studies that consider the combined effect between rare and common variants. A significant challenge, in either case, would be to adaptively adjust the strength of the regularization penalty on the observed OLS summary statistics for causal rare variants, so as to not misclassify them as spurious non-associated SNPs. Previous approaches with specific re-weighting functions for rare variants may help here [9, 28, 80] (Materials and methods). A final related extension of gene-$\varepsilon$ is to include information about standard errors when estimating $\varepsilon$-genic effects. In our analyses using the UK Biobank, some of the newly identified candidate genes contained SNPs that had large effect sizes but insignificant *P*-values in the original GWA analysis (after Bonferroni-correction; Table 2 and S19–S24 Tables). While this could be attributed to the modified SNP-level null distribution assumed by gene-$\varepsilon$, it also motivates a regularization model that accounts for the standard error of effect size estimates from GWA studies [14, 22, 29].

## Materials and methods

### Traditional association tests using summary statistics

gene-$\varepsilon$ requires two inputs: genome-wide association (GWA) marginal effect size estimates $\hat{\boldsymbol{\beta}}$, and an empirical linkage disequilibrium (LD) matrix $\Sigma$. We assumed the following generative linear model for complex traits

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \tau^2\mathbf{I}), \tag{1}$$

where $\mathbf{y}$ denotes an *N*-dimensional vector of phenotypic states for a quantitative trait of interest measured in *N* individuals; $\mathbf{X}$ is an $N \times J$ matrix of genotypes, with *J* denoting the number of single nucleotide polymorphisms (SNPs) encoded as {0, 1, 2} copies of a reference allele at each locus; $\beta$ is a *J*-dimensional vector containing the additive effect sizes for an additional copy of the reference allele at each locus on $\mathbf{y}$; $\mathbf{e}$ is a normally distributed error term with mean zero and scaled variance $\tau^2$; and $\mathbf{I}$ is an $N \times N$ identity matrix. For convenience, we assumed that the genotype matrix (column-wise) and trait of interest have been mean-centered and standardized. We also treat $\boldsymbol{\beta}$ as a fixed effect. A central step in GWA studies is to infer $\boldsymbol{\beta}$ for each SNP, given both genotypic and phenotypic measurements for each individual sample. For every SNP *j*, gene-$\varepsilon$ takes in the ordinary least squares (OLS) estimates based on Eq (1)

$$\hat{\beta}_j = (\mathbf{x}_j^{\mathsf{T}}\mathbf{x}_j)^{-1}\mathbf{x}_j^{\mathsf{T}}\mathbf{y}, \tag{2}$$

where $\mathbf{x}_j$ is the *j*-th column of the genotype matrix $\mathbf{X}$, and $\hat{\beta}_j$ is the *j*-th entry of the vector $\hat{\boldsymbol{\beta}}$. In traditional GWA studies, the null hypothesis for statistical association tests assumes $H_0$: $\beta_j = 0$ for all $j = 1, \ldots, J$ SNPs. It can be shown that two genotypic variants $\mathbf{x}_j$ and $\mathbf{x}_{j'}$ in linkage disequilibrium (LD) will produce effect size estimates $\hat{\beta}_j$ and $\hat{\beta}_{j'}$ ($j \neq j'$) that are correlated [29]. This can lead to confounded statistical tests. For the applications considered here, the LD matrix is empirically estimated from external data (e.g., directly from GWA study data, or using an LD map from a population with similar genomic ancestry to that of the samples analyzed in the GWA study).

### Regularized regression for GWA summary statistics

gene-$\varepsilon$ uses regularization on the observed GWA summary statistics to reduce inflation of SNP-level effect size estimates and increase their correlation with the assumed generative

model of complex traits. For large sample size $N$, note that the asymptotic relationship between the observed GWA effect size estimates $\hat{\boldsymbol{\beta}}$ and the true coefficient values $\boldsymbol{\beta}$ is [18, 96, 97]

$$\mathbb{E}[\hat{\beta}_j] = \sum_{j'=1}^{J} \rho(\mathbf{x}_j, \mathbf{x}_{j'})\beta_{j'} \quad \Leftrightarrow \quad \mathbb{E}[\hat{\boldsymbol{\beta}}] = \boldsymbol{\Sigma}\boldsymbol{\beta}, \tag{3}$$

where $\Sigma_{jj'} = \rho(\mathbf{x}_j, \mathbf{x}_{j'})$ denotes the correlation coefficient between SNPs $\mathbf{x}_j$ and $\mathbf{x}_{j'}$. The above mirrors a high-dimensional regression model with the misestimated OLS summary statistics as the response variables and the LD matrix as the design matrix. Theoretically, the resulting output coefficients from this model are the desired true effect size estimates. Due to the multi-collinear structure of GWA data, we cannot reuse the ordinary least squares solution reliably [98]. Thus, we derive the general regularization

$$\tilde{\boldsymbol{\beta}} = \arg\min_{\boldsymbol{\beta}} \; \| \hat{\boldsymbol{\beta}} - \boldsymbol{\Sigma}\boldsymbol{\beta} \|^2, \quad \text{subject to } (1 - \alpha) \| \boldsymbol{\beta} \|_1 + \alpha \| \boldsymbol{\beta} \|_2^2 \leq t \text{ for some } t, \tag{4}$$

where, in addition to previous notation, the solution $\tilde{\boldsymbol{\beta}}$ is used to denote the regularized solution of the observed GWA effect sizes $\hat{\boldsymbol{\beta}}$; and $\| \bullet \|_1$ and $\| \bullet \|_2^2$ denote $L_1$ and $L_2$ penalties, respectively. The free regularization parameter $t$ is chosen based off a grid $[\log t_{\min}, \log t_{\max}]$ with 100 sequential steps of size 0.01. Here, $t_{\max}$ is the minimum value such that all summary statistics are shrunk to zero. We then select the $t$ that results in a model with an $R^2$ within one standard error of the best fitted model. In other words, we choose the $t$ that (*i*) results in a more sparse solution than the best fitted model, but (*ii*) cannot be distinguished from the best fitted model in terms of overall variance explained.

The term $\alpha$ in Eq (4) distinguishes the type of regularization used, and can be chosen to induce various degrees of shrinkage on the effect size estimates. Specifically, $\alpha = 0$ corresponds to the "Least Absolute Shrinkage and Selection Operator" or LASSO solution [23], $\alpha = 1$ equates to Ridge Regression [25], while $0 < \alpha < 1$ results in the Elastic Net [24]. The LASSO solution forces some inflated coefficients to be zero; while the Ridge shrinks the magnitudes of all coefficients but does not set any of them to be exactly zero. Intuitively, the LASSO will create a regularized set of effect sizes where associated SNPs have larger effects, non-associated SNPs with spurious small-to-intermediate (or $\varepsilon$-genic) effects, and non-associated SNPs with zero-effects. It has been suggested that the $L_1$-penalty can suffer from a lack of stability [99]. Therefore, in the main text, we also highlighted gene-$\varepsilon$ using the Elastic Net (with $\alpha = 0.5$). The Elastic Net is a convex combination of the LASSO and Ridge penalties, but still produces distinguishable sets of associated, spurious, and non-associated SNPs. Note that for large GWA studies (e.g., the UK Biobank analysis in the main text), it can be impractical to construct a genome-wide LD matrix; therefore, we regularize OLS effect size estimates based on partitioned chromosome specific LD matrices. Results comparing each of the gene-$\varepsilon$ regularization implementations are given in the main text (Fig 3) and Supporting Information (S2–S24 Figs, S1–S18 and S25–S27 Tables). We will describe how we approximate the null distribution for these regularized GWA summary statistics over the next two sections.

## Estimating the SNP-level null threshold

The main innovation of gene-$\varepsilon$ is to treat spurious SNPs with $\varepsilon$-genic effects as non-associated. This leads to reformulating the GWA SNP-level null hypothesis to assume non-associated SNPs can make small-to-intermediate contributions to the phenotypic variance. Formally, we write this as

$$H_0 : \beta_j \approx 0, \quad \beta_j \sim \mathcal{N}(0, \sigma_\varepsilon^2), \quad j = 1, \dots, J \tag{5}$$

where $\sigma_\varepsilon^2$ denotes the "SNP-level null threshold" and represents the maximum proportion of phenotypic variance explained (PVE) that is contributed by spurious SNPs. Based on Eq (5), we equivalently say

$$H_0 : \mathbb{E}[\beta_j^2] \leq \sigma_\varepsilon^2. \tag{6}$$

To estimate the threshold $\sigma_\varepsilon^2$ for null SNP-level effects, we use an empirical Bayesian approach and fit a *K*-mixture of normal distributions over the (regularized) effect size estimates [18],

$$\tilde{\beta}_j \,|\, z_j = k \sim \mathcal{N}(0, \sigma_k^2), \quad \Pr[z_j = k] = \pi_k, \tag{7}$$

where $z_j \in \{1, \ldots, K\}$ is a latent variable representing the categorical membership for the *j*-th SNP. When summing over all components, Eq (7) corresponds to the following marginal distribution

$$\tilde{\beta}_j \sim \sum_{k=1}^{K} \pi_k \mathcal{N}(0, \sigma_k^2), \tag{8}$$

where $\pi_k$ is a mixture weight representing the marginal (unconditional) probability that a randomly selected SNP belongs to the *k*-th component, with $\sum_k \pi_k = 1$. The above mixture allows for distinct clusters of nonzero effects through *K* different variance components ($\sigma_k^2, k = 1, \ldots, K$) [18]. Here, we consider sequential fractions ($\pi_1, \ldots, \pi_K$) of SNPs to correspond to distinctly smaller effects ($\sigma_1^2 > \cdots > \sigma_K^2 = 0$) [18]. The goal of the mixture model is to "bin" each of the (regularized) SNP-level effects and determine an appropriate category *k* to serve as the cutoff for SNPs with null effects (i.e., choosing the threshold $\sigma_\varepsilon^2$ based on some $\sigma_k^2$). Such a threshold can be chosen based on *a priori* knowledge about the phenotype of interest. It is intuitive to assume that enriched genes will contain non-null SNPs that classify within the early-to-middle mixture components; unfortunately, the biological interpretations of the middle components may not be consistent across trait architectures. Therefore, without loss of generality in this paper, we take a conservative approach in our definition of associated SNPs within enriched genes. Here, we subjectively set the SNP-level null threshold as $\sigma_\varepsilon^2 = \sigma_2^2$. Thus, non-null SNPs are assumed to appear in the largest fraction (i.e., the alternative $H_A : \mathbb{E}[\beta_j^2] > \sigma_2^2$), while null SNPs with belong to the latter groups (i.e., the null $H_0 : \mathbb{E}[\beta_j^2] \leq \sigma_2^2$). Given Eqs (7) and (8), we write the joint log-likelihood for all *J* SNPs as the following

$$\log p(\tilde{\boldsymbol{\beta}} \,|\, \Theta) = \sum_{j=1}^{J} \log p(\tilde{\beta}_j \,|\, \Theta) = \sum_{j=1}^{J} \log \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(0, \sigma_k^2) \right\}, \tag{9}$$

where $\Theta = (\pi_1, \ldots, \pi_K, \sigma_1^2, \ldots, \sigma_K^2)$ is the complete set of parameters for the mixture model. Since there is not a closed-form solution for the maximum likelihood estimate (MLE), so we use an expectation-maximization (EM) algorithm to estimate the parameters in $\Theta$ [100–102].

**Derivation of the EM algorithm.** To derive an EM solution, we use Eqs (7) and (8) to write the joint distribution of the *J*-regularized SNP-level effect sizes and the *J*-latent random variables $\mathbf{z} = (z_1, \ldots, z_J)$, conditioned on the mixture parameters $\Theta$,

$$p(\tilde{\boldsymbol{\beta}}, \mathbf{z} \,|\, \Theta) = p(\tilde{\boldsymbol{\beta}} \,|\, \mathbf{z}, \Theta) p(\mathbf{z}) = \prod_{j=1}^{J} \prod_{k=1}^{K} [\pi_k \mathcal{N}(0, \sigma_k^2)]^{\mathbb{I}(z_j = k)}, \tag{10}$$

where $\mathbb{I}(z_j = k)$ is an indicator function and equates to one if $z_j = k$ and zero otherwise. Taking the log of this distribution yields the following

$$\log p(\tilde{\boldsymbol{\beta}}, \mathbf{z} \mid \Theta) = \sum_{j=1}^{J} \log p(\tilde{\beta}_j, z_j \mid \Theta) = \sum_{j=1}^{J} \sum_{k=1}^{K} \mathbb{I}(z_j = k)[\log \pi_k + \log \mathcal{N}(0, \sigma_k^2)]. \qquad (11)$$

As opposed to Eq (9), the augmented log-likelihood in Eq (11)) is a much simpler function for which to find a solution. The formal steps of the EM algorithm are now detailed below:

1. **E-Step: Update the probability of fraction assignment.** In the E-step of the EM algorithm, we estimate the probability that the $j$-th SNP belongs to one of the $K$ fraction groups. To begin, we use Bayes theorem to find

$$p(\mathbf{z} \mid \tilde{\boldsymbol{\beta}}, \Theta) \propto p(\tilde{\boldsymbol{\beta}} \mid \mathbf{z}, \Theta) p(\mathbf{z}) = \prod_{j=1}^{J} \prod_{k=1}^{K} [\pi_k \mathcal{N}(0, \sigma_k^2)]^{\mathbb{I}(z_j = k)}. \qquad (12)$$

Next, we take the expectation of the complete log-likelihood $\log p(\tilde{\boldsymbol{\beta}}, \mathbf{z} \mid \Theta)$, with respect to the condtional distribution $p(\mathbf{z} \mid \tilde{\boldsymbol{\beta}}, \Theta)$, under current value of the mixture parameters $\hat{\Theta}$. This yields

$$\mathbb{E}_{\mathbf{z} \mid \tilde{\boldsymbol{\beta}}, \hat{\Theta}}[\log p(\tilde{\boldsymbol{\beta}}, \mathbf{z} \mid \hat{\Theta})] = \sum_{j=1}^{J} \sum_{k=1}^{K} \hat{\gamma}_k^{(j)}[\log \pi_k + \log \mathcal{N}(0, \sigma_k^2)], \qquad (13)$$

where $\hat{\gamma}_k^{(j)}$ is referred to as the "responsibility of the $k$-th mixture component", and is given as

$$\hat{\gamma}_k^{(j)} = \Pr[z_j = k \mid \tilde{\beta}_j, \hat{\Theta}] = \frac{\hat{\pi}_k \mathcal{N}(0, \hat{\sigma}_k^2)}{\sum_{k'=1}^{K} \hat{\pi}_{k'} \mathcal{N}(0, \hat{\sigma}_{k'}^2)}. \qquad (14)$$

Intuitively, the EM algorithm uses the collection of these responsibility values to assign SNPs to one of the $K$ fraction groups. This key step may be interpreted as determining the category of SNP effects (which is determined by identifying the $k$-th component with the largest $\gamma_k^{(j)}$ for each $j$-th SNP).

2. **M-Step: Update the component variances and mixture weights.** In the M-step of the EM algorithm, we now fix the responsibility values and maximize the expectation in Eq (13), with respect to the parameters in $\hat{\Theta}$. Namely, we compute the following closed-form solutions:

$$\hat{\sigma}_k^2 = \frac{1}{J_k} \sum_{j=1}^{J} \hat{\gamma}_k^{(j)} \tilde{\beta}_j^2, \quad \hat{\pi}_k = \frac{J_k}{J} \qquad (15)$$

where $J_k = \sum_j \hat{\gamma}_k^{(j)}$ is the sum of the membership weights for the $k$-th mixture component and represents the number of SNPs assigned to that component. The $\hat{\sigma}_k^2$ estimates are used to set the SNP-level null threshold $\hat{\sigma}_\varepsilon^2$.

The gene-$\varepsilon$ software implements the above EM algorithm using the `mclust` [103] package in R. Results in the main text and Supporting Information are based on 100 iterations from 10 different parallel chains to ensure convergence. To implement the above algorithm, we use the `mclust` software package which can fit a Gaussian mixture with up to $K = 10$ distinct

components (see Software Details). Here, the function will compare the Bayesian Information Criterion (BIC) approximation to the Bayes factor for each possible $K$ [104], and produces a resulting output for the $K$ value that has the largest BIC value. Note that since the EM updates do not involve any large LD matrices, the algorithm scales to be fit efficiently over all SNPs genome-wide.

## Regularized GWA summary statistics under the null hypothesis

With an estimate of the SNP-level null threshold $\sigma_\varepsilon^2$, we now describe the probabilistic distribution of the regularized GWA summary statistics under the null hypothesis. Without loss of generality, we demonstrate this property using the general regularization approach where we fix $\alpha \in [0, 1]$ and have the following (approximate) closed form solution for the regularized effect size estimates [23–25]

$$\tilde{\boldsymbol{\beta}} \simeq \mathbf{H}\hat{\boldsymbol{\beta}}, \quad \mathbf{H} = (\boldsymbol{\Sigma} + \vartheta\mathbf{D}^{-1})^{-1} \tag{16}$$

with $\vartheta \geq 0$ being a penalization parameter that has one-to-one correspondence with $t$ in Eq (4). Here, $\mathbf{H}$ is commonly referred to as the "linear shrinkage estimator", where $\mathbf{D}$ is a diagonal weight matrix with nonzero elements dictated by the type of regularization that is being used. For example, $\mathbf{D} = \mathbf{I}$ while performing ridge regression [25], and $\mathbf{D} = \mathrm{diag}(|\tilde{\beta}_1|, \ldots, |\tilde{\beta}_p|)$ while using ridge-based approximations for the elastic net and lasso solutions [23, 24]. From Eq (16), it is clear that $\tilde{\boldsymbol{\beta}}$ may be interpreted as a marginal estimator of SNP-level effects after accounting for LD structure. Using Eqs (2) and (3), it is straightforward to show the (approximate) relationship between the regularized effect size estimates and the true coefficient values

$$\mathbb{E}[\tilde{\boldsymbol{\beta}}] \simeq \mathbf{H}\boldsymbol{\Sigma}\boldsymbol{\beta}. \tag{17}$$

As described in the main text, the accuracy of this relationship is dependent upon both the sample size and narrow-sense heritability of the trait of interest (S1 Fig). Indeed, if $\boldsymbol{\Sigma}$ is full rank and regularization is no longer implemented (i.e., $\vartheta = 0$), $\tilde{\boldsymbol{\beta}}$ is simply the ordinary least squares solution for marginal GWA summary statistics with asymptotic variance-covariance $\mathbb{V}[\tilde{\boldsymbol{\beta}}] \simeq \boldsymbol{\Sigma}$ under the null model [18, 96, 97]. In the limiting case where the number of observations in a GWA study is large (i.e., $N \to \infty$) and the trait of interest is highly heritable, $\tilde{\boldsymbol{\beta}}$ converges onto $\boldsymbol{\beta}$ in expectation; and thus is assumed to be independently and normally distributed under the null hypothesis with asymptotic variance $\sigma_\varepsilon^2\mathbf{I}$ (previously discussed in Eq (5)). As empirically demonstrated for synthetic traits in the current study, we are rarely in situations where we expect the regularized effect size estimates to have completely converged onto the true generative SNP-level coefficients (again see S1 Fig). This effectively means that we cannot expect each $\tilde{\beta}_j$ to be completely independent under the null hypothesis in practice. We accommodate this realization by assuming that under the null model

$$\mathbb{V}[\tilde{\boldsymbol{\beta}}] = \sigma_\varepsilon^2\boldsymbol{\Sigma}, \qquad \lim_{\sigma_\varepsilon^2 \to 0} \sigma_\varepsilon^2\boldsymbol{\Sigma} = \sigma_\varepsilon^2\mathbf{I}. \tag{18}$$

Our reasoning for the formulation above is that, for most quality controlled studies, SNPs in perfect LD will have been pruned such that $\rho(\mathbf{x}_j, \mathbf{x}_{j'}) < \rho(\mathbf{x}_j, \mathbf{x}_j)$ for all $j \neq j'$ variants in the data. Therefore, when traits are generated under the idealized null scenario with large sample sizes and no genetic effects, the estimate of $\sigma_\varepsilon^2 \to 0$ and the off-diagonals of $\sigma_\varepsilon^2\boldsymbol{\Sigma}$ will approach zero quicker than the diagonal elements; thus, allowing the regularized $\tilde{\boldsymbol{\beta}}$ to asymptotically

converge onto the true coefficients $\boldsymbol{\beta}$. When this scenario does not occur, we are able to appropriately deal with the remaining correlation structure (e.g., all the simulation scenarios explored in this work; see Fig 3, S2–S24 Figs, Table 1, and S1–S17 Tables).

## Using the SNP-level null threshold to detect enriched genes

We now formalize the hypothesis test for identifying significantly enriched genes conditioned on the SNP-level null threshold $\sigma_\varepsilon^2$, which we compute using the variance component estimates from the EM algorithm detailed in the previous section. The gene-$\varepsilon$ gene-level test statistic is based on a quadratic form using GWA summary statistics, which is a common approach for generating gene-level test statistics for complex traits. Let gene (or genomic region) $g$ represent a known set of SNPs $j \in \mathcal{J}_g$; for example, $\mathcal{J}_g$ may include SNPs within the boundaries of $g$ and/or within its corresponding regulatory region. Here, we conformably partition the regularized GWA effect size estimates $\tilde{\boldsymbol{\beta}}$ and define the gene-level test statistic

$$\tilde{Q}_g = \tilde{\boldsymbol{\beta}}_g^\mathsf{T} \mathbf{A} \tilde{\boldsymbol{\beta}}_g, \tag{19}$$

where $\mathbf{A}$ is an arbitrary symmetric and positive semi-definite weight matrix. We set to $\mathbf{A} = \mathbf{I}$ to be the identity matrix for all analyses in the current study; hence, $\tilde{Q}_g$ simplifies to a sum of squared SNP effects in the $g$-th gene. Indeed, similar quadratic forms have been implemented to assess the enrichment of mutations at the gene level [7, 12] and across general SNP-sets [9, 20, 28, 80]. A key feature of the gene-$\varepsilon$ framework is to assess the statistics in Eq (19) against a gene-level enrichment null hypothesis $H_0: Q_g = 0$ that is dependent on the SNP-level null threshold $\sigma_\varepsilon^2$. Due to the normality assumption for each SNP effect in Eq (5), $Q_g$ is theoretically assumed to follow a mixture of chi-square distributions,

$$Q_g \sim \sum_{j=1}^{|\mathcal{J}_g|} \lambda_j \chi_{1,j}^2, \tag{20}$$

where $|\mathcal{J}_g|$ denotes the cardinality of the set of SNPs $\mathcal{J}_g$; $\chi_{1,j}^2$ are standard chi-square random variables with one degree of freedom; and $(\lambda_1, \ldots, \lambda_{|\mathcal{J}_g|})$ are the eigenvalues of the matrix [105, 106]

$$\mathbb{V}[\tilde{\boldsymbol{\beta}}_g]^{1/2} \mathbf{A} \mathbb{V}[\tilde{\boldsymbol{\beta}}_g]^{1/2} = \sigma_\varepsilon^2 \boldsymbol{\Sigma}_g^{1/2} \mathbf{A} \boldsymbol{\Sigma}_g^{1/2}.$$

Again, in the current study, $\sigma_\varepsilon^2 = \hat{\sigma}_2^2$ from the estimates in Eq (15), and $\Sigma_g$ denotes a subset of the LD matrix only containing SNPs annotated in the $g$-th SNP-set. Again, when $\mathbf{A} = \mathbf{I}$, the eigenvalues are based on a scaled version of the local gene-specific LD matrix. Several approximate and exact methods have been suggested to obtain $P$-values under a mixture of chi-square distributions. In this study, we use Imhof's method [26] where we empirically compute an estimate of the weighted sum in Eq (20) and compare this distribution to the observed test statistic in Eq (19) (see Software Details). It is important to note here that the gene-level null hypothesis is the same for gene-$\varepsilon$ and other similar competing enrichment methods [9, 12, 20, 28, 80]; the defining characteristic that sets gene-$\varepsilon$ apart is that it assumes a different null distribution for effects on the SNP-level.

**Estimating gene specific contributions to the PVE.** In the main text, we highlight some of the additional features of the gene-$\varepsilon$ gene-level association test statistic. First, the expected enrichment for trait-associated mutations in a given gene is equal to the heritability explained by the SNPs contained in said gene. Formally, consider the expansion of Eq (19) derived from

the expectation of quadratic forms,

$$\mathbb{E}[\tilde{Q}_g] = \sum_{j=1}^{|\mathcal{J}_g|}\sum_{j'=1}^{|\mathcal{J}_g|} a_{jj'}\mathbb{E}[\tilde{\beta}_j\tilde{\beta}_{j'}] = h_g^2, \tag{21}$$

where $h_g^2$ denotes the heritability contributed by gene $g$. When $\mathbf{A} = \mathbf{I}$ (as in the current study), the gene-$\varepsilon$ hypothesis test for identifying enriched genes is based on the individual SNP contributions to the narrow-sense heritability (i.e., the sum of the expectation of squared SNP effects; see also [34])

$$\mathbb{E}[\tilde{Q}_g] = \sum_{j=1}^{|\mathcal{J}_g|}\mathbb{E}[\tilde{\beta}_j^2] = h_g^2. \tag{22}$$

Alternatively, one could choose to re-weight these contributions by specifying $\mathbf{A}$ otherwise [12, 20, 105, 107, 108]. For example, if SNP $j$ has a small effect size but is known to be functionally associated with the trait of interest, then increasing $\mathbf{A}_{jj}$ will reflect this knowledge. Specific weight functions have also been suggested for dealing with rarer variants [9, 28, 80].

## Simulation studies

We used a simulation scheme to generate SNP-level summary statistics for GWA studies. First, we randomly select a set of enriched genes and assume that complex traits (under various genetic architectures) are generated via a linear model

$$\mathbf{y} = \mathbf{Wb} + \sum_{c\in\mathcal{C}}\mathbf{x}_c\beta_c + \mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \tau^2\mathbf{I}), \tag{23}$$

where $\mathbf{y}$ is an $N$-dimensional vector containing all the phenotypes; $\mathcal{C}$ represents the set of causal SNPs contained within the associated genes; $\mathbf{x}_c$ is the genotype for the $c$-th causal SNP encoded as 0, 1, or 2 copies of a reference allele; $\beta_c$ is the additive effect size for the $c$-th SNP; $\mathbf{W}$ is an $N{\times}M$ matrix of covariates representing additional population structure (e.g., the top ten principal components from the genotype matrix) with corresponding fixed effects $\mathbf{b}$; and $\mathbf{e}$ is an $N$-dimensional vector of environmental noise. The phenotypic variance is assumed $\mathbb{V}[\mathbf{y}] = 1$. The effect sizes of SNPs in enriched genes are randomly drawn from standard normal distributions and then rescaled so they explain a fixed proportion of the narrow-sense heritability $\mathbb{V}[\sum\mathbf{x}_c\beta_c] = h^2$. The covariate coefficients are also drawn from standard normal distributions and then rescaled such that $\mathbb{V}[\mathbf{Wb}] + \mathbb{V}[\mathbf{e}] = (1 - h^2)$. GWA summary statistics are then computed by fitting a single-SNP univariate linear model via ordinary least squares (OLS): $\hat{\beta}_j = (\mathbf{x}_j^\mathsf{T}\mathbf{x}_j)^{-1}\mathbf{x}_j^\mathsf{T}\mathbf{y}$ for every SNP in the data $j = 1, \ldots J$. These effect size estimates, along with an LD matrix $\Sigma$ computed directly from the full $N{\times}J$ genotype matrix $\mathbf{X}$, are given to gene-$\varepsilon$. We also retain standard errors and $P$-values for implementation of the competing methods (VEGAS, PEGASUS, RSS, SKAT, and MAGMA). Given different model parameters, we simulate data mirroring a wide range of genetic architectures (Supporting Information).

## Software details

Source code implementing gene-$\varepsilon$ and tutorials are freely available at https://github.com/ramachandran-lab/genee and was written in R (version 3.3.3). Within this software, regularization of the OLS SNP-level effect sizes is done using the package glmnet (version 2.0-16) [109]. For large datasets, such as the UK Biobank, the software also offers regularization using the biglasso (version 1.3-6) [110] to help with memory and scalability requirements. Note

that selection of the free parameter *t* is done the same way using both the `glmnet` and `biglasso` packages. Both packages also take in an $\alpha \in [0, 1]$ to specify fitting the Ridge, Elastic Net or Lasso regularization to the OLS SNP-level effect sizes. The fitting of a *K*-mixture of Gaussian distributions for the estimation of the SNP-level null threshold $\sigma_\varepsilon^2$ is done using the package `mclust` (version 5.4.3) [103]. Lastly, the package `CompQuadForm` (version 1.4.3) was used to compute gene-$\varepsilon$ gene-level *P*-values with Imhof's method [26, 111]. Comparisons in this work were made using software for MAGMA (version 1.07b; https://ctg.cncr.nl/software/magma), PEGASUS (version 1.3.0; https://github.com/ramachandran-lab/PEGASUS), RSS (version 1.0.0; https://github.com/stephenslab/rss), SKAT (version 1.3.2.1; https://www.hsph.harvard.edu/skat), VEGAS (version 2.0.0; https://vegas2.qimrberghofer.edu.au) which are also publicly available. See all other relevant URLs below.

## URLs

gene-$\varepsilon$ software, https://github.com/ramachandran-lab/genee; UK Biobank, https://www.ukbiobank.ac.uk; Database of Genotypes and Phenotypes (dbGaP), https://www.ncbi.nlm.nih.gov/gap; NHGRI-EBI GWAS Catalog, https://www.ebi.ac.uk/gwas/; UCSC Genome Browser, https://genome.ucsc.edu/index.html; Enrichr software, http://amp.pharm.mssm.edu/Enrichr/; SNP-set (Sequence) Kernel Association Test (SKAT) software, https://www.hsph.harvard.edu/skat; Multi-marker Analysis of GenoMic Annotation (MAGMA) software, https://ctg.cncr.nl/software/magma; Precise, Efficient Gene Association Score Using SNPs (PEGASUS) software, https://github.com/ramachandran-lab/PEGASUS; Regression with Summary Statistics (RSS) enrichment software, https://github.com/stephenslab/rss; Versatile Gene-based Association Study (VEGAS) version 2, https://vegas2.qimrberghofer.edu.au.

## Supporting information

**S1 Fig. Simulation study results showing the Pearson correlation between various degrees of gene-$\varepsilon$ regularized SNP-level effect size estimates and the true effect sizes that generated the complex traits.** Assessed regularization techniques are the **(A)** LASSO [23], **(B)** Elastic Net [24], **(C)** Ridge Regression [25], and **(D)** no regularization of ordinary least squares (OLS) effect sizes which serves as a baseline. Here, we take real genotype data on chromosome 19 from *N* = 5, 000 randomly chosen individuals of European ancestry in the UK Biobank (see S1 Text). We then assumed a simple linear additive model for quantitative traits while varying the narrow-sense heritability ($h^2$ = {0.01, 0.05, 0.10, 0.15, 0.20, 0.25}). We considered two scenarios where traits are generated with and without additional population structure (colored as pink and blue lines, respectively). In the former setting, phenotypes are simulated while also using the top five principal components (PCs) of the genotype matrix as covariates to create stratification. These PCs contributed to 10% of the phenotypic variance. In both settings, GWA SNP-level effect sizes were derived via OLS without accounting for any additional structure. The y-axis shows Pearson correlation between gene-$\varepsilon$ regularized effect sizes and the truth. On the x-axis of each plot, we vary the number of causal SNPs for each trait (i.e., {1, 5, 10, 15, 20, 25}%). Results are based on ten replicates (see S1 Text), with the error bars representing standard errors across runs.
(PDF)

**S2 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations** (*N* = **5,000;** $h^2$ = **0.2**)**.** Here, the sample size *N* = 5, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.2. We compute standard GWA SNP-level effect sizes

(estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S3 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations ($N$ = 10,000; $h^2$ = 0.2).** Here, the sample size $N$ = 10, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.2. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S4 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations ($N$ = 5,000; $h^2$ = 0.6).** Here, the sample size $N$ = 5, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.6. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S5 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with population stratification ($N$ = 5,000; $h^2$ = 0.2).** Here, the sample size $N$ = 5, 000 and the narrow-

sense heritability of the simulated quantitative trait is $h^2 = 0.2$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S6 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with population stratification ($N$ = 10,000; $h^2$ = 0.2).** Here, the sample size $N = 10,000$ and the narrow-sense heritability of the simulated quantitative trait is $h^2 = 0.2$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S7 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with population stratification ($N$ = 5,000; $h^2$ = 0.6).** Here, the sample size $N = 5,000$ and the narrow-sense heritability of the simulated quantitative trait is $h^2 = 0.6$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate

for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S8 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with population stratification ($N$ = 10,000; $h^2$ = 0.6).** Here, the sample size $N$ = 10, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.6. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S9 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 5,000; $h^2$ = 0.2).** Here, the sample size $N$ = 5, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.2. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S10 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 10,000; $h^2$ = 0.2).** Here, the sample size $N$ = 10, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.2. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares).

Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S11 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 5,000; $h^2$ = 0.6).** Here, the sample size $N$ = 5, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.6. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S12 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 10,000; $h^2$ = 0.6).** Here, the sample size $N$ = 10, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.6. We compute standard GWA SNP-level effect sizes (estimated using ordinary least squares). Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S13 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 5,000; $h^2$ = 0.2).** Here, the sample size $N = 5,000$ and the narrow-sense heritability of the simulated quantitative trait is $h^2 = 0.2$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S14 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 10,000; $h^2$ = 0.2).** Here, the sample size $N = 10,000$ and the narrow-sense heritability of the simulated quantitative trait is $h^2 = 0.2$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S15 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 5,000; $h^2$ = 0.6).** Here, the sample size $N = 5,000$ and the narrow-sense heritability of the simulated quantitative trait is $h^2 = 0.6$. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with

LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S16 Fig. (A, C) Receiver operating characteristic (ROC) and (B, D) precision-recall curves comparing the performance of gene-$\varepsilon$ and competing approaches in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 10,000; $h^2$ = 0.6).** Here, the sample size $N$ = 10, 000 and the narrow-sense heritability of the simulated quantitative trait is $h^2$ = 0.6. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results for gene-$\varepsilon$ are shown with LASSO (blue), Elastic Net (EN; red), and Ridge Regression (RR; purple) regularizations. We also show the results of gene-$\varepsilon$ without regularization to illustrate the importance of the regularization step (labeled OLS; orange). We compare gene-$\varepsilon$ with five existing methods: PEGASUS (brown) [12], VEGAS (teal) [7], the Bayesian approach RSS (black) [14], SKAT (green) [20], and MAGMA (peach) [10]. Note that each was method implemented without using any covariates. **(A, C)** ROC curves show power versus false positive rate for each approach of sparse (1% enriched genes) and polygenic (10% enriched genes) architectures, respectively. Note that the upper limit of the x-axis has been truncated at 0.1. **(B, D)** Precision-Recall curves for each method applied to the simulations. Note that, in the sparse case (1% enriched genes), the top ranked genes are always true positives, and therefore the minimal recall is not 0. All results are based on 100 replicates (see S1 Text).
(PDF)

**S17 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations ($h^2$ = 0.2).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2$ = 0.2 and sample sizes are set to $N$ = 5,000 in **(A, B)** and $N$ = 10,000 in **(C, D)**. In each case, standard GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares). Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P$ = 3.55×$10^{-5}$ corrected for the 1,408 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S18 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations ($h^2 = 0.6$).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2 = 0.6$ and sample sizes are set to $N = 5,000$ in **(A, B)** and $N = 10,000$ in **(C, D)**. In each case, standard GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares). Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $Pp = 3.55 \times 10^{-5}$ corrected for the 1,408 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text). (PDF)

**S19 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with population stratification ($h^2 = 0.2$).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2 = 0.2$ and sample sizes are set to $N = 5,000$ in **(A, B)** and $N = 10,000$ in **(C, D)**. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$ corrected for the 1,408 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text). (PDF)

**S20 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with population stratification ($h^2 = 0.6$).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2 = 0.6$ and sample sizes are set to $N = 5,000$ in **(A, B)** and $N = 10,000$ in **(C, D)**. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$ corrected for the 1,408 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by

both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S21 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($h^2$ = 0.2).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2$ = 0.2 and sample sizes are set to $N$ = 5,000 in **(A, B)** and $N$ = 10,000 in **(C, D)**. In each case, standard GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares). Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P$ = 2.61×$10^{-5}$ corrected for the 1,916 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S22 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($h^2$ = 0.6).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2$ = 0.6 and sample sizes are set to $N$ = 5,000 in **(A, B)** and $N$ = 10,000 in **(C, D)**. In each case, standard GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares). Results are shown comparing the -$\log_{10}$ transformed gene-level $P$-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P$ = 2.61×$10^{-5}$ corrected for the 1,916 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S23 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($h^2$ = 0.2).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2$ = 0.2 and sample sizes are set to $N$ = 5,000 in **(A, B)** and $N$ = 10,000 in **(C, D)**. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results are shown comparing the

-log$_{10}$ transformed gene-level *P*-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$ corrected for the 1,916 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S24 Fig. Scatter plots assessing how regularization on SNP-level summary statistics affects the ability to identify enriched genes in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($h^2 = 0.6$).** Here, the narrow-sense heritability of the simulated quantitative traits is $h^2 = 0.6$ and sample sizes are set to $N = 5,000$ in **(A, B)** and $N = 10,000$ in **(C, D)**. In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. Results are shown comparing the -log$_{10}$ transformed gene-level *P*-values derived by gene-$\varepsilon$ with Elastic Net (EN) regularization on the y-axis and without regularization (labeled as OLS) on the x-axis. The horizontal and vertical dashed lines are marked at the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$ corrected for the 1,916 genes on chromosome 1 from the UK Biobank genotype data. True positive causal genes used to generate the synthetic phenotypes are colored in red, while non-causal genes are given in grey. Genes in the top right quadrant are selected by both approaches. Genes in the top left and bottom right quadrants are uniquely identified by gene-$\varepsilon$-EN and gene-$\varepsilon$-OLS, respectively. To illustrate the importance of regularization on SNP-level summary statistics, we highlight the true positive genes only identified by gene-$\varepsilon$-EN in blue. Each plot combines results from 100 simulated replicates (see S1 Text).
(PDF)

**S25 Fig. Gene-level association results from applying gene-$\varepsilon$ to body height (panels A and C) and mean platelet volume (MPV; panels B and D), assayed in European-ancestry individuals in the UK Biobank with UCSC RefSeq gene boundaries augmented by a 50 kilobase (kb) buffer.** Body height has been estimated to have a narrow-sense heritability $h^2$ in the range of 0.45 to 0.80 [6, 31–39]; while, MPV has been estimated to have $h^2$ between 0.50 and 0.70 [33, 34, 58]. Manhattan plots of gene-$\varepsilon$ gene-level association *P*-values using Elastic Net regularized effect sizes for **(A)** body height and **(B)** MPV. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 2.83 \times 10^{-6}$ correcting for 17,680 autosomal genes analyzed). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes overlapping with the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$. We highlight categories with *Q*-values (i.e., false discovery rates) less than 0.05 and annotate corresponding genes in the Manhattan plots in **(A)** and **(B)**, respectively. For height, the most enriched dbGAP category is "Body Height", with 5 of the genes identified by gene-$\varepsilon$ appearing in this category. For MPV, the four significant dbGAP categories are "Platelet Count", "Behcet Syndrome", "Psoriasis", and "Face"—all of which have been connected to trait [57, 60, 61, 112, 113].
(PDF)

**S26 Fig. Gene-level association results from applying gene-$\varepsilon$ to body mass index (BMI), assayed in European-ancestry individuals in the UK Biobank.** BMI has been estimated to have a narrow-sense heritability $h^2$ ranging from 0.25 to 0.60 [31, 33, 34, 36, 37, 39, 42–45]. Manhattan plots of gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by **(A)** using UCSC annotations directly, and **(B)** augmenting the gene boundaries by adding SNPs within a ±50kb buffer. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 3.49 \times 10^{-6}$ and $P = 2.83 \times 10^{-6}$ correcting for the 14,322 and 17,680 autosomal genes analyzed, respectively). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes previously associated with BMI in the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$ in **(A)** and **(B)**, respectively. While many of the scored categories are biologically related to BMI (e.g., "Body Mass Index", "Adiposity", and "Arteries") [66, 114–116], none of them had $Q$-values (i.e., false discovery rates) less than 0.05.
(PDF)

**S27 Fig. Gene-level association results from applying gene-$\varepsilon$ to mean corpuscular volume (MCV), assayed in European-ancestry individuals in the UK Biobank.** MCV has been estimated to have a narrow-sense heritability $h^2$ in the range of 0.20 to 0.60 [33, 34, 117, 118]. Manhattan plots of gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by **(A)** using UCSC annotations directly, and **(B)** augmenting the gene boundaries by adding SNPs within a ±50kb buffer. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 3.49 \times 10^{-6}$ and $P = 2.83 \times 10^{-6}$ correcting for the 14,322 and 17,680 autosomal genes analyzed, respectively). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes previously associated with MCV in the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$. We highlight categories with $Q$-values (i.e., false discovery rates) less than 0.05 and annotate corresponding genes in the Manhattan plots in **(A)** and **(B)**, respectively. The dbGAP categories significantly enriched for gene-level associations with MCV included "Transferrin", "Erythrocyte Indices", "Hematocrit", "Narcolepsy", and "Iron"—all of which have been connected to trait [50–57].
(PDF)

**S28 Fig. Gene-level association results from applying gene-$\varepsilon$ to platelet count (PLC), assayed in European-ancestry individuals in the UK Biobank.** PLC has been estimated to have a narrow-sense heritability $h^2$ ranging from 0.55 to 0.80 [33, 34, 58]. Manhattan plots of gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by **(A)** using UCSC annotations directly, and **(B)** augmenting the gene boundaries by adding SNPs within a ±50kb buffer. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 3.49 \times 10^{-6}$ and $P = 2.83 \times 10^{-6}$ correcting for the 14,322 and 17,680 autosomal genes analyzed, respectively). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes previously associated with PLC in the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$. We highlight categories with $Q$-values (i.e., false discovery rates) less than 0.05 and annotate corresponding genes in the Manhattan plots in **(A)** and **(B)**, respectively. The most significant dbGAP category is "Platelet Count" for

both SNP-to-gene annotation schemes. The other significant dbGAP category was "Smoking" which has been previously connected to PLC [61, 119, 120].
(PDF)

**S29 Fig. Gene-level association results from applying gene-$\varepsilon$ to waist-hip ratio (WHR), assayed in European-ancestry individuals in the UK Biobank.** WHR has been estimated to have a narrow-sense heritability $h^2$ ranging from 0.10 to 0.25 [31, 33, 35, 42, 45, 121]. Manhattan plots of gene-$\varepsilon$ gene-level association *P*-values using Elastic Net regularized effect sizes when gene boundaries are defined by **(A)** using UCSC annotations directly, and **(B)** augmenting the gene boundaries by adding SNPs within a ±50kb buffer. The purple dashed line indicates a log-transformed Bonferroni-corrected significance threshold ($P = 3.49\times10^{-6}$ and $P = 2.83\times10^{-6}$ correcting for the 14,322 and 17,680 autosomal genes analyzed, respectively). We color code all significant genes identified by gene-$\varepsilon$ in orange, and annotate genes previously associated with WHR in the database of Genotypes and Phenotypes (dbGaP). In **(C)** and **(D)**, we conduct gene set enrichment analysis using Enrichr [46, 59] to identify dbGaP categories enriched for significant gene-level associations reported by gene-$\varepsilon$ in **(A)** and **(B)**, respectively. While many of the scored categories are biologically related to WHR (e.g., "Body Mass Index", "Adiposity", and "Inflammatory Bowel Diseases") [122, 123], none of them had *Q*-values (i.e., false discovery rates) less than 0.05.
(PDF)

**S1 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations ($N = 5,000$; $h^2 = 0.2$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55\times10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S2 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations ($N = 10,000$; $h^2 = 0.2$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55\times10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based

on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S3 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations ($N = 5,000$; $h^2 = 0.6$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S4 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations ($N = 10,000$; $h^2 = 0.6$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S5 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with population stratification ($N = 5,000$; $h^2 = 0.2$).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance

gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S6 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with population stratification ($N$ = 10,000; $h^2$ = 0.2).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P$ = 3.55×10$^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S7 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with population stratification ($N$ = 5,000; $h^2$ = 0.6).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P$ = 3.55×10$^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S8 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with population stratification ($N$ = 10,000; $h^2$ = 0.6).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as

covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 3.55 \times 10^{-5}$, corrected for 1,408 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S9 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N = 5,000$; $h^2 = 0.2$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S10 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N = 10,000$; $h^2 = 0.2$).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S11 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 5,000; $h^2$ = 0.6).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue. (PDF)

**S12 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer ($N$ = 10,000; $h^2$ = 0.6).** We computed standard GWA SNP-level effect sizes (estimated using ordinary least squares) as input to each method listed. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue. (PDF)

**S13 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 5,000; $h^2$ = 0.2).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the

parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S14 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 10,000; $h^2$ = 0.2).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P$ = 2.61×10$^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S15 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 5,000; $h^2$ = 0.6).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P$ = 2.61×10$^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S16 Table. Empirical power and false discovery rates (FDR) for detecting enriched genes (genes containing at least one causal SNP) after correcting for multiple hypothesis testing in simulations with gene boundaries augmented by a 50 kilobase (kb) buffer and with population stratification ($N$ = 10,000; $h^2$ = 0.6).** In this simulation, traits were generated while using the top five principal components (PCs) of the genotype matrix as covariates. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via

ordinary least squares) without any control for the additional structure. We show the power of gene-$\varepsilon$ to identify enriched genes under the Bonferonni-corrected threshold $P = 2.61 \times 10^{-5}$, corrected for 1,916 genes simulated using chromosome 1 from the UK Biobank genotype data (see S1 Text). Results for gene-$\varepsilon$ are shown with LASSO, Elastic Net (EN), and Ridge Regression (RR) regularizations. We also show the power of gene-$\varepsilon$ without regularization to illustrate the importance of this step (OLS). Additionally, we compare the performance gene-$\varepsilon$ with five existing methods: PEGASUS [12], VEGAS [7], RSS [14], SKAT [20], and MAGMA [10]. The last is a Bayesian method and is evaluated based on the "median probability criterion" (i.e., posterior enrichment probability of a gene is greater than 0.5). All results are based on 100 replicates and standard deviations of the estimates across runs are given in the parentheses. Approaches with the greatest power are bolded in purple, while methods with the lowest FDR is bolded in blue.
(PDF)

**S17 Table. Empirical type I error estimates using different gene-$\varepsilon$ approaches.** Here, quantitative traits are simulated with just noise randomly drawn from standard normal distributions. This represents the scenario in which all SNPs are non-causal and satisfy the conventional null hypothesis $H_0: \beta_j = 0$. GWA summary statistics were computed by fitting a single-SNP univariate linear model (via ordinary least squares). Each table entry lists the mean type I error rate estimates for the four gene-$\varepsilon$ modeling approaches—which is computed as the proportion of $P$-values under some significance level $\alpha$. Empirical size for the analyses used significance levels of $\alpha = 0.05$, 0.01, 0.001, and $2.61 \times 10^{-5}$ (the Bonferonni-corrected threshold), respectively. Sample sizes of the individual-level data (used to derive the summary statistics), were set to $N = 5,000$ and 10,000 observations. These results are based on 100 simulated datasets and the standard errors across the replicated are included in the parentheses. Overall, gene-$\varepsilon$ controls the type I error rate for reasonably sized datasets, and can be slightly conservative when the sample size is small and the GWA summary statistics are less precise/more inflated.
(PDF)

**S18 Table. Characterization of the genetic architectures of six traits assayed in European-ancestry individuals in the UK Biobank.** Here, we report the way difference regularization makes when gene-$\varepsilon$ characterizes $\varepsilon$-genic effects in complex traits. Results are shown for Elastic Net (which is highlighted in the main text). We also show results when no shrinkage is applied to illustrate the importance of this step (denoted by OLS). In the former case, we regress the GWA SNP-level effect size estimates onto chromosome-specific LD matrices to derive a regularized set of summary statistics $\tilde{\boldsymbol{\beta}}$. gene-$\varepsilon$ assumes a reformulated null distribution of SNP-level effects $\tilde{\beta}_j \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, where $\sigma_\varepsilon^2$ is the SNP-level null threshold and represents the maximum proportion of phenotypic variance explained (PVE) by a spurious or non-associated SNP. We used an EM-algorithm with 100 iterations to fit $K$-mixture Gaussian models over the regularized effect sizes to estimate $\sigma_\varepsilon^2$. Here, each mixture component had distinctively smaller variances ($\sigma_1^2 > \cdots > \sigma_K^2$; with the $K$-th component fixed at $\sigma_K^2 = 0$), and the number of total mixture components $K$ was chosen based on a grid of values where the best model yielded the highest Bayesian Information Criterion (BIC). We assume associated SNPs appear in the first component, non-associated SNPs appear in the last component, and null SNPs with spurious effects fell in between (i.e., $\sigma_\varepsilon^2 = \sigma_2^2$). Thus, a SNP is considered to have some level of association with a trait if $\mathbb{E}[\beta_j^2] > \sigma_K^2 = 0$; while a SNP is considered "causal" if $\mathbb{E}[\beta_j^2] > \sigma_2^2$. Column 3 gives the $K$ used for each trait. Column 4 and 5 detail the percentage of

associated and causal SNPs, respectively. The last column gives the mean threshold for $\varepsilon$-genic effects across the chromosomes.
(PDF)

**S19 Table. Significant genes for body height in the UK Biobank analysis using gene-$\varepsilon$-EN.**
Here, we analyze 17,680 genes from $N$ = 349,468 individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our analyses, $P$ = 0.05/14322 autosomal genes = 3.49×10$^{-6}$ and $P$ = 0.05/17680 autosomal genes = 2.83×10$^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P$ = 4.67×10$^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S20 Table. Significant genes for body mass index (BMI) in the UK Biobank analysis using gene-$\varepsilon$-EN.** Here, we analyze 17,680 genes from $N$ = 349,468 individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our analyses, $P$ = 0.05/14322 autosomal genes = 3.49×10$^{-6}$ and $P$ = 0.05/17680 autosomal genes = 2.83×10$^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P$ = 4.67×10$^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S21 Table. Significant genes for mean corpuscular volume (MCV) in the UK Biobank analysis using gene-$\varepsilon$-EN.** Here, we analyze 17,680 genes from $N$ = 349,468 individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our

analyses, $P = 0.05/14322$ autosomal genes $= 3.49{\times}10^{-6}$ and $P = 0.05/17680$ autosomal genes $= 2.83{\times}10^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P = 4.67{\times}10^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S22 Table. Significant genes for mean platelet volume (MPV) in the UK Biobank analysis using gene-$\varepsilon$-EN.** Here, we analyze 17,680 genes from $N = 349{,}468$ individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our analyses, $P = 0.05/14322$ autosomal genes $= 3.49{\times}10^{-6}$ and $P = 0.05/17680$ autosomal genes $= 2.83{\times}10^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P = 4.67{\times}10^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S23 Table. Significant genes for platelet count (PLC) in the UK Biobank analysis using gene-$\varepsilon$-EN.** Here, we analyze 17,680 genes from $N = 349{,}468$ individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our analyses, $P = 0.05/14322$ autosomal genes $= 3.49{\times}10^{-6}$ and $P = 0.05/17680$ autosomal genes $= 2.83{\times}10^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant

according to a genome-wide Bonferroni-corrected threshold ($P = 4.67 \times 10^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S24 Table. Significant genes for waist-hip ratio (WHR) in the UK Biobank analysis using gene-$\varepsilon$-EN.** Here, we analyze 17,680 genes from $N$ = 349,468 individuals of European-ancestry. This file gives the gene-$\varepsilon$ gene-level association $P$-values using Elastic Net regularized effect sizes when gene boundaries are defined by (page 1) using UCSC annotations directly, and (page 2) augmenting the gene boundaries by adding SNPs within a ±50kb buffer. Significance was determined by using a Bonferroni-corrected $P$-value threshold (in our analyses, $P = 0.05/$ 14322 autosomal genes = $3.49 \times 10^{-6}$ and $P = 0.05/17680$ autosomal genes = $2.83 \times 10^{-6}$, respectively). The columns of tables on both pages provide: (1) chromosome position; (2) gene name; (3) gene-$\varepsilon$-EN gene $P$-value; (4) gene-specific heritability estimates; (5) whether or not an association between gene and trait is listed in the GWAS catalog (marked as "yes" or "no"); (6-7) the starting and ending position of the gene's genomic position; (8) number of SNPs within a gene that were included in analysis; (9) the most significant SNP according to GWA summary statistics; (10) the $P$-value of the most significant SNP; and, on the first page, (11) the corresponding gene-level posterior enrichment probability as found by RSS for comparison. Note that an "NA" in column (11) occurs wherever the MCMC for RSS failed to converge. Highlighted rows represent enriched genes whose top SNP is not marginally significant according to a genome-wide Bonferroni-corrected threshold ($P = 4.67 \times 10^{-8}$ correcting for 1,070,306 SNPs analyzed).
(XLSX)

**S25 Table. Characterization of the genetic architectures of six traits assayed in European-ancestry individuals in the UK Biobank (using un-imputed genotypes).** Here, we report the way different regularizations in gene-$\varepsilon$ characterize $\varepsilon$-genic effects in complex traits. Results are shown for Elastic Net (which is highlighted in the main text), as well as for LASSO and Ridge Regression. We also show results when no shrinkage is applied to illustrate the importance of this step (denoted by OLS). In the three former cases, we regress the GWA SNP-level effect size estimates onto chromosome-specific LD matrices to derive a regularized set of summary statistics $\tilde{\boldsymbol{\beta}}$. gene-$\varepsilon$ assumes a reformulated null distribution of SNP-level effects $\tilde{\beta}_j \sim \mathcal{N}(0, \sigma_\varepsilon^2)$, where $\sigma_\varepsilon^2$ is the SNP-level null threshold and represents the maximum proportion of phenotypic variance explained (PVE) by a spurious or non-associated SNP. We used an EM-algorithm with 100 iterations to fit $K$-mixture Gaussian models over the regularized effect sizes to estimate $\sigma_\varepsilon^2$. Here, each mixture component had distinctively smaller variances ($\sigma_1^2 > \cdots > \sigma_K^2$; with the $K$-th component fixed at $\sigma_K^2 = 0$), and the number of total mixture components $K$ was chosen based on a grid of values where the best model yielded the highest Bayesian Information Criterion (BIC). We assume associated SNPs appear in the first component, non-associated SNPs appear in the last component, and null SNPs with spurious effects fell in between (i.e., $\sigma_\varepsilon^2 = \sigma_2^2$). Thus, a SNP is considered to have some level of association with a trait if $\mathbb{E}[\beta_j^2] > \sigma_K^2 = 0$; while a SNP is considered "causal" if $\mathbb{E}[\beta_j^2] > \sigma_2^2$. Column 3 gives the $K$ used for each trait. Column 4 and 5 detail the percentage of associated and causal SNPs, respectively. The last column gives the mean threshold for $\varepsilon$-genic effects across the chromosomes.
(PDF)

**S26 Table. Comparison of the different gene-$\varepsilon$ approaches on the six quantitative traits assayed in European-ancestry individuals from the UK Biobank un-imputed genotyped**

**data.** Traits include: height; body mass index (BMI); mean corpuscular volume (MCV); mean platelet volume (MPV); platelet count (PLC); and waist-hip ratio (WHR). Here, we list the number of significant genes found when using gene-$\varepsilon$ with various regularization strategies, as well as the number of dbGAP categories enriched for significant genes identified by gene-$\varepsilon$. We also assess how well these results overlap with the gene-$\varepsilon$ -EN findings that were reported in the main text. Significant genes were determined by using a Bonferroni-corrected *P*-value threshold (in our analyses, $P = 0.05/13029$ autosomal genes = $3.84 \times 10^{-6}$). Enriched dbGAP categories were those with Enrichr *Q*-values (i.e., false discovery rates) less than 0.05.
(PDF)

**S27 Table. Comparison of the different gene-$\varepsilon$ approaches on the six quantitative traits assayed in European-ancestry individuals from the UK Biobank un-imputed genotyped data with gene boundaries augmented by a 50 kilobase (kb) buffer.** Traits include: height; body mass index (BMI); mean corpuscular volume (MCV); mean platelet volume (MPV); platelet count (PLC); and waist-hip ratio (WHR). Here, we list the number of significant genes found when using gene-$\varepsilon$ with various regularization strategies, as well as the number of dbGAP categories enriched for significant genes identified by gene-$\varepsilon$. We also assess how well these results overlap with the gene-$\varepsilon$ -EN findings that were reported in the main text. Significant genes were determined by using a Bonferroni-corrected *P*-value threshold (in our analyses, $P = 0.05/17680$ autosomal genes = $2.83 \times 10^{-6}$). Enriched dbGAP categories were those with Enrichr *Q*-values (i.e., false discovery rates) less than 0.05.
(PDF)

**S1 Text. Supplementary and background information for results mentioned in the main text.** Specifically, we give description of data quality control procedures, simulation setup and scenarios, review of other competing gene-level association methods, and additional results for the traits analyzed from the UK Biobank.
(PDF)

## Acknowledgments

## Author Contributions

**Conceptualization:** Sohini Ramachandran, Lorin Crawford.

**Data curation:** Sohini Ramachandran.

**Formal analysis:** Wei Cheng.

**Funding acquisition:** Sohini Ramachandran, Lorin Crawford.

**Investigation:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

**Methodology:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

**Project administration:** Sohini Ramachandran, Lorin Crawford.

**Resources:** Sohini Ramachandran, Lorin Crawford.

**Software:** Wei Cheng, Lorin Crawford.

**Supervision:** Sohini Ramachandran, Lorin Crawford.

**Validation:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

**Visualization:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

**Writing – original draft:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

**Writing – review & editing:** Wei Cheng, Sohini Ramachandran, Lorin Crawford.

# References

1. Visscher PM, Hill WG, Wray NR. Heritability in the genomics era–concepts and misconceptions. Nat Rev Genet. 2008; 9(4):255–266. https://doi.org/10.1038/nrg2322 PMID: 18319743

2. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. Nature. 2009; 461(7265):747–753. Available from: https://www.ncbi.nlm.nih.gov/pubmed/19812666 PMID: 19812666

3. Visscher PM, Brown MA, McCarthy MI, Yang J. Five Years of GWAS Discovery. Am J Hum Genet. 2012; 90(1):7–24. Available from: http://www.sciencedirect.com/science/article/pii/S0002929711005337 PMID: 22243964

4. Boyle EA, Li YI, Pritchard JK. An expanded view of complex traits: from polygenic to omnigenic. Cell. 2017; 169(7):1177–1186. https://doi.org/10.1016/j.cell.2017.05.038 PMID: 28622505

5. Wray NR, Wijmenga C, Sullivan PF, Yang J, Visscher PM. Common disease is more complex than implied by the core gene omnigenic model. Cell. 2018; 173(7):1573–1580. Available from: https://doi.org/10.1016/j.cell.2018.05.051 PMID: 29906445

6. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common SNPs explain a large proportion of the heritability for human height. Nat Genet. 2010; 42(7):565–569. https://doi.org/10.1038/ng.608 PMID: 20562875

7. Liu JZ, Mcrae AF, Nyholt DR, Medland SE, Wray NR, Brown KM, et al. A versatile gene-based test for genome-wide association studies. Am J Hum Genet. 2010; 87(1):139–145. https://doi.org/10.1016/j.ajhg.2010.06.009 PMID: 20598278

8. Carbonetto P, Stephens M. Integrated enrichment analysis of variants and pathways in genome-wide association studies indicates central role for IL-2 signaling genes in type 1 diabetes, and cytokine signaling genes in Crohn's disease. PLoS Genet. 2013; 9(10):e1003770–. Available from: https://doi.org/10.1371/journal.pgen.1003770 PMID: 24098138

9. Ionita-Laza I, Lee S, Makarov V, Buxbaum JD, Lin X. Sequence kernel association tests for the combined effect of rare and common variants. Am J Hum Genet. 2013; 92(6):841–853. Available from: http://www.sciencedirect.com/science/article/pii/S0002929713001766 PMID: 23684009

10. de Leeuw CA, Mooij JM, Heskes T, Posthuma D. MAGMA: generalized gene-set analysis of GWAS data. PLOS Comput Biol. 2015; 11(4):e1004219–. Available from: https://doi.org/10.1371/journal.pcbi.1004219 PMID: 25885710

11. Lamparter D, Marbach D, Rueedi R, Kutalik Z, Bergmann S. Fast and rigorous computation of gene and pathway scores from SNP-based summary statistics. PLOS Comput Biol. 2016; 12(1):e1004714–. Available from: https://doi.org/10.1371/journal.pcbi.1004714 PMID: 26808494

12. Nakka P, Raphael BJ, Ramachandran S. Gene and network analysis of common variants reveals novel associations in multiple complex diseases. Genetics. 2016; 204(2):783–798. Available from: http://www.genetics.org/content/204/2/783.abstract PMID: 27489002

13. Wang M, Huang J, Liu Y, Ma L, Potash JB, Han S. COMBAT: a combined association test for genes using summary statistics. Genetics. 2017; 207(3):883–891. https://doi.org/10.1534/genetics.117.300257 PMID: 28878002

14. Zhu X, Stephens M. Large-scale genome-wide enrichment analyses identify new trait-associated genes and pathways across 31 human phenotypes. Nat Comm. 2018; 9(1):4361. https://doi.org/10.1038/s41467-018-06805-x

**15.** Zhou X, Carbonetto P, Stephens M. Polygenic modeling with Bayesian sparse linear mixed models. PLoS Genet. 2013; 9(2):e1003264. https://doi.org/10.1371/journal.pgen.1003264 PMID: 23408905

**16.** Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in the application of mixed-model association methods. Nat Genet. 2014; 46(2):100–106. https://doi.org/10.1038/ng.2876 PMID: 24473328

**17.** Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, of the Psychiatric Genomics Consortium SWG, et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. Nat Genet. 2015; 47:291–295. Available from: http://dx.doi.org/10.1038/ng.3211 PMID: 25642630

**18.** Zhang Y, Qi G, Park JH, Chatterjee N. Estimation of complex effect-size distributions using summary-level statistics from genome-wide association studies across 32 complex traits. Nat Genet. 2018; 50 (9):1318–1326. https://doi.org/10.1038/s41588-018-0193-x PMID: 30104760

**19.** Holland D, Wang Y, Thompson WK, Schork A, Chen CH, Lo MT, et al. Estimating Effect Sizes and Expected Replication Probabilities from GWAS Summary Statistics. Front Genet. 2016; 7:15. Available from: https://www.frontiersin.org/article/10.3389/fgene.2016.00015 PMID: 26909100

**20.** Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, Hunter DJ, et al. Powerful SNP-set analysis for case-control genome-wide association studies. Am J Hum Genet. 2010; 86(6):929–942. https://doi.org/10.1016/j.ajhg.2010.05.002 PMID: 20560208

**21.** Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK Biobank resource with deep phenotyping and genomic data. Nature. 2018; 562(7726):203–209. Available from: https://doi.org/10.1038/s41586-018-0579-z PMID: 30305743

**22.** Stephens M. False discovery rates: a new deal. Biostatistics. 2017; 18(2):275–294. Available from: http://dx.doi.org/10.1093/biostatistics/kxw041 PMID: 27756721

**23.** Tibshirani R. Regression shrinkage and selection via the lasso. J R Stat Soc Series B Stat Methodol. 1996; 58(1):267–288.

**24.** Zou H, Hastie T. Regularization and variable selection via the elastic net. J R Stat Soc Series B Stat Methodol. 2005; 67(2):301–320. https://doi.org/10.1111/j.1467-9868.2005.00503.x

**25.** Hoerl AE, Kennard RW. Ridge regression: Biased estimation for nonorthogonal problems. Technometrics. 1970; 12(1):55–67. https://doi.org/10.1080/00401706.1970.10488634

**26.** Imhof JP. Computing the distribution of quadratic forms in normal variables. Biometrika. 1961; 48(3/4):419–426. Available from: http://www.jstor.org/stable/2332763

**27.** Pruitt KD, Tatusova T, Maglott DR. NCBI Reference Sequence (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. Nucleic Acids Res. 2005; 33(Database issue):D501–4. https://doi.org/10.1093/nar/gki025 PMID: 15608248

**28.** Lee S, Emond MJ, Bamshad MJ, Barnes KC, Rieder MJ, Nickerson DA, et al. Optimal unified approach for rare-variant association testing with application to small-sample case-control whole-exome sequencing studies. Am J Hum Genet. 2012; 91(2):224–237. Available from: http://www.sciencedirect.com/science/article/pii/S0002929712003163 PMID: 22863193

**29.** Zhu X, Stephens M. Bayesian large-scale multiple regression with summary statistics from genome-wide association studies. Ann Appl Stat. 2017; 11(3):1561–1592. Available from: https://projecteuclid.org:443/euclid.aoas/1507168840 PMID: 29399241

**30.** Barbieri MM, Berger JO. Optimal predictive model selection. Ann Statist. 2004; 32(3):870–897. Available from: http://projecteuclid.org/euclid.aos/1085408489

**31.** Zaitlen N, Kraft P, Patterson N, Pasaniuc B, Bhatia G, Pollack S, et al. Using extended genealogy to estimate components of heritability for 23 quantitative and dichotomous traits. PLoS Genet. 2013; 9 (5):e1003520–. Available from: https://doi.org/10.1371/journal.pgen.1003520 PMID: 23737753

**32.** Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the role of common variation in the genomic and biological architecture of adult human height. Nat Genet. 2014; 46 (11):1173–1186. https://doi.org/10.1038/ng.3097 PMID: 25282103

**33.** Heckerman D, Gurdasani D, Kadie C, Pomilla C, Carstensen T, Martin H, et al. Linear mixed model for heritability estimation that explicitly addresses environmental variation. Proc Natl Acad Sci U S A. 2016; 113(27):7377–7382. Available from: http://www.pnas.org/content/113/27/7377.abstract PMID: 27382152

**34.** Shi H, Kichaev G, Pasaniuc B. Contrasting the genetic architecture of 30 complex traits from summary association data. Am J Hum Genet. 2016; 99(1):139–153. Available from: http://www.sciencedirect.com/science/article/pii/S0002929716301483 PMID: 27346688

**35.** Xia C, Amador C, Huffman J, Trochet H, Campbell A, Porteous D, et al. Pedigree- and SNP-associated genetics and recent environment are the major contributors to anthropometric and

cardiometabolic trait variation. PLoS Genet. 2016; 12(2):e1005804–. Available from: https://doi.org/10.1371/journal.pgen.1005804 PMID: 26836320

36. Ge T, Chen CY, Neale BM, Sabuncu MR, Smoller JW. Phenome-wide heritability analysis of the UK Biobank. PLoS Genet. 2017; 13(4):e1006711–. Available from: https://doi.org/10.1371/journal.pgen.1006711 PMID: 28388634

37. Speed D, Cai N, The UCLEB Consortium, Johnson MR, Nejentsev S, Balding DJ. Reevaluation of SNP heritability in complex human traits. Nat Genet. 2017; 49:986–992. Available from: https://doi.org/10.1038/ng.3865 PMID: 28530675

38. Marouli E, Graff M, Medina-Gomez C, Lo KS, Wood AR, Kjaer TR, et al. Rare and low-frequency coding variants alter human adult height. Nature. 2017; 542(7640):186–190. https://doi.org/10.1038/nature21039 PMID: 28146470

39. Wainschtein P, Jain DP, Yengo L, Zheng Z, TOPMed Anthropometry Working Group, Trans-Omics for Precision Medicine Consortium, et al. Recovery of trait heritability from whole genome sequence data. bioRxiv. 2019;p. 588020. Available from: http://biorxiv.org/content/early/2019/03/25/588020.abstract.

40. Goldstein DB. Common genetic variation and human traits. N Engl J Med. 2009; 360(17):1696–1698. https://doi.org/10.1056/NEJMp0806284 PMID: 19369660

41. Lello L, Avery SG, Tellier L, Vazquez AI, de los Campos G, Hsu SDH. Accurate Genomic Prediction of Human Height. Genetics. 2018; 210(2):477–497. Available from: http://www.genetics.org/content/210/2/477.abstract PMID: 30150289

42. Vattikuti S, Guo J, Chow CC. Heritability and genetic correlations explained by common SNPs for metabolic syndrome traits. PLoS Genet. 2012; 8(3):e1002637. https://doi.org/10.1371/journal.pgen.1002637 PMID: 22479213

43. Yang J, Bakshi A, Zhu Z, Hemani G, Vinkhuyzen AA, Lee SH, et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. Nat Genet. 2015; 47(10):1114. https://doi.org/10.1038/ng.3390 PMID: 26323059

44. Robinson MR, English G, Moser G, Lloyd-Jones LR, Triplett MA, Zhu Z, et al. Genotype–covariate interaction effects and the heritability of adult body mass index. Nat Genet. 2017; 49(8):1174. https://doi.org/10.1038/ng.3912 PMID: 28692066

45. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, et al. Environment dominates over host genetics in shaping human gut microbiota. Nature. 2018; 555:210–215. Available from: https://doi.org/10.1038/nature25973 PMID: 29489753

46. Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. BMC Bioinform. 2013; 14(1):128. Available from: https://doi.org/10.1186/1471-2105-14-128

47. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, et al. Integrative analysis of 111 reference human epigenomes. Nature. 2015; 518(7539):317–330. Available from: https://www.ncbi.nlm.nih.gov/pubmed/25693563 PMID: 25693563

48. Eicher JD, Chami N, Kacprowski T, Nomura A, Chen MH, Yanek LR, et al. Platelet-Related Variants Identified by Exomechip Meta-analysis in 157,293 Individuals. Am J Hum Genet. 2016; 99(1):40–55. https://doi.org/10.1016/j.ajhg.2016.05.005 PMID: 27346686

49. Iotchkova V, Huang J, Morris JA, Jain D, Barbieri C, Walter K, et al. Discovery and refinement of genetic loci associated with cardiometabolic risk using dense imputation maps. Nat Genet. 2016; 48(11):1303–1312. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27668658 PMID: 27668658

50. Finberg KE, Heeney MM, Campagna DR, Aydinok Y, Pearson HA, Hartman KR, et al. Mutations in TMPRSS6 cause iron-refractory iron deficiency anemia (IRIDA). Nat Genet. 2008; 40(5):569–571. Available from: https://www.ncbi.nlm.nih.gov/pubmed/18408718 PMID: 18408718

51. Andrews NC. Genes determining blood cell traits. Nat Genet. 2009; 41:1161–1162. Available from: https://doi.org/10.1038/ng1109-1161 PMID: 19862006

52. Benyamin B, Ferreira MAR, Willemsen G, Gordon S, Middelberg RPS, McEvoy BP, et al. Common variants in TMPRSS6 are associated with iron status and erythrocyte volume. Nat Genet. 2009; 41(11):1173–1175. https://doi.org/10.1038/ng.456 PMID: 19820699

53. Chambers JC, Zhang W, Li Y, Sehmi J, Wass MN, Zabaneh D, et al. Genome-wide association study identifies variants in TMPRSS6 associated with hemoglobin levels. Nat Genet. 2009; 41(11):1170–1172. https://doi.org/10.1038/ng.462 PMID: 19820698

54. Soranzo N, Spector TD, Mangino M, Kühnel B, Rendon A, Teumer A, et al. A genome-wide meta-analysis identifies 22 loci associated with eight hematological parameters in the HaemGen consortium. Nat Genet. 2009; 41(11):1182–1190. Available from: https://www.ncbi.nlm.nih.gov/pubmed/19820697 PMID: 19820697

55. Ganesh SK, Zakai NA, van Rooij FJA, Soranzo N, Smith AV, Nalls MA, et al. Multiple loci influence erythrocyte phenotypes in the CHARGE Consortium. Nat Genet. 2009; 41(11):1191–1198. https://doi.org/10.1038/ng.466 PMID: 19862010

56. Li J, Glessner JT, Zhang H, Hou C, Wei Z, Bradfield JP, et al. GWAS of blood cell traits identifies novel associated loci and epistatic interactions in Caucasian and African-American children. Hum Mol Genet. 2013; 22(7):1457–1464. Available from: https://www.ncbi.nlm.nih.gov/pubmed/23263863 PMID: 23263863

57. Astle WJ, Elding H, Jiang T, Allen D, Ruklisa D, Mann AL, et al. The allelic landscape of human blood cell trait variation and links to common complex disease. Cell. 2016; 167(5):1415–1429. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27863252 PMID: 27863252

58. Qayyum R, Snively BM, Ziv E, Nalls MA, Liu Y, Tang W, et al. A meta-analysis and genome-wide association study of platelet count and mean platelet volume in african americans. PLoS Genet. 2012; 8(3): e1002491. https://doi.org/10.1371/journal.pgen.1002491 PMID: 22423221

59. Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. Nucleic Acids Res. 2016; 44(W1):W90–W97. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27141961 PMID: 27141961

60. Lentaigne C, Freson K, Laffan MA, Turro E, Ouwehand WH, Consortium BB, et al. Inherited platelet disorders: toward DNA-based diagnosis. Blood. 2016; 127(23):2814–2823. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27095789 PMID: 27095789

61. Mousas A, Ntritsos G, Chen MH, Song C, Huffman JE, Tzoulaki I, et al. Rare coding variants pinpoint genes that control human hematological traits. PLoS Genet. 2017; 13(8):e1006925–. Available from: https://doi.org/10.1371/journal.pgen.1006925 PMID: 28787443

62. Gibson WT, Hood RL, Zhan SH, Bulman DE, Fejes AP, Moore R, et al. Mutations in EZH2 cause Weaver syndrome. Am J Hum Genet. 2012; 90(1):110–118. Available from: https://www.cell.com/ajhg/fulltext/S0002-9297(11)00496-4 PMID: 22177091

63. Minczuk M, He J, Duch AM, Ettema TJ, Chlebowski A, Dzionek K, et al. TEFM (c17orf42) is necessary for transcription of human mtDNA. Nucleic Acids Res. 2011; 39(10):4284–4299. Available from: https://www.ncbi.nlm.nih.gov/pubmed/21278163 PMID: 21278163

64. Carel JC, Lahlou N, Roger M, Chaussain JL. Precocious puberty and statural growth. Hum Reprod. 2004; 10(2):135–147. Available from: https://academic.oup.com/humupd/article/10/2/135/617162.

65. Gong J, Schumacher F, Lim U, Hindorff LA, Haessler J, Buyske S, et al. Fine Mapping and Identification of BMI Loci in African Americans. Am J Hum Genet. 2013; 93(4):661–671. https://doi.org/10.1016/j.ajhg.2013.08.012 PMID: 24094743

66. Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR, et al. Genetic studies of body mass index yield new insights for obesity biology. Nature. 2015; 518(7538):197–206. Available from: https://www.ncbi.nlm.nih.gov/pubmed/25673413 PMID: 25673413

67. Dickinson ME, Flenniken AM, Ji X, Teboul L, Wong MD, White JK, et al. High-throughput discovery of novel developmental phenotypes. Nature. 2016; 537:508–514. Available from: https://doi.org/10.1038/nature19356 PMID: 27626380

68. Baranski TJ, Kraja AT, Fink JL, Feitosa M, Lenzini PA, Borecki IB, et al. A high throughput, functional screen of human Body Mass Index GWAS loci using tissue-specific RNAi Drosophila melanogaster crosses. PLoS Genet. 2018; 14(4):e1007222–. Available from: https://doi.org/10.1371/journal.pgen.1007222 PMID: 29608557

69. Safran M, Dalah I, Alexander J, Rosen N, Iny Stein T, Shmoish M, et al. GeneCards Version 3: the human gene integrator. Database. 2010; 2010. Available from: https://academic.oup.com/database/article/doi/10.1093/database/baq020/407450 PMID: 20689021

70. Vuillaume ML, Naudion S, Banneau G, Diene G, Cartault A, Cailley D, et al. New candidate loci identified by array-CGH in a cohort of 100 children presenting with syndromic obesity. Am J Med Genet. 2014; 164 (8):1965–1975. Available from: https://onlinelibrary.wiley.com/doi/abs/10.1002/ajmg.a.36587

71. Wheeler E, Leong A, Liu CT, Hivert MF, Strawbridge RJ, Podmore C, et al. Impact of common genetic determinants of Hemoglobin A1c on type 2 diabetes risk and diagnosis in ancestrally diverse populations: A transethnic genome-wide meta-analysis. PLoS Med. 2017; 14(9):e1002383. Available from: https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1002383 PMID: 28898252

72. Linder S, Nelson D, Weiss M, Aepfelbacher M. Wiskott-Aldrich syndrome protein regulates podosomes in primary human macrophages. Proc Natl Acad Sci U S A. 1999; 96(17):9648–9653. Available from: http://www.pnas.org/content/96/17/9648.abstract PMID: 10449748

73. Steele BM, Harper MT, Macaulay IC, Morrell CN, Perez-Tamayo A, Foy M, et al. Canonical Wnt signaling negatively regulates platelet function. Proc Natl Acad Sci U S A. 2009; 106(47):19836–19841. https://doi.org/10.1073/pnas.0906268106 PMID: 19901330

**74.** Macaulay IC, Thon JN, Tijssen MR, Steele BM, MacDonald BT, Meade G, et al. Canonical Wnt signaling in megakaryocytes regulates proplatelet formation. Blood. 2013; 121(1):188–196. Available from: http://www.bloodjournal.org/content/121/1/188 PMID: 23160460

**75.** Stocks T, Angquist L, Hager J, Charon C, Holst C, Martinez JA, et al. TFAP2B-dietary protein and glycemic index interactions and weight maintenance after weight loss in the DiOGenes trial. Hum Hered. 2013; 75(2-4):213–219. https://doi.org/10.1159/000353591

**76.** Xiang J, Yang S, Xin N, Gaertig MA, Reeves RH, Li S, et al. DYRK1A regulates Hap1–Dcaf7/WDR68 binding with implication for delayed growth in down syndrome. Proc Natl Acad Sci U S A. 2017; 114(7): E1224–E1233. Available from: https://www.pnas.org/content/114/7/E1224 PMID: 28137862

**77.** Smith CM, Finger JH, Hayamizu TF, McCright IJ, Eppig JT, Kadin JA, et al. The mouse gene expression database (GXD): 2007 update. Nucleic Acids Res. 2006; 35:D618–D623. Available from: https://academic.oup.com/nar/article/35/suppl_1/D618/1085755 PMID: 17130151

**78.** Bult CJ, Krupke DM, Begley DA, Richardson JE, Neuhauser SB, Sundberg JP, et al. Mouse Tumor Biology (MTB): a database of mouse models for human cancer. Nucleic Acids Res. 2014; 43(D1): D818–D824. Available from: https://academic.oup.com/nar/article/43/D1/D818/2439858 PMID: 25332399

**79.** Smith CL, Blake JA, Kadin JA, Richardson JE, Bult CJ, Group MGD. Mouse Genome Database (MGD)-2018: knowledgebase for the laboratory mouse. Nucleic Acids Res. 2017; 46(D1):D836–D842. Available from: https://academic.oup.com/nar/article/47/D1/D801/5165331

**80.** Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. Am J Hum Genet. 2011; 89(1):82–93. https://doi.org/10.1016/j.ajhg.2011.05.029 PMID: 21737059

**81.** Lee S, Abecasis GR, Boehnke M, Lin X. Rare-variant association analysis: study designs and statistical tests. Am J Hum Genet. 2014; 95(1):5–23. Available from: http://www.sciencedirect.com/science/article/pii/S0002929714002717 PMID: 24995866

**82.** Zuk O, Schaffner SF, Samocha K, Do R, Hechter E, Kathiresan S, et al. Searching for missing heritability: designing rare variant association studies. Proc Natl Acad Sci U S A. 2014; 111(4):E455–E464. Available from: http://www.pnas.org/content/111/4/E455.abstract PMID: 24443550

**83.** Gazal S, Loh PR, Finucane HK, Ganna A, Schoech A, Sunyaev S, et al. Functional architecture of low-frequency variants highlights strength of negative selection across coding and non-coding annotations. Nat Genet. 2018; 50(11):1600–1607. Available from: https://doi.org/10.1038/s41588-018-0231-8 PMID: 30297966

**84.** Wojcik G, Graff M, Nishimura KK, Tao R, Haessler J, Gignoux CR, et al. The PAGE Study: how genetic diversity improves our understanding of the architecture of complex traits. bioRxiv. 2018;p. 188094. Available from: http://biorxiv.org/content/early/2018/10/17/188094.abstract.

**85.** Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic risk scores may exacerbate health disparities. Nat Genet. 2019; 51(4):584–591. Available from: https://doi.org/10.1038/s41588-019-0379-x PMID: 30926966

**86.** GTEx Consortium. Genetic effects on gene expression across human tissues. Nature. 2017; 550:204–213. Available from: https://doi.org/10.1038/nature24277 PMID: 29022597

**87.** Wu Y, Zeng J, Zhang F, Zhu Z, Qi T, Zheng Z, et al. Integrative analysis of omics summary data reveals putative mechanisms underlying complex traits. Nat Comm. 2018; 9(1):918. Available from: https://doi.org/10.1038/s41467-018-03371-0

**88.** Xue A, Wu Y, Zhu Z, Zhang F, Kemper KE, Zheng Z, et al. Genome-wide association analyses identify 143 risk variants and putative regulatory mechanisms for type 2 diabetes. Nat Comm. 2018; 9 (1):2941. Available from: https://doi.org/10.1038/s41467-018-04951-w

**89.** Smemo S, Tena JJ, Kim KH, Gamazon ER, Sakabe NJ, Gomez-Marin C, et al. Obesity-associated variants within FTO form long-range functional connections with IRX3. Nature. 2014; 507(7492):371–375. https://doi.org/10.1038/nature13138 PMID: 24646999

**90.** Claussnitzer M, Dankel SN, Kim KH, Quon G, Meuleman W, Haugen C, et al. FTO Obesity Variant Circuitry and Adipocyte Browning in Humans. N Engl J Med. 2015; 373(10):895–907. Available from: https://doi.org/10.1056/NEJMoa1502214 PMID: 26287746

**91.** Lloyd-Jones LR, Zeng J, Sidorenko J, Yengo L, Moser G, Kemper KE, et al. Improved polygenic prediction by Bayesian multiple regression on summary statistics. Nat Comm. 2019; 10(1):5086. Available from: https://doi.org/10.1038/s41467-019-12653-0

**92.** Zeng P, Zhou X. Non-parametric genetic prediction of complex traits with latent Dirichlet process regression models. Nat Comm. 2017; 8:456. Available from: https://doi.org/10.1038/s41467-017-00470-2

**93.** Lee SH, Wray NR, Goddard ME, Visscher PM. Estimating missing heritability for disease from genome-wide association studies. Am J Hum Genet. 2011; 88(3):294–305. Available from: http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3059431/ PMID: 21376301

**94.** Golan D, Lander ES, Rosset S. Measuring missing heritability: inferring the contribution of common variants. Proc Natl Acad Sci U S A. 2014; 111(49):E5272–E5281. Available from: http://www.pnas.org/content/111/49/E5272.abstract PMID: 25422463

**95.** Weissbrod O, Lippert C, Geiger D, Heckerman D. Accurate liability estimation improves power in ascertained case-control studies. Nat Meth. 2015; 12:332–334. Available from: http://dx.doi.org/10.1038/nmeth.3285

**96.** Hormozdiari F, Kostem E, Kang EY, Pasaniuc B, Eskin E. Identifying causal variants at loci with multiple signals of association. Genetics. 2014; 198(2):497–508. Available from: https://pubmed.ncbi.nlm.nih.gov/25104515 PMID: 25104515

**97.** Hormozdiari F, van de Bunt M, Segrè AV, Li X, Joo JWJ, Bilow M, et al. Colocalization of GWAS and eQTL Signals Detects Target Genes. Am J Hum Genet. 2016; 99(6):1245–1260. Available from: https://doi.org/10.1016/j.ajhg.2016.10.003 PMID: 27866706

**98.** Wold S, Ruhe A, Wold H, Dunn W III. The collinearity problem in linear regression. The partial least squares (PLS) approach to generalized inverses. SIAM J Sci Comput. 1984; 5(3):735–743. https://doi.org/10.1137/0905052

**99.** Carvalho CM, Polson NG, Scott JG. The horseshoe estimator for sparse signals. Biometrika. 2010; 97(2):465–480. https://doi.org/10.1093/biomet/asq017

**100.** Dempster AP, Laird NM, Rubin DB. Maximum likelihood from incomplete data via the EM algorithm. J R Stat Soc Series B Stat Methodol. 1977; 39(1):1–22.

**101.** Benaglia T, Chauveau D, Hunter D, Young D. Mixtools: an R package for analyzing finite mixture models. J Stat Softw. 2009; 32(6):1–29. https://doi.org/10.18637/jss.v032.i06

**102.** McLachlan GJ, Lee SX, Rathnayake SI. Finite mixture models. Annual Review of Statistics and Its Application. 2019; 6(1):355–378. Available from: https://doi.org/10.1146/annurev-statistics-031017-100325

**103.** Scrucca L, Fop M, Murphy TB, Raftery AE. mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models. R J. 2016; 8(1):289–317. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27818791 PMID: 27818791

**104.** Schwarz G. Estimating the Dimension of a Model. Ann Statist. 1978; 6(2):461–464. Available from: https://projecteuclid.org:443/euclid.aos/1176344136

**105.** Zhou X. A unified framework for variance component estimation with summary statistics in genome-wide association studies. Ann Appl Stat. 2017; 11(4):2027–2051. Available from: https://projecteuclid.org:443/euclid.aoas/1514430276 PMID: 29515717

**106.** Crawford L, Zeng P, Mukherjee S, Zhou X. Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. PLoS Genet. 2017; 13(7):e1006869. Available from: https://doi.org/10.1371/journal.pgen.1006869 PMID: 28746338

**107.** Chen Z, Lin T, Wang K. A powerful variant-set association test based on chi-square distribution. Genetics. 2017; 207(3):903–910. https://doi.org/10.1534/genetics.117.300287 PMID: 28912342

**108.** Zhongxue C, Yan L, Tong L, Qingzhong L, Kai W. Gene-based genetic association test with adaptive optimal weights. Genet Epidemiol. 2017; 42(1):95–103. Available from: https://doi.org/10.1002/gepi.22098.

**109.** Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. J Stat Softw. 2010; 33(1):1. https://doi.org/10.18637/jss.v033.i01 PMID: 20808728

**110.** Zeng Y, Breheny P. The biglasso package: a memory-and computation-efficient solver for lasso model fitting with big data in R. arXiv. 2017;p. 1701.05936.

**111.** Duchesne P, Lafaye De Micheaux P. Computing the distribution of quadratic forms: Further comparisons between the Liu–Tang–Zhang approximation and exact methods. Comput Stat Data Anal. 2010; 54(4):858–862. Available from: http://www.sciencedirect.com/science/article/pii/S0167947309004381

**112.** Acikgoz N, Karincaoglu Y, Ermis N, Yagmur J, Atas H, Kurtoglu E, et al. Increased mean platelet volume in Behcet's disease with thrombotic tendency. Tohoku J Exp Med. 2010; 221(2):119–123. https://doi.org/10.1620/tjem.221.119 PMID: 20484842

**113.** Canpolat F, Akpinar H, Eskioglu F. Mean platelet volume in psoriasis and psoriatic arthritis. Clin Rheumatol. 2010; 29(3):325–328. https://doi.org/10.1007/s10067-009-1323-8 PMID: 20012663

**114.** Faeh D, Braun J, Bopp M. Body mass index vs cholesterol in cardiovascular disease risk prediction models. JAMA Intern Med. 2012; 172(22):1766–1768. https://doi.org/10.1001/2013.jamainternmed.327

115. Kurth T, Gaziano JM, Berger K, Kase CS, Rexrode KM, Cook NR, et al. Body mass index and the risk of stroke in men. JAMA Intern Med. 2002; 162(22):2557–2562. https://doi.org/10.1001/archinte.162.22.2557

116. Speakman JR, Loos RJF, O'Rahilly S, Hirschhorn JN, Allison DB. GWAS for BMI: a treasure trove of fundamental insights into the genetic basis of obesity. Int J Obes (Lond). 2018; 42(8):1524–1531. https://doi.org/10.1038/s41366-018-0147-5

117. Garner C, Tatu T, Reittie J, Littlewood T, Darley J, Cervino S, et al. Genetic influences on F cells and other hematologic variables: a twin heritability study. Blood. 2000; 95(1):342–346. https://doi.org/10.1182/blood.V95.1.342.001k33_342_346 PMID: 10607722

118. Van't Erve TJ, Wagner BA, Martin SM, Knudson CM, Blendowski R, Keaton M, et al. The heritability of hemolysis in stored human red blood cells. Transfusion. 2015; 55(6):1178–1185. https://doi.org/10.1111/trf.12992

119. Guerrero JA, Rivera J, Quiroga T, Martinez-Perez A, Antón AI, Martínez C, et al. Novel loci involved in platelet function and platelet count identified by a genome-wide study performed in children. Haematologica. 2011; 96(9):1335–1343. Available from: https://www.ncbi.nlm.nih.gov/pubmed/21546496 PMID: 21546496

120. Justice AE, Winkler TW, Feitosa MF, Graff M, Fisher VA, Young K, et al. Genome-wide meta-analysis of 241,258 adults accounting for smoking behaviour identifies novel loci for obesity traits. Nat Comm. 2017; 8:14977 EP –. Available from: https://doi.org/10.1038/ncomms14977

121. Loh PR, Kichaev G, Gazal S, Schoech AP, Price AL. Mixed-model association for biobank-scale datasets. Nat Genet. 2018; 50(7):906–908. Available from: https://doi.org/10.1038/s41588-018-0144-6 PMID: 29892013

122. Shungin D, Winkler TW, Croteau-Chonka DC, Ferreira T, Locke AE, Mägi R, et al. New genetic loci link adipose and insulin biology to body fat distribution. Nature. 2015; 518(7538):187–196. Available from: https://www.ncbi.nlm.nih.gov/pubmed/25673412 PMID: 25673412

123. Emdin CA, Khera AV, Natarajan P, Klarin D, Zekavat SM, Hsiao AJ, et al. Genetic association of waist-to-hip ratio with cardiometabolic traits, type 2 diabetes, and coronary heart disease. JAMA. 2017; 317(6):626–634. Available from: https://doi.org/10.1001/jama.2016.21042 PMID: 28196256