

1 **Improved analyses of GWAS summary statistics by reducing data heterogeneity and**  
2 **errors**

3

4 Wenhan Chen<sup>1</sup>, Yang Wu<sup>1</sup>, Zhili Zheng<sup>1</sup>, Ting Qi<sup>1,2</sup>, Peter M Visscher<sup>1</sup>, Zhihong Zhu<sup>1</sup>, Jian Yang<sup>1,\*</sup>

5

6 <sup>1</sup>Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,  
7 Australia

8 <sup>2</sup>Present address: School of Life Sciences, Westlake University, Hangzhou, Zhejiang 310024,  
9 China

10 \*Correspondence: Jian Yang ([jian.yang.qt@gmail.com](mailto:jian.yang.qt@gmail.com))

11

12 **Abstract**

13 Summary statistics from genome-wide association studies (GWAS) have facilitated the  
14 development of various summary data-based methods, which typically require a reference  
15 sample for linkage disequilibrium (LD) estimation. Analyses using these methods may be biased  
16 by errors in GWAS summary data and heterogeneity between GWAS and LD reference. Here we  
17 propose a quality control method, DENTIST, that leverages LD among genetic variants to detect  
18 and eliminate errors in GWAS or LD reference and heterogeneity between the two. Through  
19 simulations, we demonstrate that DENTIST substantially reduces false-positive rate (FPR) in  
20 detecting secondary signals in the summary-data-based conditional and joint (COJO) association  
21 analysis, especially for imputed rare variants (FPR reduced from >28% to <2% in the presence  
22 of ancestral difference between GWAS and LD reference). We further show that DENTIST can  
23 improve other summary-data-based analyses such as LD score regression analysis, and  
24 integrative analysis of GWAS and expression quantitative trait locus data.

25

## 26 **Introduction**

27 Genome-wide association studies (GWASs) have been extraordinarily successful in uncovering  
28 genetic variants associated with complex human traits and diseases<sup>1,2</sup>. Summary statistics  
29 available from GWASs have facilitated the development of various summary-data-based  
30 methods<sup>3</sup> such as those for fine-mapping<sup>4-9</sup>, imputing summary statistics at untyped variants<sup>10,11</sup>,  
31 estimating SNP-based heritability<sup>12-14</sup>, assessing causal or genetic relationship between traits<sup>15-  
32 17</sup>, prioritizing candidate causal genes for a trait<sup>18-21</sup>, and polygenic risk prediction<sup>8,22,23</sup>. Most of  
33 the summary-data-based methods require linkage disequilibrium (LD) structure of the variants  
34 used, which are not available in the summary data but can be estimated from a reference cohort  
35 with individual-level genotypes assuming a homogeneous LD structure between the GWAS and  
36 reference cohorts. Hence, summary-data-based analyses can be affected by not only errors in the  
37 GWAS and LD reference data sets but also differences between them for the following reasons.  
38 First, there are often errors in GWAS summary statistics resulting from the data generation and  
39 analysis processes (e.g., genotyping/imputation errors and genetic variants with mis-specified  
40 effect alleles)<sup>24,25</sup>, some of which are not easy to detect, even if individual-level data are  
41 available. Second, there is often heterogeneity between data sets (e.g., between the discovery  
42 GWAS and LD reference) because of differences in ancestry, and genotyping platform, analysis  
43 pipeline. Although the recommended practice is to use an ancestry-matched reference cohort,  
44 samples with similar ancestries, such as populations of European ancestry, can still have  
45 discernable differences in LD structure<sup>26</sup>, and the effects of such differences on summary-data-  
46 based analyses are largely unexplored. To the best of our knowledge, there is no existing method  
47 specifically designed to detect data heterogeneity that affect summary data-based analyses.

48

49 In this study, we propose a quality control (QC) method to identify errors in GWAS summary  
50 data and heterogeneity between summary data and LD reference by testing the difference  
51 between the observed z-score of each variant and its predicted value from the surrounding  
52 variants. The method has been implemented in a software tool named DENTIST (detecting  
53 errors in analyses of summary statistics). We show by simulation that DENTIST can effectively  
54 detect simulated errors of several kinds. We then demonstrate the utility of DENTIST as a QC  
55 step for multiple, frequently-used, summary data-based methods, including the conditional and  
56 joint analysis (COJO)<sup>6</sup> of summary statistics, LD score regression<sup>12</sup>, and heterogeneity in  
57 dependent instruments (HEIDI) test<sup>21</sup>.

58

## 59 **Results**

### 60 **Overview of the DENTIST method**

61 Details of the methodology can be found in the Methods section. In brief, we first use a sliding  
62 window approach to divide the variants into 2Mb segments with a 500kb overlap between two  
63 adjacent segments. Within each segment, we randomly partition variants into two subsets, S1  
64 and S2, with an equal number of variants, and apply the statistic below to test the difference  
65 between the observed z-score of a variant  $i$  ( $z_i$ ) in S1 and its predicted value ( $\tilde{z}_i$ ) based on z-  
66 scores of an array of variants  $\mathbf{t}$  in S2 (**Methods**).

$$67 \quad T_{d(i)} = \frac{(z_i - \tilde{z}_i)^2}{1 - \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{R}'_{it}} \text{ with } \tilde{z}_i = \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t \quad (1)$$

68 where  $\mathbf{z}_t$  is a vector of z-scores of variants  $\mathbf{t}$  in S2, and  $\mathbf{R}$  is the LD correlation matrix calculated  
69 from a reference sample with  $\mathbf{R}_{tt}$  to denote the LD between variants  $\mathbf{t}$  and  $\mathbf{R}_{it}$  to denote the LD  
70 between variant  $i$  and variants  $\mathbf{t}$ .  $T_d$  follows approximately a  $\chi^2$  distribution with 1 degree of  
71 freedom. Note that methods that leverage LD to predict GWAS test-statistic of a variant (i.e.,  $\tilde{z}_i$ )  
72 from test-statistics of its adjacent variants (i.e.,  $\mathbf{z}_t$ ) have been developed in prior work<sup>10,11</sup>. A  
73 significant difference between the observed and predicted z-scores indicates errors in the  
74 discovery GWAS or LD reference, or heterogeneity between them. If the difference between  $z_i$   
75 and  $\tilde{z}_i$  is due to error in  $z_i$ , the power of  $T_d$  depends on how  $z_i$  deviates from its true value and  
76 how well variant  $i$  is tagged by variants  $\mathbf{t}$ . We conduct a truncated singular value decomposition  
77 (SVD) on  $\mathbf{R}_{tt}$  to mitigate the sampling noise in LD estimated from the reference that is often  
78 independent from the discovery GWAS and to perform a pseudo inverse when  $\mathbf{R}_{tt}$  is singular<sup>13</sup>  
79 (see **Methods**).

80  
81 One challenge for the DENTIST method is that errors can be present in both S1 and S2, and  
82 errors in S2 can inflate  $T_d$  statistics of the variants in S1. To mitigate this issue, we propose an  
83 iterative partitioning approach. In each iteration, we partition the variants at random into two  
84 subsets (S1 and S2) and remove variants with  $P_{DENTIST} < 5 \times 10^{-8}$  (capped at 0.5% variants with the  
85 smallest  $P$ -values). This step is to create a more reliable set of variants for the next iteration. The  
86 problematic variants are prioritized and filtered out in the first few iterations so that the  
87 prediction of  $\tilde{z}_i$  (Equation 1) becomes more accurate in the following iterations. We set the  
88 number of iterations to 10 in practice. All variants with  $P_{DENTIST} < 5 \times 10^{-8}$  are removed in the final  
89 step.

90

### 91 **Detecting simulated errors in GWAS data**

92 To assess the performance of DENTIST in detecting errors, we simulated GWAS data with  
93 genotyping errors and allelic errors (i.e., variants with the effect allele mis-labelled) using whole  
94 genome sequence (WGS) data of chromosome 22 on 3,642 unrelated individuals from the  
95 UK10K project<sup>27,28</sup> (denoted by UK10K-WGS). A descriptive summary of all the data sets used in

96 this study can be found in **Supplementary Table 1** and the Methods section. We simulated a  
97 trait affected by 50 common, causal variants with effects drawn from  $N(0,1)$ , which together  
98 explained 20% of the phenotypic variation (proportion of variance explained by a causal variant,  
99 denoted by  $q^2$ , was 0.4%, on average). Prior to the simulations with errors, we showed by a  
100 simulation under the null (i.e., simulating a scenario without errors and applying DENTIST using  
101 the discovery GWAS as the reference) that the DENTIST test-statistics were well calibrated,  
102 meaning that DENTIST will only remove a very small proportion of variants if there are no  
103 errors and heterogeneity in the data (**Supplementary Figure 1**). We then simulated genotyping  
104 and allelic errors at 0.5% randomly selected variants respectively. Genotyping errors of each of  
105 these variants were simulated by altering the genotypes of a certain proportion ( $f_{error} = 0.05, 0.1$   
106 or 0.15) of randomly selected individuals, and allelic error of each of the variants was introduced  
107 by swapping the effect allele by the other allele. The simulation was repeated 200 times with the  
108 causal and erroneous variants re-sampled in each simulation. We then ran DENTIST using  
109 UK10K-WGS or an independent sample (UKB-8K-1KGP) as the LD reference after standard QCs  
110 of the discovery GWAS: removing variants with a Hardy-Weinberg Equilibrium (HWE) P-value  $<$   
111  $10^{-6}$  using the individual-level data or  $\Delta AF > 0.1$  with  $\Delta AF$  being the difference in allele frequency  
112 (AF) between the summary data and reference sample. The independent sample UKB-8K-1KGP  
113 is referred to as a set of 8000 unrelated individuals from the UK Biobank<sup>29</sup> (UKB) with variants  
114 imputed from the 1000 Genomes Project (1KGP). The statistical power (sensitivity) was  
115 measured by the proportion of erroneous variants in the data that can be detected from QC. We  
116 also computed the fold enrichment in probability of an erroneous variant being detected from  
117 QC compared to a random guess (i.e., the ratio of the percentage of true erroneous variants in  
118 the variants detected by DENTIST to that in all variants).

119

120 When using UKB-8K-1KGP as the reference,  $\sim 45\%$  of the genotyping and  $\sim 95\%$  of the allelic  
121 errors could be removed by the standard QCs (**Supplementary Table 2**). However, the  $\Delta AF$   
122 approach performed poorly for the very common variants, e.g., the power was  $\sim 16\%$  for  
123 variants with  $MAF > 0.45$ . After the standard QCs, DENTIST was able to detect  $\sim 42\%$  of the  
124 remaining genotyping and  $\sim 78\%$  of the remaining allelic errors (**Figure 1** and **Supplementary**  
125 **Table 3**), with only  $\sim 0.3\%$  variants being removed in total (**Supplementary Table 4**). The fold  
126 enrichment was 212 for allelic errors and of 112 for genotyping errors (**Supplementary Table**  
127 **5**), showing good specificity of DENTIST in detecting the simulated errors. Notably, the power to  
128 detect allelic errors was  $\sim 78\%$  for variants with  $MAF > 0.45$ , compensating the low power of the  
129  $\Delta AF$  approach in this MAF range (note: another shortcoming of the  $\Delta AF$  approach is that the  
130 threshold is heavily sample size dependent, and currently there is no consensus guidance on the  
131 choice of a  $\Delta AF$  threshold in practice). When restricted to variants passing the genome-wide

132 significance level (i.e.,  $p < 5 \times 10^{-8}$ ), the DENTIST detection power increased to  $\sim 87\%$  for the  
133 genotyping errors and  $\sim 84\%$  for the allelic errors (**Supplementary Table 3**). The power also  
134 varied with the genotyping error rate ( $f_{error}$ ), e.g., the power in the  $f_{error}=0.05$  scenario was, on  
135 average, lower than that for  $f_{error}=0.15$  (**Figure 1c** and **Supplementary Table 3**). When using  
136 UK10K-WGS as the reference (mimicking the application of DENTIST in a scenario where  
137 individual-level data of the discovery GWAS are available), the power remained similar, but the  
138 fold enrichment was much higher compared to that using UKB-8K-1KGP (**Supplementary**  
139 **Tables 5** and **6**). In addition, using this same simulation setting, we explored the choice of the  
140 parameter  $\theta_k$  (i.e., the proportion of eigenvectors retained in SVD; see Methods for details) and  
141 reference sample size ( $n_{ref}$ ), and the results suggested a choice of  $\theta_k=0.5$  and  $n_{ref} \geq 5000$  in  
142 practice (**Supplementary Figure 2**). Together, these results demonstrate the power of DENTIST  
143 to identify allelic and genotyping errors even after the standard QCs, suggesting that DENTIST  
144 can complement existing QC filters for either individual- or summary-level GWAS data. On the  
145 other hand, DENTIST was parsimonious in data filtering, with  $\sim 0.3\%$  of the variants being  
146 removed in total across all the simulation scenarios (**Supplementary Table 4**). Nevertheless,  
147 we acknowledge that this simulation did not cover the full complexity of real case scenarios,  
148 which may involve multiple independent samples with heterogeneous LD structures caused by  
149 several factors, such as imputation errors or ancestry mismatches (**Supplementary Figure 3**).  
150 These cases are difficult to mimic in this simulation but will be assessed in the following  
151 analyses.

152

### 153 **Applying DENTIST to COJO with simulated phenotypes**

154 COJO<sup>6</sup> is a method that uses summary data from a GWAS or meta-analysis and LD data from a  
155 reference sample to run a conditional and joint multi-SNP regression analysis. We used  
156 simulations to assess the performance of COJO in the presence of heterogeneity between  
157 discovery GWAS and LD reference before and after DENTIST filtering. To mimic the reality that  
158 causal signals are often not perfectly captured by imputed variants, we simulated a phenotype  
159 affected by one or two sequenced variants using WGS data (i.e., UK10K-WGS) and performed  
160 association analyses using imputed data of the same individuals (imputing 312,264 variants, in  
161 common with those on an SNP array, to the 1KGP<sup>28,30</sup>; denoted by UK10K-1KG). More  
162 specifically, we first randomly selected one or two variants from two MAF bins as causal  
163 variants, i.e., variants with  $MAF \geq 0.01$  (denoted by common-causal) and  $0.01 > MAF \geq 0.001$   
164 (denoted by rare-causal) to generate a phenotype (note:  $MAF > 0.001$  is equivalent to minor  
165 allele count  $> 7$  in this sample). The causal variant  $q^2$  was set to 2% to achieve similar power to a  
166 scenario with  $q^2 = 0.03\%$  and  $n = 250,000$  (note: the mean  $q^2$  of 697 height variants discovered  
167 in Wood et al.<sup>31</sup> is 0.03%) because the power of GWAS is determined by  $nq^2/(1 - q^2)$ . Then, we

168 ran a GWAS using UK10K-1KGP and performed COJO analyses using multiple LD references,  
169 including the discovery GWAS sample, UKB-8K-1KGP, the Health Retirement Study (HRS)<sup>32</sup>, and  
170 the Atherosclerosis Risk In Communities (ARIC) study<sup>33</sup>, with different degrees of ancestral  
171 differences with UK10K-1KGP (**Supplementary Figure 4**). For a fair comparison, only the  
172 variants shared between these reference samples were included. We repeated the simulation  
173 100 times for each autosome and computed the false positive rate (FPR, i.e., the frequency of  
174 observing two COJO signals in the scenario where there was only one causal variant) and power  
175 (the frequency of observing two COJO signals in the scenario where there were two distinct  
176 causal variants with LD  $r^2 < 0.1$  between them). It should be noted that the false positive COJO  
177 signals defined here are not false associations but falsely claimed as jointly associated (also  
178 known as quasi-independent) signals.

179  
180 When using the discovery GWAS sample as the LD reference, the FPR of COJO was 0.1% for  
181 common-causal and 0.2% for rare-causal (**Figure 2; Table 1**), which can be regarded as a  
182 baseline for comparison as there was no data heterogeneity in this case. The FPRs were higher  
183 than the expected values because the causal variants were not perfectly tagged by the imputed  
184 variants (**Supplementary Table 7**). When using UKB-8K-1KGP (i.e., 1KGP-imputed data of 8000  
185 UKB participants with similar ancestry to the UK10K participants as shown in **Supplementary**  
186 **Figure 4**) as the LD reference, the FPR was close to the benchmark for common-causal (1%) and  
187 slightly inflated for rare-causal (2.7%) (**Figure 2**). After DENTIST filtering (using UKB-8K-1KGP  
188 as the LD reference), the FPR for rare-causal decreased to 1.3%. Moreover, when using LD  
189 computed from European-American individuals in HRS or ARIC, the FPR of COJO was strongly  
190 inflated in the whole MAF range: >7% for common-causal and >28% for rare-causal, likely  
191 because of the difference in ancestry between HRS/ARIC and UK10K-1KGP. DENTIST can  
192 effectively control the FPR of COJO to <1% for common-causal and <2% for rare-causal (**Figure**  
193 **2**). Taken together, the FPR of COJO was reasonably well controlled for common variants but  
194 substantially inflated for rare variants especially when there was a difference in ancestry  
195 between the GWAS and LD reference samples, and most of the false positive COJO signals could  
196 be removed by DENTIST.

197  
198 The power of COJO (without DENTIST) using in-sample LD from UK10K-1KGP or out-of-sample  
199 LD from the other references were similar: 77-81% for common-causal and 26-30% for rare-  
200 causal (**Table 1**). The low power for rare-causal was because they were poorly captured by  
201 imputation (**Supplementary Table 7**). DENTIST filtering caused a <2% loss of power for  
202 common-causal, and 5-10% for rare-causal (**Table 1**). Hence, the control of FPR of COJO by

203 DENTIST was to some extent at the expense of power although the reduction in FPR was larger  
204 than that in power.

205

206 We also examined the effect of imputation INFO score-based QC on the FPR and power of COJO.  
207 Take the analysis with HRS as an example. By removing variants with INFO-scores  $< 0.9$  from the  
208 HRS data, the FPR of COJO decreased to 2.3% for common-causal and 6% for rare-causal  
209 (**Supplementary Table 8**), both of which were higher than those using DENTIST (FPR = 0.5%  
210 for common-causal and 1.7% for rare-causal) (**Table 1**). Meanwhile, the power of COJO after the  
211 INFO score-based QC decreased to 70% for common-causal and 12% for rare-causal  
212 (**Supplementary Table 8**), both of which were lower than those using DENTIST (power = 81%  
213 for common-causal and 27% rare-causal). The other less stringent INFO-score threshold were  
214 even less effective, and the results from analyses using the other references were similar  
215 (**Supplementary Table 8**). These results suggest that filtering variants by imputation INFO is  
216 less effectively than that by DENTIST.

217

### 218 **Applying DENTIST to COJO for real phenotypes**

219 Having assessed the performance of DENTIST in COJO analyses by simulation, we then applied it  
220 to COJO analyses for height in the UKB. The height GWAS summary statistics ( $n = 328,577$ ) were  
221 generated from a GWAS analysis of all the unrelated individuals of European ancestry in the UKB  
222 (denoted by UKBv3-329K) except 20,000 individuals (denoted by UKBv3-20K), which were used  
223 as a non-overlapping LD reference. Genotype imputation of the UKB data was performed by the  
224 UKB team with most of the variants imputed from the Haplotype Reference Consortium (HRC)<sup>34</sup>.  
225 We performed COJO analyses with a host of references: overlapping in-sample references with  
226 sample sizes ( $n_{\text{ref}}$ ) varying from 10,000 to 150,000, non-overlapping in-sample references  
227 including UKBv3-8K ( $n = 8,000$ ) and UKBv3-20K (containing UKBv3-8k), and out-of-sample  
228 references including ARIC and HRS (**Supplementary Table 1**). We excluded from the analysis  
229 variants with  $\text{MAF} < 0.001$  to ensure sufficient number of minor alleles for rare variants in  
230 reference samples with  $n_{\text{ref}} < 10\text{k}$ . We first performed a COJO analysis using the actual GWAS  
231 sample as the reference and identified 1,279 signals from variants with  $\text{MAFs} > 0.01$ , and 1310  
232 signals from variants with  $\text{MAFs} > 0.001$  (**Table 2**). These results can be regarded as a  
233 benchmark. When using the overlapping in-sample LD references, the number of COJO signals  
234 first decreased as  $n_{\text{ref}}$  increased and then started to stabilize when  $n_{\text{ref}}$  exceeded 30,000  
235 (**Supplementary Figure 5**). The results from using the two non-overlapping in-sample  
236 references (UKBv3-8K and UKBv3-20K) were comparable to those from using the overlapping  
237 in-sample references with similar sample sizes (**Table 2** and **Supplementary Table 9**) because  
238 the non-overlapping in-sample references, despite being excluded from the GWAS, were

239 consistent with the GWAS sample with respect to ancestry, data collection, and analysis  
240 procedures.

241  
242 When using LD from an out-of-sample reference (either HRS or ARIC), there was substantial  
243 inflation in the number of COJO signals compared to the benchmark (by 15.5-16.1% for common  
244 variants and 18.7-25.6% for all variants), with a few variants in weak LD with those identified  
245 from the benchmark analysis (**Supplementary Figure 6**). The results from using the two out-of-  
246 sample references became more consistent with the benchmark after DENTIST filtering, with the  
247 inflation reduced to 4.6-5.8% for common variants and 2.7-6.7% for all variants, comparable to  
248 the results using an in-sample LD reference with a similar sample size (**Table 2**). Polygenic score  
249 analysis shows that the reduction in the number of COJO signals owing to DENTIST QC had  
250 almost no effect on the accuracy of using the COJO signals to predict height in HRS  
251 (**Supplementary Table 10**), suggesting the redundancy of the COJO signals removed by  
252 DENTIST. We further found that compared to using the imputed data from ARIC or HRS, using  
253 UK10K-WGS (n=3,642) as the reference showed lower inflation (10% for common variants and  
254 <12.4% for all variants) before DENTIST QC but larger inflation after DENTIST QC (**Table 2**),  
255 suggesting a large reference sample size is essential even for WGS data. In all the DENTIST  
256 analyses above, the total number of removed variants varied from 0.05% to 0.98%  
257 (**Supplementary Table 11**). All these results are consistent with what we observed from  
258 simulations, demonstrating the effectiveness of DENTIST in eliminating heterogeneity between  
259 GWAS and LD reference samples.

260  
261 We further applied DENTIST to Educational Attainment (EA), Coronary Artery Disease (CAD),  
262 Type 2 Diabetes (T2D), Crohn's Disease (CD), Major Depressive Disorder (MDD), Schizophrenia  
263 (SCZ), Ovarian Cancer (OC), Breast Cancer (BC), Height and Body Mass Index (BMI) using GWAS  
264 summary data from the public domain<sup>35-43</sup> (**Supplementary Table 12**) and three LD reference  
265 samples (i.e., ARIC, HRS, and UKBv3-8K). Since these published studies focus on common  
266 variants (rare variants are not available in most of the data sets), we used a MAF threshold of  
267 0.01 in this analysis. When using ARIC as the LD reference, the proportion of variants removed  
268 by DENTIST QC ranged from 0.02% (BMI) to 0.94% (CAD) with a median of 0.28%  
269 (**Supplementary Table 13**), and the reduction in the number of COJO signals for common  
270 variants ranged from 0% (OC and MDD) to 11.9% (CAD) with a median of 1.5%  
271 (**Supplementary Table 14**). The results from using the other two references are similar  
272 (**Supplementary Table 13 and 14**).

273

274 **Improved HEIDI test with DENTIST**



275 The summary data-based Mendelian randomization (SMR) is a method that integrates summary-  
276 level data from a GWAS and an expression quantitative trait loci (eQTL) study to test pleiotropic  
277 associations between a trait and expression levels of genes<sup>21</sup>. It features the HEIDI test that  
278 utilizes multiple cis-eQTL variants at a locus to distinguish pleiotropy (the trait and expression  
279 level of a gene are affected by the same causal variants) from linkage (causal variants for the  
280 trait are in LD with a distinct set of causal variants affecting gene expression). The HEIDI test  
281 uses summary data from two studies and LD from a reference so that any errors in and  
282 heterogeneity between the GWAS, eQTL and reference samples can cause inflated HEIDI test  
283 statistics, giving rise to more SMR associations being rejected than expected by chance<sup>21,44</sup>. Here,  
284 we performed simulations to assess the effect of data heterogeneity on HEIDI and sought to  
285 mitigate it using DENTIST. We first generated a trait based on a causal variant ( $q^2 = 1\%$ )  
286 randomly sampled from the variants on chromosome 22 in the ARIC data. To simulate a  
287 pleiotropic model, we used the same causal variant to simulate the expression level of a gene in a  
288 subset of the HRS data ( $n = 3,000$ ; denoted by HRS-3K) with  $q^2$  for the gene expression level  
289 randomly sampled from the eQTL  $q^2$  distribution reported by the Consortium for the  
290 Architecture of Gene Expression (CAGE)<sup>45</sup>. To simulate a linkage model, a second causal variant  
291 in LD ( $r^2 > 0.25$ ) with the trait causal variant was selected to generate the gene expression level,  
292 again with the eQTL  $q^2$  value sampled from the CAGE. In addition to the two-sample scenario  
293 above, we also simulated a one-sample scenario in which both the trait and gene expression  
294 level were generated in the HRS-3K sample. The UKB-8K-1KGP sample was used as the LD  
295 reference for both the SMR-HEIDI and DENTIST analyses. For each scenario, the simulation was  
296 repeated 4000 times with the causal variants re-sampled in each replicate. The FPR of the HEIDI  
297 test was calculated as the proportion of pleiotropic models detected with  $P_{\text{HEIDI}} < 0.05$ , and the  
298 power was defined as the proportion of linkage models detected with  $P_{\text{HEIDI}} < 0.05$ .

299  
300 We found that the FPR of HEIDI was close to the expected value (5%) in the one-sample scenario  
301 (5.8%) but inflated (9.8%) in the two-sample scenario (**Figures 3a and 3b**, and **Supplementary**  
302 **Table 15**). To mitigate the inflation, we performed DENTIST in both the GWAS and eQTL  
303 summary data using UKB-8K-1KGP as the reference. After DENTIST filtering, the FPR of HEIDI in  
304 the two-sample scenario decreased to 7.6%; the decrease was small but statistically significant  
305 ( $P_{\text{difference}} = 0.002$ ). The results remained similar when the discovery GWAS sample (i.e., HRS-3K)  
306 was used as the reference (**Supplementary Table 15**). These results suggest that the HEIDI test  
307 statistic was inflated in the two-sample scenario likely because of LD heterogeneity between the  
308 GWAS and eQTL samples. The power of HEIDI to detect the linkage model remained almost the  
309 same before and after DENTIST filtering (**Figure 3c and Supplementary Table 15**). To further  
310 validate if the inflation of HEIDI FPR was due to heterogeneity between the GWAS and eQTL

311 samples, we increased the difference in ancestry between the two discovery samples by  
312 simulating GWAS and eQTL data from UKB-8K-1KGP and HRS-3K, respectively, and performed  
313 the HEIDI analysis using ARIC as the reference. In this case, the FPR of HEIDI increased to 12.0%  
314 before and to 9.1% after DENTIST filtering (**Supplementary Table 15**). All these results suggest  
315 that DENTIST slightly improved the FPR of HEIDI in the presence of data heterogeneity at almost  
316 no expense of power and that in the presence of ancestry difference between the GWAS and  
317 eQTL samples, HEIDI tends to be conservative (rejecting more SMR associations than expected  
318 by chance) even after DENTIST filtering.

319

### 320 **Improved LD score regression analysis with DENTIST**

321 LD score regression (LDSC) is an approach which was originally developed to distinguish  
322 polygenicity from population stratification in GWAS summary data set by a weighted regression  
323 of GWAS  $\chi^2$  statistics against LD scores computed from a reference<sup>12</sup> but has often been used to  
324 estimate the SNP-based heritability ( $h_{SNP}^2$ ). We investigated the impact of DENTIST on LDSC  
325 using the height GWAS summary data generated using the UKBv3-329K sample along with  
326 several reference samples including four imputation-based samples (i.e., the discovery GWAS  
327 sample, HRS, ARIC and UKB-8K-1KGP) and two WGS-based samples (i.e., UK10K-WGS and the  
328 European individuals from the 1KGP (1KGP-EUR)). We performed the one- and two-step LDSC  
329 analyses using LD scores of the variants, in common with those in the HapMap3, computed from  
330 each of the references before and after DENTIST-based QC. Note that DENTIST was performed  
331 for all common variants but only those overlapped with HapMap3 were included in the LDSC  
332 analyses.

333

334 Using the discovery GWAS sample as the reference, the estimates of  $h_{SNP}^2$  and regression  
335 intercept from the one-step LDSC were 46% (SE = 0.02) and 1.13 (SE = 0.04) respectively (**Table**  
336 **3**). When using the other reference samples, the results were very close to the benchmark except  
337 for a noticeably larger estimate of regression intercept using HRS (1.24, SE = 0.04). After  
338 DENTIST filtering, the intercept estimate using HRS decreased to 1.15 (SE = 0.04) with little  
339 difference in  $\hat{h}_{SNP}^2$  (increased from 45% to 46%) (**Table 3**). To better understand the effect of  
340 DENTIST QC on LDSC using HRS, we plotted the mean  $\chi^2$ -statistic against the mean LD score  
341 across the LD score bins. We found that the GWAS mean  $\chi^2$ -statistic in the bin with the smallest  
342 mean LD score deviated from the value expected from a linear relationship between the LD  
343 score and  $\chi^2$ -statistic (**Figure 4a**), and the deviation was removed by filtering out a small  
344 proportion of variants with very small LD scores but large  $\chi^2$ -statistic by DENTIST  
345 (**Supplementary Figure 7**). These results show how the quality of LD reference can impact the  
346 LDSC analysis, but such effect is typically unknown *a priori* and varied across different reference

347 samples. We also re-ran the analyses using the two-step LDSC, where the intercept was  
348 estimated using the variants with  $\chi^2$  values  $< 30$  in the first step and constrained in the second  
349 step to estimate  $h_{SNP}^2$  using all the variants. Compared to the one-step approach, the two-step  
350 approach provides larger estimates of the intercepts and smaller estimates of  $h_{SNP}^2$  either before  
351 or after DENTIST filtering. It is noteworthy that when using 1KGP-EUR as the LD reference,  
352 DENTIST suggested many more variants for removal compared to that using the other  
353 references, which caused a substantially smaller estimate of  $h_{SNP}^2$  (**Table 3**). This is because LD  
354 correlations computed from references of small sample size are noisy due to sampling variation,  
355 which cause inflated test-statistic and thereby elevated FPR of DENTIST (**Supplementary**  
356 **Figure 2**). This result cautions the use of DENTIST with LD references with small sample sizes  
357 (e.g.,  $n < 5000$ ). In addition, we applied LDSC to the 10 published GWAS data sets mentioned  
358 above (**Supplementary Table 12**) using ARIC and UKBv3-8K as the reference. The  
359 improvement of LDSC by DENTIST QC was small particularly for traits whose LDSC intercepts  
360 were close to 1 before DENTIST QC (**Supplementary Table 16**), demonstrating the robustness  
361 of LDSC to data heterogeneity and errors.

362

## 363 Discussion

364 In this study, we developed DENTIST, an QC tool for summary data-based analyses, which  
365 leverages LD from a reference sample to detect and filter out problematic variants by testing the  
366 difference between the observed z-score of a variant and a predicted z-score from the  
367 neighboring variants. From simulations and real data analyses, we show that some of the  
368 commonly-used analyses including the COJO, SMR-HEIDI and LDSC, can be biased to various  
369 extents in the presence of data heterogeneity, e.g., inflated number of COJO signals, elevated rate  
370 of rejecting pleiotropic models for SMR-HEIDI, or biased estimates of regression intercept and  
371  $h_{SNP}^2$  for LDSC. For most of these analyses, DENTIST-based QC can substantially mitigate the  
372 biases.

373

374 Our results suggest that summary-data-based analyses are generally well calibrated in the  
375 absence of data heterogeneity but biased otherwise. For example, we showed that the mismatch  
376 in ancestry between the discovery GWAS and LD reference (e.g., European vs. British ancestry)  
377 caused inflated FPRs of both COJO and SMR-HEIDI analyses (**Table 2** and **Supplementary Table**  
378 **15**). Also, we found that the FPR of COJO for rare variants was much higher than that for  
379 common variants likely because rare variants are more difficult to impute such that they are  
380 more likely to have discrepancy in LD between two imputed data sets. It should be clarified that  
381 the false-positive COJO signals as defined here are not false associations but falsely claimed as  
382 jointly significant associations. DENTIST substantially reduced false-positive detections from

383 COJO analyses especially for rare variants even when there was a difference in ancestry between  
384 the GWAS and LD reference samples. This extends the utility of COJO, which was originally  
385 developed for common variants, to rare variants. This message is important for the field because  
386 more and more GWASs and meta-analyses have started included rare variants from imputation.  
387 The FPR of HEIDI was only marginally reduced by DENTIST but at almost no cost of power. It  
388 should also be clarified again that the inflated HEIDI test-statistics would not lead to false  
389 discoveries because higher HEIDI test-statistics correspond to higher probability of rejecting  
390 SMR associations rather than tending to claim more significant associations. Among all the  
391 methods tested, LDSC was least affected by errors or data heterogeneity, but in one case where  
392 HRS was used as the LD reference for the analysis of the UKB height summary data, the  
393 estimates were biased but could be corrected by DENTIST (**Figure 4a**). DENTIST has the unique  
394 feature to detect heterogeneity between a GWAS summary data set and an LD reference. The  
395 benefit of using DENTIST as a QC tool has been demonstrated in the three case studies above,  
396 but we believe that it can potentially be applied to all GWAS summary data-based analyses that  
397 require a LD reference such as fine mapping methods<sup>7,9,46</sup> and joint modeling of all variants for  
398 polygenic risk prediction<sup>22,23,47</sup>. We have also shown by simulation that DENTIST can even add  
399 value to the standard GWAS QC process in a single-cohort-based GWAS to detect  
400 genotyping/imputation errors.

401  
402 Given that a QC step can potentially remove true signals, we make sure that DENTIST is  
403 conservative in filtering variants. We show by simulation that in the absence of errors and data  
404 heterogeneity, the DENTIST test-statistics were not inflated, and on average, only <0.05%  
405 variants were filtered out by DENTIST (**Supplementary Figure 1**). In practice, we implemented  
406 two strategies to avoid widespread inflation of the DENTIST statistics in the presence of data  
407 errors or heterogeneity: 1) we used the SVD truncation method to control for sampling variation  
408 in LD estimated from the reference; 2) we applied an iterative approach to prioritize the  
409 elimination of larger outliers in earlier iterations (Methods). We optimized the parameters  
410 related to these two steps through simulations (**Supplementary Figure 2**). Throughout all the  
411 analyses performed in this study, we found in no cases DENTIST degraded the results, and  
412 DENTIST often only needed to remove a very small proportion of the variants to correct or  
413 alleviate the biases (**Table 3** and **Supplementary Tables 4, 11, 13**). To the contrary, filtering  
414 variants based on imputation INFO score<sup>48</sup> caused significant loss of power in the COJO analysis  
415 when a stringent INFO score threshold was used to achieve a similar level of FPR as that using  
416 DENTIST.

417

418 Our method is an early attempt at QC for summary-data-based analyses. To avoid misuse, we  
419 summarize the usages and limitations, in addition to the features mentioned above. Firstly,  
420 DENTIST is a QC method for detecting not only errors in summary-data but also heterogeneity  
421 between discovery and reference data. As shown from our simulation, DENTIST does not  
422 guarantee the filtering of all the errors but most of them with large GWAS z-scores and a large  
423 proportion of them with small z-scores. Secondly, DENTIST can identify errors that passed the  
424 standard QC approaches (such as HWE test and allelic frequency checking), which makes it a  
425 good complementary method to existing QC filters. We suggest that DENTIST-based QC should  
426 be applied after the standard QC as DENTIST is more powerful when the proportion of errors is  
427 smaller (**Supplementary Table 3**). DENTIST can also be used as a method for checking  
428 summary data sanity by running it with a reliable LD reference sample. Thirdly, regarding the  
429 choice of a reference sample, DENTIST expects unrelated individuals from a closely matched  
430 ancestry, with a large sample size ( $n > 5,000$ ). From simulations, we found that small reference  
431 sample size biased the DENTIST test statistics leading to significantly elevated FPR  
432 (**Supplementary Figure 2**). Lastly, DENTIST assumes the test statistics of different variants  
433 have similar sample sizes; violation of this assumption will lead to variants with significantly  
434 smaller or larger sample sizes being mistakenly recognized as problematic variants by DENTIST.  
435

436 In summary, we have proposed a new QC approach to improve summary-data-based analyses  
437 that are potentially affected by the errors in summary data or heterogeneity between data sets.  
438 This method has been implemented in a user-friendly software tool DENTIST. The software tool  
439 is multi-threaded so that it is computationally efficient when enough computing resources are  
440 available. For example, when running each chromosome in parallel, it took <1h to run DENTIST  
441 for all variants with  $MAF > 1\%$  and <5h for all variants with  $MAF > 0.01\%$  on each chromosome  
442 (**Supplementary Table 16**).  
443

## 444 **Methods**

### 445 **The DENTIST test-statistic**

446 Given an ancestrally homogeneous sample and a genotype matrix  $\mathbf{X}$  consisting of  $n$  unrelated  
 447 individuals genotyped/imputed at  $m$  variants, an association study is carried out at each variant  
 448 by performing a linear regression between the variant and a phenotype of interest. This  
 449 provides a set of summary data that include the estimate of variant effect, the corresponding  
 450 standard error, and thereby the z-statistic. Under the null hypothesis of no association, the z-  
 451 scores of  $m$  variants follow a multivariate normal distribution,  $\mathbf{Z} \sim MVN(\mathbf{0}, \Sigma)$  with  $\mathbf{Z} =$   
 452  $(Z_1, Z_2, \dots, Z_m)$ , with  $\Sigma$  a LD correlation matrix of the variants.

453  
 454 The aim of this method is to test the difference between the z-statistic of a variant and that  
 455 predicted from adjacent variants. To do this, we use a sliding window approach to divide  
 456 genome into 2Mb segments with a 500kb overlap between one another and randomly partition  
 457 the variants in a segment into two subsets, S1 and S2, with similar numbers of variants. We then  
 458 use variants in S2 to predict those in S1 and vice versa. According to previous studies<sup>10,11</sup>, the  
 459 distribution of z-statistic of a variant  $i$  from S1, conditional on the observed z-scores of a set of  
 460 variants from S2 is

$$461 \quad Z_i | \mathbf{z}_t = \mathbf{z}_t \sim N(\Sigma_{it} \Sigma_{tt}^{-1} \mathbf{z}_t, \Sigma_{ii} - \Sigma_{it} \Sigma_{tt}^{-1} \Sigma'_{it}) \quad (1),$$

462 where  $\Sigma_{it}$  denotes the correlation of z-scores between variant  $i$  from S1 and variants  $t$  from S2,  
 463 and  $\Sigma_{tt}$  is the correlation matrix of variants  $t$ . We use the correlation matrix calculated from an  
 464 ancestry-matched reference sample (denoted by  $\mathbf{R}$ ) to replace that in the discovery sample if  
 465 individual-level genotypes of in the discovery GWAS are unavailable. In this case, **Equation 1**  
 466 can be rewritten as

$$467 \quad Z_i | \mathbf{z}_t \sim N(\mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t, \mathbf{1} - \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{R}'_{it}) \quad (2).$$

468 We can use  $E(Z_i | \mathbf{z}_t)$  as a predictor of  $Z_i$ , i.e.,  $\tilde{z}_i = \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t$ , and can therefore use the test-  
 469 statistic below to test the difference between the observed and predicted z-scores

$$470 \quad T_{d(i)} = (z_i - \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t)^2 / (\mathbf{1} - \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{R}'_{it}) \quad (3)$$

471 which approximately follows a  $\chi^2$  distribution with 1 degree of freedom. A deviation of  $T_{d(i)}$  from  
 472  $\chi^2_1$  can be attributed to 1) errors in the summary data; 2) errors in the reference data; or 3)  
 473 heterogeneity between the two data sets. Using **Equation 3**, the test statistic  $T_d$  can be  
 474 calculated for each variant in S1 given z-scores from S2. As in previous studies<sup>10,11</sup>, the method is  
 475 derived under the null hypothesis of no association, but the test-statistics are well calibrated in  
 476 the presence of true association signals (**Supplementary Figure 2**).

477

### 478 **The iterative partitioning approach**

479 One challenge of using **Equation 3** is that errors in  $\mathbf{z}_t$  or discrepancy between  $\mathbf{R}_{it}$  and  $\Sigma_{it}$  can  
480 affect the accuracy of predicting  $\tilde{z}_i$ . To mitigate this, we use an iterative partitioning approach.  
481 That is, in each iteration, we randomly partition the variants into two sets, S1 and S2, predict the  
482 z-statistic of each variant in S1 using its adjacent variants in S2 and vice versa, and run the  $T_d$   
483 test to remove a small fraction of variants with  $P_{\text{DENTIST}} < 5 \times 10^{-8}$  (capped at 0.5% variants with  
484 the smallest  $P_{\text{DENTIST}}$  if more than 0.5% of the variants exceeding this threshold). The default  
485 number of iterations is set to 10. In this iterative process, variants with very large errors or LD  
486 heterogeneity between the discovery and LD reference samples are prioritized for removal in  
487 the first few iterations so that the prediction accuracy increases in the following iterations. After  
488 the iterations are completed, any SNPs with  $P_{\text{DENTIST}} < 5 \times 10^{-8}$  are removed in the final step.

489

### 490 Accounting for sampling noise in LD

491 A simple replacement of the LD correlation matrix  $\Sigma$  by  $\mathbf{R}$  introduces additional noises, which  
492 can inflate  $T_d$ , because the sampling variations in  $\mathbf{R}$  differ from those  $\Sigma$ . Therefore, we adopt a  
493 truncated singular value decomposition (SVD) approach used in a previous study<sup>13</sup> to suppress  
494 the sampling noises. The essential idea was to remove variance components of  $\mathbf{R}_{tt}$  that  
495 corresponded to the smallest singular values in SVD, as these variance components were likely  
496 to be induced by sampling noises. Given the equivalence between SVD and eigen decomposition  
497 of  $\mathbf{R}_{tt}$ , we perform pseudoinverse of  $\mathbf{R}_{tt}$  using eigen decomposition, set small eigen values to 0,  
498 and retain only  $k$  components with large eigen values.

$$499 \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t = \mathbf{R}_{it} \mathbf{R}_{tt}^+ \mathbf{z}_t = \sum_{1..k} 1/w_k (\mathbf{R}_{it} \mathbf{v}_k) (\mathbf{v}_k' \mathbf{z}_t) \quad (4)$$

$$500 \mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{R}'_{it} = \mathbf{R}_{it} \mathbf{R}_{tt}^+ \mathbf{R}'_{it} = \sum_{1..k} 1/w_k (\mathbf{R}_{it} \mathbf{v}_k)^2 \quad (5)$$

501  $\mathbf{R}_{tt}^+$  denotes the pseudo inversion of  $\mathbf{R}_{tt}$ . The scalars  $w_1, \dots, w_k$  correspond to the largest  $k$   
502 eigenvalues, and vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are the corresponding  $k$  eigenvectors. Given  $q = \text{rank}(\mathbf{R}_{tt})$ , the  
503 suggested value of  $k$  is  $k \ll q$ . Let  $\theta_k = k/q$ . We show by simulation that  $\theta_k = 0.5$  appears to be a  
504 good choice meanwhile a large reference sample size (e.g.,  $n_{\text{ref}} \geq 5000$ ) is need  
505 (**Supplementary Figure 2**). This pseudoinverse also prevents the problem of rank deficiency  
506 due to strongly correlated variants when computing  $\mathbf{R}_{tt}^{-1}$ .

507

508 According to the term  $\mathbf{R}_{it} \mathbf{R}_{tt}^{-1} \mathbf{z}_t$  in **Equation 3**, which is a weighted sum of multiple z-scores, a  
509 variant displaying a strong correlation with  $i$  can overrule the information from the rest of the  
510 variants in S2. This would affect the robustness of our method. Therefore, we prune the variants  
511 for LD with an  $r^2$  threshold of 0.95 (note: we do not actually remove variants in this case). For a  
512 set of variants in high LD ( $r^2 > 0.95$ ), variants pruned out by this pruning process are assigned  
513 with the same  $T_d$  value as that of the variant retained.

514

515

## 516 **Genotype data sets**

517 This study is approved by the University of Queensland Human Research Ethics Committee  
518 (approval number: 2011001173). A summary of the genotype data sets used in this study as well  
519 as their relevant information can be found in **Supplementary Table 1**. These data are from four  
520 GWAS cohorts of European descendants, including the Health Retirement Study (HRS)<sup>32</sup>,  
521 Atherosclerosis Risk in Communities (ARIC) study<sup>33</sup>, UK10K<sup>27</sup>, and UK Biobank (UKB)<sup>29</sup>. The  
522 samples were genotyped using either WGS or SNP array technology (**Supplementary Table 1**).  
523 Imputation of the UKB data had been performed in a previous study<sup>49</sup> using the Haplotype  
524 Reference Consortium (HRC)<sup>34</sup> and UK10K reference panels<sup>29,50</sup>. We used different subsets of the  
525 imputed UKB data as the LD reference in this study, denoted with the prefix “UKBv3”, such as  
526 UKBv3-unrel (all the unrelated individuals of European ancestry,  $n = 348,577$ ), UKBv3-329K (a  
527 subset of 328,577 individuals of UKBv3-unrel), UKBv3-20K (another subset of 20,000 individuals  
528 of UKBv3-unrel, independent of UKBv3-392K) and UKBv3-8K (a subset of 8,000 individuals of  
529 UKBv3-20K). HRS, ARIC and UK10K cohorts were imputed to the 1KGP reference panel in prior  
530 studies<sup>28,30</sup>, and a subset of 8000 unrelated individuals from UKB were imputed to the 1KGP  
531 reference panel in this study (referred to as UKB-8K-1KGP). The UK10K variants in common with  
532 those on an Illumina CoreExome array were used for 1KGP imputation<sup>30</sup>. The imputation dosage  
533 values were converted to best-guess genotypes in all the data sets except for UKBv3-all, in which  
534 the hard-called genotypes were converted from the imputation dosage values using PLINK2 --  
535 hard-call-threshold 0.1 (Ref<sup>51</sup>). For all the data set, standard QCs were performed to remove  
536 variants with HWE test  $P$ -value  $< 10^{-6}$ , imputation INFO score  $< 0.3$ , or MAF  $< 0.001$ . Since the hard-  
537 called genotypes had missing values, in UKBv3 and its subsets, we further removed variants with  
538 missingness rate  $> 0.05$ .

539

## 540 **Genome-wide association analysis for height using the UKB data**

541 We performed a genome-wide association analysis for height using the genotype data of UKBv3-  
542 329K, i.e., all the unrelated individuals of European ancestry in the UKB ( $n=328,577$ ) except for  
543 20000 individuals randomly selected to create a non-overlapping reference sample (i.e., UKBv3-  
544 20K). The height phenotype was pre-adjusted for sex and age. We conducted the association  
545 analysis using the simple linear regression model in fastGWA<sup>52</sup> with the first 10 principle  
546 components (PCs) fitted as covariates.

547

## 548 **Data availability**

549 All the data sets used in this study are available in the public domain (**Supplementary Table 12**).

550



## 551 **Code availability**

552 The software tool DENTIST was written in C++ as a command-line tool. The source code and  
553 pre-compiled executable for 64-bit Linux distributions are available at [https://github.com/Yves-](https://github.com/Yves-CHEN/DENTIST/)  
554 CHEN/DENTIST/.

555

## 556 **Acknowledgements**

557 We are very grateful for constructive comments from Naomi Wray, Loic Yengo, Ying Wang and  
558 Jian Zeng and technical supports from Allan McRae, Julia Sidorenko, and Futao Zhang. This  
559 research was supported by the Australian Research Council (FT180100186, FL180100072), the  
560 Australian National Health and Medical Research Council (1113400 1107258), and the Sylvia &  
561 Charles Viertel Charitable Foundation.

562

## 563 **Author Contributions**

564 JY conceived and supervised the study. WC, ZZh and JY developed the method. WC, YW, ZZh, and  
565 JY designed the experiment. WC performed the simulations and data analyses under the  
566 assistance and guidance from YW, ZZl, TQ, PMV, ZZh and JY. WC developed the software tool.  
567 PMV and JY contributed funding and resources. WC and JY wrote the manuscript with the  
568 participation of all authors. All authors reviewed and approved the final manuscript.

569

## 570 **Competing Interests**

571 The authors declare no competing interests.

572

## 573 **References**

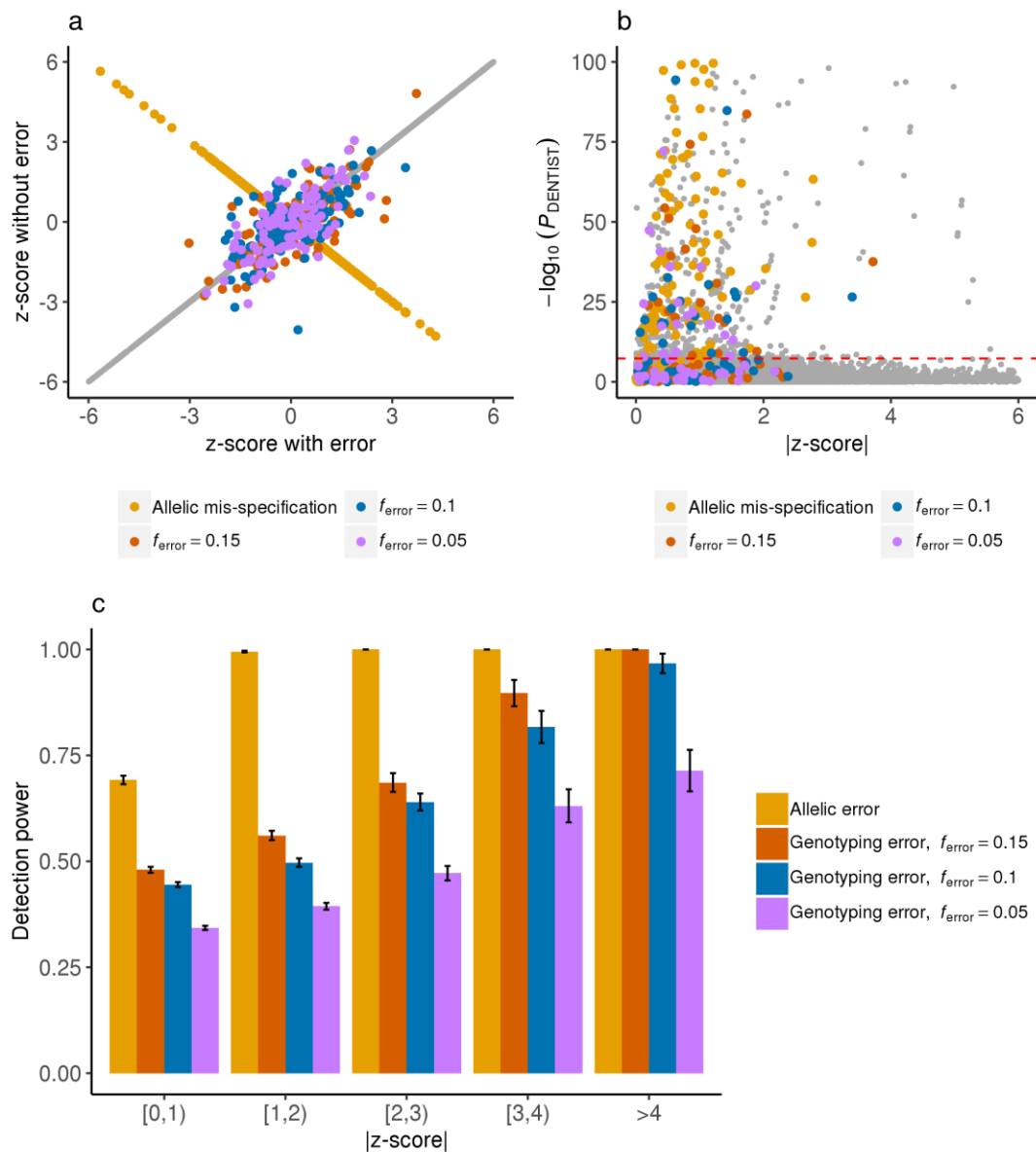
- 574 1. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association  
575 studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* **47**, D1005-D1012  
576 (2019).
- 577 2. Visscher, P.M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am J*  
578 *Hum Genet* **101**, 5-22 (2017).
- 579 3. Pasaniuc, B. & Price, A.L. Dissecting the genetics of complex traits using summary  
580 association statistics. *Nat Rev Genet* **18**, 117-127 (2017).
- 581 4. Schaid, D.J., Chen, W. & Larson, N.B. From genome-wide associations to candidate causal  
582 variants by statistical fine-mapping. *Nat Rev Genet* **19**, 491-504 (2018).
- 583 5. Dadaev, T. *et al.* Fine-mapping of prostate cancer susceptibility loci in a large meta-  
584 analysis identifies candidate causal variants. *Nat Commun* **9**, 2256 (2018).
- 585 6. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics  
586 identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-75, S1-3  
587 (2012).
- 588 7. Chen, W. *et al.* Fine mapping causal variants with an approximate Bayesian method using  
589 marginal test statistics. *Genetics* **200**, 719-736 (2015).
- 590 8. Zhu, X. & Stephens, M. Bayesian large-scale multiple regression with summary statistics  
591 from genome-wide association studies. *Ann Appl Stat* **11**, 1561 (2017).

- 592 9. Wang, G., Sarkar, A.K., Carbonetto, P. & Stephens, M. A simple new approach to variable  
593 selection in regression, with application to genetic fine-mapping. *bioRxiv*, 501114 (2019).
- 594 10. Pasaniuc, B. *et al.* Fast and accurate imputation of summary statistics enhances evidence  
595 of functional enrichment. *Bioinformatics* **30**, 2906-14 (2014).
- 596 11. Lee, D., Bigdeli, T.B., Riley, B.P., Fanous, A.H. & Bacanu, S.A. DIST: direct imputation of  
597 summary statistics for unmeasured SNPs. *Bioinformatics* **29**, 2925-7 (2013).
- 598 12. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding from polygenicity  
599 in genome-wide association studies. *Nat Genet* **47**, 291-5 (2015).
- 600 13. Shi, H., Kichaev, G. & Pasaniuc, B. Contrasting the Genetic Architecture of 30 Complex  
601 Traits from Summary Association Data. *Am J Hum Genet* **99**, 139-153 (2016).
- 602 14. Finucane, H.K. *et al.* Partitioning heritability by functional annotation using genome-wide  
603 association summary statistics. *Nat Genet* **47**, 1228 (2015).
- 604 15. Zhu, Z. *et al.* Causal associations between risk factors and common diseases inferred from  
605 GWAS summary data. *Nat Commun* **9**, 224 (2018).
- 606 16. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits.  
607 *Nat Genet* **47**, 1236-41 (2015).
- 608 17. Hartwig, F.P., Davies, N.M., Hemani, G. & Davey Smith, G. Two-sample Mendelian  
609 randomization: avoiding the downsides of a powerful, widely applicable but potentially  
610 fallible technique. *Int J Epidemiol* **45**, 1717-1726 (2016).
- 611 18. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic  
612 association studies using summary statistics. *PLoS Genet* **10**, e1004383 (2014).
- 613 19. Gamazon, E.R. *et al.* A gene-based association method for mapping traits using reference  
614 transcriptome data. *Nat Genet* **47**, 1091-8 (2015).
- 615 20. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association  
616 studies. *Nat Genet* **48**, 245-52 (2016).
- 617 21. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex  
618 trait gene targets. *Nat Genet* **48**, 481-7 (2016).
- 619 22. Vilhjalmsón, B.J. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of Polygenic  
620 Risk Scores. *Am J Hum Genet* **97**, 576-92 (2015).
- 621 23. Lloyd-Jones, L.R. *et al.* Improved polygenic prediction by Bayesian multiple regression on  
622 summary statistics. *Nat Commun* **10**, 5086 (2019).
- 623 24. Johnson, E.O. *et al.* Imputation across genotyping arrays for genome-wide association  
624 studies: assessment of bias and a correction strategy. *Hum Genet* **132**, 509-22 (2013).
- 625 25. Winkler, T.W. *et al.* Quality control and conduct of genome-wide association meta-  
626 analyses. *Nat Protoc* **9**, 1192-212 (2014).
- 627 26. Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98 (2008).
- 628 27. UK10K consortium. The UK10K project identifies rare variants in health and disease.  
629 *Nature* **526**, 82-90 (2015).
- 630 28. Yang, J. *et al.* Genetic variance estimation with imputed variants finds negligible missing  
631 heritability for human height and body mass index. *Nat Genet* **47**, 1114-20 (2015).
- 632 29. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide  
633 range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779 (2015).
- 634 30. Wu, Y., Zheng, Z., Visscher, P.M. & Yang, J. Quantifying the mapping precision of genome-  
635 wide association studies using whole-genome sequencing data. *Genome Biol* **18**, 86  
636 (2017).
- 637 31. Wood, A.R. *et al.* Defining the role of common variation in the genomic and biological  
638 architecture of adult human height. *Nat Genet* **46**, 1173-86 (2014).
- 639 32. Sonnega, A. *et al.* Cohort profile: the health and retirement study (HRS). *Int J Epidemiol* **43**,  
640 576-585 (2014).
- 641 33. ARIC INVESTIGATORS. The atherosclerosis risk in community (aric) study: Design and  
642 objectives. *American Journal of Epidemiology* **129**, 687-702 (1989).
- 643 34. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nature*  
644 *genetics* **48**, 1279 (2016).

- 645 35. Lee, J.J. *et al.* Gene discovery and polygenic prediction from a genome-wide association  
646 study of educational attainment in 1.1 million individuals. *Nat Genet* **50**, 1112-1121  
647 (2018).
- 648 36. van der Harst, P. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an  
649 Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circ Res* **122**, 433-  
650 443 (2018).
- 651 37. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-  
652 density imputation and islet-specific epigenome maps. *Nat Genet* **50**, 1505-1513 (2018).
- 653 38. Liu, J.Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel  
654 disease and highlight shared genetic risk across populations. *Nat Genet* **47**, 979-986  
655 (2015).
- 656 39. Wray, N.R. *et al.* Genome-wide association analyses identify 44 risk variants and refine the  
657 genetic architecture of major depression. *Nat Genet* **50**, 668-681 (2018).
- 658 40. Pardini, A.F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant  
659 genes and in regions under strong background selection. *Nat Genet* **50**, 381-389 (2018).
- 660 41. Phelan, C.M. *et al.* Identification of 12 new susceptibility loci for different histotypes of  
661 epithelial ovarian cancer. *Nat Genet* **49**, 680-691 (2017).
- 662 42. Michailidou, K. *et al.* Association analysis identifies 65 new breast cancer risk loci. *Nature*  
663 **551**, 92-94 (2017).
- 664 43. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass  
665 index in approximately 700000 individuals of European ancestry. *Hum Mol Genet* **27**,  
666 3641-3649 (2018).
- 667 44. Wu, Y. *et al.* Integrative analysis of omics summary data reveals putative mechanisms  
668 underlying complex traits. *Nat Commun* **9**, 918 (2018).
- 669 45. Lloyd-Jones, L.R. *et al.* The Genetic Architecture of Gene Expression in Peripheral Blood.  
670 *Am J Hum Genet* **100**, 228-237 (2017).
- 671 46. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-  
672 wide association studies. *Bioinformatics* **32**, 1493-1501 (2016).
- 673 47. Robinson, M.R. *et al.* Genetic evidence of assortative mating in humans. *Nature Human*  
674 *Behaviour* **1**(2017).
- 675 48. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method for  
676 genome-wide association studies by imputation of genotypes. *Nat Genet* **39**, 906-13  
677 (2007).
- 678 49. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.  
679 *Nature* **562**, 203-209 (2018).
- 680 50. Huang, J. *et al.* Improved imputation of low-frequency and rare variants using the UK10K  
681 haplotype reference panel. *Nature communications* **6**, 1-9 (2015).
- 682 51. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based  
683 linkage analyses. *The American journal of human genetics* **81**, 559-575 (2007).
- 684 52. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale  
685 data. *Nature Genetics* **51**, 1749-1755 (2019).

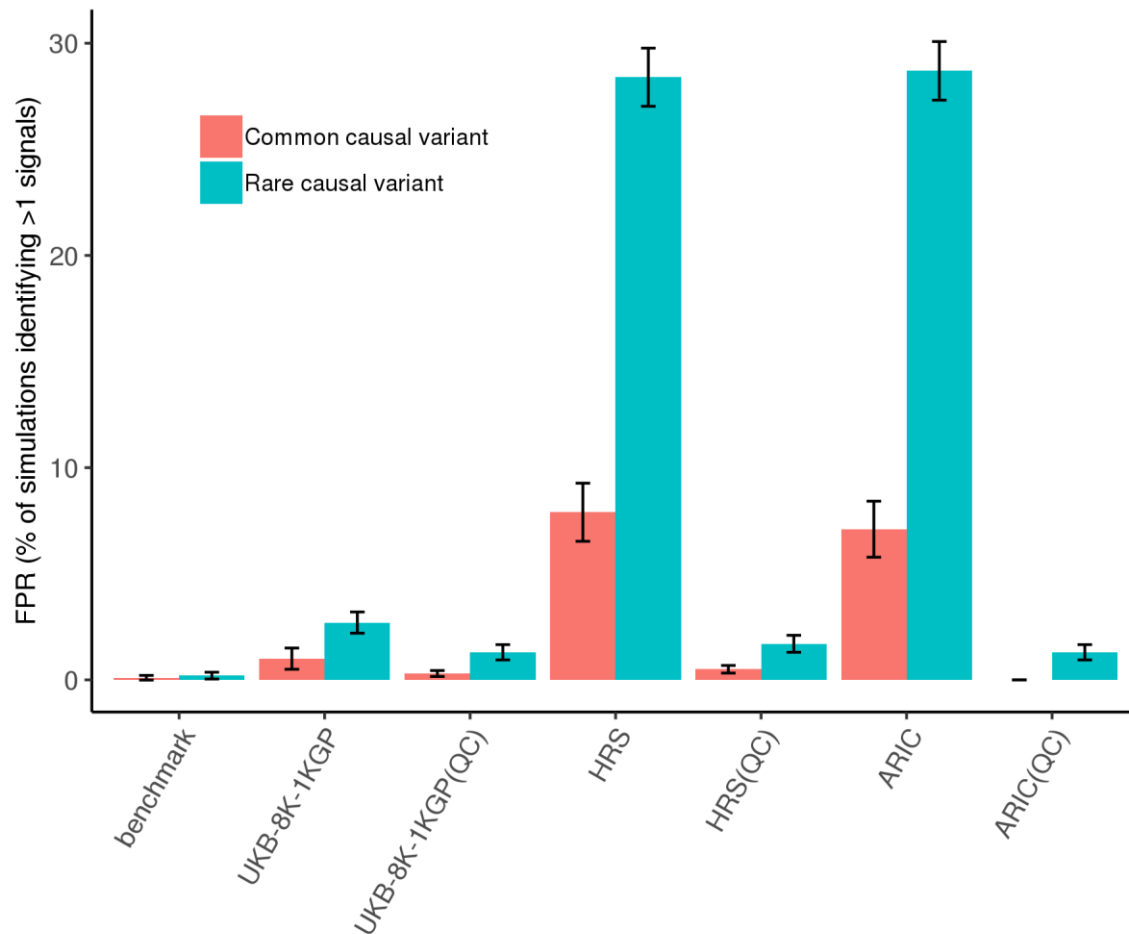
686

687



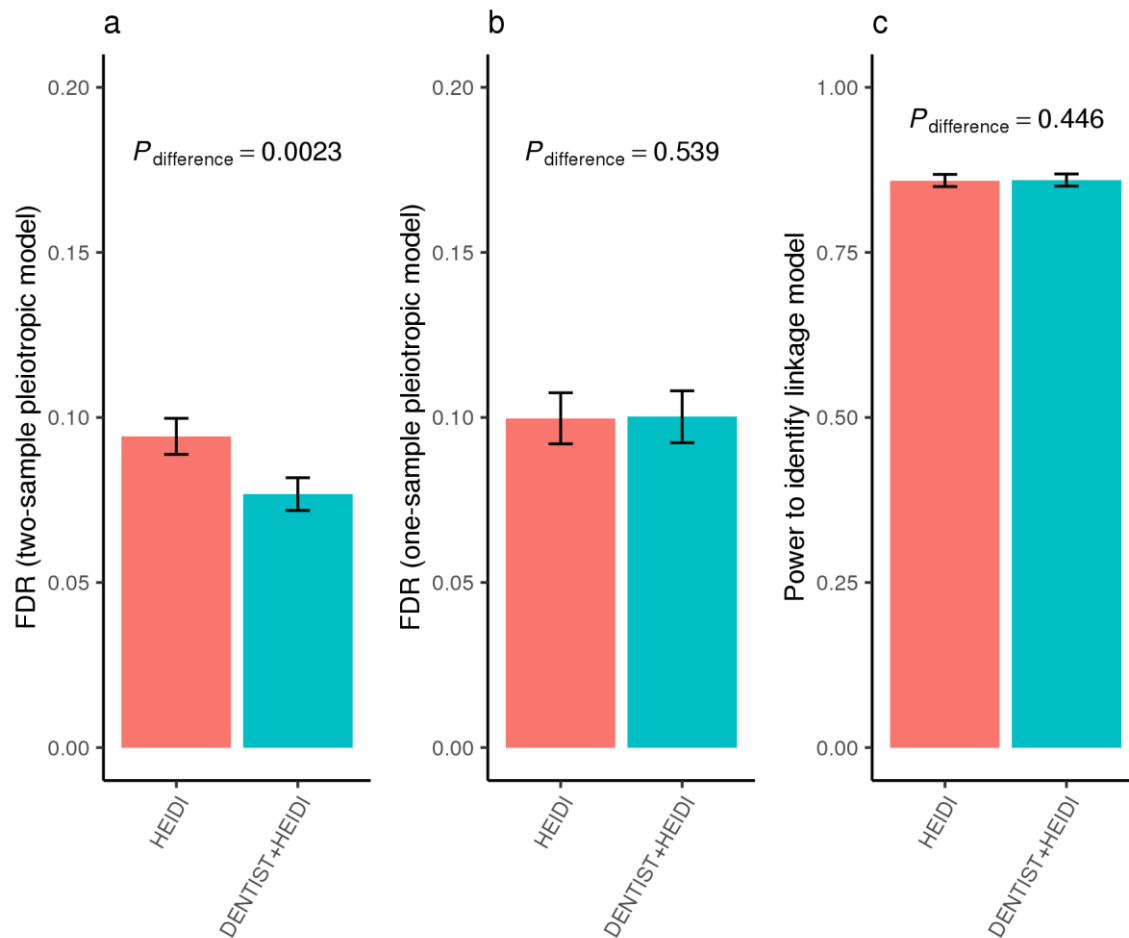
688

689 **Figure 1. Detecting simulated allelic and genotyping errors using DENTIST.** We assessed  
 690 the power of DENTIST in detecting allelic and genotyping errors. There are three levels of  
 691 genotyping error rate ( $f_{\text{error}}=0.15, 0.1$  or  $0.05$ ), defined as the proportion of individuals with  
 692 erroneous genotypes for a variant. Panel a) is a plot of the GWAS z-scores from data with  
 693 simulated errors against those from data without such errors. The gray dots in the diagonal  
 694 represents z-scores of variants without errors. In panel b), the DENTIST  $P$ -values are plotted  
 695 against the absolute values of the GWAS z-scores for all the variants. The horizontal dashed line  
 696 corresponds to  $P = 5 \times 10^{-8}$ . In panel c), to demonstrate the power in each  $|z\text{-score}|$  bin, we  
 697 pooled the results from 200 simulations. Each bar plot (+/- s.e.) represents power computed  
 698 from the pooled results.



699

700 **Figure 2. FPRs of COJO with and without DENTIST.** Based on simulations with one causal  
701 signal, we assessed the FPRs of COJO analyses when performed with and without DENTIST-  
702 based QC (FPR is defined as the frequency of observing more than one COJO signals in the  
703 scenario in which only one causal variant was simulated). The x-axis labels indicate the LD  
704 reference samples used in the COJO analyses, and those performed after DENTIST QC are labeled  
705 with “QC” in the parentheses. The error bars correspond to the standard error of FPRs calculated  
706 from 2200 replications, each with a re-sampled causal variant.



707

708 **Figure 3. FPRs and power of the HEIDI test with and without DENTIST.** Shown are the

709 results from simulations to quantify the FPR of HEIDI under a pleiotropic model (panels **a** and **b**)

710 and the power of HEIDI under a linkage model (panel **c**). The two-sample pleiotropic model in

711 panel **a** refers to the scenario where the eQTL and GWAS summary data were simulated based

712 on two different samples (HRS-3K and ARIC). The one-sample scenario in panel **b** refers to the

713 scenario where the GWAS and eQTL data were simulated using the same sample (HRS-3K). In

714 both scenarios, an independent sample (UK10K-1KGP) was used as the LD reference.  $P_{\text{difference}}$  is

715 to test if the FPR after QC is significantly different that without QC.  $P_{\text{difference}}$  is calculated from a

716 posterior distribution of  $k_1 \sim \text{Binomial}(n, p)$  with  $p$  from a prior distribution of  $p \sim \text{Beta}(k_2, n -$

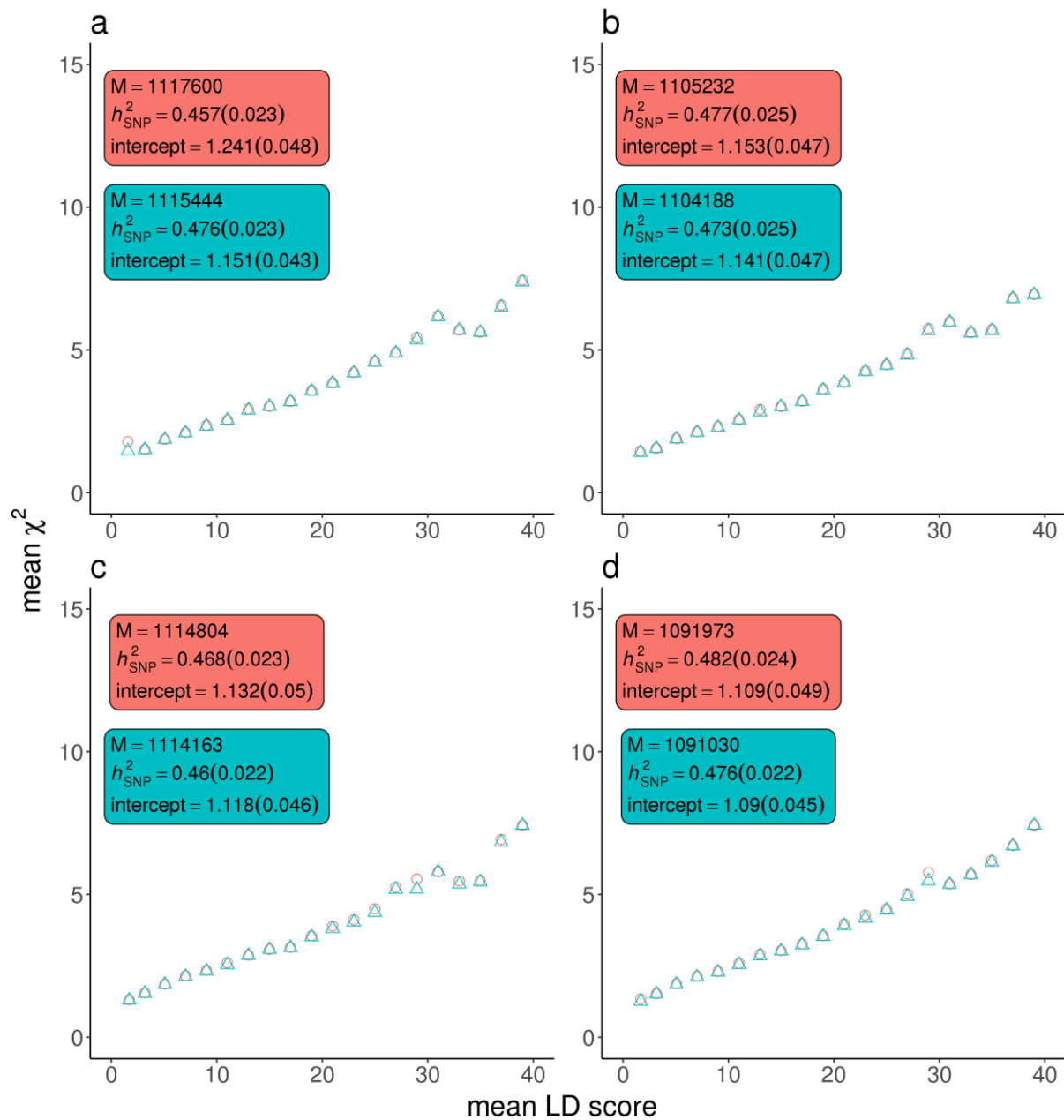
717  $k_2)$ , where  $n$  is the number of simulation replicates, and  $k_1$  and  $k_2$  are the numbers of simulation

718 replicates in which the HEIDI test correctly identified the right model with and without the

719 DENTIST-based QC, respectively. The error bars correspond to the standard errors of the

720 correspond metrics calculated from 4000 replications with re-sampled causal variants.

721



722

723

724

725

726

727

728

729

730

731

732

**Figure 4. The effect of DENTIST-based QC on LDSC analysis of the UKB height summary**

**data.** We assessed the effect of DENTIST on LDSC when different LD references were used, including a) HRS, b) ARIC, c) UKB-8K-1KGP, and d) UK10K-WGS. For each reference sample, LDSC was performed before and after DENTIST-based QC, and the corresponding results are shown in the red and cyan text boxes, respectively, on each plot. The variants are binned by their LD scores. Each dot on the plots represents the mean LD score value of each bin on the x-axis and the mean  $\chi^2$  value on the y-axis, with those before and after DENTIST-based QC in red and cyan colors respectively. In the textbox, “M” represents the number of variants, “ $h^2_{SNP}$ ” represents the estimate of SNP-based heritability, and “intercept” represents the LDSC intercept, with the corresponding standard errors given in the parentheses.

733 **Table 1.** FPRs and power of COJO before and after DENTIST-based QC in simulations.

Analysis method (LD reference)	FPR (%)		Power (%)	
	Common-causal	Rare-causal	Common-causal	Rare-causal
Benchmark	0.1 ± 0.11	0.2 ± 0.16	78.8 ± 0.9	30.6±1.4
COJO without DENTIST (UKB-8K-1KGP)	1.0 ± 0.50	2.7 ± 0.50	79.0 ± 0.9	28.6±1.9
COJO with DENTIST (UKB-8K-1KGP)	0.3 ± 0.14	1.3 ± 0.36	77.3 ± 0.9	22.2±1.6
COJO without DENTIST (HRS)	7.9 ± 1.37	28.4 ± 1.37	81.8±3.9	27.7±1.4
COJO with DENTIST (HRS)	0.5 ± 0.18	1.7 ± 0.40	81.2±1.3	16.8±1.1
COJO without DENTIST (ARIC)	7.1 ± 1.32	28.7 ± 1.38	77.6±0.9	26.7±1.7
COJO with DENTIST (ARIC)	0 ± 0.00	1.3 ± 0.36	75.8 ± 0.9	17.1±1.4

734 Benchmark: COJO analysis using the discovery GWAS as the reference without DENTIST. Shown  
735 are mean ± standard error.

736

737



738 **Table 2.** Numbers of COJO signals from analyses of the UKB height summary data using different  
 739 LD reference samples with and without DENTIST-based QC.

	<b>Benchmark</b>	<b>UKBv3-20K (n = 20,000)</b>	<b>UKBv3-8K (n = 8,000)</b>	<b>HRS (n = 8,557)</b>	<b>ARIC (n = 7,703)</b>	<b>UK10K-WGS (n = 3,642)</b>
MAF > 0.01 Without DENTIST	1279	1296 (1.3%)	1337 (4.5%)	1477 (15.5%)	1485 (16.1%)	1417 (10.8%)
MAF > 0.01 With DENTIST	/	1300 (1.6%)	1319 (3.1%)	1338 (4.6%)	1353 (5.8%)	1413 (10.5%)
MAF > 0.001 Without DENTIST	1310	1313 (0.2%)	1337 (2.0%)	1555 (18.7%)	1645 (25.6%)	1473 (12.4%)
MAF > 0.001 With DENTIST	/	1326 (1.2%)	1326 (1.2%)	1346 (2.7%)	1398 (6.7%)	1421 (8.5%)

740 Benchmark: COJO analysis using the discovery GWAS (UKBv3-329K) as the reference without  
 741 DENTIST. The inflation rate as compared to the benchmark is shown in the parentheses.

742

743 **Table 3.** Estimates from LDSC analyses of the UKB height GWAS summary data with and without  
 744 DENTIST-based QC.

	Number of variants	One-step approach		Two-step approach	
		$h_{SNP}^2$	Intercept	$h_{SNP}^2$	Intercept
<b>Reference = the discovery GWAS sample</b>					
Benchmark	1114780	0.46 (0.023)	1.13 (0.049)	0.41(0.017)	1.34 (0.030)
<b>Reference = HRS</b>					
Without DENTIST	1117600	0.45 (0.022)	1.24 (0.047)	0.42 (0.018)	1.36 (0.028)
With DENTIST	1116249	0.46 (0.022)	1.15 (0.040)	0.41 (0.017)	1.35 (0.027)
<b>Reference = ARIC</b>					
Without DENTIST	1105232	0.47 (0.025)	1.15 (0.047)	0.42 (0.018)	1.34 (0.030)
With DENTIST	1102229	0.47 (0.024)	1.14 (0.047)	0.41 (0.018)	1.34 (0.030)
<b>Reference = UKB-8K-1KGP</b>					
Without DENTIST	1114804	0.46(0.023)	1.13(0.050)	0.41 (0.017)	1.33 (0.030)
With DENTIST	1113260	0.46(0.022)	1.12(0.048)	0.40 (0.016)	1.33 (0.030)
<b>Reference = UK10K-WGS</b>					
Without DENTIST	1091973	0.48 (0.023)	1.10(0.048)	0.42 (0.018)	1.31 (0.029)
With DENTIST	1074399	0.47 (0.022)	1.07 (0.042)	0.41 (0.016)	1.30(0.028)
<b>Reference = 1KGP-EUR</b>					
Without DENTIST	1133151	0.48 (0.024)	1.15 (0.047)	0.43 (0.018)	1.33 (0.028)
With DENTIST	1071447	0.40 (0.018)	1.03 (0.030)	0.35 (0.014)	1.22 (0.024)

745 Standard errors are given in the parentheses.