

Exposing flaws in S-LDSC; reply to Gazal *et al.*

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.

In our recent publication,¹ we examined the two heritability models most widely used when estimating SNP heritability: the GCTA Model, which is used by the software GCTA² and upon which LD Score regression (LDSC) is based,³ and the LDAK Model, which is used by our software LDAK.⁴ First we demonstrated the importance of choosing an appropriate heritability model, by showing that estimates of SNP heritability can be highly sensitive to which model is assumed. Then we empirically tested the GCTA and LDAK Models on GWAS data for a wide variety of complex traits. We found that the LDAK Model fits real data both significantly and substantially better than the GCTA Model, indicating that LDAK estimates more accurately describe the genetic architecture of complex traits than those from GCTA or LDSC.

Some of our most striking results were our revised estimates of functional enrichments (the heritability enrichments of SNP categories defined by functional annotations). In general, estimates from LDAK were substantially more modest than previous estimates based on the GCTA Model. For example, we estimated that DNase I hypersensitive sites (DHS) were 1.4-fold (SD 0.1) enriched, whereas a study using GCTA had found they were 5.1-fold (SD 0.5) enriched,⁵ and we estimated that conserved SNPs were 1.3-fold (SD 0.3) enriched, whereas a study using S-LDSC (stratified LDSC) had found they were 13.3-fold (SD 1.5) enriched.⁶

In their correspondence, Gazal *et al.* dispute our findings. They assert that the heritability model assumed by LDSC is more realistic than the LDAK Model, and that estimates of enrichment from S-LDSC⁷ are more accurate than those from LDAK. Here, we explain why their justification for preferring the model used by LDSC is incorrect, and provide a simple demonstration that S-LDSC produces unreliable estimates of enrichment.

The GCTA and LDAK Models.

Let h_j^2 denote the heritability contributed by SNP j , defined so that $h_{\text{SNP}}^2 = \sum_j h_j^2$ is the SNP heritability of the trait. The GCTA Model assumes a prior distribution for effect sizes such that each SNP is expected to contribute equal heritability: $\mathbb{E}[h_j^2] \propto 1$.^{1,2} By contrast, the LDAK Model assumes

$$\mathbb{E}[h_j^2] \propto (f_j(1 - f_j))^{0.75} w_j = q_j,$$

where f_j is the minor allele frequency (MAF) of SNP j and w_j is its LDAK weight (SNPs in regions of high LD tend to have lower w_j , and vice versa).^{1,4} If r_{jl}^2 denotes the squared correlation between SNPs j and l , then $v_j^2 = \sum_l r_{jl}^2 h_l^2$ is the total heritability tagged by SNP j (in theory, the summation is across all SNPs, but in practice³ we consider only those within 1 cM). Under the GCTA Model, $\mathbb{E}[v_j^2] = l_j h_{\text{SNP}}^2/m$, where $l_j = \sum_l r_{jl}^2$ is the LD score of SNP j ,³ whereas under the LDAK Model, $\mathbb{E}[v_j^2] = l'_j h_{\text{SNP}}^2 / \sum_j q_j$, where $l'_j = \sum_l r_{jl}^2 q_j$ is the “LDAK score” of SNP j .

LDSC is based on the GCTA Model

Suppose we have a GWAS on n individuals and m SNPs. The $\chi^2(1)$ additive association test statistic for SNP j has value⁸

$$S_j = nc_j^2 = n(v_j^2 + a_j + e_j), \quad (1)$$

where c_j^2 is the phenotypic variance explained by SNP j , which can be partitioned into v_j^2 , a_j and e_j , components corresponding to causal variation, confounding and noise, respectively. e_j has expectation $1/n$; LDSC seeks to estimate the expected values of v_j^2 and a_j . For this it assumes the model³

$$\mathbb{E}[S_j] = nh_{\text{SNP}}^2 l_j/m + na + 1.$$

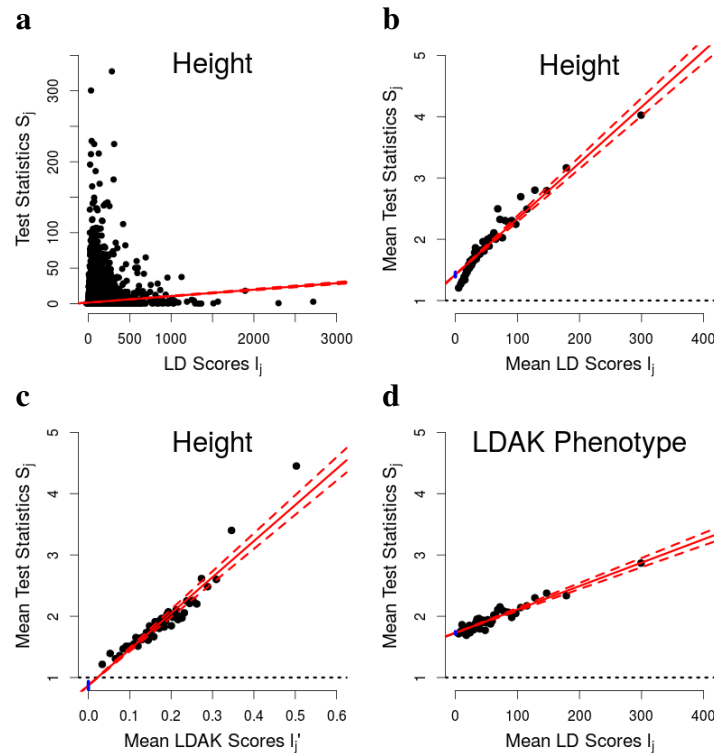


Figure 1: Test statistics are correlated with both LD and LDAK Scores. (a) Test statistics versus LD scores from the most recent Giant Consortium meta-analysis for height;¹² to avoid correlated datapoints, we restrict to a subset of 121 310 SNPs with $MAF > 0.01$ in approximate linkage equilibrium (obtained by pruning so that no two SNPs within 1 cM have $r_{jl}^2 > 0.2$). (b) The correlation can be magnified by first dividing SNPs into 50 bins based on LD Scores, then plotting mean test statistic versus mean LD score for each bin.³ (c) The same as (b), except we consider LDAK scores instead of LD Scores. (d) The same as (b), except that instead of using the test statistics for height, we generate new ones based on the LDAK Model. In each plot, the solid red line is the line of best fit from least-squares regression; the dashed red lines and solid blue segments indicate, respectively, 95% confidence intervals for the slope and intercept from this regression.

30 We can see that this follows from Equation (1) if we assume the GCTA Model (as then v_j^2 has expected value $l_j h_{SNP}^2/m$) and that a_j is
31 constant across the genome.

32 Evidence for the LDAK Model

33 In our previous work,¹ we performed a careful evaluation of the GCTA and LDAK Models. We collected GWAS data for 42 different
34 traits, both binary and quantitative, then performed stringent quality control, checking that any confounding due to population structure
35 or cryptic relatedness was at most slight.^{9,10} We demonstrated that it was valid to compare models using the REML likelihood, then
36 used this approach to show that the LDAK Model was both significantly and substantially more realistic than the GCTA Model; it fit
37 better for 37 of the 42 traits ($P < 10^{-7}$) and resulted in an average increase in log likelihood of 9.8 per trait. We also investigated
38 attempts to improve the accuracy of the GCTA Model by partitioning (we focused on GCTA-LDMS,¹¹ but the same arguments apply to
39 S-LDSC⁷). While partitioning allowed GCTA to achieve log likelihoods comparable to those from LDAK, this came at the cost of 19
40 extra parameters which were arbitrarily defined, added little to model interpretation and reduced the precision of heritability estimates.

41 Evidence for the GCTA Model

42 In their correspondence, Gazal *et al.* make no mention of the evidence we provided in support of the LDAK Model. Instead, their
43 rationale for preferring the GCTA Model is the observation that for many traits *the marginal effect size of a SNP has been shown to have*
44 *a strong linear dependency on its LD score* (in our notation, that there is a significant correlation between l_j and S_j). We do not dispute
45 that these correlations exist; for example, Figures 1a & 1b demonstrate that l_j and S_j are correlated for human height, using data from

46 the most recent Giant Consortium meta-analysis.¹² However, we disagree with the reasoning that because the GCTA Model predicts
47 $v_j^2 \propto l_j$, an observed correlation between l_j and S_j is evidence to prefer the GCTA Model. Firstly, it does not immediately follow from
48 Equation (1) that all correlation between l_j and S_j is driven by correlation between l_j and v_j^2 . While this would be true if $a_j = a$ (or
49 more generally, if a_j is orthogonal to l_j), no empirical evidence was provided to support this assumption.³ Considering that l_j correlates
50 with factors such as MAF, genotyping certainty and population axes (Supplementary Figure 1), it seems plausible that a_j does correlate
51 with l_j . It was to avoid uncertainty regarding a_j , that when comparing the GCTA and LDAK Models,¹ we restricted ourselves to GWAS
52 where we were confident that $a_j \approx 0$.

53 Secondly, a significant correlation between l_j and v_j^2 only proves that the GCTA Model fits better than the model $\mathbb{E}[h_j^2] = 0$,
54 not that it fits better than the LDAK Model. The LDAK Model predicts $v_j^2 \propto l'_j$; while Figure 1c shows that for height there is also
55 significant correlation between l'_j and S_j , it would be equally absurd of us to claim that the LDAK Model was superior to the GCTA
56 Model based on this evidence alone. For Figure 1d, we generate test statistics under the LDAK Model (assuming no confounding);
57 specifically, we sample S_j from a χ_1^2 distribution with non-centrality parameter $5.2 l'_j$ (we chose 5.2 so that the mean test statistic is
58 2.29, matching that observed for height). This simulation demonstrates that l_j and S_j will also be correlated for LDAK phenotypes, on
59 account of the strong correlation between l_j and l'_j (for these data, their correlation is 0.51). Moreover, it highlights the dangers of using
60 (S-)LDSC when the GCTA Model is not appropriate. The model used by LDSC makes strong predictions about how v_j^2 , and therefore
61 S_j , vary across the genome; for example, the 95% range of l_j is 38 to 228, and the 10% (1%) of SNPs with highest l_j are on average
62 expected to tag 2.8 (5.8) times as much heritability as the average SNP. When the data do not align with these predictions, LDSC will
63 compensate by under-estimating h_{SNP}^2 (the slope of the line) and over-estimating a (the intercept).

64 **Demonstrating problems with S-LDSC**

65 The original S-LDSC model contained 53 categories: 28 functional annotations (which include coding, conserved and DHS regions),
66 24 buffers and the base category containing all SNPs.⁶ Recently, this was expanded to 75 categories, by adding 3 more functional
67 annotations, 3 extra buffers, 10 MAF tranches and 6 continuous LD-related annotations.⁷ We now construct an additional category of
68 “thinned SNPs”, by pruning so that no two SNPs within 1 cM have $r_{jl}^2 > 0.2$, and also the corresponding buffer (all thinned SNPs
69 and those within 500 bp). Table 1 and Supplementary Table 1 report average estimates of enrichment for coding, conserved, DHS and
70 thinned SNPs, estimated using six versions of LDSC (which vary according to choice of category), as well as GCTA and LDAK. We use
71 two sources of data: LDSC requires only summary statistics, so we first analyze published results from 24 large-scale GWAS (12 binary
72 traits, 12 quantitative, average sample size 121 000; see Supplementary Table 2); GCTA and LDAK need raw data, so we also perform
73 25 GWAS using data from the Wellcome Trust Case Control Consortium¹³ and the eMerge Network¹⁴ (18 binary traits, 7 quantitative,
74 average sample size 9 700; see Supplementary Table 3).

75 Table 1 highlights two shortcomings with using S-LDSC to estimate enrichments. Firstly, there are many arbitrary choices
76 underlying S-LDSC, such as which functional categories and LD annotations to include, the size and number of buffer regions and
77 how to partition SNPs by MAF; we see that estimates from S-LDSC vary substantially depending on these choices. Secondly, both
78 old and new S-LDSC find thinned SNPs to be highly enriched for heritability: 20.3-fold (SD 0.4) and 14.5-fold (SD 0.5), respectively.
79 Considering that we selected thinned SNPs simply by pruning, and not based on biological criteria, we see no reason why they should
80 be many-fold enriched for heritability. By contrast, LDAK estimates their enrichment to be 0.86-fold (SD 0.06), indicating that the high
81 estimates from S-LDSC are a consequence of the GCTA Model not accounting for LD.

82 In summary, Gazal *et al.* have argued that the heritability model used by LDSC better reflects real data than the LDAK Model,
83 and that high estimates of functional enrichment from S-LDSC should be preferred to those from LDAK. We do not agree with their first
84 claim; whereas we provided rigorous evidence to support the LDAK Model,¹ Gazal *et al.* rely on the observation that for many traits,
85 association test statistics correlate with LD scores, something we have shown is also to be expected under the LDAK Model. Nor do we
86 agree with their second claim; we have shown that S-LDSC estimates can be highly sensitive to category choice, without it being clear
87 which choice to prefer, and that simply by thinning SNPs, we can construct a category which S-LDSC finds to be over ten-fold enriched
88 for heritability.

Summary Statistics	2-Part LDSC	10.4 (0.5)	18.1 (0.5)	5.3 (0.1)	24.2 (0.4)
	3-Part LDSC	7.5 (0.5)	10.6 (0.6)	3.2 (0.1)	24.6 (0.4)
	Old S-LDSC	6.2 (0.5)	12.0 (0.5)	1.7 (0.2)	
	Old S-LDSC+	4.6 (0.3)	7.9 (0.3)	1.5 (0.1)	20.3 (0.4)
	New S-LDSC	4.5 (0.4)	7.6 (0.4)	1.4 (0.1)	
	New S-LDSC+	4.0 (0.3)	6.3 (0.3)	1.4 (0.1)	14.5 (0.5)
Raw Data	2-Part LDSC	18.3 (1.7)	18.3 (1.4)	8.2 (0.2)	28.7 (0.8)
	2-GRM GCTA	15.3 (1.5)	15.8 (1.3)	7.6 (0.2)	22.3 (0.9)
	2-GRM LDAK	2.9 (0.4)	1.9 (0.3)	1.3 (0.1)	0.9 (0.1)

Table 1: Enrichment of coding, conserved, DHS and thinned SNPs. For each of the four annotations, values report average estimates of enrichment based on either summary statistics from 24 published GWAS, or analysis of 25 GWAS for which we have raw genotype and phenotype data. We use six versions of LDSC: 2-part (the annotation SNPs and the base category containing all SNPs); 3-part (the annotation SNPs, the corresponding 500 bp buffer and the base category); old S-LDSC (53 categories, including coding, conserved and DHS SNPs); old S-LDSC+ (the 53 categories, plus thinned SNPs and the corresponding buffer); new S-LDSC (75 categories); new S-LDSC+ (75 categories, plus thinned SNPs and the corresponding buffer). We also estimate enrichments using GCTA and LDAK, each time constructing two genomic similarity matrices (GSMs), the first corresponding to the annotation SNPs, the second to all other SNPs.

89 URLs

90 LDAK, <http://ldak.org>; GCTA, <http://cnsgenomics.com/software/gcta>; LDSC, <http://github.com/bulik/>
91 ldsc.

92 Methods

93 Full details for repeating our analyses are provided in the Supplementary Note.

94 Acknowledgments

95 Access to Wellcome Trust Case Control Consortium data was authorized as work related to the project “Genome-wide association study
96 of susceptibility and clinical phenotypes in epilepsy,” access to eMerge Network data was granted under dbGaP Project 14422, “Compre-
97 hensive testing of SNP-based prediction models.” D.S. is funded by the UK Medical Research Council under grant MR/L012561/1, by
98 the European Unions Horizon 2020 Research and Innovation Programme under the Marie Skłodowska-Curie grant agreement number
99 754513, and by Aarhus University Research Foundation (AUFF). The eMERGE Network was initiated and funded by NHGRI through
100 the following grants: U01HG006828 (Cincinnati Childrens Hospital Medical Center/Boston Childrens Hospital); U01HG006830 (Chil-
101 drens Hospital of Philadelphia); U01HG006389 (Essentia Institute of Rural Health, Marshfield Clinic Research Foundation and Pennsyl-
102 vania State University); U01HG006382 (Geisinger Clinic); U01HG006375 (Group Health Cooperative); U01HG006379 (Mayo Clinic);
103 U01HG006380 (Icahn School of Medicine at Mount Sinai); U01HG006388 (Northwestern University); U01HG006378 (Vanderbilt
104 University Medical Center); and U01HG006385 (Vanderbilt University Medical Center serving as the Coordinating Center).

105 Author contributions

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

107 **Competing financial interests**

108 The authors declare no competing financial interests.

109 **Code availability.** Step-by-step code for performing the analyses underlying Figure 1 and Table 1 are provided in the Supplementary
110 Note.

111 **Data availability.** WTCCC¹³ and eMerge Network¹⁴ were applied for and downloaded from the European Genome-phenome Archive
112 and dbGaP, respectively (see Supplementary Table 3 for accession codes). Results for each of the 24 summary GWAS are available to
113 download from the websites of the corresponding study (see Supplementary Table 2 for references).

114 **References**

- 116 1. Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.* **49**, 986–992 (2017).
- 117 2. Yang, J., Lee, S., Goddard, M. & Visscher, P. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82
118 (2011).
- 119 3. Bulik-Sullivan, B. *et al.* LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat.*
120 *Genet.* **47**, 291–295 (2014).
- 121 4. Speed, D., Hemani, G., Johnson, M. & Balding, D. Improved heritability estimation from genome-wide SNP data. *Am. J. Hum.*
122 *Genet.* **91**, 1011–1021 (2012).
- 123 5. Gusev, A. *et al.* Partitioning Heritability of Regulatory and Cell-Type-Specific Variants across 11 Common Diseases. *Am. J. Hum.*
124 *Genet.* **95**, 535–552 (2014).
- 125 6. Finucane, H. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.*
126 **47**, 1228–1235 (2015).
- 127 7. Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits shows action of negative selection. *Nat.*
128 *Genet.* **49** (2017).
- 129 8. Yang, J. *et al.* Genomic inflation factors under polygenic inheritance. *Eur. J. Hum. Genet.* **19**, 807–812 (2011).
- 130 9. Yang, J. *et al.* Genomic partitioning of genetic variation for complex traits using common SNPs. *Nat. Genet.* **43**, 519–525 (2011).
- 131 10. Speed, D. *et al.* Describing the genetic architecture of epilepsy through heritability analysis. *Brain* **137**, 26802689 (2014).
- 132 11. Yang, J. *et al.* Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body
133 mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
- 134 12. Wood, A. *et al.* Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat. Genet.*
135 **46**, 1173–1186 (2014).
- 136 13. The Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and
137 3,000 shared controls. *Nature* **447**, 661–678 (2007).
- 138 14. Verma, S. *et al.* Imputation and quality control steps for combining multiple genome-wide datasets. *Front. Genet.* **5**, 370 (2015).