

Better estimation of SNP heritability from summary statistics provides a new understanding of the genetic architecture of complex traits

Doug Speed^{1,2,*} and David J Balding^{2,3}

*Corresponding author: doug@aias.au.dk

¹Aarhus Institute of Advanced Studies (AIAS), Aarhus University, Denmark.

²UCL Genetics Institute, University College London, United Kingdom.

³Melbourne Integrative Genomics, School of BioSciences and School of Mathematics & Statistics, University of Melbourne, Australia.

LD Score Regression (LDSC) has been widely applied to the results of genome-wide association studies. However, its estimates of SNP heritability are derived from an unrealistic model in which each SNP is expected to contribute equal heritability. As a consequence, LDSC tends to over-estimate confounding bias, under-estimate the total phenotypic variation explained by SNPs, and provide misleading estimates of the heritability enrichment of SNP categories. Therefore, we present SumHer, software for estimating SNP heritability from summary statistics using more realistic heritability models. After demonstrating its superiority over LDSC, we apply SumHer to the results of 24 large-scale association studies (average sample size 121 000). First we show that these studies have tended to substantially over-correct for confounding, and as a result the number of genome-wide significant loci has under-reported by about 20%. Next we estimate enrichment for 24 categories of SNPs defined by functional annotations. A previous study using LDSC reported that conserved regions were 13-fold enriched, and found a further twelve categories with above 2-fold enrichment. By contrast, our analysis using SumHer finds that conserved regions are only 1.6-fold (SD 0.06) enriched, and that no category has enrichment above 1.7-fold. SumHer provides an improved understanding of the genetic architecture of complex traits, which enables more efficient analysis of future genetic data.

LD Score Regression (LDSC) has been frequently used to analyze summary statistics from genome-wide association studies (GWAS).¹⁻⁴ It has four main uses: to estimate the average bias due to confounding, to estimate the “SNP heritability” of a trait (the proportion of phenotypic variance explained by all SNPs), to estimate the heritability enrichments of SNP categories, and to estimate the genetic correlations between pairs of traits. LDSC estimates are derived from a specific heritability model in which each SNP in the genome is expected to contribute equally.¹ Although this model is widely used in statistical genetics, we recently showed that across a range of human traits, it poorly reflects reality.⁵ In particular, it fails to appreciate that in regions of high linkage disequilibrium (LD), the average heritability of each SNP tends to be lower due to multiple tagging of causal variation.⁶ As a result of this model misspecification, LDSC tends to over-estimate confounding bias, under-estimate SNP heritability and produce exaggerated estimates of enrichment.

We propose SumHer, software for estimating SNP heritability from summary statistics that allows the user to specify the heritability model. We apply SumHer to publicly-available GWAS results for 24 disease and quantitative traits,⁷ using a heritability model that we have previously shown to perform well.^{5,6} We first show that these GWAS tended to over-correct for confounding; when we adjust their results using SumHer, the total number of genome-wide significant loci increases from 1 760 to 2 190. A previous study by Finucane *et al.* used LDSC to estimate enrichments for 24 categories of SNPs defined by functional annotations.³ The authors concluded that heritability is highly concentrated in specific functional categories; most notably, they estimated that across 17 diseases, conserved regions contribute 35% of SNP heritability, 13-fold higher than their expected contribution. When we repeat this analysis using SumHer and our 24 traits, the estimated enrichments are more modest: for example, conserved regions are estimated to contribute only 5.8% of SNP heritability and the highest enrichment is only 1.7-fold (transcription start sites), consistent with an omnigenic model of genetic architecture.⁸ We finish by providing an example of how results from SumHer can enable more efficient analysis of genetic data. We show how for body mass index, height, HDL & LDL cholesterol and triglyceride, we are able to significantly improve the predictive performance of polygenic risk scores by incorporating our preferred heritability model and estimates of enrichments. We make SumHer freely available within our software package LDAK (www.ldak.org).⁶

41 Results

42 SumHer

43 SumHer has the same four aims as LDSC;¹⁻³ we outline them here, with methodological details provided in Online Methods. Suppose
44 we are provided with summary statistics from a GWAS where each of m SNPs has been tested individually for association with a
45 particular trait. Suppose also that we have access to a well-matched reference panel, from which we can reliably estimate r_{jl}^2 , the
46 squared correlation between SNPs j and l . Let S_1, S_2, \dots, S_m denote the $\chi^2(1)$ test statistics from single-SNP analysis; the first aim is
47 to estimate the average inflation of these test statistics due to confounding. Let h_j^2 denote the heritability directly contributed by SNP j ;
48 the second aim is to estimate $h_{\text{SNP}}^2 = \sum_{j=1}^m h_j^2$, the SNP heritability of the trait.⁹ Let \mathbb{C} index a category of SNPs; the third aim is
49 to estimate $(\sum_{j \in \mathbb{C}} h_j^2)/h_{\text{SNP}}^2$, the proportion of SNP heritability contributed by SNPs in \mathbb{C} (we can then estimate the enrichment of
50 the category by dividing its estimated proportion of SNP heritability by its expected proportion). Finally, if we are also provided with
51 summary statistics from a second trait, the fourth aim is to estimate the correlation between SNP effect sizes for the two traits.¹⁰

52 In order to achieve these four aims, we must specify a heritability model, which describes how h_j^2 is expected to vary across the
53 genome. Suppose this heritability model takes the form $\mathbb{E}[h_j^2] \propto q_j$. The main difference between LDSC and SumHer is that SumHer
54 allows for any heritability model (i.e., the user can specify arbitrary q_j), whereas LDSC assumes all q_j are the same. We recommend
55 using SumHer with the “LDAK Model”: $q_j = [f_j(1 - f_j)]^{0.75} w_j$, where f_j is the minor allele frequency (MAF) of SNP j and w_j is
56 a weighting based on local levels of LD.^{5,6} In this model, a SNP with high MAF is expected to contribute more heritability than one
57 with low MAF, while a SNP in a region of low LD is expected to contribute more than one in a region of high LD. By contrast, LDSC
58 estimates are obtained by setting $q_j = 1$, which corresponds to the assumption that all SNPs are expected to contribute equally.¹ We
59 refer to this as the “GCTA Model” as this is a core assumption of the software GCTA.^{5,9}

60 A second difference between LDSC and SumHer is how they estimate confounding bias. In a GWAS with no confounding,
61 $\mathbb{E}[S_j] \approx 1 + nv_j^2$, where n is the sample size and $v_j^2 = h_j^2 + \sum_l \text{near } j r_{jl}^2 h_l^2$ is the heritability tagged by SNP j (a working definition
62 of “near” is within 1 Mb¹). Both LDSC and SumHer estimate the deviation of test statistics from their expected values assuming no
63 confounding. LDSC uses the model $\mathbb{E}[S_j] \approx 1 + A + nv_j^2$, where A indicates the average amount each test statistic is inflated additively
64 due to confounding (LDSC then reports $1 + A$, which it refers to as “the intercept”). By contrast, we recommend using the model
65 $\mathbb{E}[S_j] \approx C(1 + nv_j^2)$, where C reflects how much each test statistic is inflated multiplicatively. There are two reasons why we prefer our
66 approach. Firstly, it is standard practice to correct test statistics by scaling (i.e., divide each by C);¹⁰ in theory, one could instead shift
67 test statistics (i.e., subtract A from each), but this would result in some negative values. Secondly, the SumHer model for estimating
68 bias accommodates test statistics that have been subjected to genomic control.¹¹ Although genomic control is intended to reduce bias,
69 we find that it is often the biggest source of bias in GWAS results and a major hindrance when using summary statistics to interrogate
70 genetic architecture (see below).

71 In total we use six versions of SumHer, which differ according to their assumed heritability model and allowance for confounding
72 bias. **LDSC-Zero** assumes the GCTA Model and that there is no bias ($A = 0, C = 1$); this is equivalent to using the LDSC software¹
73 with the option `--intercept-h2 1`. **LDSC** assumes the GCTA Model and allows for additive bias (A free to vary, $C = 1$); this
74 is equivalent to using the LDSC software¹ with default options. **SumHer-Zero** assumes the LDAK Model and that there is no bias
75 ($A = 0, C = 1$); this is our recommended version when estimating h_{SNP}^2 or enrichments and confident that confounding is negligible.
76 **SumHer-GC** assumes the LDAK Model and allows for multiplicative bias ($A = 0, C$ free to vary); this is our recommended version
77 when estimating confounding or genetic correlations, or when estimating h_{SNP}^2 or enrichments and it is likely that test statistics are
78 biased due to population structure or relatedness, or were obtained using genomic control or mixed-model association analysis (see
79 below). **Hybrid-Zero** assumes the heritability model

$$q_j = (1 - p) \times 1/m + p \times [f_j(1 - f_j)]^{0.75} w_j / Q' \quad \text{where} \quad Q' = \sum_j [f_j(1 - f_j)]^{0.75} w_j, \quad (1)$$

80 and that there is no bias ($A = 0, C = 1$), while **Hybrid-GC** assumes the same heritability model but allows for multiplicative bias
81 ($A = 0, C$ free to vary). Model (1) is a linear combination of the GCTA and LDAK models, where p indicates the weight assigned to

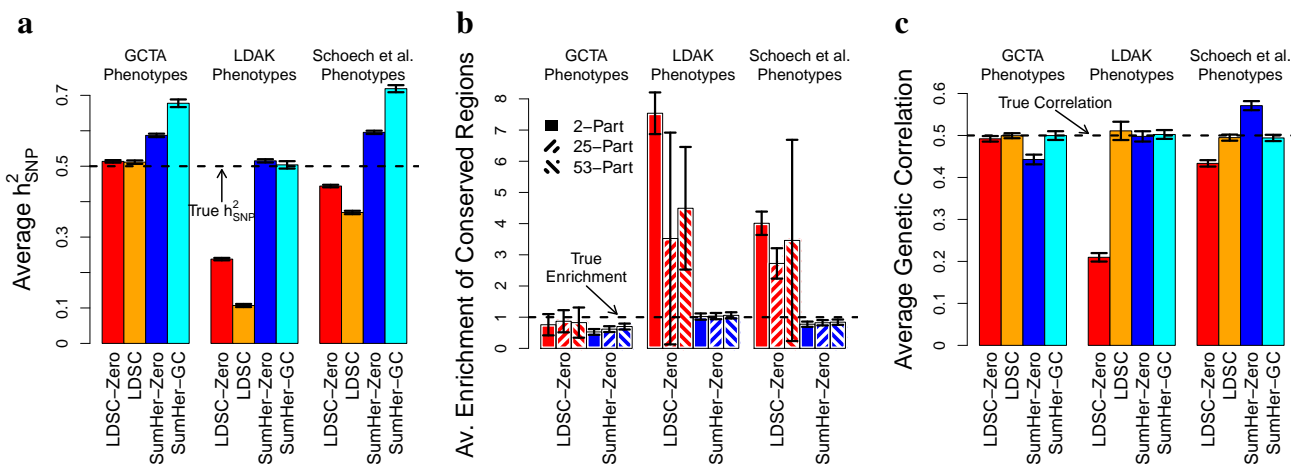


Figure 1: **Estimates can depend sensitively on the heritability model.** (a) We generate 500 phenotypes with $h^2_{\text{SNP}} = 0.5$ for each of three heritability models, GCTA, LDAK and Schoech *et al.*¹³ (see main text for details of each model). Bars report average estimates of h^2_{SNP} from LDSC-Zero, LDSC, SumHer-Zero and SumHer-GC. (b) For the same phenotypes, bars now report average estimates of the enrichment of heritability in conserved regions from LDSC-Zero and SumHer-Zero (true enrichment is 1). For each method, the three bars correspond to estimating enrichment using a 2-part, 25-part or 53-part model (see main text). (c) For each of the three heritability models, we now generate 500 additional pairs of phenotypes with genetic correlation 0.5. Bars report average estimates of genetic correlation from LDSC-Zero, LDSC, SumHer-Zero and SumHer-GC. In all plots, vertical line segments mark 95% confidence intervals for the average estimates.

82 the LDAK model; using this heritability model allows us to compare the fit of the GCTA and LDAK models on real data (see below).

83 For all analyses, our reference panel is 8 850 unrelated Caucasian individuals from the Health and Retirement Study¹² (HRS).
 84 When estimating enrichments of SNP categories, we use the 24 functional annotations used by Finucane *et al.*,³ which include coding,
 85 conserved, enhancer and promoter regions (see Supplementary Table 1 for a full list). When analyzing real data, we exclude SNPs within
 86 the major histocompatibility complex (Chromosome 6: 25-34 Mb), as well as SNPs which individually explain $>1\%$ of phenotypic
 87 variation, and SNPs in LD with these (within 1 cM and $r^2_{jl} > 0.1$).

88 Simulated phenotypes

89 We first demonstrate the importance of choosing an appropriate heritability model via simulations. For this we use 7 548 unrelated
 90 Caucasian individuals from the three control cohorts of the Wellcome Trust Case Control Consortium (WTCCC), recorded for 3 280 768
 91 common SNPs (MAF >0.01). We first generate 1 000 phenotypes each with 2 000 causal SNPs and $h^2_{\text{SNP}} = 0.5$; for half the phenotypes,
 92 we sample causal SNP effect sizes according to the GCTA Model, for the other half, according to the LDAK Model. We then analyze
 93 each phenotype using LDSC-Zero, LDSC, SumHer-Zero and SumHer-GC.

94 Figure 1a shows that, as expected, accurate estimates of h^2_{SNP} are returned when phenotypes are analyzed assuming the matching
 95 heritability model (i.e., when GCTA phenotypes are analyzed using LDSC-Zero or when LDAK phenotypes are analyzed using SumHer-
 96 Zero), but that using a different heritability model can result in very poor estimates; SumHer-Zero tends to over-estimate h^2_{SNP} by
 97 about 20% when applied to GCTA phenotypes, while LDSC-Zero tends to under-estimate h^2_{SNP} by about 50% when applied to LDAK
 98 phenotypes. Supplementary Figure 1a shows LDSC correctly infers that there is no confounding when used on GCTA phenotypes
 99 (average intercept 1.000, SD 0.0003), and therefore its estimates of h^2_{SNP} closely match those from LDSC-Zero. However, when used
 100 on LDAK phenotypes, LDSC wrongly infers that much of the causal signal is in fact confounding (average intercept 1.033, SD 0.0003),
 101 and as a result, its estimates of h^2_{SNP} are on average about half those from LDSC-Zero and about 75% lower than the true value.

102 Figure 1b reports the estimated enrichment of SNPs in conserved regions. As causal SNPs were picked at random from across the

103 genome, the true enrichment is 1. Again, we see that assuming the correct heritability model produces reliable estimates, but assuming
104 the wrong model can lead to misleading conclusions. In particular, we find that when LDSC-Zero is used to analyze LDAK phenotypes,
105 it infers that conserved regions are at least 3-fold enriched for heritability. We chose to focus on conserved regions as this was the
106 category that, by applying LDSC to real data, Finucane *et al.*³ found to be most enriched. We have previously shown that the LDAK
107 Model better reflects real data than the GCTA Model⁵ (and provide further evidence below), and thus our simulations suggest that a
108 substantial portion of the enrichment observed by Finucane *et al.* is an artifact of misspecifying the heritability model.

109 Figure 1b also examines how estimates of enrichment are affected by the way the genome is divided into SNP categories (Sup-
110 plementary Fig. 1c provides a zoomed-in version). The simplest approach is to use a 2-part model, in this case partitioning SNPs into
111 those inside and those outside conserved regions. Next we use a 25-part model, dividing the genome into the 24 functional categories
112 (of which conserved regions are one), plus a category containing all SNPs. Finally, we use a 53-part model, constructed by adding 28
113 “buffer regions”, the approach recommended by Finucane *et al.* We find that when the correct heritability model is assumed, results ap-
114 pear insensitive to the approach used. However, when the wrong model is assumed, the three approaches can give substantially different
115 estimates.

116 For Figure 1c, we generate 500 additional pairs of phenotypes, each with genetic correlation 0.5 (again, half the phenotypes are
117 generated under the GCTA Model, half under the LDAK Model). As expected, using LDSC-Zero to analyze GCTA phenotypes or
118 SumHer-Zero to analyze LDAK phenotypes produces accurate estimates of genetic correlation, whereas using LDSC-Zero on LDAK
119 phenotypes or SumHer-Zero on GCTA phenotypes results in biased estimates. Surprisingly, it appears that biases can be minimized by
120 allowing for confounding bias in the analysis (Supplementary Fig. 1f shows why this is the case). However, even if both LDSC and
121 SumHer-GC produce unbiased estimates of genetic correlation, it remains that highest precision is achieved when the correct heritability
122 model is assumed; the SD of LDSC estimates is about half that of SumHer-GC estimates when analyzing GCTA phenotypes, but about
123 twice as high when analyzing LDAK phenotypes.

124 As well as the GCTA and LDAK phenotypes, for each analysis we additionally generate phenotypes according to a heritability
125 model recently proposed by Schoech *et al.*¹³ We formally define this model in the Discussion, however, loosely speaking, it is interme-
126 diate between the GCTA and LDAK models. This is reflected by the estimates of h^2_{SNP} in Figure 1a: on average those from LDSC-Zero
127 are about 5% too low, while those from SumHer-Zero are about 10% too high. We also consider changing the reference panel. While
128 we prefer using the HRS dataset, as its larger sample size enables more accurate estimation of r^2_{jl} , Supplementary Figure 2 shows that
129 estimates are similar if we instead use the 404 non-Finnish Europeans from the 1000 Genomes Project.¹⁴

130 Real phenotypes

131 We now show that the choice of heritability model is also important when analyzing real data. We use “25 raw GWAS” (18 binary traits,
132 7 quantitative, average sample size 9 700; see Supplementary Table 2), for which we have individual-level genotype and phenotype
133 data from either the WTCCC¹⁵ or the eMERGE Network.¹⁶ The 13 WTCCC GWAS all examine diseases (e.g., Bipolar Disorder,
134 Ischaemic Stroke and Parkinson’s Disease), while the 12 eMERGE GWAS consider a mixture of diseases (e.g., age-related macular
135 degeneration, heart failure and peripheral artery disease) and clinical measurements (e.g., blood pressure, height and lipid levels). After
136 imputation, strict quality control, and excluding SNPs not present in our reference panel, the WTCCC data contain on average 2.3 M
137 SNPs, while the eMERGE data contain 2 972 162 SNPs. When performing single-SNP analysis, we include as covariates sex and ten
138 principal components (five calculated from the data, five derived from the 1000 Genomes Project¹⁴). As we have access to raw data,
139 we can use REML to estimate how much of the total phenotypic variance explained by SNPs is inflation due to population structure
140 or relatedness.^{17,18} We estimate that on average 3.6% of the variance explained is inflation (range -0.9% to 6.8%), indicating that
141 confounding due to population structure and relatedness is modest (Supplementary Fig. 3).

142 Figure 2a and Supplementary Table 2 show that across the 25 traits, estimates of h^2_{SNP} from SumHer-Zero are on average 2.0
143 times larger (SD 0.05) than those from LDSC-Zero. Figure 2b and Supplementary Table 3 report estimates of enrichment for the 24
144 functional categories, averaged across the 25 traits. First we estimate enrichment using LDSC-Zero with a 53-part model (the approach

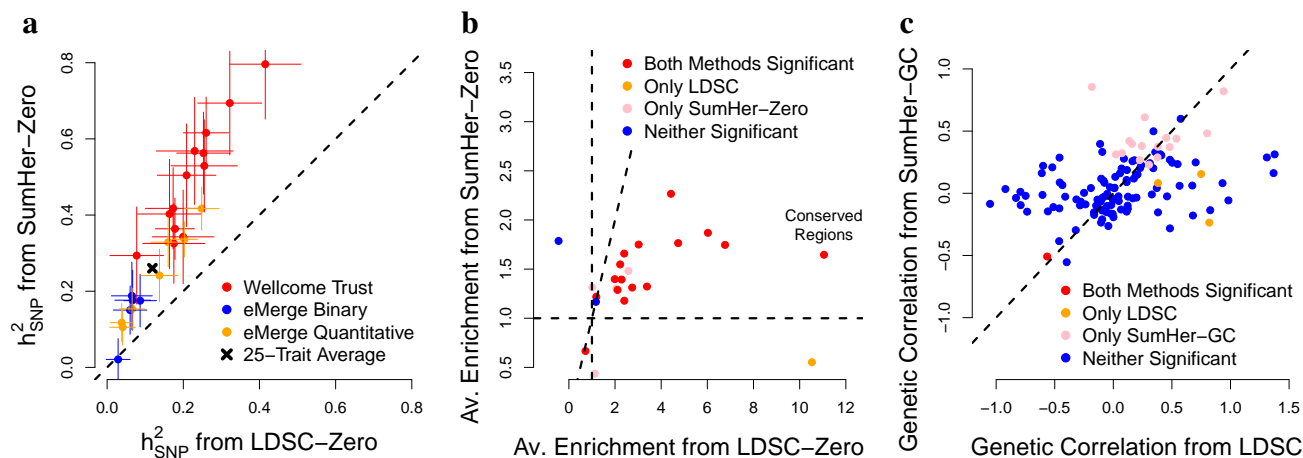


Figure 2: **Importance of the heritability model for the 25 raw GWAS.** (a) Estimates of h^2_{SNP} from LDSC-Zero (x axis) and SumHer-Zero (y axis). Colors distinguish between the 13 WTCCC, the 5 binary eMERGE and 7 quantitative eMERGE traits (black denotes the 25-trait average). Horizontal and vertical line segments mark 95% confidence intervals. Estimates for binary traits are on the observed scale. (b) Average estimates of enrichment for the 24 functional categories from LDSC-Zero (using a 53-part model) and SumHer-Zero (using a 25-part model). Colors indicate significant enrichments ($P < 0.05$) from one or both methods. (c) Estimates of genetic correlations between pairs of eMERGE traits using LDSC (x axis) and SumHer-GC (y axis). Colors indicate significant correlations ($P < 0.05$) from one or both methods.

145 taken by Finucane *et al.*³), then using SumHer-Zero with a 24-part model (our recommended approach). While there is reasonable
 146 agreement between which categories SumHer-Zero and LDSC-Zero declare to be significant, their estimates of enrichment are very
 147 different. Notably, LDSC-Zero estimates that conserved regions on average contribute 28% (SD 4) of SNP heritability (corresponding to
 148 11-fold enrichment), whereas SumHer-Zero estimates that they contribute only 5.1% (SD 0.7) of SNP heritability (1.6-fold enrichment).
 149 Supplementary Figure 4 compares estimates of enrichment from 2-part, 25-part and 53-part models. When we use LDSC-Zero, we find
 150 substantial differences between the results of the three approaches; in particular, not one of the 24 estimates from the 2-part model is
 151 consistent ($P > 0.05/24$) with either of the corresponding estimates from the 25-way or 53-way models. By contrast, when we use
 152 SumHer-Zero, there is much stronger concordance between the three sets of results; for example, all 24 (19 out of 24) of the 2-part
 153 model estimates are consistent with those from the 25-way (53-way) model.

154 Figure 2c and Supplementary Table 4 report genetic correlations between pairs of traits. In general, it is only possible to get a
 155 meaningful estimate when both traits have substantial h^2_{SNP} , so we restrict to the 18 traits for which both LDSC-Zero and SumHer-Zero
 156 find significant h^2_{SNP} ($P < 0.05/25$). As predicted by our earlier simulations, we find good concordance between the genetic correlations
 157 reported by LDSC and SumHer-GC, but that the latter produces more precise estimates: on average the SumHer-GC estimates have SD
 158 about half that of the LDSC estimates, and SumHer-GC finds 18 pairs of traits with significant correlation ($P < 0.05$), whereas LDSC
 159 finds only 4.

160 Comparing heritability models

161 Previously, we compared different heritability models based on the likelihood from REML analysis; we showed that across 42 human
 162 traits (average sample size 7 400), log likelihood was on average 9.8 higher if we assumed the LDK Model rather than the GCTA
 163 Model.⁵ Those 42 traits included the 13 WTCCC traits we study here. If we repeat that analysis using the 12 eMERGE traits (average
 164 sample size 13 000), we find that log likelihood is on average 17 higher under the LDK Model (Supplementary Table 5). In Sup-
 165 plementary Table 6, we show that it remains possible to compare models based on likelihood if only summary statistics are available.
 166 However, an alternative, and easier to visualize, method is to fit both the GCTA and LDK Models simultaneously and allow the data to
 167 decide the relative weighting of each. Specifically, we use Hybrid-Zero, with the focus on estimating p in Model (1), the “proportion of
 168 LDK” in the heritability model. Figure 3a demonstrates proof of principle: we see that Hybrid-Zero correctly estimates $p = 0$ when

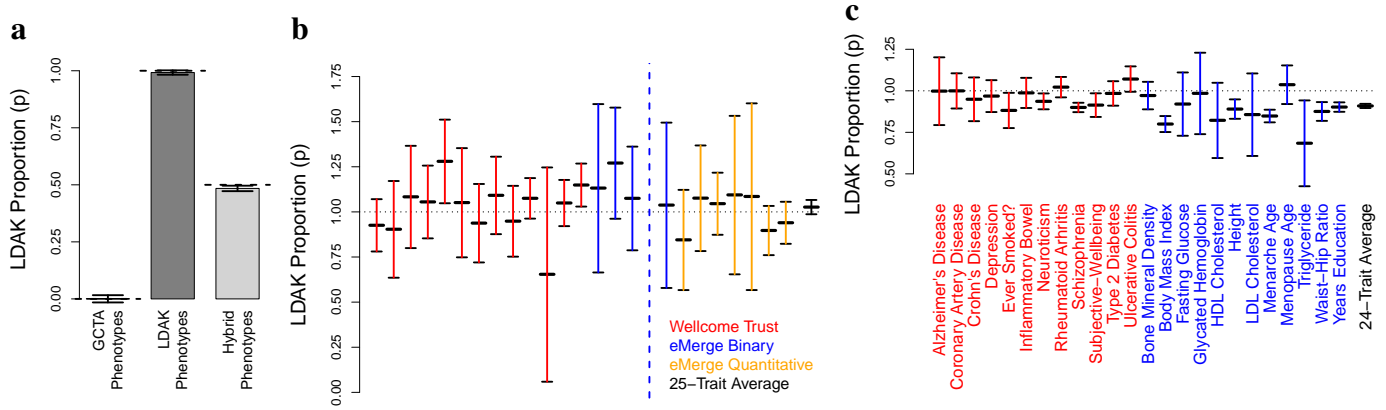


Figure 3: **Comparing the GCTA and LDAK Models.** These analyses use Hybrid-Zero and Hybrid-GC, versions of SumHer which assign weights $1-p$ and p to the GCTA and LDAK heritability models, respectively. (a) Estimates of p from Hybrid-Zero for GCTA phenotypes (true $p = 0$), LDAK phenotypes (true $p = 1$) and hybrid phenotypes (true $p = 0.5$). (b) Estimates of p from Hybrid-Zero for the 25 raw GWAS. Colors distinguish between the 13 WTCCC, the 5 binary eMERGE and 7 quantitative eMERGE traits (black denotes the 25-trait average). A precise estimate of p was not possible for Shingles (dashed vertical line), due to the trait having very low h^2_{SNP} . (c) Estimates of p from Hybrid-GC for the 24 summary GWAS and the 24-trait average. In all plots, vertical line segments mark 95% confidence intervals.

169 applied to GCTA phenotypes, $p = 1$ for LDAK phenotypes, and $p = 0.5$ for “Hybrid Phenotypes” (each created by summing a GCTA
 170 and an LDAK phenotype). We then apply Hybrid-Zero to our 25 raw GWAS (Figure 3b and Supplementary Table 6), finding on average
 171 $p = 1.03$ (SD 0.02), indicating that the data overwhelmingly support the LDAK Model over the GCTA Model.

172 Population structure and relatedness

173 Until recently, it was standard to estimate confounding bias via the genomic inflation factor¹¹ (GIF). However, the GIF tends to over-
 174 estimate bias, because it makes the assumption that all observed inflation of test statistics is due to confounding.¹⁹ LDSC provides a
 175 method for estimating confounding bias which appreciates that substantial inflation can instead be due to causal variation.¹ However,
 176 our above simulations indicate that LDSC also tends to over-estimate bias, due to the poor fit of the GCTA heritability model. This
 177 is supported by our analysis of the 25 raw GWAS (Supplementary Table 2), for which we performed very careful quality control and
 178 verified using REML that confounding due to population structure and relatedness is modest (Supplementary Fig. 3). LDSC finds
 179 substantial bias; its average estimate of the intercept ($1 + A$) is 1.031 (SD 0.002), which is similar to the average GIF, 1.035. By contrast,
 180 SumHer-GC finds that bias is slight; its average estimate of the scaling factor (C) is 1.001 (SD 0.002).

181 We now construct GWAS when there is substantial confounding. For each of the 13 WTCCC GWAS, we replace 2504 of the
 182 controls (on average 67%) with 2504 individuals from POBI²⁰ (People of the British Isles); this generates population structure because,
 183 although both WTCCC and POBI individuals were recruited from the UK, the latter predominately came from isolated, rural regions
 184 (Supplementary Figure 5). Supplementary Table 7 reports estimates of confounding bias and h^2_{SNP} for each GWAS. SumHer-GC now
 185 finds substantial bias; its average estimate of the scaling factor is 1.049 (SD 0.003). For comparison, the average GIF is 1.098, while
 186 the average LDSC intercept is 1.088 (SD 0.002). SumHer-Zero estimates of h^2_{SNP} after switching controls are on average 95% (SD 4)
 187 higher than SumHer-Zero estimates prior to switching, demonstrating how population structure can lead to substantial inflation when
 188 estimating SNP heritability. However, SumHer-GC estimates of h^2_{SNP} after switching controls are on average only 4% (SD 7) higher,
 189 consistent with SumHer-GC being able to reliably estimate SNP heritability despite the population structure.

190 Genomic control and mixed-model association analysis

191 Although its use is declining, it remains that the majority of published GWAS have performed genomic control (divided test statistics by
 192 the GIF) at least once in their analyses.¹¹ As the GIF tends to over-estimate confounding,¹⁹ genomic control tends to produce negatively

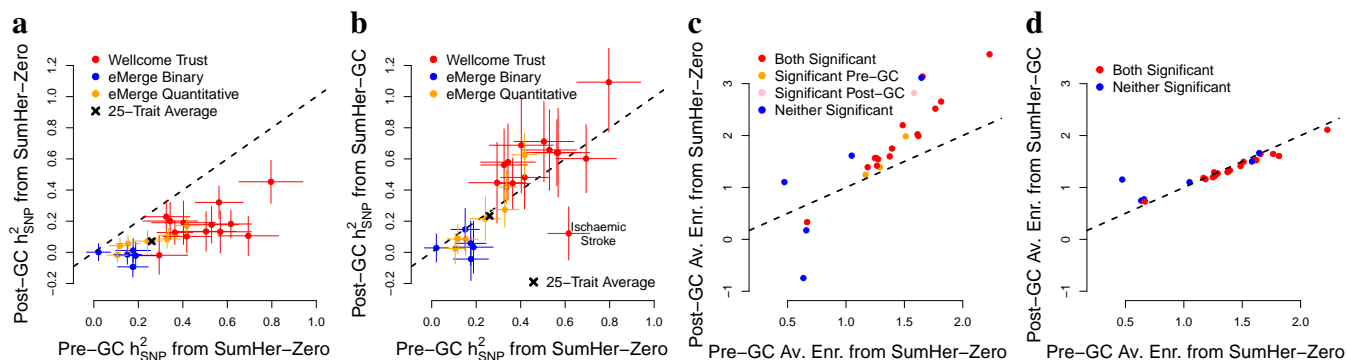


Figure 4: **Correcting for genomic control.** For the 25 raw GWAS, each plot compares estimates from SumHer using raw test statistics (x -axis) with those using test statistics subjected to genomic control (y -axis). Horizontal and vertical line segments mark 95% confidence intervals. **(a)** Estimates of h^2_{SNP} using SumHer-Zero both pre and post genomic control. **(b)** Estimates of h^2_{SNP} using SumHer-Zero pre genomic control and SumHer-GC post genomic control. **(c)** Estimates of average enrichment for the 24 functional categories using SumHer-Zero both pre and post genomic control. **(d)** Estimates of average enrichment for the 24 functional categories using SumHer-Zero pre genomic control and SumHer-GC post genomic control.

193 biased test statistics. For this reason, we recommend using SumHer-GC to estimate h^2_{SNP} and enrichment, instead of SumHer-Zero, if it
 194 is likely that genomic control has been applied. By way of demonstration, we perform genomic control for each of the 25 raw GWAS.
 195 Figure 4a shows that if we continue to use SumHer-Zero, estimates of h^2_{SNP} are centered on zero. However, if we instead use SumHer-
 196 GC, estimates of h^2_{SNP} are in general consistent with estimates from SumHer-Zero prior to genomic control (Figure 4b); an exception is
 197 stroke, indicative of the LDAK Model being sub-optimal for this trait (Supplementary Figure 6). Similarly, Figures 4c & d show that if
 198 genomic control has been performed, it is beneficial to use SumHer-GC instead of SumHer-Zero when estimating enrichments.

199 SumHer, like LDSC, is designed to be used with test statistics from classical regression (see Online Methods). However, with the
 200 development of software such as Fast-LMM, GCTA-LOCO and Bolt-LMM,^{21–23} a popular alternative is for GWAS to use mixed-model
 201 association analysis.^{24,25} Supplementary Figures 6 & 7 show that when estimating h^2_{SNP} and enrichment, the impact of mixed-model
 202 analysis is similar to, albeit less severe than, that of genomic control. For example, SumHer-Zero estimates of h^2_{SNP} based on mixed-
 203 model test statistics are on average about half those based on classical test statistics. However, as with genomic control, reliable estimates
 204 can be obtained by using SumHer-GC instead of SumHer-Zero.

205 Analysis of published GWAS results

206 For our main analysis, we use “24 summary GWAS” (12 binary traits, 12 quantitative, average sample size 121 000; see Table 1). These
 207 are GWAS for which we do not have individual-level data, but have downloaded summary statistics from previously-published analyses.⁷
 208 After excluding SNPs not present in our reference panel, on average there are summary statistics for 2.2M SNPs per trait. All of the
 209 GWAS either performed genomic control or used mixed-model association analysis, so in general we use SumHer-GC. However, we
 210 also repeat all analyses using SumHer-Zero restricted to the 11 traits that did not use mixed-model analysis and for which the impact of
 211 genomic control was lowest (Supplementary Table 8).

212 Figure 3c reports estimates of p , the proportion of LDAK in the heritability model, obtained using Hybrid-GC. Across the 24
 213 traits, we estimate $p = 0.91$ (SD 0.01), indicating that again the data strongly support the LDAK Model over the GCTA Model, although
 214 not as strongly as for the 25 raw GWAS (see Discussion). Supplementary Table 9 shows that we reach the same conclusion if we instead
 215 use Hybrid-Zero restricted to the 11 traits least impacted by genomic control ($p = 0.90$; SD 0.01), or if we compare the GCTA and
 216 LDAK Models based on likelihood (the log likelihood is on average 103 higher under the LDAK Model than the GCTA Model).

217 Confounding

218 Table 1 reports the LDSC intercept ($1 + A$) and the SumHer-GC scaling factor (C) for each trait. LDSC finds that test statistics are on

Trait	n	GIF	LDSC		SumHer-GC		Number of significant loci after dividing test statistics by			
			h_{SNP}^2 (SD)	$1 + A$ (SD)	h_{SNP}^2 (SD)	C (SD)	1	GIF	$1 + A$	C
Alzheimer's Diseases ²⁶	54 000	1.07	0.02 (0.01) [†]	1.07 (0.01)	0.09 (0.02)	1.03 (0.01)	19	17	17	18
Coronary Artery ²⁷	80 000	1.07	0.03 (0.01)	1.07 (0.01)	0.11 (0.01)	0.99 (0.01)	11	7	7	11
Crohn's Disease ²⁸	21 000	1.08	0.18 (0.05) [†]	1.10 (0.02)	0.67 (0.08)	0.97 (0.02)	55	51	50	56
Depression ²⁹	161 000	1.12	0.01 (0.00)	1.05 (0.01)	0.06 (0.01)	0.97 (0.01)	1	0	1	1
Ever Smoked? ³⁰	74 000	1.05	0.04 (0.00) [†]	1.03 (0.01)	0.12 (0.01)	0.96 (0.01)	0	0	0	0
Inflammatory Bowel ²⁸	35 000	1.13	0.12 (0.02) [†]	1.14 (0.02)	0.49 (0.05)	0.98 (0.02)	60	51	51	62
Neuroticism ²⁹	171 000	1.24	0.04 (0.00)	1.09 (0.01)	0.16 (0.02)	0.88 (0.02)	14	5	12	24
Rheumatoid Arthritis ³¹	58 000	1.00	0.06 (0.01)	1.00 (0.01)	0.26 (0.04)	0.89 (0.01)	71	71	71	81
Schizophrenia ³²	81 000	1.35	0.22 (0.01) [†]	1.24 (0.01)	0.65 (0.03)	0.90 (0.01)	82	35	44	108
Subjective-Wellbeing ²⁹	298 000	1.12	0.01 (0.00)	1.09 (0.01)	0.05 (0.01)	0.93 (0.02)	1	0	0	2
Type 2 Diabetes ³³	157 000	1.08	0.03 (0.00)	1.08 (0.01)	0.11 (0.01)	0.95 (0.01)	31	27	26	33
Ulcerative Colitis ²⁸	27 000	1.09	0.07 (0.02) [†]	1.11 (0.01)	0.38 (0.05)	1.00 (0.01)	30	28	27	30
Bone Mineral Density ³⁴	33 000	1.08	0.08 (0.01) [†]	1.07 (0.01)	0.27 (0.03)	1.00 (0.01)	16	16	16	16
Body Mass Index ³⁵	229 000	0.92	0.07 (0.01)	0.82 (0.01)	0.28 (0.02)	0.56 (0.02)	60	70	92	239
Fasting Glucose ³⁶	58 000	1.04	0.04 (0.01)	1.05 (0.01)	0.11 (0.02)	0.99 (0.01)	20	18	18	21
Glycated Hemoglobin ³⁷	46 000	1.02	0.02 (0.01) [†]	1.03 (0.01)	0.07 (0.02)	1.00 (0.01)	8	8	8	8
HDL Cholesterol ³⁸	95 000	0.95	0.08 (0.02)	1.00 (0.04)	0.35 (0.04)	0.74 (0.02)	99	106	99	142
Height ³⁹	245 000	1.69	0.16 (0.01) [†]	1.73 (0.05)	0.38 (0.03)	1.02 (0.03)	580	242	234	563
LDL Cholesterol ³⁸	90 000	0.96	0.07 (0.03)	1.00 (0.04)	0.36 (0.08)	0.74 (0.04)	83	89	83	117
Menarche Age ⁴⁰	253 000	1.36	0.12 (0.01) [†]	1.25 (0.02)	0.28 (0.02)	0.91 (0.02)	244	134	160	293
Menopause Age ⁴¹	69 000	1.05	0.05 (0.01)	1.07 (0.02)	0.21 (0.03)	0.93 (0.02)	36	31	30	42
Triglyceride Levels ³⁸	92 000	0.96	0.11 (0.03)	0.94 (0.03)	0.34 (0.07)	0.73 (0.03)	67	71	72	109
Waist-Hip Ratio ⁴²	141 000	0.95	0.05 (0.00)	0.93 (0.01)	0.17 (0.01)	0.77 (0.01)	21	25	28	48
Years Education ⁴³	329 000	1.18	0.06 (0.00) [†]	1.14 (0.01)	0.17 (0.01)	0.84 (0.01)	151	88	88	166
Average	121 000	1.11		1.06 (0.00)		0.93 (0.00)	73	50	51	91
Total							1760	1190	1234	2190

Table 1: **Estimates of confounding bias for the 24 summary GWAS.** Columns 2 & 3 report the average sample size (n) and genomic inflation factor (GIF) for each trait. Columns 4-7 report estimates of h_{SNP}^2 and confounding from both LDSC and SumHer-GC (LDSC measures confounding via the intercept $1 + A$, while SumHer-GC uses the scaling factor C). Columns 8-11 report the number of significant loci based on the published test statistics, then after correction via genomic control, LDSC and SumHer-GC (dividing test statistics by the GIF, $1 + A$ and C , respectively). When estimating h_{SNP}^2 , LDSC requires uncorrected test statistics from classical regression. However the test statistics from all 24 GWAS were calculated either using genomic control or mixed-model association analysis; [†] indicates the 11 traits for which classical regression was used and the impact of genomic control was lowest.

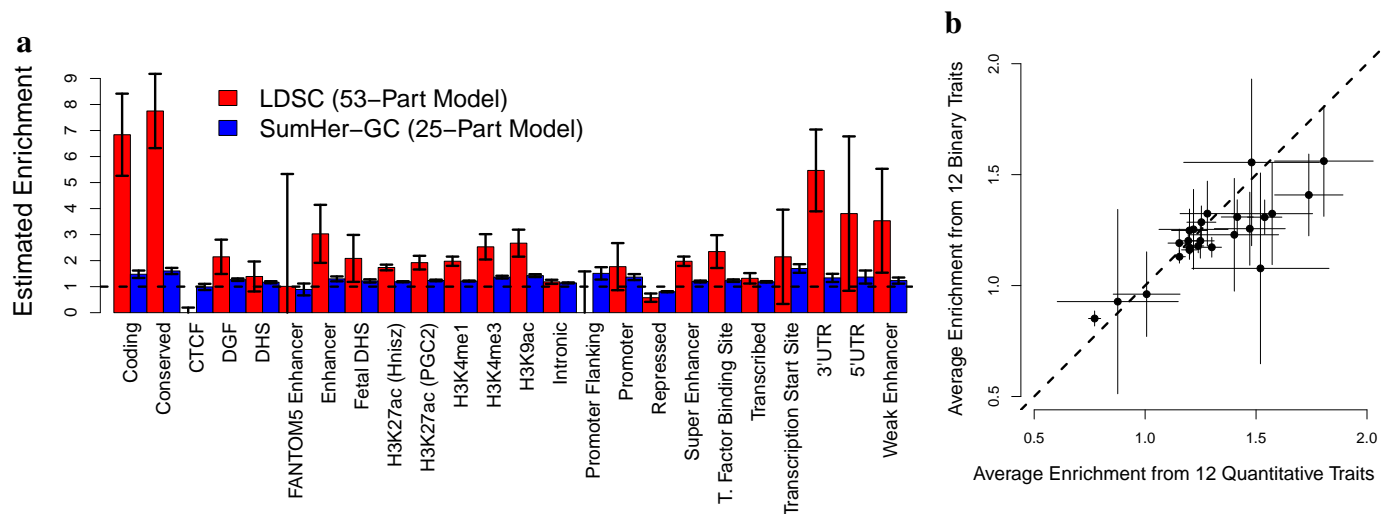


Figure 5: **Functional enrichments across the 24 summary GWAS.** Both plots show average estimated enrichments for the 24 functional categories. Horizontal and vertical line segments mark 95% confidence intervals. **(a)** Estimates from LDSC using a 53-part model (red bars) and from SumHer-GC using a 25-part model (blue bars). **(b)** Estimates from SumHer-GC using a 25-part model, based either on the 12 quantitative traits (x -axis) or on the 12 binary traits (y -axis).

219 average inflated by 5.7% (SD 0.2), implying that the GWAS tended to under-correct for confounding. By contrast, SumHer-GC finds that
 220 test statistics are on average deflated by 7.4% (SD 0.3), indicating that the studies tended to over-correct. We note that the four GWAS
 221 with lowest C' (0.56 for body mass index, 0.74 for HDL and LDL cholesterol, and 0.73 for triglyceride levels), are all meta-analyses that
 222 used genomic control both before and after combining results across cohorts.^{35,38} Table 1 also reports the number of independent loci
 223 with $P < 5 \times 10^{-8}$. Without adjustment (i.e., using the published test statistics), there are on average 73 loci per trait. If we correct
 224 using LDSC, this number is reduced to 51, but if we correct using SumHer-GC, it increases to 91. For these counts, we defined two
 225 significant SNPs as dependent if they are within 1 cM and have $r_{jl}^2 > 0.05$, but we find that results are very similar if instead we increase
 226 the window size to 3 cM, or use $r_{jl}^2 > 0.2$ (Supplementary Table 10).

227 Functional enrichments

228 Figure 5a and Supplementary Table 11 report estimates of enrichment for the 24 functional categories, averaged across the 24 traits.
 229 We again see striking differences between the estimates from LDSC (using a 53-part model) and those from SumHer-GC (using a 25-
 230 part model). For example, LDSC estimates of enrichment range from -1.5 to 7.8, whereas SumHer-GC estimates range from 0.80 to
 231 1.7. Supplementary Table 12 shows that large differences remain if we instead estimate enrichment using LDSC-Zero and SumHer-
 232 Zero restricted to the 11 traits least impacted by genomic control. Based on the results from SumHer-GC, we conclude that conserved
 233 regions^{44,45} and transcription start sites⁴⁶ are most enriched for heritability; they are estimated to contribute 5.8% (SD 0.2) and 3.0%
 234 (SD 0.2) of h_{SNP}^2 , respectively, 1.6 (SD 0.06) and 1.7 (SD 0.08) times higher than their expected contributions under the LDK Model.
 235 Repressed regions⁴⁶ are the only category significantly depleted; although they are estimated to contribute 36% (SD 0.4) of h_{SNP}^2 , this is
 236 0.80 (SD 0.01) times lower than expected. Further results are provided in Supplementary Figure 8, which confirm that LDSC estimates
 237 are similar if we change the SNP sets (LDSC recommends that the reference panel contains as many SNPs as possible, but that only
 238 HapMap 3⁴⁷ SNPs with MAF > 0.05 are used when performing the regression^{1,3}), or if we use the 75-part model used by Gazal *et al.*⁴⁸

239 Genetic correlations

240 Supplementary Figure 9 and Supplementary Table 13 provide estimates of genetic correlations for the 276 pairs of traits. As expected,
 241 there is strong concordance between estimates from LDSC and SumHer-GC, but the SumHer-GC estimates are more precise; for exam-
 242 ple, across the 41 pairs of traits that both methods find to be significantly correlated ($P < 0.05/276$), the SD of SumHer-GC estimates
 243 is on average a third lower than the SD of LDSC estimates.

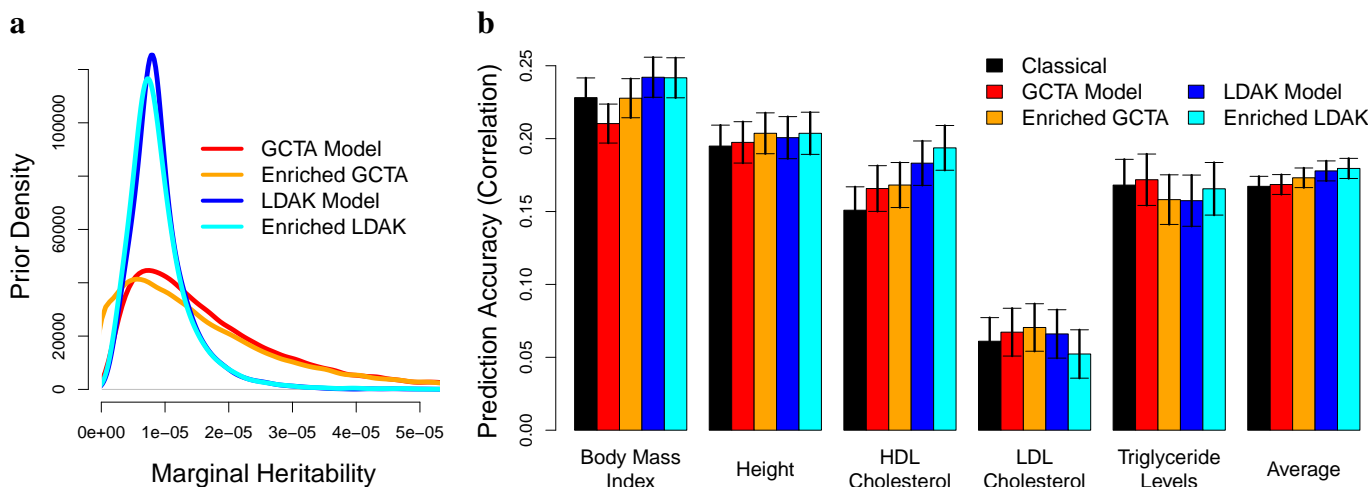


Figure 6: **Prediction of five quantitative traits.** For each trait, we use data from the 24 summary GWAS to construct Bayesian polygenic risk scores (PRS) corresponding to four heritability models: GCTA, Enriched-GCTA, LDAK and Enriched-LDAK (see main text for details of each model). (a) The prior distribution of v_j^2 , the heritability tagged by SNP j , corresponding to each heritability model. (b) Prediction accuracy, measured as correlation between observed and predicted phenotypes in the (independent) eMERGE data, for the Classical PRS (effect sizes are frequentist estimates from single-SNP analysis) and for each of the four Bayesian PRS (effect sizes are posterior means). Vertical line segments mark 95% confidence intervals.

244 Improving the efficiency of future analyses

245 Figure 5b shows that there is strong concordance between the average estimates of enrichment obtained from the 12 binary traits and those
 246 from the 12 quantitative traits. This suggests broad similarities between the genetic architectures of different traits, which in turn implies
 247 that it should be possible to use information from existing GWAS to improve the efficiency of future analyses. As a demonstration, we
 248 consider prediction using polygenic risk scores (PRS). We focus on body mass index, height, HDL & LDL cholesterol and triglyceride
 249 levels, as for these five traits we can train prediction models using the 24 summary GWAS, then measure how well these perform on the
 250 independent eMERGE data.

251 To construct a PRS, we need estimates of SNP effect sizes. The current standard is to use estimates from single-SNP analysis
 252 (“Classical PRS”). However, in Online Methods, we explain how, given a heritability model, we can obtain a prior distribution for v_j^2 ,
 253 the heritability tagged by SNP j , then calculate a “Bayesian PRS” using the posterior mean effect sizes. For each trait, we construct four
 254 Bayesian PRS corresponding to four heritability models. First we use the GCTA heritability model ($q_j = 1$). Next we use the “Enriched
 255 GCTA Model”, obtained by scaling the q_j based on the (53-part) estimates of enrichment (e.g., if a category was estimated to have 2-fold
 256 enrichment, then the SNPs it contains would have average $q_j = 2$). We similarly construct PRS based on the LDAK Model, then the
 257 “Enriched LDAK Model”, where the q_j are scaled according to the SumHer-GC (25-part) estimates of enrichment.

258 Figure 6a compares the four prior distributions for v_j^2 . We see the two priors derived from the GCTA Model are more diffuse than
 259 the two from the LDAK Model. As explained in the Online Methods, this is because the GCTA Model predicts that v_j^2 scales with local
 260 levels of LD, which vary considerably across the genome, whereas the LDAK Model predicts that v_j^2 scales with local MAFs, which
 261 vary less. Figure 6b and Supplementary Table 14 report the performance of each PRS, measured as correlation between predicted and
 262 observed phenotypes for the eMERGE individuals. Averaged across the five traits, the Bayesian PRS constructed from the GCTA and
 263 Enriched GCTA Model are, respectively, 1.1% (SD 2.0) worse and 2.3% (SD 2.0) better than the Classical PRS, whereas the Bayesian
 264 PRS constructed from the LDAK and Enriched LDAK Model are, respectively, 5.8% (SD 2.0) and 7.5% (SD 2.0) better. The fact that
 265 the PRS based on the Enriched GCTA Model outperforms the PRS based on the GCTA Model, is because introducing partitions relaxes
 266 the unrealistic assumption that v_j^2 scales with local levels of LD; however, the fact that the PRS based on the LDAK Model outperforms
 267 the PRS based on the Enriched GCTA Model, reflects that it is better to start with a more realistic model, than try to fix a sub-optimal
 268 model by partitioning.⁵

269 Discussion

270 We have presented SumHer, software for estimating confounding bias, SNP heritability, enrichments of heritability and genetic corre-
271 lations from GWAS results. While the aims of SumHer are the same as those of LDSC, the key difference is that SumHer allows the
272 user to specify the heritability model. If SumHer is run using the GCTA Model, its estimates will match those from LDSC. However,
273 we instead recommend using the LDAK Model, which we have shown better reflects real data, and therefore produces more accurate
274 estimates. We have analyzed GWAS results for tens of traits, showing that the impact of using an improved heritability model is often
275 substantial, and overall provides a very different description of the genetic architecture of complex traits than has to date been obtained
276 from LDSC analyses.

277 While the GCTA Model has been used almost exclusively when estimating SNP heritability from summary statistics, alternative
278 models have been used when analyzing individual-level data. A recent submission by Schoech *et al.*,¹³ which has three authors in
279 common with the original LDSC publication,¹ uses the following model: $q_j = [f_j(1 - f_j)]^{0.62}(1 - 0.3LLD_j)$, where LLD_j is a local
280 measure of tagging (obtained by computing LD Scores, binning by MAF, quantile normalizing, then truncating). Like the LDAK Model,
281 the model of Schoech *et al.* assumes that a SNP with high MAF contributes more heritability than one with low MAF, and that a SNP in
282 a region of low LD contributes more than one in a region of high LD. Schoech *et al.* claim that this model is more realistic than both the
283 GCTA Model and the LDAK Model.¹³ However, while they compared their model with the GCTA Model on real data (comparing the
284 two models according to REML likelihood), their only comparison with the LDAK Model was using phenotypes simulated according
285 to their heritability model.¹³ Supplementary Table 5 compares the model of Schoech *et al.* to the GCTA and LDAK Models for the 25
286 raw GWAS. While we agree with Schoech *et al.* that their model is superior to the GCTA Model (on average, its log likelihood is 9.4
287 higher), we find it to be inferior to the LDAK Model (on average, its log likelihood is 9.5 lower). Moreover, Figure 1 shows that even
288 if the model of Schoech *et al.* does accurately reflect the genetic architecture of complex traits, it would remain the case that LDSC is
289 producing inaccurate estimates of confounding bias, h_{SNP}^2 and enrichment.

290 We note that, despite remaining superior to the GCTA Model, the LDAK Model does not perform as well for the 24 summary
291 GWAS (average LDAK Proportion 0.91) as for the 25 raw GWAS (average 1.02). One possible reason for this is that for the 24
292 summary GWAS we had to rely on the quality control performed by the original analysts, which was generally much less strict than we
293 would recommend for heritability analysis.¹⁸ For example, when analyzing the 25 raw GWAS, we restricted to SNPs with imputation
294 information score > 0.99 , whereas for the 24 summary GWAS, the most common thresholds were 0.5 and 0.8. There will be a correlation
295 between the heritability of a SNP and its genotyping certainty (a SNP genotyped with error will tag less causal variation than were it
296 perfectly typed),⁵ and similarly, there will be a correlation between genotyping certainty and local levels of LD (low-LD regions tend
297 to contain more low-MAF SNPs, which are often hard to genotype reliably, while imputation is easier in high-LD regions). Therefore,
298 including lower-certainty SNPs in a GWAS will generate correlation between the heritability of each SNP and levels of LD, which will
299 result in traits appearing more “GCTA-like”. Alternatively, the results for the 24 summary traits might simply reflect that the LDAK
300 Model is not perfect, and will fit some traits better than others. Therefore, we encourage readers to find ways to improve the LDAK
301 Model, either generally or on a per-trait level, which they can do using the tools provided by SumHer for testing and comparing different
302 heritability models on large-scale GWAS data.

303 Two practical advantages of SumHer over LDSC are evidenced by the results for height in Table 1 (for these we used summary
304 statistics from the most recent GIANT Consortium meta-analysis,³⁹ which has average sample size 245 000). First we focus on estimates
305 of h_{SNP}^2 . Accurate estimates of h_{SNP}^2 are important not only because they improve our understanding of genetic architecture, but also
306 because they are now being incorporated in software for analyzing complex traits (e.g., the mixed-model association software Bolt-
307 LMM²³ and the prediction software LDPred⁴⁹). Multiple studies have estimated that common SNPs explain at least 40% of the variation
308 in height;^{6,9,17,39,50} this indicates that the SumHer-GC estimate of h_{SNP}^2 (0.38, SD 0.03) is closer to the truth than the LDSC estimate
309 (0.16, SD 0.01). Complementary to estimates of h_{SNP}^2 are estimates of confounding bias. These are used both to assess the quality of a
310 GWAS and when correcting test statistics. LDSC estimates the intercept $(1 + A)$ to be 1.73 (SD 0.05), implying that there is substantial
311 bias and that strong correction is required; were we to divide test statistics by 1.73, the number of significant loci drops by over half.

312 By contrast, SumHer-GC estimates the scaling factor (C) to be 1.02 (SD 0.03), indicating that confounding is slight and that almost no
313 correction of test statistics is needed (dividing the test statistics by 1.02 reduces the number of significant loci by only 3%).

314 The most striking differences between LDSC and SumHer are observed when estimating heritability enrichments, as highlighted
315 in Figure 5a. Whereas analyses using LDSC have found that heritability is highly focused in specific genomic regions, SumHer instead
316 shows that heritability is spread far more diffusely across the genome, supporting an omnigenic view of genetic architecture.⁸ While the
317 realization that complex traits are even more complicated than previously thought is daunting, as our prediction example demonstrates,
318 it is only by properly understanding their complexity that we can develop more efficient tools for analyzing genetic data.

319 URLs

320 LDAK, <http://www.ldak.org/>; LDSC, <http://www.github.com/bulik/ldsc>; 24 functional annotations, [https://](https://data.broadinstitute.org/alkesgroup/LDSCORE/baseline_bedfiles.tgz)
321 data.broadinstitute.org/alkesgroup/LDSCORE/baseline_bedfiles.tgz.

322 Methods

323 Methods, including statements of data availability and any associated accession codes and references, are available in the online version
324 of the paper.

325 Acknowledgments

326 We thank Alkes Price, Hilary Finucane, Paul O'Reilly and Maria Speed for helpful discussions. Access to Wellcome Trust Case Control
327 Consortium data was authorized as work related to the project "Genome-wide association study of susceptibility and clinical phenotypes
328 in epilepsy," access to eMERGE Network data was granted under dbGaP Project 14422, "Comprehensive testing of SNP-based predic-
329 tion models," while access to the Health and Retirement Study was granted under dbGaP Project 15139, "Developing summary-statistic
330 tools for analysing genetic association study data." D.S. is funded by the UK Medical Research Council under grant MR/L012561/1, by
331 the European Unions Horizon 2020 Research and Innovation Programme under the Marie Skłodowska-Curie grant agreement number
332 754513, and by Aarhus University Research Foundation (AUFF). The eMERGE Network was initiated and funded by NHGRI through
333 the following grants: U01HG006828 (Cincinnati Childrens Hospital Medical Center/Boston Childrens Hospital); U01HG006830 (Chil-
334 drens Hospital of Philadelphia); U01HG006389 (Essentia Institute of Rural Health, Marshfield Clinic Research Foundation and Pennsyl-
335 vania State University); U01HG006382 (Geisinger Clinic); U01HG006375 (Group Health Cooperative); U01HG006379 (Mayo Clinic);
336 U01HG006380 (Icahn School of Medicine at Mount Sinai); U01HG006388 (Northwestern University); U01HG006378 (Vanderbilt Uni-
337 versity Medical Center); and U01HG006385 (Vanderbilt University Medical Center serving as the Coordinating Center). The Health
338 and Retirement Study genetic data is sponsored by the National Institute on Aging (grant numbers U01AG009740, RC2AG036495,
339 and RC4AG039029) and was conducted by the University of Michigan. Analyses were performed with the use of the UCL Computer
340 Science Cluster and the help of the CS Technical Support Group, as well as the use of the UCL Legion High-Performance Computing
341 Facility (Legion@UCL) and associated support services.

342 Author contributions

343 D.S. performed the analysis, D.S. and D.J.B. wrote the manuscript.

344 Competing financial interests

345 The authors declare no competing financial interests.

- 347 1. Bulik-Sullivan, B. *et al.* LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat.*
348 *Genet.* **47**, 291–295 (2014).
- 349 2. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* **47**, 1236–1241 (2015).
- 350 3. Finucane, H. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.*
351 **47**, 1228–1235 (2015).
- 352 4. Zheng, J. *et al.* LD Hub: a centralized database and web interface to perform ld score regression that maximizes the potential of
353 summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272–279 (2016).
- 354 5. Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.* **49**, 986–992 (2017).
- 355 6. Speed, D., Hemani, G., Johnson, M. & Balding, D. Improved heritability estimation from genome-wide SNP data. *Am. J. Hum.*
356 *Genet.* **91**, 1011–1021 (2012).
- 357 7. Pasaniuc, B. & Price, A. Dissecting the genetics of complex traits using summary association statistics. *Nat. Rev. Genet.* **18**,
358 117–127 (2017).
- 359 8. Boyle, E., Li, Y. & Pritchard, J. An expanded view of complex traits: From polygenic to omnigenic. *Cell* **169**, 1177–1186 (2017).
- 360 9. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* **42**, 565–569 (2010).
- 361 10. Cross-Disorder Group of the Psychiatric Genomics Consortium. Genetic relationship between five psychiatric disorders estimated
362 from genome-wide SNPs. *Nat. Genet.* **45**, 984–994 (2013).
- 363 11. Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997–1004 (1999).
- 364 12. Juster, F. & Suzman, R. An overview of the Health and Retirement Study. *J. Hum. Resources* **30**, S7–S56 (1995).
- 365 13. Schoech, A. *et al.* Quantification of frequency-dependent genetic architectures and action of negative selection in 25 UK Biobank
366 traits (2017). Preprint available on BioRxiv.
- 367 14. The 1000 Genomes Project Consortium. A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–
368 1073 (2010).
- 369 15. The Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and
370 3,000 shared controls. *Nature* **447**, 661–678 (2007).
- 371 16. Verma, S. *et al.* Imputation and quality control steps for combining multiple genome-wide datasets. *Front. Genet.* **5**, 370 (2015).
- 372 17. Yang, J. *et al.* Genomic partitioning of genetic variation for complex traits using common SNPs. *Nat. Genet.* **43**, 519–525 (2011).
- 373 18. Speed, D. *et al.* Describing the genetic architecture of epilepsy through heritability analysis. *Brain* **137**, 26802689 (2014).
- 374 19. Yang, J. *et al.* Genomic inflation factors under polygenic inheritance. *Eur. J. Hum. Genet.* **19**, 807–812 (2011).
- 375 20. Leslie, S. *et al.* The fine-scale genetic structure of the British population. *Nature* **519**, 309–314 (2015).
- 376 21. Lippert, C. *et al.* FaST linear mixed models for genome-wide association studies. *Nat. Methods* **8**, 833–835 (2011).
- 377 22. Yang, J., Zaitlen, N., Goddard, M., Visscher, P. & Price, A. Advantages and pitfalls in the application of mixed-model association
378 methods. *Nat. Genet.* **46**, 100–106 (2014).
- 379 23. Loh, P. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* **47**, 284–290 (2015).

- 380 24. Yu, J. *et al.* A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat. Genet.* **38**,
381 203–208 (2006).
- 382 25. The International Multiple Sclerosis Genetics Consortium *et al.* Genetic risk and a primary role for cell-mediated immune mecha-
383 nisms in multiple sclerosis. *Nature* **476**, 214–219 (2011).
- 384 26. Lambert, J. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer’s disease. *Nat. Genet.* **45**,
385 1452–1458 (2013).
- 386 27. Schunkert, H. *et al.* Large-scale association analysis identifies 13 new susceptibility loci for coronary artery disease. *Nat. Genet.*
387 **43**, 333–338 (2011).
- 388 28. Liu, J. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk
389 across populations. *Nat. Genet.* **47**, 979–986 (2015).
- 390 29. Okbay, A. *et al.* Genetic variants associated with subjective well-being, depressive symptoms, and neuroticism identified through
391 genome-wide analyses. *Nat. Genet.* **48**, 626–633 (2016).
- 392 30. The Tobacco and Genetics Consortium. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nat.*
393 *Genet.* **42**, 441–447 (2010).
- 394 31. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* **506**, 376–381 (2014).
- 395 32. Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights from 108 schizophrenia-associated
396 genetic loci. *Nature* **511**, 421–427 (2014).
- 397 33. Scott, R. *et al.* An expanded genome-wide association study of type 2 diabetes in Europeans. *Diabetes* **66**, 2888–2902 (2017).
- 398 34. Zheng, H. *et al.* Whole-genome sequencing identifies EN1 as a determinant of bone density and fracture. *Nature* **526**, 112–117
399 (2015).
- 400 35. Locke, A. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197–206 (2015).
- 401 36. Manning, A. *et al.* A genome-wide approach accounting for body mass index identifies genetic variants influencing fasting glyce-
402 mic traits and insulin resistance. *Nat. Genet.* **44**, 659–669 (2012).
- 403 37. Soranzo, N. *et al.* Common variants at 10 genomic loci influence hemoglobin a(c) levels via glyce-
404 mic and nonglyce-
405 mic pathway. *Diabetes* **59**, 3229–3239 (2010).
- 405 38. Global Lipids Genetics Consortium. Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* **45**, 1274–1283
406 (2013).
- 407 39. Wood, A. *et al.* Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat. Genet.*
408 **46**, 1173–1186 (2014).
- 409 40. Perry, J. *et al.* Parent-of-origin-specific allelic associations among 106 genomic loci for age at menarche. *Nature* **514**, 92–97 (2014).
- 410 41. Day, F. *et al.* Large-scale genomic analyses link reproductive aging to hypothalamic signaling, breast cancer susceptibility and
411 brca1-mediated dna repair. *Nat. Genet.* **47**, 1294–1303 (2015).
- 412 42. Shungin, D. *et al.* New genetic loci link adipose and insulin biology to body fat distribution. *Nat. Genet.* **518**, 187–196 (2015).
- 413 43. Okbay, A. *et al.* Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* **533**, 539–542
414 (2016).
- 415 44. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* **478**, 476–482 (2011).
- 416 45. Ward, L. & Kellis, M. Evidence of abundant purifying selection in humans for recently acquired regulatory functions. *Science* **337**,
417 1675–1678 (2012).

- 418 46. Hoffman, M. *et al.* Integrative annotation of chromatin elements from encode data. *Nucleic Acids Res.* **41**, 827–841 (2013).
- 419 47. The International HapMap 3 Consortium. Integrating common and rare genetic variation in diverse human populations. *Nature* **467**,
420 52–58 (2010).
- 421 48. Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits shows action of negative selection. *Nat.*
422 *Genet.* **49** (2017).
- 423 49. Vilhjálmsson, B. *et al.* Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *Am. J. Hum. Genet.* **97**, 576–592
424 (2015).
- 425 50. Zhou, X., Carbonetto, P. & Stephens, M. Polygenic modeling with Bayesian sparse linear mixed models. *PLoS Genet.* **9**, e1003264
426 (2013).
- 427 51. Dempster, E. & Lerner, I. Heritability of threshold characters. *Genetics* **35**, 212–236 (1950).
- 428 52. Lee, S., Wray, N., Goddard, M. & Visscher, P. Estimating missing heritability for disease from genome-wide association studies.
429 *Am. J. Hum. Genet.* **88**, 294–305 (2011).
- 430 53. Wakefield, J. Bayes factors for genome-wide association studies: comparison with p -values. *Genet. Epidemiol.* **33**, 79–86 (2009).
- 431 54. Euesden, J., Lewis, C. & O’Reilly, P. PRSice: polygenic risk score software. *Bioinformatics* **31**, 1466–1468 (2015).
- 432 55. Yang, J., Lee, S., Goddard, M. & Visscher, P. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82
433 (2011).
- 434 56. Delaneau, O., Zagury, J. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat.*
435 *Methods* **10**, 5–6 (2013).
- 436 57. Howie, B., Marchini, J. & Stephens, M. Genotype imputation with thousands of genomes. *G3* **1**, 457–470 (2011).

437 **Online Methods**

438 The Supplementary Note provides step-by-step code for using SumHer to estimate SNP heritability, confounding bias, heritability
439 enrichments and genetic correlations from summary statistics.

440 **Estimating SNP heritability.** Suppose that we have summary statistics from a GWAS on n individuals and m SNPs; let S_j denote the
441 $\chi^2(1)$ test statistic from regressing the phenotype on X_j , the vector of additively-coded genotypes for SNP j , and let $n_j \leq n$ denote the
442 number of individuals used in this regression (note that in the main text, for simplicity, we assumed $n_j = n$). If S_j was obtained using
443 classical (i.e., least-squares) linear regression, then¹⁹

$$\mathbb{E}[S_j] \approx 1 + n_j v_j^2 \quad \text{with} \quad v_j^2 = h_j^2 + \sum_{l \neq j} \text{Cor}(X_j, X_l)^2 h_l^2, \quad (2)$$

where v_j^2 is the total amount of heritability tagged by SNP j . In the main text, we referred to h_j^2 as the heritability “directly contributed”
by SNP j , to emphasize that while a causal variant can contribute to multiple v_j^2 (i.e., be tagged by multiple SNPs), it can only contribute
to one h_j^2 . More formally, the h_j^2 represent a partitioning of h_{SNP}^2 , the total heritability tagged by the m SNPs genotyped by the GWAS;
this formal definition appreciates that a causal variant need not be typed to contribute towards h_{SNP}^2 , provided it is tagged by one or more
SNPs that have been typed (in which case its heritability will be shared across the h_j^2 of the tagging SNPs, even though none of these
“directly contribute” this heritability). If there is no population structure or cryptic relatedness, then $\text{Cor}(X_j, X_l)^2$ will be negligible for
distant SNPs, while for local SNPs, an unbiased estimate of $\text{Cor}(X_j, X_l)^2$ is ¹

$$r_{jl}^2 = \text{Cor}(X_j', X_l')^2 - \frac{1 - \text{Cor}(X_j', X_l')^2}{n' - 2},$$

444 where X'_j is the vector of SNP j genotypes for the n' individuals in a reference panel (for accurate estimates of r_{jl} , these individuals
445 should have similar ancestry to those used in the GWAS). Therefore, in place of Equation (2) we use

$$\mathbb{E}[S_j] \approx 1 + n_j v_j^2 \quad \text{with} \quad v_j^2 = h_j^2 + \sum_{l \in N_j} r_{jl}^2 h_l^2, \quad (3)$$

446 where the set N_j indexes those SNPs “near” SNP j ; a working definition of near is two SNPs within 1 cM (Supplementary Figure 2).
447 Finally, given a heritability model defined as $\mathbb{E}[h_j^2] \propto q_j$, we replace h_j^2 by its expected value $q_j h_{\text{SNP}}^2 / Q$, where $h_{\text{SNP}}^2 = \sum h_j^2$ and
448 $Q = \sum_j q_j$, resulting in

$$\mathbb{E}[S_j] \approx 1 + u_j h_{\text{SNP}}^2 \quad \text{with} \quad u_j = (q_j + \sum_{l \in N_j} q_l r_{jl}^2) n_j / Q, \quad (4)$$

449 which allows us to estimate h_{SNP}^2 by regressing $(S_j - 1)$ on u_j . To account for correlated datapoints and heteroscedasticity we use
450 weighted least squares regression. Specifically, if D is a diagonal matrix whose non-zero entries are the regression weights, then the
451 estimate of h_{SNP}^2 would be $(u^T D u)^{-1} u^T D (S - 1)$, where u and S are vectors containing the m values for u_j and S_j , respectively.
452 Following LDSC,¹ we use $1/D_{jj} = (\sum_{l \in N_j} r_{jl}^2)(1 + u_j h_{\text{SNP}}^2)$. Since h_{SNP}^2 is unknown, we proceed iteratively starting at $h_{\text{SNP}}^2 = 0$,
453 then successively updating D_{jj} and h_{SNP}^2 until convergence. Again following LDSC, we estimate standard errors via block jackknifing
454 (by default we use 200 blocks).

455 **Comparing heritability models.** The (weighted) log likelihood is

$$L(S_j | h_{\text{SNP}}^2, D) = -\frac{d}{2} (\log(2\pi\gamma/d) + 1) \quad \text{with} \quad d = \sum_j D_{jj} \quad \text{and} \quad \gamma = \sum_j D_{jj} (S_j - 1 - u_j h_{\text{SNP}}^2)^2. \quad (5)$$

456 To evaluate $L(S_j | h_{\text{SNP}}^2, D)$ requires values for h_{SNP}^2 and D . While the natural choice is to use the values after the final iteration, in order
457 to compare different heritability models based on $L(S_j | h_{\text{SNP}}^2, D)$, we must use the same D for each (else models leading to lower D_{jj}
458 would have an unfair advantage). Therefore, SumHer reports $L(S_j | \widehat{h_{\text{SNP}}^2}, D^0)$, where $\widehat{h_{\text{SNP}}^2}$ is the final estimate of h_{SNP}^2 , but D^0 is the
459 initial weight matrix (obtained by setting $1/D_{jj} = \sum_{l \in N_j} r_{jl}^2$); Supplementary Table 6 shows that for the 25 raw GWAS, comparisons
460 of heritability models based on $L(S_j | \widehat{h_{\text{SNP}}^2}, D^0)$ align closely with those based on REML likelihood.

461 **Estimating confounding bias.** We recommend replacing Equation (3) with $\mathbb{E}[S_j] \approx C(1 + u_j h_{\text{SNP}}^2)$, where C denotes the multiplicative
462 inflation of test statistics due to confounding; we can then estimate C and $C h_{\text{SNP}}^2$ by jointly regressing S_j on 1 and u_j . To instead copy
463 LDSC, which considers the additive inflation of test statistics, we replace Equation (4) with $\mathbb{E}[S_j] \approx 1 + A + u_j h_{\text{SNP}}^2$, then estimate A
464 and h_{SNP}^2 by jointly regressing $(S_j - 1)$ on 1 and u_j .

465 **Estimating enrichments.** Suppose we have K categories; let $I_{jk} \in \{0, 1\}$ indicate whether SNP j belongs to Category k . We wish to
466 estimate $h_{\text{Cat } k}^2 = \sum_j I_{jk} h_j^2$, the heritability contributed by SNPs in Category k . We now use the heritability model

$$\mathbb{E}[h_j^2] = q_j \sum_k I_{jk} h_{\text{Part } k}^2 / Q_k \quad \text{with} \quad Q_k = \sum_j q_j I_{jk}, \quad (6)$$

467 where $q_j h_{\text{Part } k}^2 / Q_k$ indicates the contribution of Category k to the expected heritability of SNP j , and $h_{\text{SNP}}^2 = h_{\text{Part } 1}^2 + \dots + h_{\text{Part } K}^2$.
468 We consider two scenarios: in Scenario 1, the K categories partition the genome ($\sum_k I_{jk} = 1$); in Scenario 2, the first $K - 1$ categories
469 correspond to annotations, and the K th is the base category containing all SNPs ($I_{jK} = 1$). Using Equation (6) ensures that when there
470 is no enrichment, $\mathbb{E}[h_j^2] = q_j h_{\text{SNP}}^2 / Q$. To appreciate why, consider that for Scenario 1, $h_{\text{Part } k}^2 = h_{\text{Cat } k}^2$, so that when no categories are
471 enriched, $h_{\text{Part } k}^2 = Q_k h_{\text{SNP}}^2 / Q$; for Scenario 2, $h_{\text{Part } 1}^2, \dots, h_{\text{Part } K-1}^2$ indicate how much each annotation increases the SNP heritability
472 from $h_{\text{Part } K}^2$, so that when there is no enrichment, $h_{\text{Part } 1}^2 = \dots = h_{\text{Part } K-1}^2 = 0$ and $h_{\text{Part } K}^2 = h_{\text{SNP}}^2$. Now when we replace h_j^2 in
473 Equation (3) by its expected value, we obtain

$$\mathbb{E}[S_j] \approx 1 + \sum_k u_{jk} h_{\text{Part } k}^2 \quad \text{with} \quad u_{jk} = (q_j I_{jk} + \sum_{l \in N_j} q_l I_{lk} r_{jl}^2) n_j / Q_k, \quad (7)$$

474 and therefore we can estimate $h_{\text{Part } 1}^2, \dots, h_{\text{Part } K}^2$ by jointly regressing $(S_j - 1)$ on $u_{j,1}, \dots, u_{j,K}$. Given these, our estimate of $h_{\text{Cat } k}^2$ is
475 $\sum_{k'} (\sum_j I_{jk} q_j I_{jk'}) h_{\text{Part } k'}^2 / Q_{k'}$, which we then divide by $Q_k h_{\text{SNP}}^2 / Q$ to get an estimate of the enrichment of Category k .

476 **Estimating genetic correlation.** Suppose we have summary statistics from two GWAS. Instead of $\chi^2(1)$ test statistics, we now use
 477 (signed) Z -statistics. Let Z_{Aj} and Z_{Bj} denote the two Z -statistics for SNP j , computed using n_{Aj} and n_{Bj} individuals, respectively, of
 478 which n_{Cj} were common to both GWAS (if the two GWAS were independent, $n_{Cj} = 0$). We assume

$$\mathbb{E}[Z_{Aj}Z_{Bj}] \approx \frac{c_{AB}n_{Sj}}{\sqrt{n_{Aj}n_{Bj}}} + u'_j h_{AB}^2 \quad \text{with} \quad u'_j = (q_j + \sum_{l \in N_j} q_l r_{jl}^2) \sqrt{n_{Aj}n_{Bj}}/Q, \quad (8)$$

479 where c_{AB} is the phenotypic correlation between the two traits and h_{AB}^2 is their genetic covariance. This equation matches that used by
 480 LDSC,² except we have replaced r_{jl}^2/m by $q_j r_{jl}^2/Q$. By regressing $Z_{Aj}Z_{Bj}$ on u'_j , we obtain an estimate of h_{AB}^2 , which we then divide
 481 by estimates of $\sqrt{h_{\text{SNP}}^2}$ for each trait to produce an estimate of their genetic correlation.

482 **Regions of extreme LD and large-effect SNPs.** Estimates of SNP heritability can be unduly affected by regions of extreme LD and
 483 by SNPs with disproportionately large effect size.^{4,5} Therefore, for all analyses, we exclude SNPs within the major histocompatibility
 484 complex (Chromosome 6: 25-34 Mb), as well as SNPs which individually explain $>1\%$ of phenotypic variation ($S_j > n_j/99$), and SNPs
 485 in LD with these (within 1 cM and $r_{jl}^2 > 0.1$); the latter resulted in the exclusion of SNPs for 8 of the 25 raw traits and for 6 of the 24
 486 summary traits (Supplementary Table 15).

487 **Binary phenotypes.** So far, we have implicitly assumed the trait is quantitative. If instead it is binary (i.e., for case-control GWAS),
 488 there are two considerations. Firstly, heritability estimates now correspond to the ‘‘observed scale,’’ SumHer will convert them to the
 489 ‘‘liability scale’’ if provided with the prevalence and ascertainment.^{51,52} Secondly, it is likely the p -values came from (classical) logistic
 490 regression instead of linear regression; however, Supplementary Figure 10 shows that, because linear and logistic p -values closely match
 491 for SNPs with small or moderate effect, this tends to have limited impact.

492 **Polygenic risk scoring.** To predict phenotypes, we construct PRS of the form $\sum_j \beta_j X'_j$, where the vector X'_j contains standardized
 493 genotypes for SNP j (obtained by centering X_j , then scaling to have variance 1). Without loss of generality, we assume phenotypes
 494 have also been standardized, in which case the estimate of β_j from classical linear regression is ρ_j , the correlation observed between
 495 SNP j and the phenotype, and has variance $(1 - \rho_j^2)/n$. Note that given summary statistics, we can recover ρ_j by appreciating that the
 496 Wald test Z -statistic is $\rho_j/\sqrt{(1 - \rho_j^2)/n}$. For the Classical PRS, our estimate of β_j is ρ_j . For the Bayesian PRS, we must specify a prior
 497 distribution for β_j . Given a heritability model, we use $\beta_j \sim \mathbb{N}(0, \sigma_j^2)$, where $\sigma_j^2 = (q_j + \sum_{l \in N_j} q_l r_{jl}^2) h_{\text{SNP}}^2/Q$; this prior is motivated
 498 by recognizing that for the ‘‘true PRS’’, $\beta_j^2 = v_j^2$, the heritability tagged by SNP j . We then estimate β_j by its posterior mean, a shrunken
 499 version of the classical estimate. To calculate this, we approximate the likelihood distribution by $\mathbb{N}(\rho_j, (1 - \rho_j^2)/n)$,⁵³ then the posterior
 500 mean equals $\rho_j \sigma_j^2 / (\sigma_j^2 + (1 - \rho_j^2)/n_j)$. Two technicalities. Firstly, to construct each Bayesian PRS requires a value for h_{SNP}^2 , so we
 501 used the corresponding estimate from LDSC-Zero or SumHer-Zero (for Enriched GCTA we used the 53-part model, for Enriched LDAK
 502 we used the 25-part model). Supplementary Table 14 confirms that the ranking of PRS remains the same if we instead agnostically set
 503 $h_{\text{SNP}}^2 = 0.5$. Secondly, for the Bayesian PRS, we performed clumping⁵⁴ (thinned SNPs to ensure no pair within 1 cM had $r_{jl}^2 > 0.05$),
 504 whereas for the Classical PRS we did not; this was because we found clumping benefited all Bayesian PRS, but was detrimental to the
 505 Classical PRS. Supplementary Table 14 shows that the ranking of Bayesian PRS remains the same if we do not clump, but then they are
 506 inferior to the Classical PRS.

507 **Choosing a heritability model.** Although the heritability model is defined in terms of h_j^2 , as is clear from Equation (3), its fit depends
 508 only on how well it models $v_j^2 = h_j^2 + \sum_{l \in N_j} r_{jl}^2 h_l^2$. Therefore, when specifying q_j , the focus should be on accurately describing how
 509 v_j^2 is expected to vary across the genome. LDSC assumes $\mathbb{E}[h_j^2]$ is constant ($q_j = 1$), and therefore that v_j^2 correlates with local levels
 510 of LD; we refer to this as the GCTA Model⁵ because the same assumption is made by GCTA, software for estimating SNP heritability
 511 via REML.⁵⁵ The GCTA Model is widely used in statistical genetics. For example, it is implicitly assumed by any regression method
 512 where SNPs are standardized and the same penalty function / prior distribution is applied to each, and likewise by any simulation study
 513 where causal SNPs are picked at random and their standardized effect sizes are sampled from a common distribution. However, as we
 514 have shown, the GCTA Model poorly reflects real data.⁵ Even for traits where a correlation is observed between v_j^2 and S_j (those with
 515 $p < 1$ in Figures 3b & 3c), the correlation is substantially weaker than predicted by the GCTA Model (hence $p \gg 0$), and as we
 516 discussed above, might be a consequence of including lower-certainty SNPs in the GWAS. We currently recommend using the LDAK

517 Model. Originally, this took the form $q_j = w_j$, where the weights w_j are calculated based on the assumption that $\mathbb{E}[v_j^2]$ is constant.⁶ We
518 recently updated it to $q_j = w_j [f_j(1 - f_j)]^{0.75} r_j$, reflecting that v_j^2 tends to depend on MAF, f_j , and genotype certainty, r_j (the latter
519 is not relevant in this study as either $r_j \approx 1$ or was unknown), and we recognize that further improvements will be possible.

520 **Quality control and association analysis.** To prepare the 13 WTCCC GWAS, we used our previously described protocol.⁵ In summary,
521 after excluding apparent population outliers, samples with extreme missingness or heterozygosity and SNPs with MAF < 0.01 , call rate
522 < 0.95 or Hardy-Weinberg $P < 1 \times 10^{-6}$, we phased using SHAPEIT⁵⁶, then imputed using IMPUTE2⁵⁷ with the 1000 Genomes
523 Project Phase 3 (2014) reference panel.¹⁴ When merging case and control datasets, we converted genotype probabilities to hard calls
524 using a certainty threshold of 0.95, then retained only autosomal SNPs that in all cohorts had MAF ≥ 0.01 and info score ≥ 0.99
525 (using IMPUTE2 `r2_type2` for directly genotyped SNPs). Finally, we thinned individuals, so no pair remained with estimated
526 relatedness > 0.05 . We performed the same quality control for the our reference panel, the Health and Retirement Study.¹² The emerge
527 data were provided post-imputation; this was also performed using SHAPEIT and IMPUTE2, but used the 1000 Genomes Project Phase
528 2 reference panel.¹⁶ We converted genotype probabilities to hard calls using a certainty threshold of 0.95, then retained only biallelic
529 SNPs with MAF ≥ 0.01 , call rate ≥ 0.95 , info score ≥ 0.99 and whose genomic position matched that in the 1000 Genomes Project.
530 Finally, we excluded individuals ancestrally inconsistent ($P < 0.05$) with non-Finnish Europeans from the 1000 Genomes Project (see
531 Supplementary Figure 11) and those whose ethnicity was reported as “Hispanic or Latino”, then filtered until no pair remained with
532 estimated relatedness > 0.05 (which left 25 875 individuals).

533 For the 25 raw GWAS, we performed the association analysis using linear regression (regardless of whether the trait was quan-
534 titative or binary), including as covariates sex and ten principal components (five derived from the reference panel, five from the 1000
535 Genomes Project¹⁴). For the 24 summary GWAS, we used publicly-available summary statistics (Supplementary Table 8). For each SNP,
536 SumHer requires the two alleles, the $\chi^2(1)$ test statistic, n_j and, when estimating genetic correlations, the direction of effect (relative to
537 the first allele). Given summary statistics, we generally used all SNPs in common with our reference panel (and with consistent alleles),
538 however, for the five GWAS which provided info scores, we also excluded SNPs with score < 0.95 . Per-SNP sample sizes were only
539 available for eight of the summary GWAS, so for the remainder, we set $n_j = n$, the total sample size.

540 **Run times.** Given summary statistics and a reference panel, a SumHer analysis has two steps: the first generates a “tagfile”, which
541 contains $q_j + \sum_{l \in N_j} q_l r_{jl}^2$ for each SNP; the second performs the regression. The latter is trivial, and typically finishes within minutes.
542 The time to compute the tagfile depends mainly on the size of the reference panel; using our preferred reference panel (8 850 individuals)
543 takes ~ 1 day on a single CPU, whereas using the non-Finnish Europeans from the 1000 Genomes Project¹⁴ (404 individuals) takes a
544 couple of hours. While we recommend users compute tagfiles from scratch, starting with the subset of SNPs common to both the GWAS
545 and reference panel, we alternatively provide a pre-computed tagfile, which should suffice when the GWAS coverage is high.