

The power of pathway-based polygenic risk scores

Shing Wan Choi

Icahn School of Medicine at Mount Sinai

Judit Garcia-Gonzalez

Icahn School of Medicine at Mount Sinai

Yunfeng Ruan

Broad Institute

Hei Man Wu

<https://orcid.org/0000-0003-1559-7586>

Jessica Johnson

Clive Hoggart

Icahn School of Medicine at Mount Sinai

Paul O'Reilly (✉ paul.oreilly@mssm.edu)

Icahn School of Medicine at Mount Sinai <https://orcid.org/0000-0001-7515-0845>

Biological Sciences - Article

Keywords: Genetic Liability, Functional Genome Sub-structure, GWAS Signal Enrichment, Heterogeneity

Posted Date: June 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-643696/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

The power of pathway-based polygenic risk scores

Shing Wan Choi^{1*}, Judit García-González^{1*}, Yunfeng Ruan², Hei Man Wu¹, Jessica Johnson¹, Clive J Hoggart¹, Paul F. O'Reilly¹

¹ Department of Genetics and Genomic Sciences, Icahn School of Medicine, Mount Sinai, New York City, NY 10029, USA

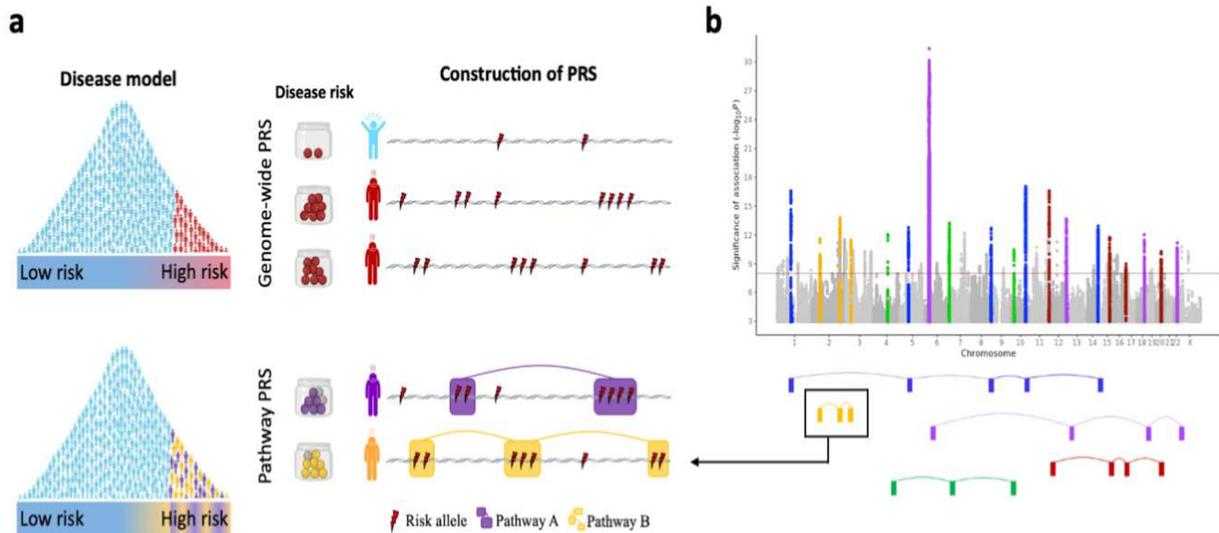
² The Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

* These authors contributed equally to this work

Polygenic risk scores (PRSs) have been among the leading advances in biomedicine in recent years. As a proxy of genetic liability¹, PRSs are utilised across multiple fields and applications². While numerous statistical and machine learning methods have been developed to optimise their predictive accuracy³⁻⁵, all of these distil genetic liability to a single number based on aggregation of an individual's genome-wide alleles. This results in a key loss of information about an individual's genetic profile, which could be critical given the functional sub-structure of the genome and the heterogeneity of complex disease. Here we evaluate the performance of pathway-based PRSs, in which polygenic scores are calculated across genomic pathways for each individual, and we introduce a software, PRSet, for computing and analysing pathway PRSs. We find that pathway PRSs have similar power for evaluating pathway enrichment of GWAS signal as the leading methods^{6,7}, with the distinct advantage of providing estimates of pathway genetic liability at the individual-level. Exemplifying their utility, we demonstrate that pathway PRSs can stratify diseases into subtypes in the UK Biobank with substantially greater power than genome-wide PRSs. Compared to genome-wide PRSs, we expect pathway-based PRSs to offer greater insights into the heterogeneity of complex disease and treatment response, generate more biologically tractable therapeutic targets, and provide a more powerful path to precision medicine.

29 As proxies for genetic liability to human traits or diseases ¹, polygenic risk scores (PRSs) have
30 been applied in numerous applications, including prediction of disease risk ^{8–13}, patient
31 stratification ¹⁴, investigation of treatment response ^{15–18} and genetically-informed experimental
32 perturbation ^{19,20}. Most leading PRS methods, including those that incorporate functional
33 annotation ^{21,22}, are based on the classical polygenic model of disease, which assumes that
34 individuals lie on a linear spectrum from low to high genetic risk and that summarises an
35 individual's genetic profile to a single value estimate of liability. This incurs substantial loss of
36 information about an individual's genetic profile, such as how the burden of genetic risk varies
37 across biological processes and pathways. This information may be more informative for many
38 applications of PRS, such as patient stratification and prediction of treatment response.

39 Here we evaluate a new polygenic risk score approach that accounts for genomic sub-structure,
40 constitutes an extension to the classic polygenic model of disease, and may better reflect disease
41 heterogeneity (**Fig. 1a**). Instead of aggregating the estimated effects of risk alleles across the
42 entire genome, pathway-based PRSs aggregate risk alleles across k pathways (or gene sets)
43 separately. Therefore, rather than a single genome-wide PRS, each individual has k PRSs
44 corresponding to k pathways across the genome. Pathways can be defined in a variety of ways,
45 including by existing canonical databases (e.g. KEGG, REACTOME ^{23,24}), by analytically derived
46 modules of e.g. gene co-expression or protein-protein interactions, or from functional output of
47 experimental perturbations ^{25–27}. Ideally, pathways reflect the encoding of different biological
48 functions, separable in the same way that different environmental risk factors, such as smoking
49 or dietary factors, are considered separately in epidemiological prediction models. From this
50 perspective, GWAS results can be considered a composite of signal corresponding to function
51 encoded by different genomic pathways (**Fig. 1b**).



52

53 **Fig. 1. The pathway polygenic risk score approach.** Coloured boxes represent genes, lines link genes
 54 that are within the same genomic pathway. **a**, Upper model: Classical polygenic model of disease, in which
 55 individuals lie on a linear spectrum from low to high risk and genome-wide PRSs are constructed as the
 56 sum of risk alleles across the genome. Disease risk depicted by the Jar model ²⁸. Lower model: Pathway
 57 polygenic model of disease, in which multiple liabilities are considered and PRSs are constructed by
 58 aggregating risk alleles over different genomic pathways. **b**, GWAS results Manhattan plot illustrated as a
 59 hypothetical composite of signals that each correspond to alternative functional routes to disease. Pathways
 60 that only make a small contribution to disease risk across the population, or a contribution in a small fraction
 61 of individuals (e.g. nicotine receptor pathway in those individuals who smoke), are likely to harbour risk
 62 variants of relatively small effect.

63

64 A key concern in application of PRS computed over relatively short genomic regions is whether
 65 they are sufficiently powered to be useful. Here we show, for the first time, that polygenic risk
 66 scores capture genetic signal at the pathway-level with comparable power to that of leading
 67 pathway enrichment methods MAGMA ⁶ and LD score regression (LDSC) when applied to large
 68 target sample sizes ^{29,7}. Therefore, pathway PRS may be powered for a range of other

69 applications for which genome-wide PRS are presently used. To test this premise using real data,
70 we perform a head-to-head performance comparison of pathway-specific versus genome-wide
71 PRS for patient stratification, utilising the large number of individuals in the UK Biobank with
72 common diseases such as type 2 diabetes (T2D), coronary artery disease (CAD),
73 hypercholesterolemia (HC) and their comorbidities. We show that our pathway PRS method
74 outperforms both genome-wide PRS alternatives, lassosum and PRSice, in 22 out of 23
75 supervised stratification scenarios tested, often by a wide margin. We expect the power of
76 pathway PRS to improve substantially in the future with improved functional annotation and more
77 context-specific definition of pathways. Our new method and accompanying software (called
78 PRSet) builds on the popular PRSice genome-wide PRS tool ^{3,30} and is likewise user-friendly, fast,
79 intuitive and openly available.

80

81 **PRSet model overview**

82 Our PRSet method for calculating pathway-based PRSs leverages the classical genome-wide
83 PRS method¹ - clumping + thresholding (C+T) - to calculate k PRSs corresponding to k genomic
84 pathways for an individual i , as follows:

$$85 \quad PRS_{ik} = \sum_{j=1}^{m_k} \beta_j G_{ij}$$

86 where m_k is the number of clumped SNPs in pathway k , β_j is the SNP effect size estimated from
87 a GWAS on the studied phenotype, G_{ij} is the genotype of individual i . In contrast to the C+T
88 method, where SNPs are clumped across the whole genome, PRSet performs clumping on each
89 pathway independently. This is done to retain pathway signal and to account for correlation
90 between SNPs in nearby genes of the same pathway (**Supplementary Fig. 1**). The use of the P -
91 value thresholding procedure is dependent on the use-case. For example, while P -value

92 thresholding is not performed in pathway enrichment analyses, it is performed in each pathway
93 independently in the disease sub-typing application of this study (see **Methods**). Each pathway
94 PRS can be tested for association with a phenotype of interest in a target sample by regressing
95 the phenotype on the PRS, as in standard PRS analyses. Pathway enrichment is evaluated by
96 computing an empirical “competitive” *P*-value, which accounts for pathway size using permutation
97 (see **Methods**). Many applications of standard genome-wide PRS can be adapted to pathway
98 PRS, the analyses of which can be evaluated and reported similarly.

99 In the next section, we benchmark the power of pathway PRS for assessing pathway enrichment,
100 versus MAGMA and LDSC, using a range of comparisons that can be separated into (i) those
101 that use canonical pathways, and (ii) those that define pathways by tissue and cell-type specific
102 gene expression. These analyses are performed only to benchmark the methods’ relative power
103 as pathway enrichment tools, in particular to assess how well pathway PRS capture GWAS signal
104 and, thus, their potential for wider use, and are not optimized to produce novel findings, e.g. for
105 cell-type specificity.

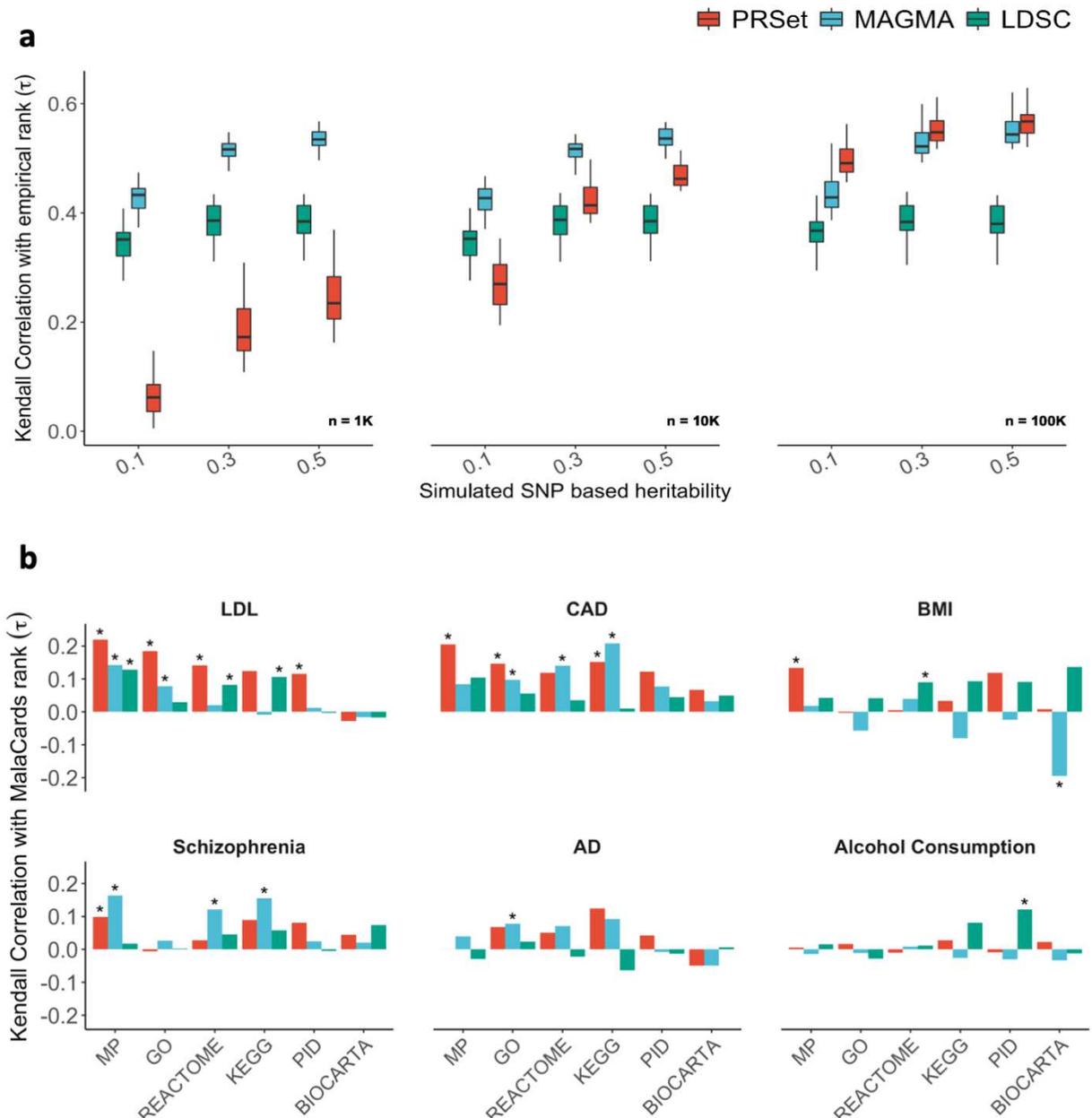
106

107 **Evaluating pathway enrichment**

108 **Canonical pathways:** In this sub-section, pathways are defined using six publicly available
109 databases (Biocarta ³¹, Pathway Interaction Database ³², Reactome ²⁴, Mouse Genome Database
110 ³³, KEGG ²³ and GO ^{34,35}) and pathway enrichment of genetic signal is tested by (i) a simulation
111 study, (ii) real data using *MalaCards* gene scores (**Methods**). First, we performed simulations in
112 which quantitative traits of varying heritability ($h^2 = 0.1, 0.3, 0.5$) were simulated from real
113 genotype data of UK Biobank individuals, with causal variants randomly spread across 50 (real)
114 pathways that were randomly selected from the pathway databases (**Methods**). GWAS were then
115 performed on 250k individuals and their simulated traits, and an additional 1k, 10k and 100k

116 individuals were selected as target data. PRSet, MAGMA and LDSC were applied to the data
117 (**Methods**) to test for pathway enrichment and the pathways were subsequently ranked by their
118 inferred enrichment and compared to those of the known simulated enrichment via Kendall's
119 correlation. This process was repeated 20 times. **Figure 2a** displays the results, showing best
120 overall performance for MAGMA (Kendall $\tau^2 = 26.0\%$), then PRSet (Kendall $\tau^2 = 15.2\%$) and then
121 LDSC (Kendall $\tau^2 = 8.62\%$). All methods perform better with larger h^2 and target sample sizes,
122 with PRSet the best-performing method in large (100k) target data.

123 Next, we apply the three methods to the real data of the UK Biobank, and that of publicly available
124 GWAS, across six traits: low-density lipoproteins (LDL), CAD, schizophrenia, BMI, Alzheimer's
125 disease (proxy status) and alcohol consumption. Since the true GWAS pathway enrichment of
126 each pathway is unknown, we produce a *disease relevance score* for each pathway by summing
127 *MalaCards* gene scores (**Methods**), which assign values to genes based on systematic
128 phenotype-specific text-mining of the literature (note that most genes are assigned a value of 0).
129 In **Figure 2b**, we report the Kendall's correlations between the pathway enrichment ranked by the
130 three methods versus the *MalaCards disease relevance scores*. While performance varies widely
131 depending on pathway resource (**Fig. 2b**) and trait (**Supplementary Fig. 2**), the three methods
132 show broadly similar results, with PRSet having the highest mean correlation ($\tau = 0.078$) between
133 its pathway enrichment ranks and those of the *MalaCards* scores, followed by MAGMA ($\tau = 0.050$)
134 and LDSC ($\tau = 0.043$). We also repeated the analysis removing all genes with *MalaCards* scores
135 greater than 0 to examine evidence of pathway enrichment among genes not yet highlighted in
136 the literature and found that the correlations were eliminated (**Supplementary Figs. 3, 4**). This
137 may indicate that the methods have limited power to identify weak effects across pathways or that
138 only a modest fraction of genes in pathways influence disease contribution to risk.



139

140 **Fig. 2. Evaluating pathway enrichment using canonical pathways. a**, Simulation study. Performance
 141 was defined as the Kendall correlation between the pathway ranks based on competitive P -values of
 142 enrichment computed by each software and the empirical pathway ranks based on the true (simulated)
 143 effects across the pathways. Boxplots illustrate the values of Kendall rank correlation coefficients (τ) for
 144 PRSet, MAGMA and LDSC for each combination of heritability ($h^2 = 0.1, 0.3, 0.5$) and target sample size n
 145 $= (1K, 10K, 100K)$. **b**, Real data study. Kendall correlation coefficients (τ) between pathway ranks based on

146 competitive *P*-values of enrichment computed by each software and pathway ranks based on *MalaCards*
147 disease relevance scores. *Empirical *P*-value < 0.05. MP, Mouse Genome Database; GO, Gene Ontology
148 database; KEGG, Kyoto Encyclopaedia of Genes and Genomes; PID, Pathway Interaction Database; AD,
149 Alzheimer's disease (proxy status).

150

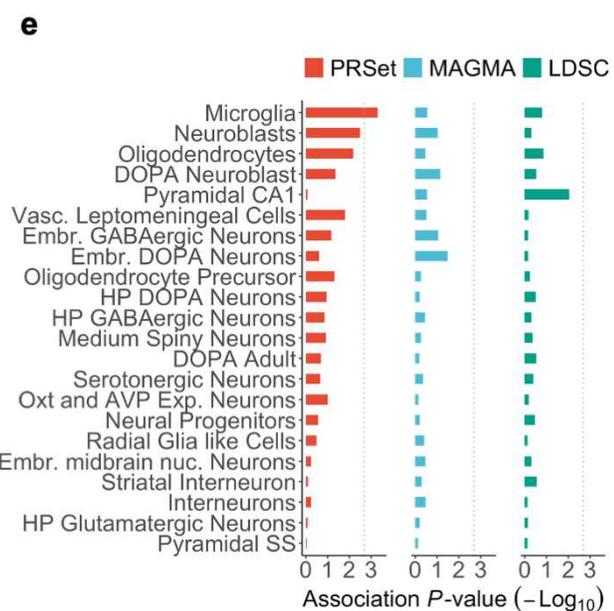
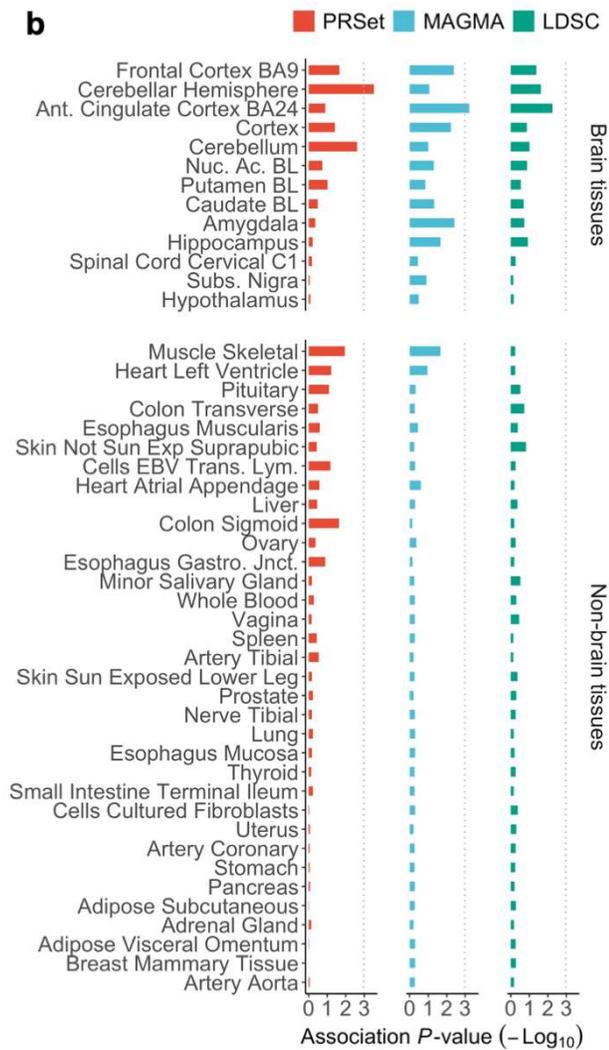
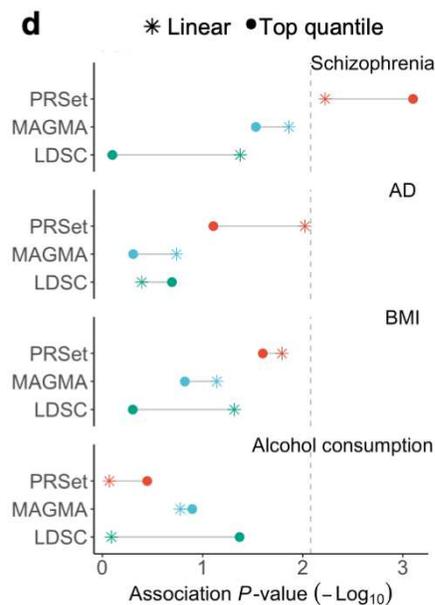
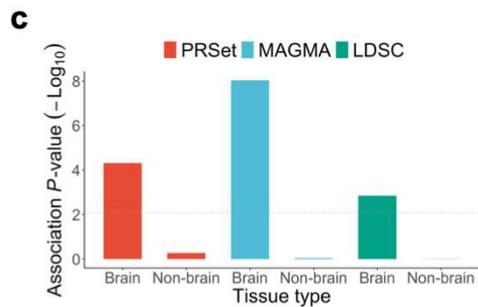
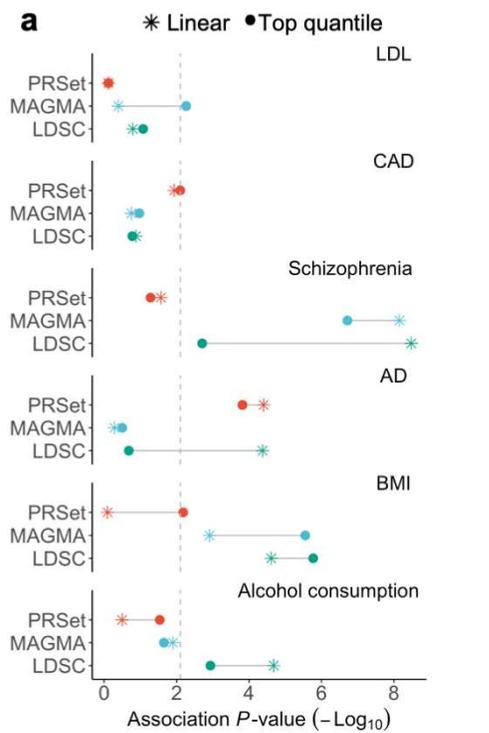
151 **Pathways defined using tissue/cell-type specificity:** To further interrogate the power of PRS
152 to capture genetic signal at the pathway-level compared to MAGMA and LDSC, we compared
153 their performance in tissue/cell-type expression specificity analyses using the approach
154 introduced in *Skene et al 2018*³⁶. In this approach, genes are grouped into 11 quantiles of
155 increasing expression specificity based on expression reported across 47 bulk-tissues and 24
156 brain cell types (**Methods**). Tissue and cell types with greatest GWAS signal – as evaluated by
157 MAGMA and LDSC (and here PRSet) – among quantiles with the most specific expression, are
158 thus implicated in disease etiology. The pathway enrichment of the genes in the top quantile is
159 evaluated as standard for canonical pathways, which we refer to as the 'top quantile' test model,
160 while the linear trend of enrichment is also assessed and referred to as the 'linear' test model (see
161 **Extended Data Fig. 1, Methods**).

162 Here we perform these analyses in the same data and traits used in the previous section. In the
163 absence of well-established roles for individual tissue and cell types in these outcomes, we sought
164 *a priori* candidates from two domain experts for each outcome to provide an agnostic way to
165 evaluate the performance of the different methods in this setting (**Methods**). We observed
166 significant associations between expert opinion and the tissue-type specificity results, although
167 results varied substantially depending on the pathway method and test model used (**Fig. 3a**).
168 Associations were strongest for schizophrenia (**Fig. 3b-c, Extended Data Fig. 2**) and BMI
169 (**Extended Data Fig. 3**), in which MAGMA and LDSC had a higher correlation with expert's
170 opinion than PRSet, while PRSet had best performance in Alzheimer's disease (**Extended Data**

171 **Fig. 4)** and CAD (**Extended Data Fig. 5**). The associations relating to the cell-type specific
172 analyses were relatively weak (**Fig. 3d**), with significant results only observed for MAGMA and
173 PRSet in relation to schizophrenia. For AD, the strongest and only significant result was that of
174 PRSet implicating microglia using the top quantile test model (**Fig. 3e**), which is notable since
175 microglia has been extensively linked to AD etiology in the literature³⁷. However, individual results
176 reported here should be treated with caution, since they appear highly sensitive to the test model
177 (top quantile / linear) and the number of quantiles used (**Extended Data Fig. 2-7**). Moreover,
178 there have been several extensions of the *Skene et al* approach, including an extension of
179 MAGMA designed specifically for tissue/cell-type analyses that likely has substantially higher
180 power than the standard MAGMA enrichment tool used here³⁸; the basic version of MAGMA as
181 an enrichment tool was used here to enable like-for-like comparisons with PRSet and LDSC
182 regarding power to capture pathway signal.

183 Our results benchmarking these pathway enrichment tools in multiple settings suggest that PRSet
184 has broadly comparable power to capture genetic signal in pathways as MAGMA and LDSC, with
185 the distinct advantage of providing individual-level estimates of pathway liability, which could be
186 useful in a wide-range of applications. Below, we test the power of pathway PRS for one such
187 application, that of disease stratification.

188



190 **Fig 3. Performance of PRSet, MAGMA and LDSC for ranking of pathways defined by tissue and cell-**
191 **type expression specificity. a,** Association between pathway enrichment *P*-value and expert opinion of
192 tissue relevance for each software and six diseases. Colours indicate the software used to calculate
193 enrichment: Red, *PRSet*; Blue, *MAGMA*; Green, *LDSC*. Results are shown for both the *top quantile* and
194 *linear* specificity test methods (**Methods**). Dashed line corresponds to Bonferroni significance threshold of
195 0.05 for 6 tests (3 methods x 2 test models). **b,** Pathway enrichment results for schizophrenia under the top
196 quantile test model. Bar plots show enrichment *P*-value for each tissue and pathway method. Dashed line
197 corresponds to Bonferroni significance threshold for 47 tissues ($-\text{Log}_{10}(0.05/47) = 2.97$). Ant, Anterior; Nuc.
198 Ac., Nucleus Accumbens; BL, Basal Ganglia; Subs, Substantia; Exp, Exposed; EBV Trans. Lym., Epstein-
199 Barr virus transformed lymphocytes; Gastro. Jnct, Gastroesophageal Junction. **c,** Enrichment of
200 schizophrenia signal is higher in brain tissues vs non-brain tissues under the top quantile test model. Bar
201 plots show the meta-analysis enrichment *P*-value using the Fisher's method for brain vs non-brain tissue
202 and method. Dashed line corresponds to Bonferroni significance threshold for the 6 tests conducted. **d,**
203 Associations between pathway enrichment *P*-value and expert opinion on cell type relevance for each
204 software and four diseases. Colours indicate the software used to calculate tissue type enrichment: Red,
205 *PRSet*; Blue, *MAGMA*; Green, *LDSC*. Dashed line indicates Bonferroni significance threshold of 0.05 for 6
206 tests (3 methods x 2 models). **e,** Pathway enrichment results for AD under the top quantile test model. Bar
207 plots show enrichment *P*-values for each cell type and method. Dashed line corresponds to Bonferroni
208 significance threshold for 22 cell types. DOPA, Dopaminergic, Vasc, Vascular; Emb, Embryonic; HP,
209 hypothalamic; Oxt and AVP Exp Neurons, Oxytocin and Vasopressin Expressing Neurons; Nuc, Nucleus;
210 SS, Somatostatin.

211

212

213

214 **Pathway PRSs for disease stratification**

215 To evaluate the power of pathway PRSs, versus genome-wide PRSs, for performing disease
216 stratification, or *sub-typing*, we leveraged the large number of individuals in the UK Biobank with
217 major diseases: T2D (N = 19,668), CAD (N = 22,388), HC (N = 26,561), obesity (N = 92,818),
218 and additionally “extreme height” (top 5th percentile; N = 20,512) because height GWASs are
219 presently the most powerful for PRS analyses (**Methods**). We compared the power of PRSet for
220 sub-typing with two genome-wide PRS approaches, lassosum and PRSice (C+T), using two
221 strategies. First, in the absence of well-established subtypes for these outcomes we produced
222 “pseudo subtypes” by combining the 5 outcomes into pairs. We then meta-analysed the two
223 GWAS of each pair for use as base data and applied each of the PRS methods to stratify the set
224 of pooled UK Biobank patients (pooled across each pair, removing comorbid individuals),
225 comparing inferred subtypes to the original patient groups (**Methods**). Supervised stratification
226 was performed with 5-fold cross-validation, using lasso regression applied to the most enriched k
227 pathway PRS predictors from PRSet (competitive P -value < 0.05), or standard regression with
228 genome-wide PRS predictors generated by lassosum or PRSice (**Methods**).

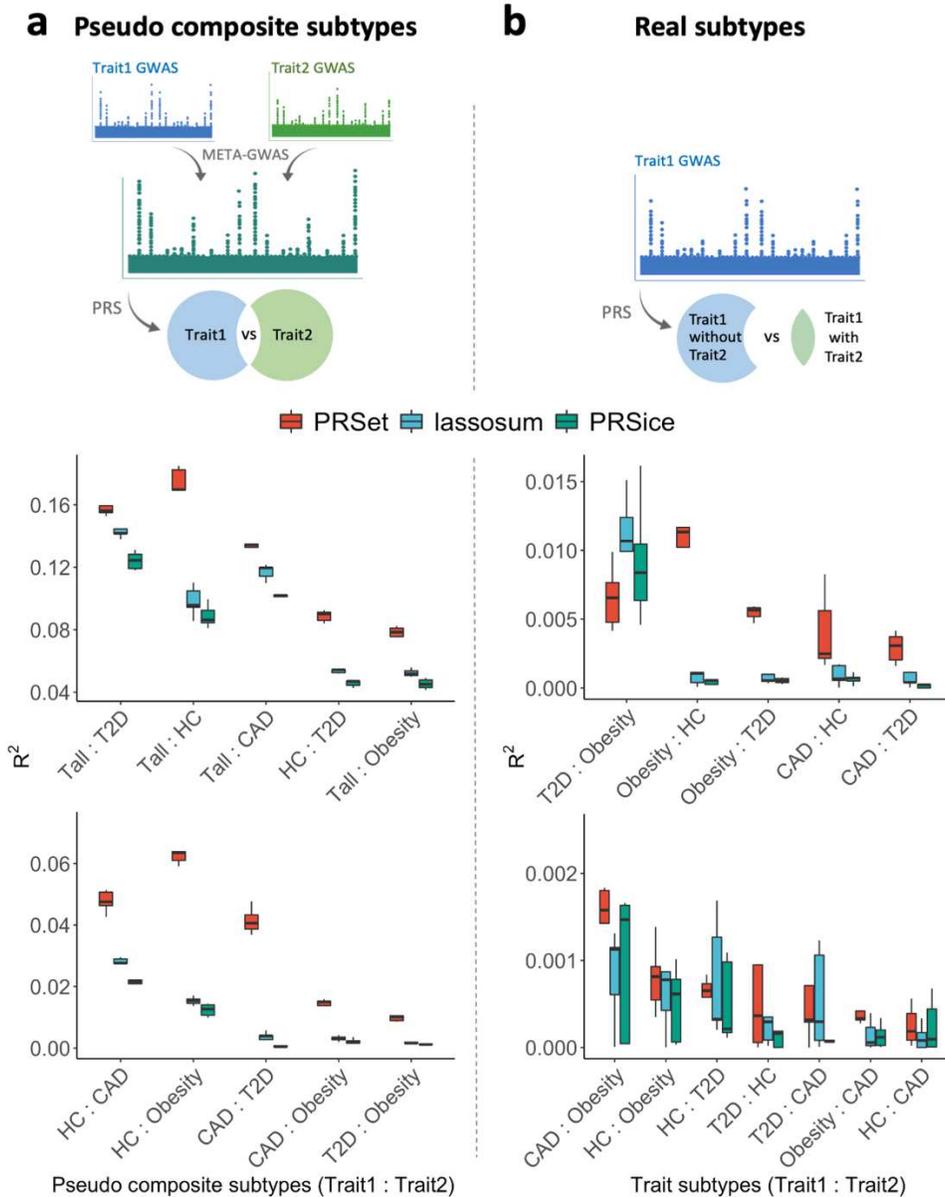
229 PRSet had the highest subtype classification performance across all ten comparisons (**Fig. 4a**).
230 Most of these showed strikingly higher discriminatory power for PRSet, notable since different
231 PRS methods typically have similar predictive power^{39–41}. We performed several sensitivity
232 analyses to explore the causes of the relatively high power of PRSet here, repeating the analyses:
233 i) with genes/pathways shifted by 5Mb to reduce their biological meaning (**Supplementary Note**
234 **2**), ii) using single diseases and testing case/control status instead (**Supplementary Note 3**), (iii)
235 extending lassosum and PRSice to also utilize a lasso regression framework in their optimisation
236 (**Supplementary Note 4**). Lower relative predictive power was observed for PRSet in each of
237 these scenarios, indicating that the improvement in PRSet performance is partially due to (i) the

238 enrichment of signal in genes and/or pathways, (ii) the application to sub-typing, rather than
239 case/control prediction, and (iii) the use of a regularisation framework (see **Discussion**).

240 Next, we defined subtypes of CAD, HC, T2D and obesity as the presence/absence of comorbidity
241 within each pair of these four diseases (**Fig. 4b**). PRSs were calculated based on one of the
242 disease GWASs before performing sub-typing. For example, PRSs based on CAD GWAS were
243 used to discriminate between CAD patients with or without T2D, with or without HC, and with or
244 without obesity. Classification performance estimates for these “real subtype” analyses were
245 lower than for the pseudo subtypes, with R^2 estimates < 0.016 (**Fig. 4b**). In 11 of the 12
246 comparisons PRSet outperformed lassosum and PRSice, while for T2D with/without obesity,
247 lassosum and PRSice outperformed PRSet.

248 Finally, we tested stratification performance in a disease with a relatively small number of patients
249 in the UK Biobank but with well-established subtypes: inflammatory bowel disease (IBD) (N =
250 6,008) and its subtypes Crohn’s disease (N = 2,101) and ulcerative colitis (N = 3,681). Supervised
251 sub-typing was performed as above, with PRSet again showing higher discriminatory power (R^2
252 = 0.007, P -value = 1.2×10^{-6}) than PRSice ($R^2 = 0.005$, P -value = 0.01). Given the availability of a
253 known number of subtypes for validation, we also performed unsupervised stratification using K -
254 means clustering ($K = 2$; **Methods**) and found that neither PRSet (P -value = 0.69) or PRSice (P -
255 value = 0.98) produced clusters that were associated with the two true subtypes.

256



257

258 **Figure 4. Performance of PRSs to classify trait subtypes.** Boxplots illustrate model R^2 from a 5-fold
 259 cross validation approach. Pathway specific PRSs for 4,079 pathways were calculated using *PRSet*.
 260 Genome-wide PRSs were calculated using *lassosum* and *PRSice*. **a**, Pseudo composite subtype analysis.
 261 PRSs were calculated using effect sizes derived from meta-analysis of Trait1 and Trait2 GWASs. The target
 262 sample was defined as pooled cases of Trait1 and Trait2, removing comorbid individuals. **b**, Real subtype
 263 analysis. PRSs were calculated based on the GWAS of Trait1. The target sample was defined as cases of
 264 Trait1 with/without comorbidity with Trait2.

265 **Discussion**

266 Here we have introduced a novel, pathway-based, polygenic risk score approach and software
267 tool, PRSet, for performing pathway PRS analyses. We demonstrated that pathway PRSs can
268 capture genetic signal across pathways with similar power to MAGMA and LDSC, with the distinct
269 advantage of providing individual-level estimates of pathway liability. We do not presently
270 recommend PRSet as an enrichment tool over these established methods, given its lower power
271 under simulation in small target sample sizes (**Fig. 2a**). However, its capacity to capture genetic
272 signal over pathways highlights the promise of pathway PRSs as higher-resolution, more
273 biologically interpretable, alternatives to genome-wide PRSs.

274 Next, we directly compared the performance of the two PRS approaches in an application for
275 which much hope is placed in genome-wide PRSs: disease stratification. In a series of supervised
276 stratification analyses we found that PRSet outperformed genome-wide PRS methods, lassosum
277 and PRSice, in almost all scenarios (22 out of 23) and often by striking margins. While we
278 observed no significant unsupervised stratification of known IBD subtypes, we expect higher
279 power when the corresponding pathway PRSs are trained in large case/control data. In contrast,
280 genome-wide PRSs may be limited-by-design in their application to disease stratification because
281 they are dominated by variants that affect multiple disease subtypes, while pathway PRSs allow
282 for the independent influence of pathways that differentiate between subtypes.

283 Pathway PRSs, as outlined here, have two major limitations: (i) pathways are not well-defined
284 and so are likely a poor proxy of biological function, (ii) variants are linked to pathways by location
285 (i.e. inside genes) only (**Methods**). However, the rapid advances being made in functional
286 genomics, via the integration of rich resources of multi-omics data now available, can help to
287 address both issues. For example, future pathway PRSs could be enhanced so that pathways
288 are defined according to robust differential gene co-expression or protein-protein interaction

289 networks and annotated by regulatory elements active in tissue and cell types relevant to the
290 disease under study.

291 We believe that pathway-based PRSs may ultimately offer greater promise of delivering stratified
292 medicine for complex diseases than genome-wide PRSs, which aggregate disparate forms of risk
293 into a single number. However, despite promising early results for pathway PRSs reported here,
294 they have several limitations that need addressing, some of which rely on field-level advances,
295 before their potential can be realised. A better understanding of how genetics leads to biological
296 function, and the role of pivotal genes in signalling and mechanistic cascades, will contribute to
297 more accurate and powerful modelling of the multiple genetic liabilities underlying complex
298 disease.

299 Our new method, PRSet, provides a novel approach to computing and analysing polygenic risk
300 scores, motivated by the functional sub-structure of the genome and the heterogeneity of disease.
301 Unlike genome-wide PRSs, pathway-based PRSs provide high-resolution information about an
302 individual's genetic risk profile aligned to biological function, and thus have the potential to offer
303 greater insights into disease and a more direct route to precision medicine.

304

305

306

307 **Methods**

308 **Participants**

309 *UK Biobank*

310 The UK Biobank (UKB) is a prospective multi-ethnic cohort of 502,493 participants, aged 40-69
311 years, initially recruited across the United Kingdom between 2006 and 2010, with follow up since.
312 The UK Biobank genetic data used in this study included 488,377 samples and 805,426 SNPs.

313 Standard quality controls were performed, removing SNPs with genotype missingness > 0.02 ,
314 minor allele frequency (MAF) < 0.01 and with Hardy Weinberg Equilibrium (HWE) P -value $< 1 \times 10^{-8}$.
315 We removed all individuals who had withdrawn consent, who had a high degree of missingness
316 or heterozygosity and who had mismatching genetically inferred and self-reported sex as reported
317 by the UK biobank data processing team. We also removed individuals who were not of European
318 ancestry based on a 4-mean clustering on the first two principal components, and related samples
319 with kinship coefficient > 0.044 using a greedy algorithm, since present PRS methods have been
320 shown to have relatively poor portability between global ancestries. A total of 387,392 individuals
321 and 557,369 SNPs remained after quality control.

322 *Sweden-Schizophrenia Population-Based cohort*

323 Samples from the Sweden-Schizophrenia Population-Based cohort are a subset of the samples
324 of the Psychiatric Genomics Consortium (PGC) Schizophrenia Working Group and were obtained
325 from the database of Genotypes and Phenotypes (Study Accession: phs000473.v2.p2). Data
326 processing and quality controls performed on these data are described elsewhere⁴². A total of
327 4,834 individuals diagnosed with schizophrenia and 6,128 controls were included.

328 **Phenotypes**

329 In order to optimise statistical power for benchmarking the performance of the methods tested in
330 the study, we selected complex phenotypes with high SNP-heritability estimates, with publicly
331 available summary statistics from large GWASs and that were measured in the UK Biobank or
332 the Sweden-Schizophrenia Population-Based cohort (**Extended Data Table 1 and 2**). As such,
333 we extracted data from the UK Biobank on the following phenotypes: BMI, LDL, CAD, alcohol
334 consumption, standing height, T2D, and a proxy of AD based on parental history of the disease
335 (**Supplementary Information**). Schizophrenia cases and controls were extracted from the
336 Sweden-Schizophrenia Population-Based cohort.

337 **GWAS Data Sets**

338 GWAS data sets were downloaded on BMI ⁴³, LDL ⁴⁴, AD ⁴⁵, CAD ⁴⁶, T2D ⁴⁷, alcohol consumption
339 ⁴⁸ and standing height ⁴⁹ from public online databases and used without modification. Due to
340 sample overlap between UKB and the Sweden-Schizophrenia Population-Based cohort, we used
341 a version of the Sweden-Schizophrenia GWAS ⁴² with UKB data removed to prevent inflation of
342 results.

343 **Definition of pathways**

344 KEGG ²³, BioCarta ⁵⁰, Pathway Interaction Database (PID) ⁵¹ and Reactome ⁵² canonical
345 pathways were obtained from the Molecular Signatures Database (MSigDB v7.0) ⁵³. Pathways
346 from the Gene Ontology database (GO, accessed on 2021-03-17) ^{34,35} and Mouse Genome
347 Database (MGD, accessed on 2021-03-17) ⁵⁴ were also included. For MGD pathways, we i) used
348 the human-mouse homolog list provided by MGD to convert the mouse gene names to their
349 human counterpart and ii) restricted our analyses to pathways with ontology level > 4 to avoid
350 inclusion of pathways that are extremely specific. To avoid performing analyses on under-
351 powered or overly genetic pathways, as standard in pathway analyses ⁶, we removed pathways

352 with fewer than 10 genes or more than 2000 genes to exclude over specific or too broad pathways.
353 A total of 4,079 pathways across the six pathway database resources were included in the
354 analyses.

355 **Pathway enrichment analyses**

356 *PRSet*

357 Pathway specific PRS analyses were performed using PRSice-2 (v2.3.5) on genotype data. The
358 Major histocompatibility complex region (MHC, chr6:25Mb-34Mb) was removed for all the
359 diseases assessed and the APOE region (chr19:44Mb-46Mb) was also removed for AD. SNPs
360 were annotated to genes and pathways based on GTF files obtained from ENSEMBL
361 (GRCh37.75). We extended the gene coordinates 35 kilobases (kb) upstream and 10 kb
362 downstream of each gene to include potential regulatory elements. Ambiguous SNPs (A/T and
363 G/C) and SNPs not present in both GWAS summary statistics and genotype data were excluded.
364 10,000 permutations were performed to obtain empirical “competitive” *P*-values, which account
365 for pathway size, as follows: first, a “background” pathway containing all genic SNPs is
366 constructed, and clumping is performed within this pathway. For pathways with *m* SNPs, *N* null
367 pathways are generated by randomly selecting *m* “independent” SNPs from the “background”
368 pathway. The competitive *P*-value can then be calculated as

$$369 \quad \text{competitive } P - \text{value} = \frac{\sum_{n=1}^N I(P_n < P_o) + 1}{N + 1}$$

370 where *I*(.) is an indicator function, taking a value of 1 if the association *P*-value of the observed
371 gene set (*P*_o) is larger than the one obtained from the *n*th null set (*P*_{*n*}), and 0 otherwise. A pseudo-
372 count of 1 is added to the numerator and denominator to avoid competitive *P*-values of 0 and
373 conservatively counting the observed gene set as 1 potential null set⁵⁵. One consideration of this

374 permutation procedure is that the smallest achievable competitive P -value is $1/(N + 1)$, which
375 can lead to difficulties in ranking highly significant gene sets.

376 *MAGMA*

377 To directly compare the performance of PRSet vs MAGMA (v1.07b) given identical input data, we
378 removed all ambiguous SNPs and non-overlapping SNPs prior to MAGMA analyses. It is
379 important to note that this step is unnecessary for MAGMA and might negatively impact its
380 performance. After filtering, gene-based analyses were performed on summary statistics and
381 genotype data independently. Same as for PRSet analyses, a 35kb window upstream and a 10kb
382 window downstream were added to gene coordinates, the MHC region was excluded for all traits,
383 and the APOE region was excluded for AD. Gene-based results were meta-analysed using the
384 inbuilt `--meta` function and were subsequently used as input to the pathway analysis.

385 *LDSC*

386 Partitioned LD scores were calculated using the 1000 Genomes European genotype data ⁵⁶.
387 Similar to PRSet, SNPs were annotated to genes and pathway with 35kb upstream and 10kb
388 downstream extension prior to calculation of LD scores. Ambiguous SNPs and non-overlapping
389 SNPs were removed prior to LDSC analyses to allow for direct comparison between PRSet and
390 LDSC. GWAS were performed on the target genotype data using PLINK v1.90b6.7 ⁵⁷, and were
391 meta-analysed with the external GWAS summary statistics using METAL (2011-03-25)⁵⁸.
392 Partitioned LD score regression was then performed using LDSC v1.01 ^{7,29}, with the MHC (all
393 traits) and APOE (AD only) regions excluded.

394 **Evaluating pathway enrichment using canonical pathway definitions**

395 **Assessment of pathway enrichment by simulation**

396 *Generation of causal pathways*

397 Out of 4,079 empirical pathways extracted from six publicly available collections (see “definition
398 of pathways” section), we randomly select 50 pathways. These pathways were simulated with
399 different levels of enrichment, ranging from 5-50% randomly assigned causal SNPs per pathway,
400 with step size of 5%. The empirical enrichment level for all pathways included in our studies were
401 calculated as the proportion of causal SNPs included. The simulation process was repeated 20
402 times and results over the 20 sets of simulations provided.

403 *Phenotype simulation and sample selection*

404 Simulation was performed using UKB genotype data to assess the performance of PRSet in
405 enrichment analyses comparing to MAGMA and LDSC. Quantitative traits (Y) with SNP-based
406 heritability (h^2) of 0.1, 0.3 or 0.5 were simulated as $Y=X\beta+\varepsilon$, where X is the standardized genotype
407 matrix, ε is the random error defined as $\varepsilon\sim N(\text{mean} = 0, \text{sd} = \sqrt{\text{var}(X\beta)(1 - h^2)})$, and β is a
408 vector of SNP effect sizes which follows a point-normal distribution $\beta\sim N(\text{mean} = 0, \text{sd} = \sqrt{h^2})$,
409 with non-causal SNPs assigned with $\beta = 0$.

410 For each trait, 250,000 individuals from European ancestry were randomly selected to generate
411 the GWAS summary statistics using PLINK v1.90.b6.7. An independent set of either 1k, 10k or
412 100k individuals were then randomly selected as the target samples. Pathway analyses were
413 performed as described in the previous sections.

414

415 **Assessment of pathway enrichment using *MalaCards* relevance scores**

416 To assess whether pathway enrichment results were in line with previous biological knowledge
417 on the phenotypes of interest, disease-associated relevance scores for each pathway were
418 constructed using information from the *MalaCards* database ⁵⁹. The *MalaCards* database
419 provides a disease relevance score *for each gene* based on experimental evidence and co-
420 citation in the literature. For the six diseases included in this analysis (schizophrenia, AD, alcohol
421 consumption, LDL, CAD and BMI), we downloaded the *MalaCards* disease-associated relevance
422 scores (Accessed on 2020-11-27, see **Extended Data Table 3** for disease terms used and
423 number of genes). Next, we performed a rank normalization of the scores where, assuming that
424 a disease has n genes with *MalaCards* scores, a score of $(r+1)/(n+1)$ were assigned to each gene,
425 with r being the inverse ranking of the gene with *MalaCards* score. Genes without a *MalaCards*
426 score are assigned a score of 0. *MalaCards* provide gene information as gene symbols, which
427 were transformed to ENSEMBL gene names.

428 To obtain disease-associated relevance scores *for each pathway*, we calculated the sum of the
429 rank transformed *MalaCards* scores for the genes included in a pathway and divided by the
430 number of genes in the pathway to account for pathway size.

431 Agreement between pathway enrichment results for PRSet, MAGMA and LDSC and the
432 *Malacards disease relevance scores* was assessed by calculating the Kendall correlations
433 between the $-\text{Log}_{10}$ competitive P -value generated by each pathway enrichment tool, and the
434 *MalaCards* relevance score for each pathway.

435 **Evaluating pathway enrichment using tissue/cell-type defined pathways**

436 *Calculation of tissue specificity from bulk-tissue RNA-sequencing data*

437 To calculate tissue specificity across pathways, we obtained bulk-tissue RNA-sequencing gene
438 expression data from 55 tissues from the GTEx consortium ⁶⁰(v8, median across samples).

439 Tissues with less than 100 individuals, cancer related tissue types (e.g. EBV-transformed
440 lymphocytes and Leukemia cell line), and testis (which were considered as an outlier ⁶¹) were
441 removed, retaining a total of 47 tissues. We filtered out all non-protein-coding genes and genes
442 not expressed in any tissue.

443 Gene expression specificity was calculated by dividing the expression of each gene by its total
444 expression across tissues ⁶¹. The resulting gene expression specificity ranged from 0 (gene is not
445 expressed) to 1 (gene is exclusively expressed in this tissue). Next, expression specificity of each
446 tissue was divided into 11 quantiles following the approach introduced in *Skene et al 2018* ³⁶,
447 where the first quantile contained all non-expressed genes in a given tissue, and the 11th quantile
448 contained the most specifically expressed genes. Genes within each quantile were grouped into
449 a single pathway.

450 *Defining the cell type specificity sets*

451 Cell type specificity data were obtained from supplementary materials of Skene et al (2018) ³⁶
452 which includes gene expression specificity information for 24 brain cell types obtained from single
453 cell RNA-sequencing data. Again, expression specificity of each brain cell type was divided into
454 11 quantiles with the first quantile containing all non-expressed genes in a given cell type. Genes
455 within each quantile were grouped into a single pathway.

456 *Ranking the importance of cell type / tissue*

457 To provide an objective estimate of tissue / cell-type importance for each phenotype, we invited
458 two experts (per phenotype) who were blind to our experiment and algorithm design to provide
459 their opinion on what cell-type(s) and tissue(s) are expected to be implicated for each disease
460 context (**Extended Data Table 4**). The expert response was coded as “none” (both experts think
461 tissue/cell type is not implicated), “single” (only one expert thinks a tissue/cell type is important)
462 and “both” (both experts agree about the importance of a tissue/cell type).

463 *Cell type and tissue specificity analyses*

464 We used two testing strategies to assess the relationship between disease GWAS signals and
465 tissue(s) / cell type(s) specificity. For the *Top quantile enrichment strategy* GWAS signals are
466 enriched in the most specifically expressed genes^{36,61,62}; whereas for the *Linear enrichment*
467 *strategy* GWAS signals increase linearly with expression specificity^{6,38}. The top quantile strategy
468 reports the competitive *P*-value of the pathway defined by those genes in the top expression
469 specificity quantile for each software and tissue/cell type. The linear enrichment strategy fits a
470 linear regression with the $-\text{Log}_{10}$ competitive *P*-value for each of the pathways defined by the
471 expression specificity quantiles as dependent variable, and the quantile ranks as the predictor
472 variable, and reports the one-sided *P*-value for a positive association.

473 The concurrence of the methods' ranking of the tissues / cell types with that of the experts within
474 each disease for both the top quantile and linear enrichment strategies was measured by
475 regressing the inverse normalized $-\text{Log}_{10}$ *P*-value for the top quantile / linear enrichment
476 strategies for each cell type / tissue against the expert opinion, coded as factor.

477 MAGMA has a specific model which accounts for expression specificity (`--gene-covar`). However,
478 in favour of a more consistent analysis between the three software methods, this model was not
479 used. It is thus possible that MAGMA can provide more powerful results using the dedicated
480 model.

481 Results from the regressions against the expert confidence score assessed the association of the
482 gene expression specificity and GWAS signal with the expert opinion for each pathway
483 enrichment software under each of the two hypotheses.

484 **Disease Stratification**

485 *Pseudo subtype analysis*

486 Composite phenotypes were generated by combining two distinct phenotypes, including T2D,
487 CAD, Obesity (defined as BMI > 30), extreme height (defined as top 5% each sex) and
488 hypercholesterolemia (defined as LDL > 4.9 mmol/L) from the UK Biobank. To simplify the
489 analysis, we removed controls and co-morbid cases so that the analysis was only performed in
490 participants presenting either case outcome. Phenotypes were encoded mimicking sub-
491 phenotypes of a given disease, for example, for the phenotype CAD-Obesity, samples with CAD
492 (and not Obesity) were coded as 0 and those with Obesity (and not CAD) were coded as 1.

493 *Comorbid subtype analysis*

494 For the analysis of subtypes with presence/absence of comorbid diseases, we used T2D, CAD,
495 Obesity and hypercholesterolemia, as these diseases present high comorbidity between them.
496 For each disease, three traits were derived mimicking the subtypes of the disease as the
497 presence/absence of the other disorders.

498 For computational expediency, sex, age, age of diagnosis (for CAD and T2D), genotyping batch,
499 recruitment centre and first 15 principal components were adjusted for by using pseudo residuals
500 from logistic regression analyses as the outcome variable in the PRS analyses.

501 *Meta-analysis of GWASs*

502 Previously published GWASs from each pair of traits were combined into a meta-analysis to
503 obtain joint summary statistics for each composite phenotype. Meta-analyses were performed
504 using METAL⁵⁸ with the sample-size weighted fixed-effects algorithm. To truly mimic a composite
505 phenotype GWAS, only variants included in both GWAS summary statistics were retained.

506 *Calculation of PRSs and evaluation of model performance*

507 PRS analyses were performed using a 5-fold cross validation. For each iteration, samples were
508 split into 80:20 training:validation subsets. Summary statistics of the “composite” phenotype were
509 used as base sample, and UK Biobank samples with the residualised phenotypes were used as
510 target sample (**Extended Data Tables 5 and 6**).

511 Pathway-specific PRSs for 4,079 pathways (see definition of pathways) were calculated using
512 PRSet. Competitive P -values were calculated using 10,000 permutations and pathways with
513 competitive P -value < 0.05 were defined as enriched. PRSs for the enriched pathways were
514 recalculated using P -value thresholding, such that the predictive power of each PRS was
515 maximized. The PRSs with the “best” predictive performance for each enriched pathway were
516 then included in a generalized linear model with lasso regularization using the `cv.glmnet` function
517 from the glmnet package (v4.0-2) in R with 5-fold cross-validation to select the lambda parameter
518 that generates the smallest out-of-sample mean squared error (MSE). The resultant best fitting
519 model was applied to the validation sample to calculate the model R^2 .

520 We also performed genome-wide PRS analyses using lassosum and PRSice-2. Parameters (P -
521 value thresholds for PRSice; penalty factor λ and soft-thresholding parameter s for lassosum)
522 were optimised in the training samples and were applied to the validation samples to calculate
523 the model R^2 .

524 *Stratification of inflammatory bowel disease subtypes*

525 For the stratification of inflammatory bowel disease (IBD) subtypes, pathway and genome-wide
526 PRS were calculated using IBD summary statistics as base sample, and cases of Crohn’s disease
527 vs cases of ulcerative colitis as target sample. Comorbid cases were excluded from the analysis.
528 To adjust for sex, genotyping batch, recruitment centre and first 15 principal components, logistic

529 regressions were performed and the pseudo residuals were used as the outcome variable in the
530 PRS analyses.

531 PRS analyses were performed using a 5-fold cross validation approach. For each iteration,
532 samples were split into 80:20 training:validation subsets. For the training sample, we performed
533 a similar approach to that used for supervised classification of subtypes (See *Calculation of PRS*
534 *and evaluation of model performance*) but with the PRS standardised to have mean 0 and
535 standard deviation of 1. In this case, only PRSet and PRSice were used to calculate the PRSs,
536 and pathways with non-zero coefficients after the lasso regularisation were used in the
537 unsupervised algorithm step.

538 PRS for the pathways selected in the training step were calculated for the validation sample using
539 PRSet (for pathway PRS) and PRSice (for genome-wide PRS). The silhouette method was
540 applied to select the optimal number of centroids for K-mean clustering, and we used the optimal
541 number of centroids and also 2 centroids (when the optimal is not 2) for the down-stream k-mean
542 clustering. K-mean cluster groups were regressed against the subtypes to assess the
543 classification performance of the K-means algorithm.

544

545

546

547

548

549 **Code Availability**

550 The scripts used to perform quality control on UK Biobank data are available at
551 https://gitlab.com/choishingwan/ukb_process. The scripts used in the current study are available
552 at https://gitlab.com/choishingwan/prset_analyses. PRSet is a module within PRSice and is
553 available on github repository [<https://github.com/choishingwan/PRSice>].

554 **Acknowledgments**

555 We thank the participants in the UK Biobank and the scientists involved in the construction of this
556 resource. We thank Dr Kristen Brennand, Dr Jason Kovacic, Professor Alison Goate, Professor
557 Ruth Loos, Dr Edoardo Marcora, Dr Alexander Charney, Dr Manav Kapoor and Dr Jacqueline
558 Meyers for providing their expert knowledge for each specific disease. We thank Dr Conrad
559 Iyegbe, Laura Sloofman, Collin Spencer, Chris Porras, Dr Zhe Wang and Dr Jiayi Xu for useful
560 discussions and feedback. This research has been conducted using the UK Biobank Resource
561 under application 18177 (P.F.O) and was supported by grants from the UK Medical Research
562 Council (MR/N015746/1) and the National Institute of Health (R01MH122866) to PFO. This work
563 was supported in part through the computational resources and staff expertise provided by
564 Scientific Computing at the Icahn School of Medicine at Mount Sinai, specifically the Minerva
565 Supercomputer and the Mount Sinai Data Ark data commons. Research reported in this paper
566 was supported by the Office of Research Infrastructure of the National Institutes of Health under
567 award number S10OD026880. The content is solely the responsibility of the authors and does
568 not necessarily represent the official views of the National Institutes of Health, NHS, the NIHR or
569 the Department of Health. Figure 1 was partially created using the resource BioRender.com.

570

571 **Author Information**

572 These authors contributed equally: Shing Wan Choi, Judit García-González

573 **Corresponding authors**

574 Correspondence to Shing Wan Choi and Paul O'Reilly.

575 **Contributions**

576 P.F.O conceived the design of the study, with S.W.C and Y.R providing critical contributions.

577 S.W.C developed the computational methods and the PRSet software, with J.G.G, Y.R, H.M.W

578 and P.F.O providing critical contributions. S.W.C, J.G.G and Y.R performed the analyses. J.J

579 performed quality control of the Sweden-Schizophrenia Population-Based cohort data. P.F.O

580 supervised the project, with C.J.H providing critical contributions. S.W.C, J.G.G and P.F.O wrote

581 the manuscript, with critical feedback from all other authors.

582

583

584

585

586

587

588

589 **References**

- 590 1. Choi, S. W., Mak, T. S.-H. & O'Reilly, P. F. Tutorial: a guide to performing polygenic risk
591 score analyses. *Nat. Protoc.* **15**, 2759–2772 (2020).
- 592 2. Torkamani, A., Wineinger, N. E. & Topol, E. J. The personal and clinical utility of
593 polygenic risk scores. *Nat. Rev. Genet.* **19**, 581–590 (2018).
- 594 3. Euesden, J., Lewis, C. M. & O'Reilly, P. F. PRSice: Polygenic Risk Score software.
595 *Bioinforma. Oxf. Engl.* **31**, 1466–1468 (2015).
- 596 4. Vilhjálmsson, B. J. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of
597 Polygenic Risk Scores. *Am. J. Hum. Genet.* **97**, 576–592 (2015).
- 598 5. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A. & Smoller, J. W. Polygenic prediction via
599 Bayesian regression and continuous shrinkage priors. *Nat. Commun.* **10**, 1776 (2019).
- 600 6. Leeuw, C. A. de, Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: Generalized Gene-
601 Set Analysis of GWAS Data. *PLOS Comput. Biol.* **11**, e1004219 (2015).
- 602 7. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide
603 association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 604 8. International Schizophrenia Consortium *et al.* Common polygenic variation contributes to
605 risk of schizophrenia and bipolar disorder. *Nature* **460**, 748–752 (2009).
- 606 9. Musliner, K. L. *et al.* Association of Polygenic Liabilities for Major Depression, Bipolar
607 Disorder, and Schizophrenia With Risk for Depression in the Danish Population. *JAMA*
608 *Psychiatry* **76**, 516–525 (2019).

- 609 10. Zheutlin, A. B. *et al.* Penetrance and Pleiotropy of Polygenic Risk Scores for
610 Schizophrenia in 106,160 Patients Across Four Health Care Systems. *Am. J. Psychiatry* **176**,
611 846–855 (2019).
- 612 11. Khera, A. V. *et al.* Genome-wide polygenic scores for common diseases identify
613 individuals with risk equivalent to monogenic mutations. *Nat. Genet.* **50**, 1219–1224 (2018).
- 614 12. Aung, N. *et al.* Genome-Wide Analysis of Left Ventricular Image-Derived Phenotypes
615 Identifies Fourteen Loci Associated With Cardiac Morphogenesis and Heart Failure
616 Development. *Circulation* **140**, 1318–1330 (2019).
- 617 13. Haas, M. E. *et al.* Genetic Association of Albuminuria with Cardiometabolic Disease and
618 Blood Pressure. *Am. J. Hum. Genet.* **103**, 461–473 (2018).
- 619 14. Mavaddat, N. *et al.* Polygenic Risk Scores for Prediction of Breast Cancer and Breast
620 Cancer Subtypes. *Am. J. Hum. Genet.* **104**, 21–34 (2019).
- 621 15. Zhang, J.-P. *et al.* Schizophrenia Polygenic Risk Score as a Predictor of Antipsychotic
622 Efficacy in First-Episode Psychosis. *Am. J. Psychiatry* **176**, 21–28 (2019).
- 623 16. Natarajan, P. *et al.* Polygenic Risk Score Identifies Subgroup With Higher Burden of
624 Atherosclerosis and Greater Relative Benefit From Statin Therapy in the Primary Prevention
625 Setting. *Circulation* **135**, 2091–2101 (2017).
- 626 17. Mega, J. L. *et al.* Genetic risk, coronary heart disease events, and the clinical benefit of
627 statin therapy: an analysis of primary and secondary prevention trials. *Lancet Lond. Engl.*
628 **385**, 2264–2271 (2015).

- 629 18. Pain, O. *et al.* Antidepressant Response in Major Depressive Disorder: A Genome-wide
630 Association Study. *medRxiv* 2020.12.11.20245035 (2020)
631 doi:10.1101/2020.12.11.20245035.
- 632 19. Hoekstra, S. D., Stringer, S., Heine, V. M. & Posthuma, D. Genetically-Informed Patient
633 Selection for iPSC Studies of Complex Diseases May Aid in Reducing Cellular
634 Heterogeneity. *Front. Cell. Neurosci.* **11**, 164 (2017).
- 635 20. Dobrindt, K. *et al.* Publicly Available hiPSC Lines with Extreme Polygenic Risk Scores
636 for Modeling Schizophrenia. *Complex Psychiatry* **6**, 68–82 (2020).
- 637 21. C, M.-L. *et al.* LDpred-funct: incorporating functional priors improves polygenic prediction
638 accuracy in UK Biobank and 23andMe data sets. *bioRxiv* (2018) doi:10.1101/375337.
- 639 22. Hu, Y. *et al.* Leveraging functional annotations in genetic risk prediction for human
640 complex diseases. *PLoS Comput. Biol.* **13**, e1005589 (2017).
- 641 23. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*
642 *Acids Res.* **28**, 27–30 (2000).
- 643 24. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–
644 D503 (2020).
- 645 25. Saelens, W., Cannoodt, R. & Saeys, Y. A comprehensive evaluation of module detection
646 methods for gene expression data. *Nat. Commun.* **9**, 1090 (2018).
- 647 26. Szklarczyk, D. *et al.* STRING v10: protein–protein interaction networks, integrated over
648 the tree of life. *Nucleic Acids Res.* **43**, D447–D452 (2015).
- 649 27. Markowetz, F. How to Understand the Cell by Breaking It: Network Analysis of Gene
650 Perturbation Screens. *PLOS Comput. Biol.* **6**, e1000655 (2010).

- 651 28. Austin, J. C. & Honer, W. G. Psychiatric genetic counselling for parents of individuals
652 affected with psychotic disorders: a pilot study. *Early Interv. Psychiatry* **2**, 80–89 (2008).
- 653 29. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from
654 polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
- 655 30. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-scale
656 data. *GigaScience* **8**, (2019).
- 657 31. Nishimura, D. BioCarta. *Biotech Softw. Internet Rep.* **2**, 117–120 (2001).
- 658 32. Schaefer, C. F. *et al.* PID: the Pathway Interaction Database. *Nucleic Acids Res.* **37**,
659 D674–679 (2009).
- 660 33. Bult, C. J. *et al.* Mouse Genome Database (MGD) 2019. *Nucleic Acids Res.* **47**, D801–
661 D806 (2019).
- 662 34. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nat. Genet.* **25**,
663 25–29 (2000).
- 664 35. The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Res.* **47**,
665 D330–D338 (2019).
- 666 36. Skene, N. G. *et al.* Genetic identification of brain cell types underlying schizophrenia.
667 *Nat. Genet.* **50**, 825–833 (2018).
- 668 37. Hemonnot, A.-L., Hua, J., Ulmann, L. & Hirbec, H. Microglia in Alzheimer Disease: Well-
669 Known Targets and New Opportunities. *Front. Aging Neurosci.* **11**, (2019).
- 670 38. Watanabe, K., Umićević Mirkov, M., de Leeuw, C. A., van den Heuvel, M. P. &
671 Posthuma, D. Genetic mapping of cell type specificity for complex traits. *Nat. Commun.* **10**,
672 3222 (2019).

- 673 39. Pain, O. *et al.* Evaluation of polygenic prediction methodology within a reference-
674 standardized framework. *PLoS Genet.* **17**, e1009021 (2021).
- 675 40. Lloyd-Jones, L. R. *et al.* Improved polygenic prediction by Bayesian multiple regression
676 on summary statistics. *Nat. Commun.* **10**, 5086 (2019).
- 677 41. Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics*
678 **36**, 5424–5431 (2020).
- 679 42. Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological
680 insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421–427 (2014).
- 681 43. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity
682 biology. *Nature* **518**, 197–206 (2015).
- 683 44. Willer, C. J. *et al.* Discovery and refinement of loci associated with lipid levels. *Nat.*
684 *Genet.* **45**, 1274–1283 (2013).
- 685 45. Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer’s disease identifies
686 new risk loci and implicates A β , tau, immunity and lipid processing. *Nat. Genet.* **51**, 414–430
687 (2019).
- 688 46. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association
689 meta-analysis of coronary artery disease. *Nat. Genet.* **47**, 1121–1130 (2015).
- 690 47. Scott, R. A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in
691 Europeans. *Diabetes* **66**, 2888–2902 (2017).
- 692 48. Sanchez-Roige, S. *et al.* Genome-Wide Association Study Meta-Analysis of the Alcohol
693 Use Disorders Identification Test (AUDIT) in Two Population-Based Cohorts. *Am. J.*
694 *Psychiatry* **176**, 107–118 (2019).

- 695 49. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological
696 architecture of adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- 697 50. Nishimura, D. BioCarta. *Biotech Softw. Internet Rep.* **2**, 117–120 (2001).
- 698 51. Schaefer, C. F. *et al.* PID: the Pathway Interaction Database. *Nucleic Acids Res.* **37**,
699 D674–D679 (2009).
- 700 52. Fabregat, A. *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids Res.* **46**,
701 D649–D655 (2018).
- 702 53. Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**,
703 1739–1740 (2011).
- 704 54. Bult, C. J. *et al.* Mouse Genome Database (MGD) 2019. *Nucleic Acids Res.* **47**, D801–
705 D806 (2019).
- 706 55. North, B. V., Curtis, D. & Sham, P. C. A Note on the Calculation of Empirical P Values
707 from Monte Carlo Procedures. *Am. J. Hum. Genet.* **71**, 439–441 (2002).
- 708 56. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 709 57. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
710 datasets. *GigaScience* **4**, (2015).
- 711 58. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of
712 genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
- 713 59. Espe, S. Malacards: The Human Disease Database. *J. Med. Libr. Assoc. JMLA* **106**,
714 140–141 (2018).

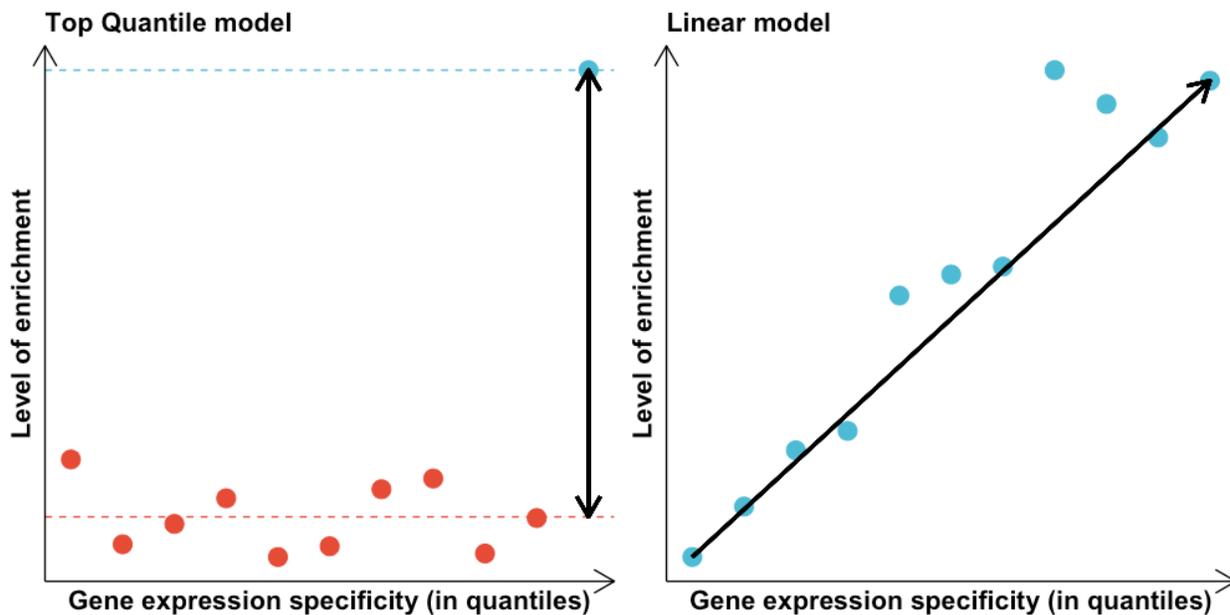
- 715 60. Aguet, F. *et al.* The GTEx Consortium atlas of genetic regulatory effects across human
716 tissues. *bioRxiv* 787903 (2019) doi:10.1101/787903.
- 717 61. Bryois, J. *et al.* Genetic identification of cell types underlying brain complex traits yields
718 insights into the etiology of Parkinson's disease. *Nat. Genet.* **52**, 482–493 (2020).
- 719 62. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies
720 disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
- 721 63. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body
722 mass index in ~700000 individuals of European ancestry. *Hum. Mol. Genet.* **27**, 3641–3649
723 (2018).
- 724 64. Sanchez-Roige, S. *et al.* Genome-Wide Association Study Meta-Analysis of the Alcohol
725 Use Disorders Identification Test (AUDIT) in Two Population-Based Cohorts. *Am. J.*
726 *Psychiatry* **176**, 107–118 (2018).
- 727 65. Sinnott-Armstrong, N. *et al.* Genetics of 35 blood and urine biomarkers in the UK
728 Biobank. *Nat. Genet.* **53**, 185–194 (2021).
- 729 66. Global Lipids Genetics Consortium *et al.* Discovery and refinement of loci associated
730 with lipid levels. *Nat. Genet.* **45**, 1274–1283 (2013).
- 731 67. Inouye, M. *et al.* Genomic Risk Prediction of Coronary Artery Disease in 480,000 Adults:
732 Implications for Primary Prevention. *J. Am. Coll. Cardiol.* **72**, 1883–1893 (2018).
- 733 68. Marioni, R. E. *et al.* GWAS on family history of Alzheimer's disease. *Transl. Psychiatry* **8**,
734 1–7 (2018).

735

736

737 **Extended Data Figures**

738

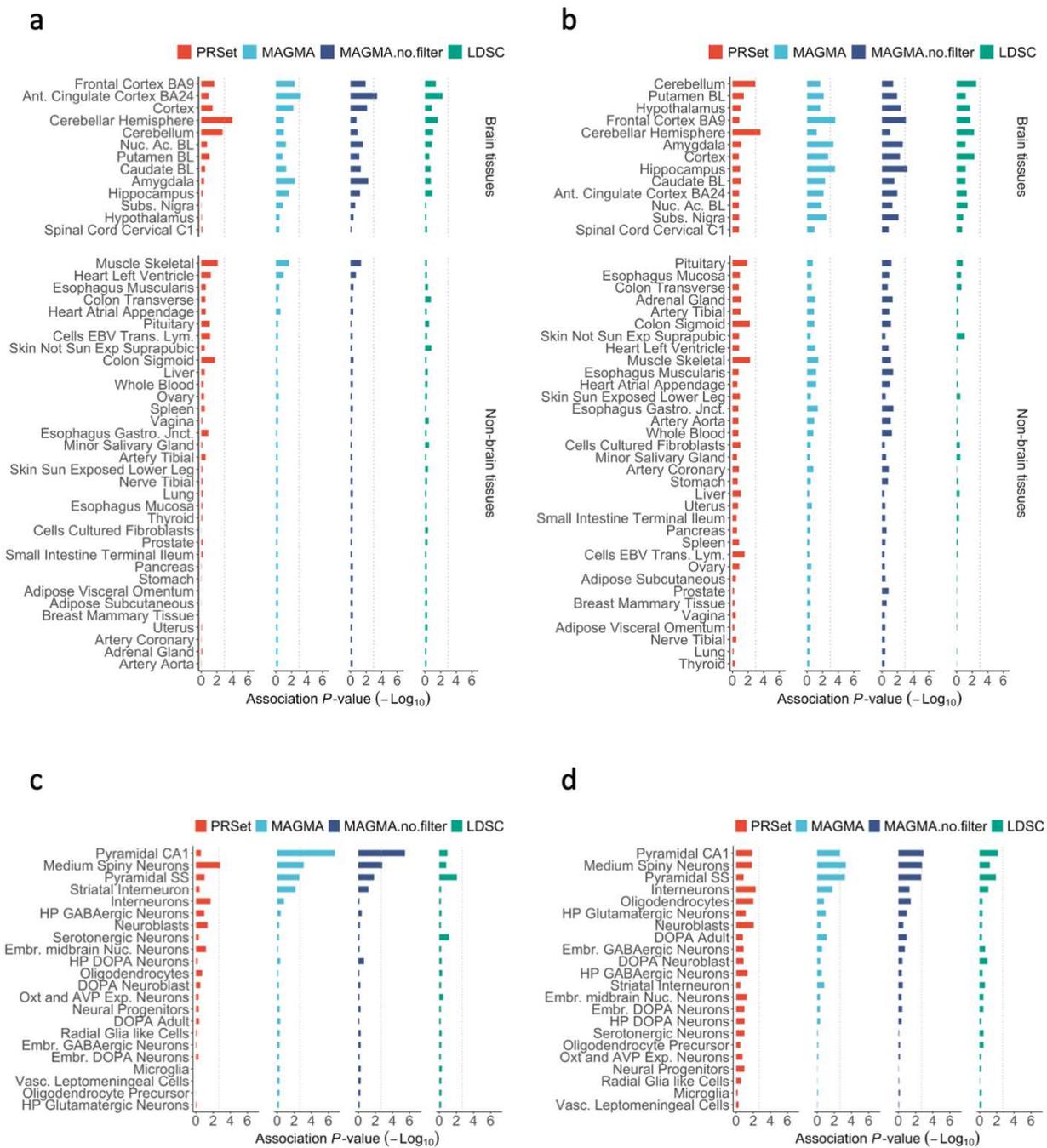


739

740 **Extended Data Figure 1.** Illustration of the test models used to assess cell type and tissue
741 specificity. Left panel: illustrates the “top quantile” test model, which assumes that GWAS signal
742 enrichment is concentrated in the most specifically expressed genes. Right panel: illustrates the
743 “linear” test model, which assumes that enrichment of GWAS signal increases linearly with
744 expression specificity.

745

746



747

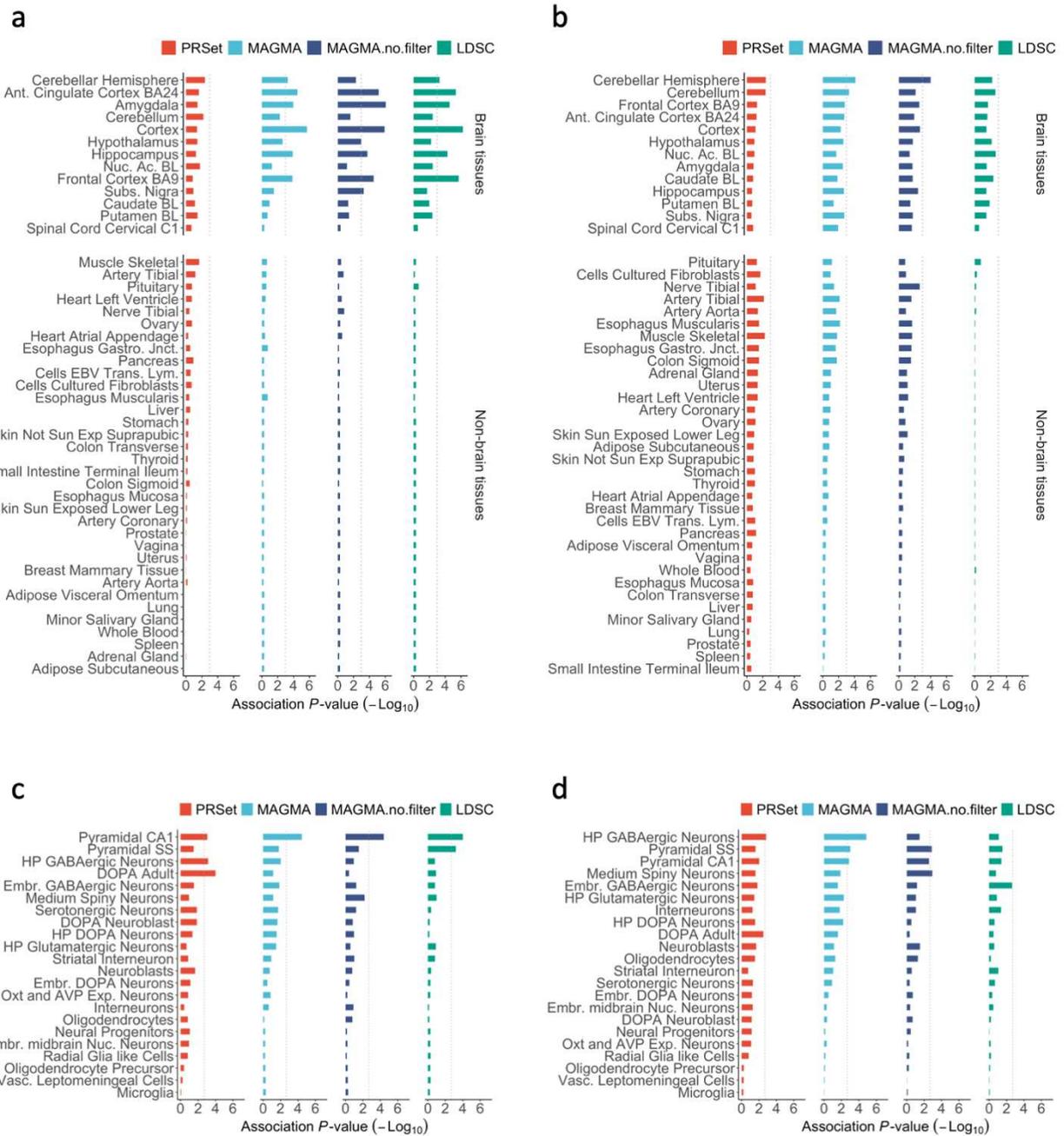
748 **Extended Data Figure 2.** Association between pathway enrichment P -value and expert opinion

749 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software

750 and schizophrenia. Results for 'MAGMA' (light blue) were obtained after removing ambiguous

751 and non-overlapping SNPs in the analysis to compare software performance given identical input

752 data. Results for 'MAGMA.no.filter' (dark blue) show results including all available SNPs.



753

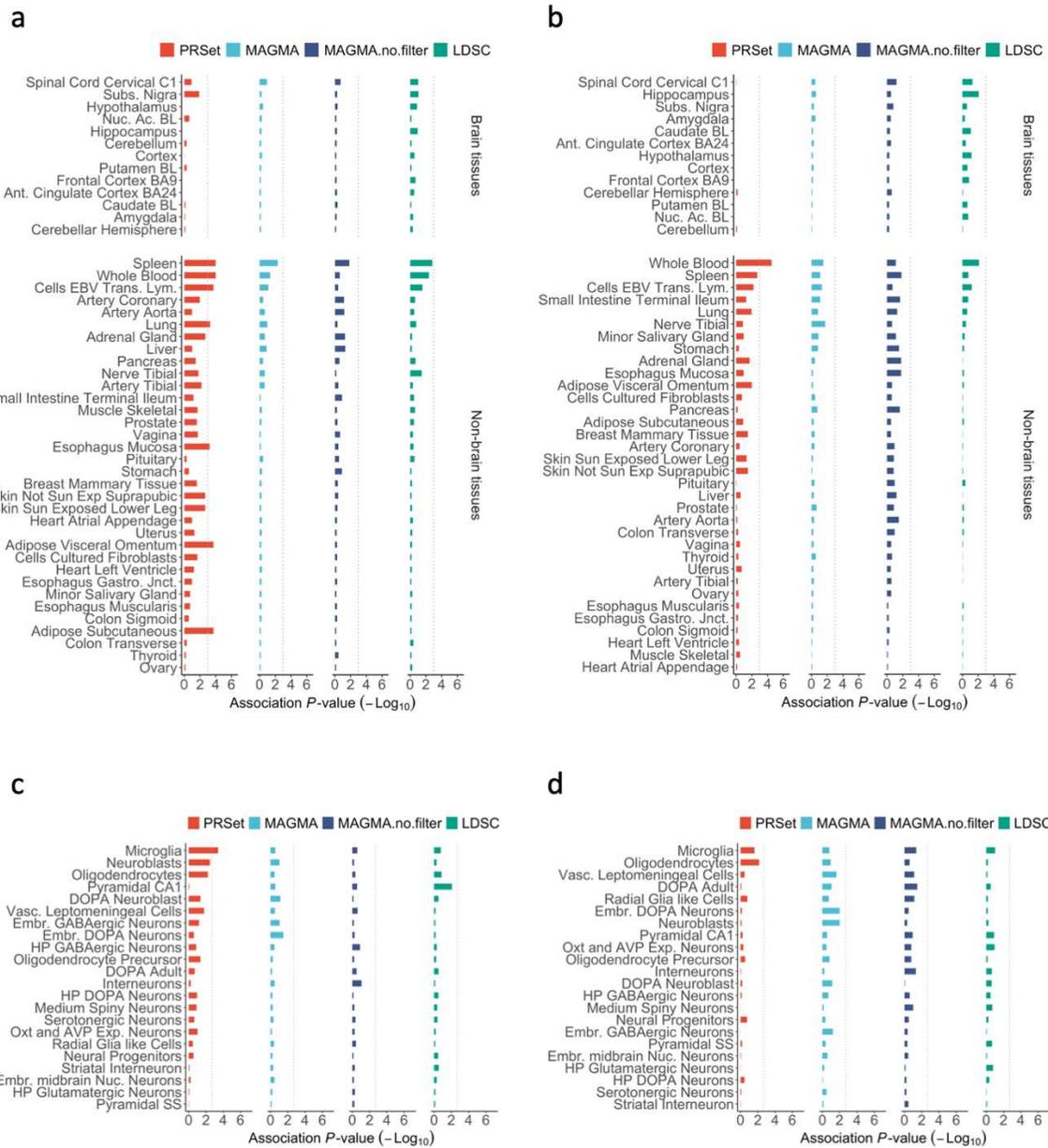
754 **Extended Data Figure 3.** Association between pathway enrichment P -value and expert opinion

755 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software

756 and BMI. Results for 'MAGMA' (light blue) were obtained after removing ambiguous and non-

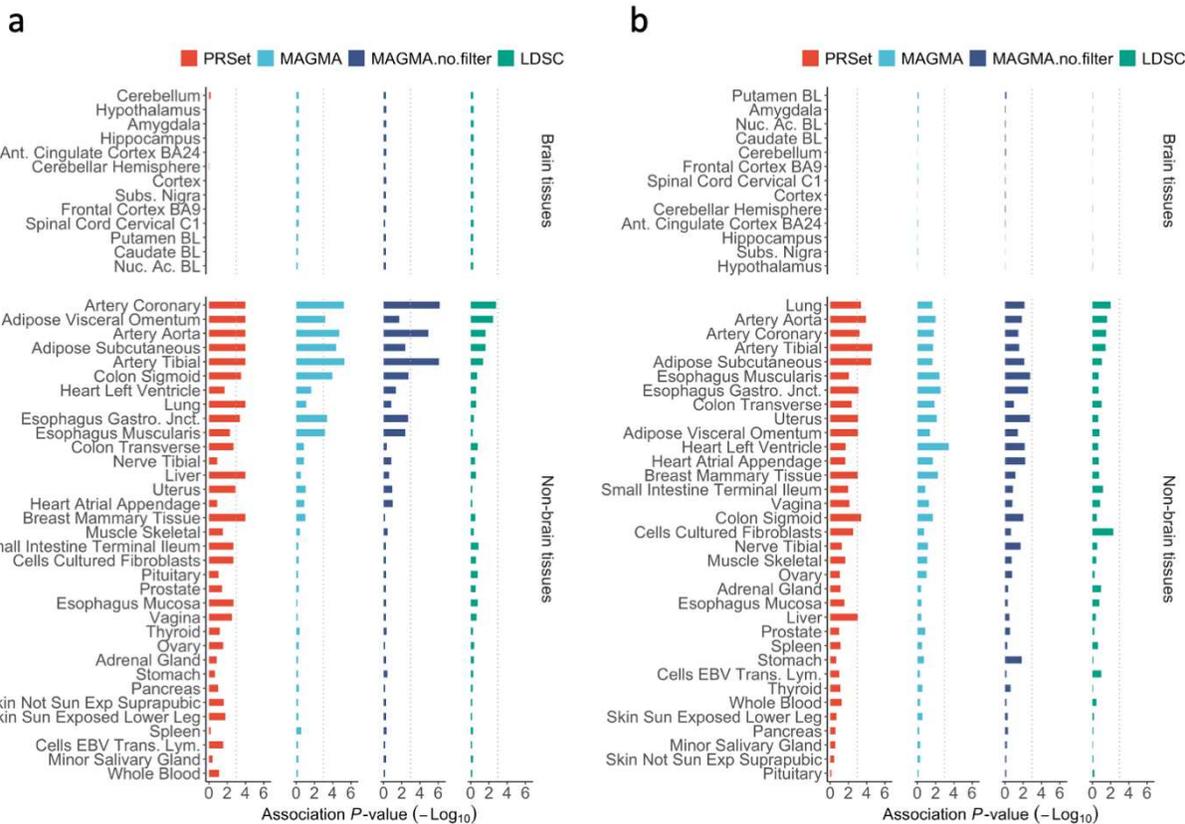
757 overlapping SNPs in the analysis to compare software performance given identical input data.

758 Results for 'MAGMA.no.filter' (dark blue) show results including all available SNPs.



759

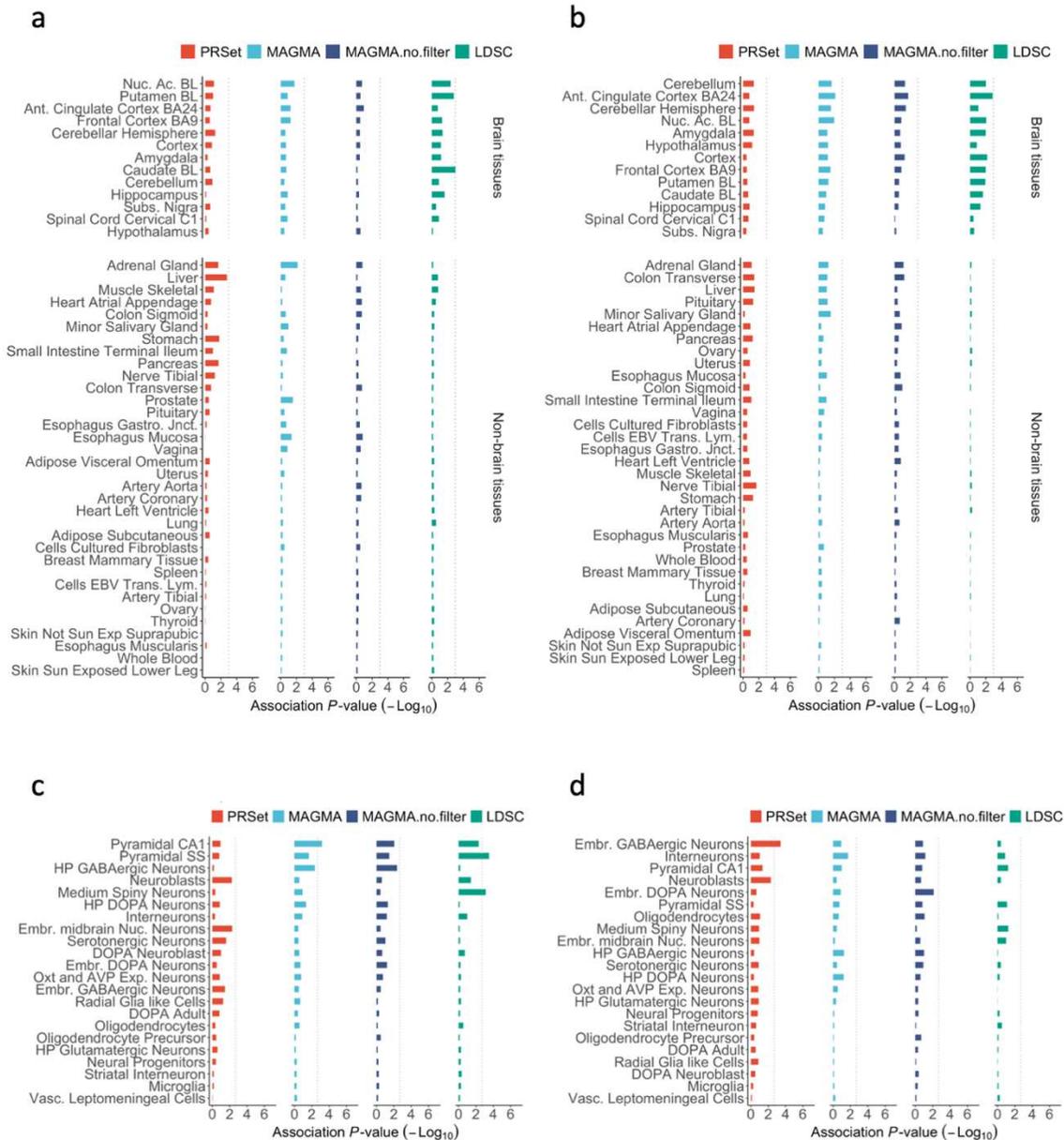
760 **Extended Data Figure 4.** Association between pathway enrichment P -value and expert opinion
 761 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software
 762 and Alzheimer's disease. Results for 'MAGMA' (light blue) were obtained after removing
 763 ambiguous and non-overlapping SNPs in the analysis to compare software performance given
 764 identical input data. Results for 'MAGMA.no.filter' (dark blue) show results including all available
 765 SNPs.



766

767 **Extended Data Figure 5.** Association between pathway enrichment P -value and expert opinion
 768 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software
 769 and Coronary Artery Disease. Results for ‘MAGMA’ (light blue) were obtained after removing
 770 ambiguous and non-overlapping SNPs in the analysis to compare software performance given
 771 identical input data. Results for ‘MAGMA.no.filter’ (dark blue) show results including all available
 772 SNPs.

773



774

775 **Extended Data Figure 6.** Association between pathway enrichment *P*-value and expert opinion

776 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software

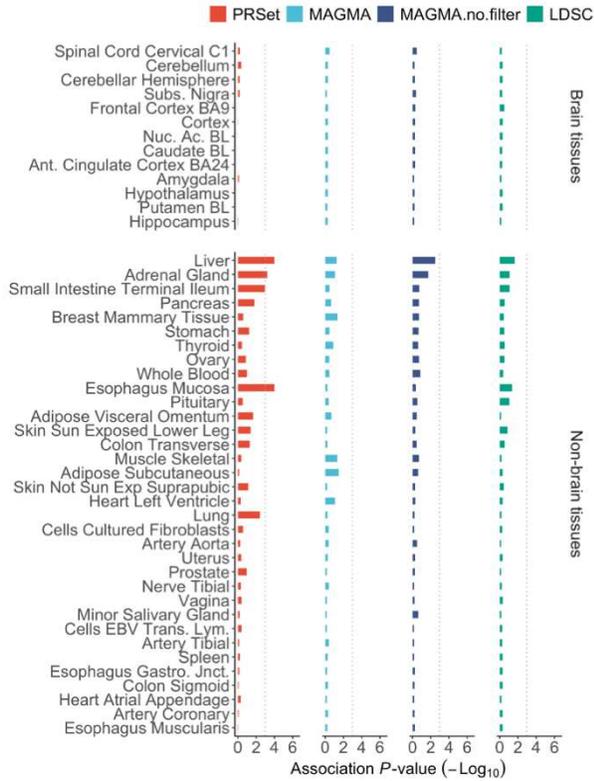
777 and alcohol consumption. Results for 'MAGMA' (light blue) were obtained after removing

778 ambiguous and non-overlapping SNPs in the analysis to compare software performance given

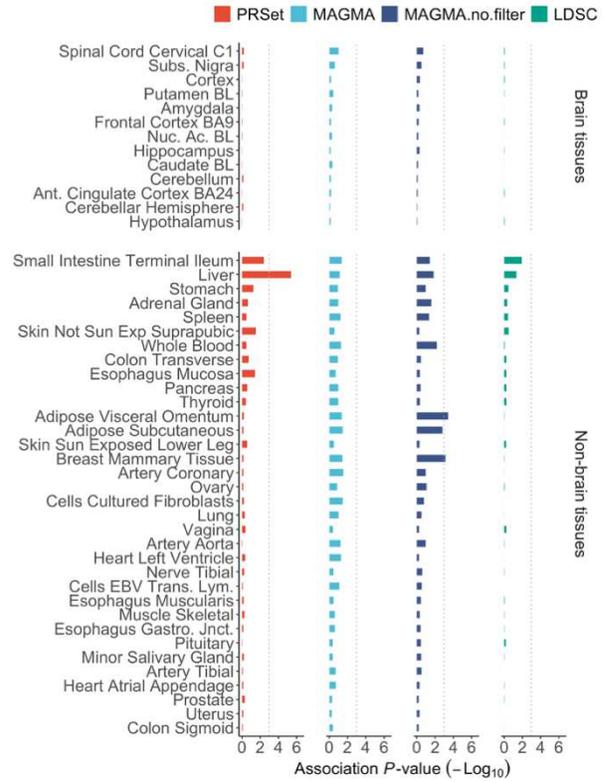
779 identical input data. Results for 'MAGMA.no.filter' (dark blue) show results including all available

780 SNPs.

a



b



781

782 **Extended Data Figure 7.** Association between pathway enrichment P -value and expert opinion
 783 of tissue relevance (panels a and b) and cell type relevance (panels c and d) for each software
 784 and low-density lipoproteins. Results for 'MAGMA' (light blue) were obtained after removing
 785 ambiguous and non-overlapping SNPs in the analysis to compare software performance given
 786 identical input data. Results for 'MAGMA.no.filter' (dark blue) show results including all available
 787 SNPs.

788

789 **Supplementary information for:**

790 **The power of pathway-based polygenic risk scores**

791 Shing Wan Choi^{1*}, Judit García-González^{1*}, Yunfeng Ruan², Hei Man Wu¹, Jessica Johnson¹,

792 Clive J Hoggart¹, Paul F. O'Reilly¹

793

794 ¹ Department of Genetics and Genomic Sciences, Icahn School of Medicine, Mount Sinai, New York City, NY 10029, USA

795 ²The Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

796 * These authors contributed equally to this work

797

798 **Supplementary Information Contents:**

799

800 **Supplementary Methods**

801

802 **Ascertainment of phenotypes in the UK Biobank cohort**

803 **PRSet implementation**

804 **Gene set based clumping**

805 **Permutation optimization**

806

807 **Supplementary Notes**

808

809 **Supplementary Note 1.** Kendall correlations between pathway enrichment and *MalaCards*
810 scores – Sensitivity analysis excluding genes in *MalaCards* database.

811 **Supplementary Note 2.** Performance of PRSet with a 5Mbp shift in gene boundaries.

812 **Supplementary Note 3.** Performance of PRSet vs genome-wide PRS methods for single traits.

813 **Supplementary Note 4.** Performance of PRSet vs genome-wide with lasso regression.

814

815 **Supplementary References**

816

817

818 **Supplementary Methods**

819 **Ascertainment of phenotypes in the UK Biobank cohort**

820 **Body Mass Index (BMI):** BMI information was extracted from Field ID 21001. For the analysis of
821 pathway enrichment and comparison with *MalaCards* disease relevance scores, residuals were
822 calculated for men and women separately using linear regression, and were adjusted for age,
823 recruitment centre, genotyping batch and 15 first principal components. Inverse normal
824 transformation was carried out on the residuals⁶³.

825 **Alcohol Consumption:** Alcohol consumption scores were extracted using the consumption
826 component of the Alcohol Use Disorder Identification Test (AUDIT-C), which was included in the
827 mental health questionnaire⁶⁴. AUDIT-C was based on three questions: frequency of drinking
828 alcohol (Field ID 20414), amount of alcohol drunk on a typical drinking day (Field ID 20403), and
829 frequency of consuming six or more units of alcohol (Field ID 20416). We used linear regression
830 to obtain residuals of the AUDIT-C score adjusted for age, sex, recruitment centre, genotyping
831 batch and 15 principal components⁶⁴.

832 **Low-Density Lipoprotein (LDL):** Levels of LDL across participants (Field ID 30780) were
833 adjusted for statin use, following the procedure from Sinnott-Armstrong et al⁶⁵. Medication use
834 for 14 statins was extracted from Field ID 20003 (**Extended Data Table 2**). This information was
835 then used to identify 1,382 individuals with LDL measurements that were not taking statins upon
836 enrolment (years 2006-2010) but were taking statins at the time of first repeat assessment (years
837 2012-2013). For these individuals, we applied a statin correction factor on LDL measurements
838 (*corrected LDL measurement = LDL upon enrolment / statin correction factor*). The statin
839 correction factor was calculated as the mean of each individual's LDL measurement upon
840 enrolment divided by LDL measurement at the time of first repeat assessment. Residuals of the
841 LDL measurements were calculated by adjusting for age, sex, recruitment centre, genotyping

842 batch, fasting status (Field ID 74), dilution factor (Field ID 30897) and 15 PCs⁶⁶ using linear
843 regression.

844 **Coronary Artery Disease (CAD):** CAD was defined as in Inouye et al (2018)⁶⁷. CAD cases were
845 ascertained as individuals who had suffered fatal or nonfatal myocardial infarction as indicated by
846 their hospital records (Field IDs 41270, 41202, 41204, 41203, 41205 and 41271), death records
847 (40001 and 40002) and medical history (Field IDs 6150 and 20002), and patients who have
848 undergone percutaneous transluminal coronary angioplasty, or coronary artery bypass grafting
849 (Field IDs 20004, 41200 and 41272). The age of event in cases was determined as the self-
850 reported age and calculated age based on the earliest hospital record of the event or based on
851 the death records; if more than one age were available, the smaller value was used. Age of the
852 controls was determined as the latest self-reported age. For the analysis of pathway enrichment
853 and comparison with MalaCards relevance scores, residuals of the CAD case control status were
854 calculated by adjusting for age of event, sex, recruitment centre, genotyping batch and 15
855 principal components using logistic regression.

856 **Alzheimer's Disease (AD):** We generated a proxy phenotype for AD case control status based
857 on family history of AD, similar to the procedure used in Marioni et al (2018)⁶⁸. First, we removed
858 participants diagnosed with AD as indicated by their hospital records (Field IDs 41202, 41204),
859 death records (Field IDs 40001 and 40002) and primary care data, participants who were adopted
860 (Field ID 1767), as well as participants whose parents were aged under 60 years (Field IDs 1845
861 and 2946), dead before age 60 years, or without age information (Field IDs 3526 and 1807). After
862 merging with the genetic data, 41,164 participants had at least one parent who was affected by
863 AD and 223,253 participants with parents non-affected by AD. Parental AD status was then
864 defined as the number of parents affected by AD. Residuals of the parental AD status were
865 calculated by adjusting for the sex of the participant, paternal age, maternal age, recruitment
866 centre, genotyping batch and 15 principal components using linear regression.

867 **Height:** Standing height information was extracted from Field ID 50. Residuals were calculated
868 by adjusting for sex, age, recruitment centre, genotyping batch and 15 principal components using
869 linear regression.

870 **Type 2 diabetes:** Type 2 diabetes cases were ascertained as individuals diagnosed with non-
871 insulin-dependent diabetes mellitus as indicated by their hospital records and primary care data
872 (Field IDs 41202, 41204 and 41270) as well as death records (Field ID 40001 and 40002).

873 **Inflammatory Bowel Disease:** Cases of Inflammatory Bowel Disease were ascertained as
874 individuals diagnosed with Crohn's disease or Ulcerative colitis as indicated by the self-reported
875 questionnaire (Field ID 20002), their hospital records and primary care data (Field IDs 41202,
876 41204 and 41270) as well as death records (Field ID 40001 and 40002). Cases diagnosed with
877 both Crohn's disease and Ulcerative Colitis were excluded. Residuals were calculated by
878 adjusting for sex, recruitment centre, genotyping batch and 15 principal components using logistic
879 regression.

880 **PRSet implementation**

881 **Gene set based clumping**

882 Genome-wide clumping can be suboptimal for gene set PRS calculation, as SNPs inside the gene
883 set may be clumped out by SNPs outside the gene set. Therefore, clumping should be performed
884 within each gene set to maximize signal retention. However, independently performing clumping
885 on each gene set can be computationally intensive. PRSet optimizes the per-gene set clumping
886 by utilizing a bit-flag system, where gene set membership is represented as a bit-flag
887 (**Supplementary Fig. 1**). If the SNP is a member of the k^{th} gene set, then the k^{th} bit of the bit-flag
888 will be set. During the clumping procedure, the index SNP will "remove" gene set membership
889 from the clumped SNP if and only if they fall within the same gene set. At the end, SNPs without

890 any gene set memberships are removed from the analysis. This allows PRSet to perform the gene
 891 set based clumping without repeating the entire clumping procedure.

892

	Set A	Set B	Set C		Set A	Set B	Set C
Index SNP	1	0	1		1	0	1
Clumped SNP 1	1	1	0	→	0	1	0
Clumped SNP 2	1	0	1		0	0	0

893

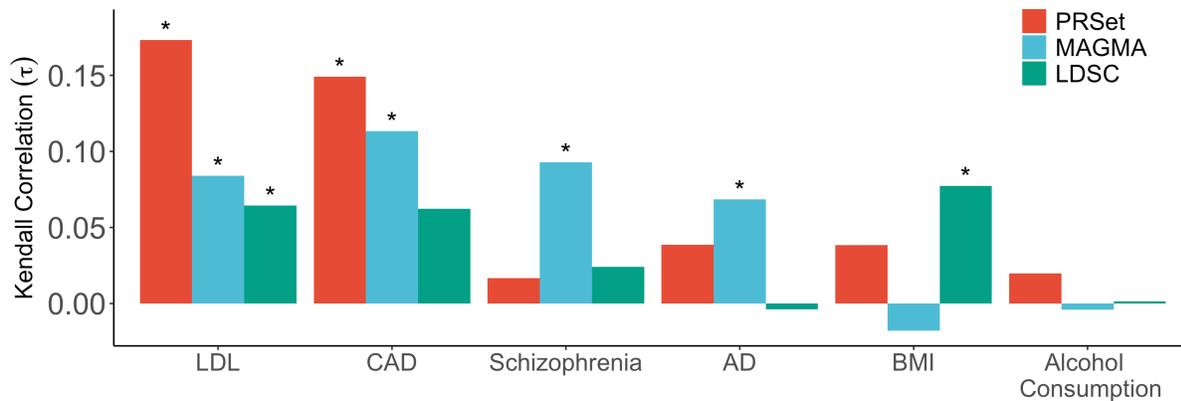
894 **Supplementary Fig 1.** Illustration of bit operation involved in PRSet clumping. The index SNP
 895 will “remove” gene set memberships from the clumped SNPs if and only if they fall within the same
 896 gene set. Clumped SNP without any gene set membership will be removed at the end of clumping.
 897 Here, clumped SNP 2 will be removed.

898 **Permutation optimisation**

899 PRSet calculates the competitive *P*-value via a permutation procedure, which is computationally
 900 expensive. To speed-up the permutation process, PRSet employes multiple techniques to
 901 optimize the computation procedure. First, the genotype of all “background” SNPs were loaded
 902 into the memory (using the --ultra parameter), such that the cost of repeated I/O during
 903 permutation is reduced. Then, gene sets with the same number of independent SNPs are grouped
 904 together and share the same null. In addition, during set-based permutation, an inverse
 905 regression is performed with *X* and *Y* switched. Inverted *X* and *Y* in the permutation results in the
 906 same z-score, but this allows pre-decomposition of the matrix and reuse for all permutations,
 907 drastically speeding up the performance.

908 **Supplementary Notes**

909 **Supplementary Note 1. Kendall correlations between pathway enrichment**
910 **and *MalaCards* scores**



911

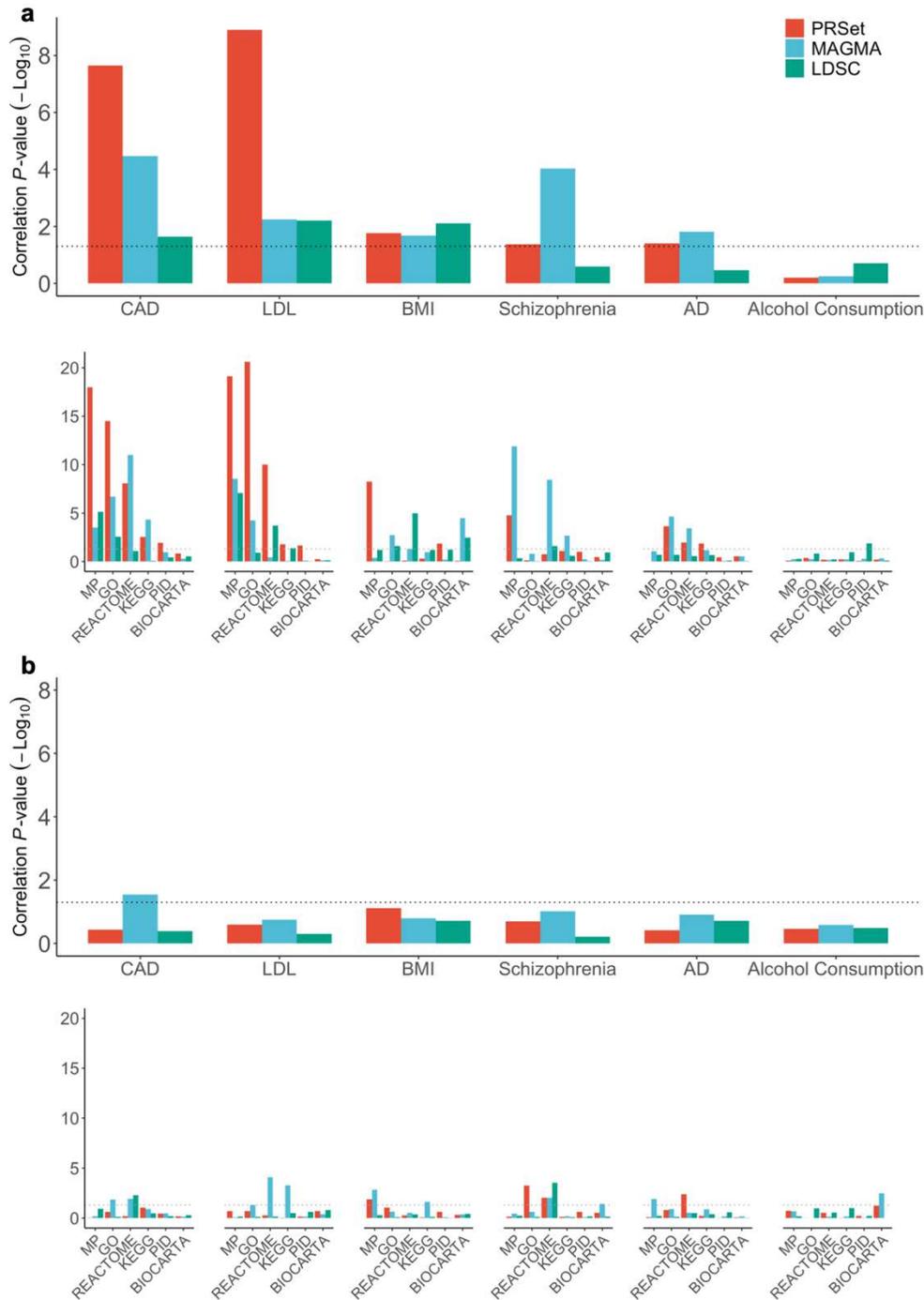
912 **Supp Fig 2.** Kendall correlation coefficients (τ) between pathway enrichment analyses and
913 *MalaCards* relevance scores. Bar plots illustrate joint results of the six databases used to define
914 pathways. *empirical P -value < 0.05 .

915 **Sensitivity analysis excluding genes in *MalaCards* database**

916 One concern with using *MalaCards* as a proxy to empirical evidence is that the *MalaCards* scores
917 may have inflated scores for genes with large effect sizes or genes that are well captured by
918 GWAS because such genes may have particularly high experimental follow-up for this reason. To
919 test whether potential biases in *MalaCards* scores could have a different impact across the
920 enrichment algorithms, we ran a sensitivity analysis in which we removed all genes that have a
921 *MalaCards* score above 0 from the pathway definitions and repeated the pathway enrichment
922 analyses.

923 For the three methods, the correlation between pathway enrichment and the *MalaCards* relevance
924 scores decreased substantially (**Supplementary Fig. 3 and 4**), confirming that *MalaCards*

925 pathway relevance scores include part of the disease signal captured by the three pathway
926 enrichment tools. The effect was stronger for PRSet and LDSC than for MAGMA. For PRSet and
927 LDSC correlations were no longer significant (**Supplementary Fig. 3b**), and correlation
928 coefficients across the six diseases decreased from $\tau = 0.077$ to $\tau = -0.0015$ for PRSet, and from
929 $\tau = 0.043$ to $\tau = 0.0076$ for LDSC (**Supplementary Fig. 4**); whereas for MAGMA, results were still
930 correlated with the original pathway relevance scores, albeit with reduced correlation and
931 significance ($\tau = 0.029$; P -value = 7.35×10^{-11}).



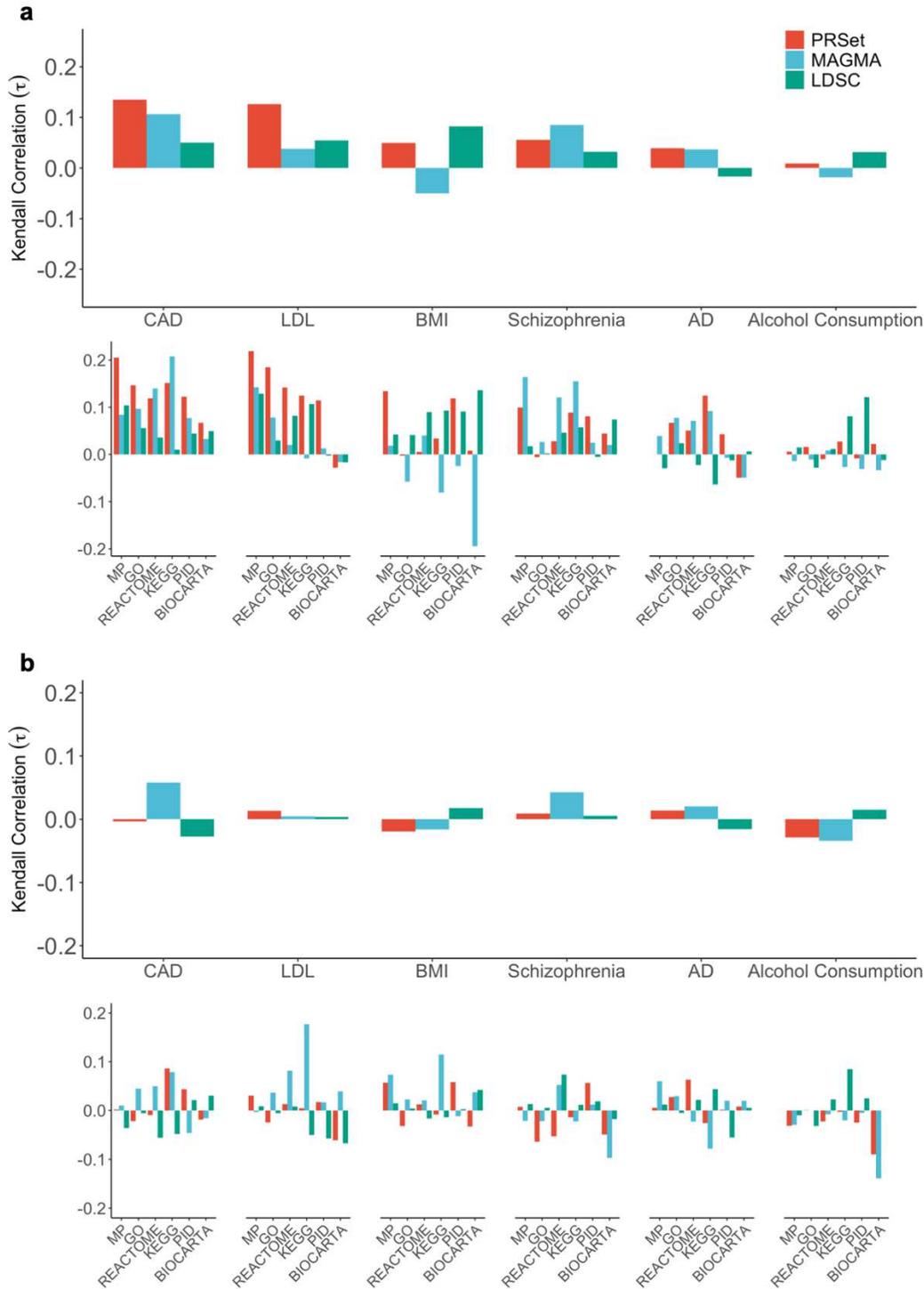
932

933 **Supplementary Fig 3.** Kendall correlation P -values between pathway enrichment analyses and

934 *MalaCards* relevance scores. **a**, Enrichment analyses use pathways that contain genes with

935 *MalaCards* scores **b**, Enrichment analyses use pathways where genes that have *MalaCards*

936 scores have been removed.



937

938 **Supplementary Fig 4.** Kendall correlation coefficients (τ) between pathway enrichment analyses

939 and *MalaCards* relevance scores. **a**, Enrichment analyses use pathways that contain genes with

940 *MalaCards* scores. **b**, Enrichment analyses use pathways where genes that have *MalaCards*
941 scores are removed.

942 **Supplementary Note 2. Performance of PRSet with a 5Mb shift in gene**
943 **boundaries.**

944 *Rationale:* Genome-wide PRS methods use the most associated SNPs (smallest *P*-value) to
945 calculate PRSs. For heterogeneous diseases, the most associated SNPs are likely to contain risk
946 alleles that are common across disease subtypes, as those will have more power to be detected.
947 Whereas prioritizing SNPs that are common across subtypes can be useful for the prediction of
948 the disease risk, those SNPs may not be useful for disease subtype classification.

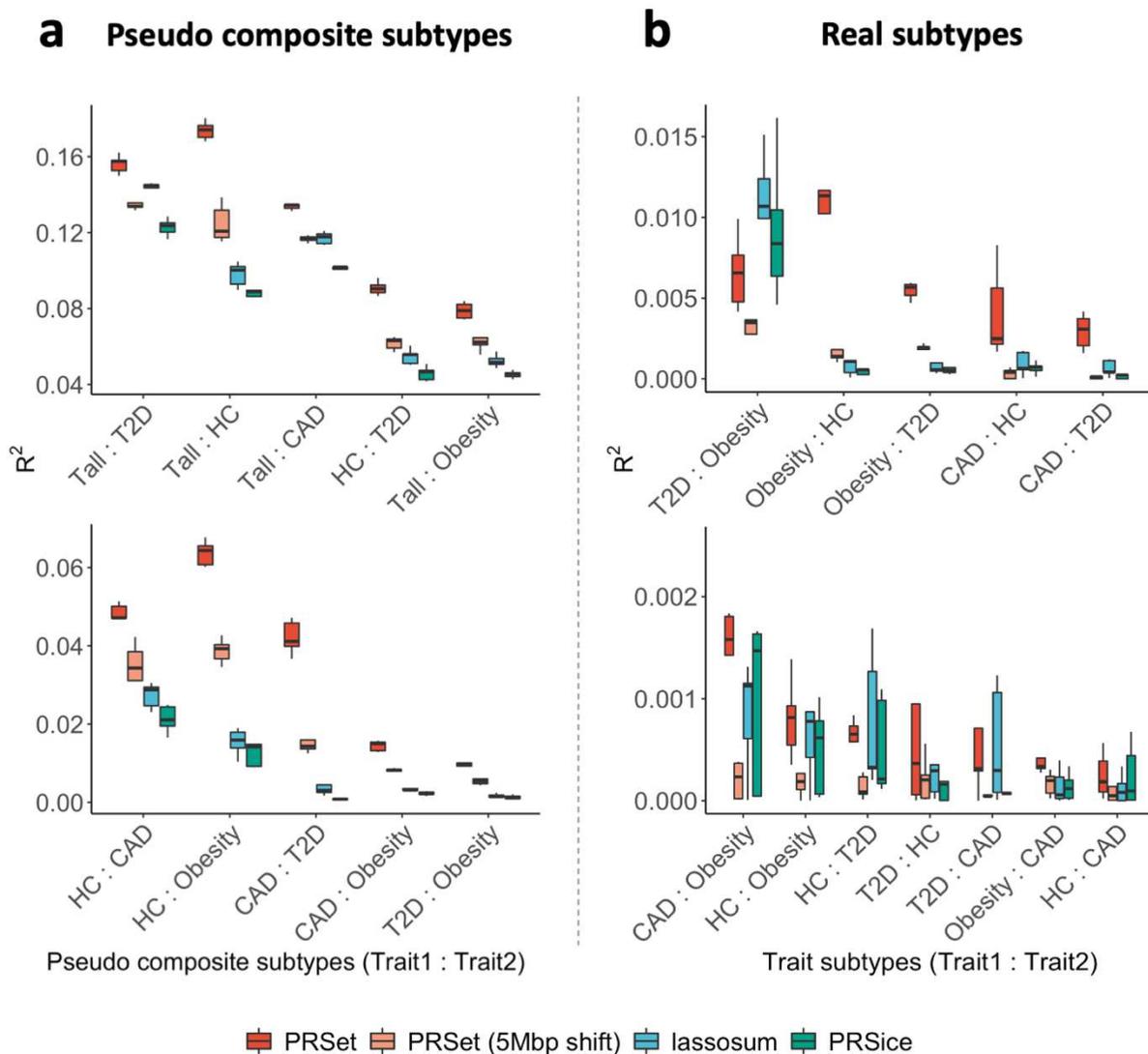
949 PRSet splits the genome into *k* pathways or ‘chunks’ and calculates PRSs for each pathway after
950 clumping and thresholding. It is possible that splitting the genome into chunks and optimising
951 parameters for each chunk improves classification: under this scenario, the most associated
952 SNPs that are common across subtypes are less likely to mask other SNPs that -although less
953 strongly associated- may be better subtype classifiers.

954 One important consideration is that the improvement in classification after splitting the genome
955 into chunks may be true regardless of a chunk being part of a gene/biological pathway. To test
956 whether the pathway structure is sufficient to explain the improvement in classification
957 performance, or whether a *biologically informed* pathway structure is needed, we assessed the
958 ability of PRSet to classify disease subtypes when we have reduced biological relevance of the
959 pathway regions while retaining the pathway structures.

960 *Methods:* To keep the same pathway structure and overlap as the original analyses, we modified
961 the gtf by shifting the gene boundaries 5Mbp, and we re-run the clustering analyses as previously
962 described (Methods). Shifting was performed using the R packages ‘genomation’ (v.1.14.0) and

963 'GenomicRanges' (v.1.34.0). It is important to note that ~5% genes fall outside of chromosomal
 964 boundary after the shift and they will be excluded from the analysis.

965 *Results:* After shifting gene boundaries 5Mb, 77.6% of the SNPs were out of genic regions. For
 966 all the composite phenotypes, performance of PRSet with the 5Mb shift decreased substantially
 967 as expected (**Supplementary Fig. 5**), suggesting that the biological information contained in the
 968 pathway structure is useful for subtype classification.



970 **Supplementary Fig 5.** Classification performance of PRSet with a 5Mb shift in the gene
971 coordinate boundaries versus standard version of PRSet and genome-wide PRS methods
972 (lassosum and PRSice). **a**, pseudo composite subtype classification results. **b**, real subtype
973 classification results.

974

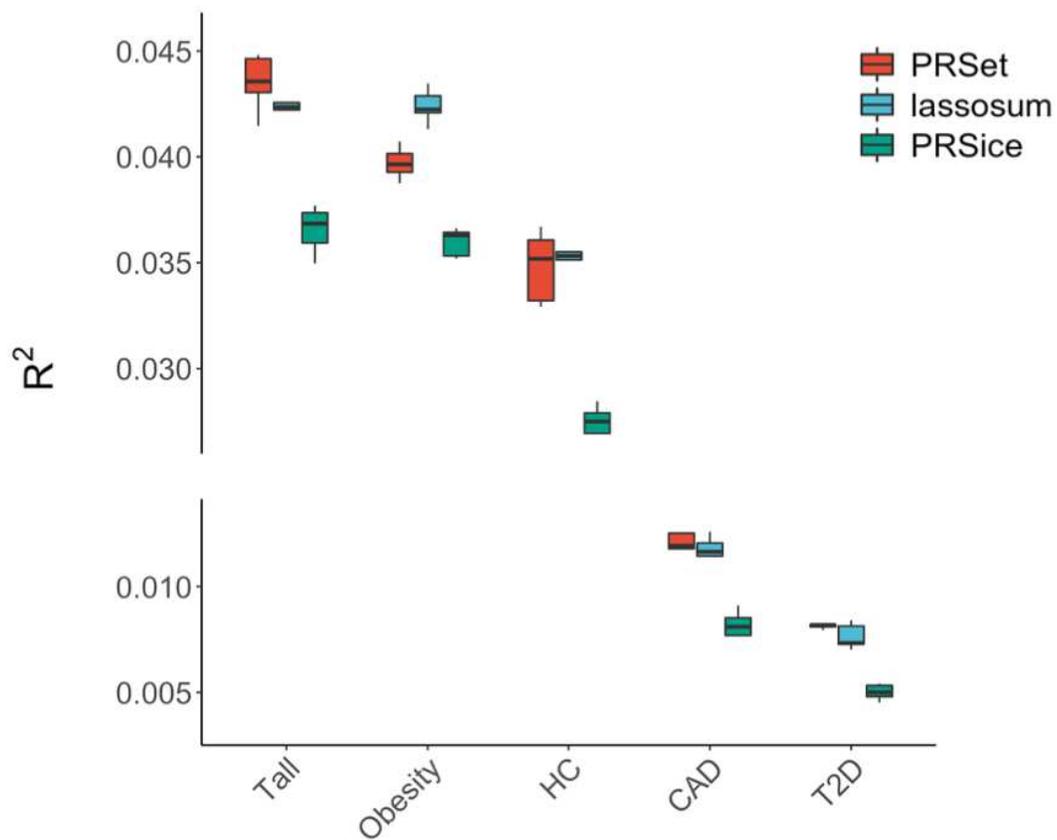
975 **Supplementary Note 3. Performance of PRSet vs genome-wide PRS methods**
976 **for single traits**

977 *Rationale:* Following results from Supplementary Note 2, we hypothesise that PRSet outperforms
978 genome-wide PRS methods at classifying disorders due to the use of pathways with biological
979 information. PRSet calculates pathway specific PRSs, which allows SNPs with effect size in only
980 one subtype to remain in the PRS, whereas PRSice and lassosum are more likely to optimize
981 SNP signal that affects both subtypes -and therefore may not be so useful as classifiers.

982 If this hypothesis is true, PRSet would outperform genome-wide PRS tools for classification tasks,
983 but its performance will be more limited when PRSs are used for the prediction of single, more
984 homogeneous traits. In those cases, we expect genome-wide PRS to do a better job at trait
985 prediction as they have more power to optimize genetic signal associated with disease. To test
986 whether PRSet results are due to the use of pathways or due to methodological aspects that could
987 have contributed to the improvement in performance, we assessed pathway and whole-genome
988 PRSs for single traits (as opposed to than the composite phenotypes that were generated in the
989 main analyses).

990 *Methods:* Whole-genome and pathway specific PRS were calculated for the same five
991 traits/diseases that were used for the sub-phenotype classification in real data: type 2 diabetes,
992 CAD, obesity (defined as BMI > 30), LDL and extreme height. We calculated PRS for these traits
993 using publicly available GWAS data for individuals from the UK Biobank cohort (**Methods and**
994 **Extended Data Tables 5 and 6**).

995 *Results:* For single trait analyses, the improvement in performance for PRSet vs the whole-
996 genome methods were reduced, and in some cases (BMI and LDL) lassosum performed better
997 than PRSet. For all the traits assessed, phenotypic variance explained by PRSice performance
998 was the lowest (**Supplementary Fig. 6**).



999

1000 **Supplementary Fig 6.** Performance of PRSet vs genome-wide PRS methods for prediction of
 1001 single traits.

1002

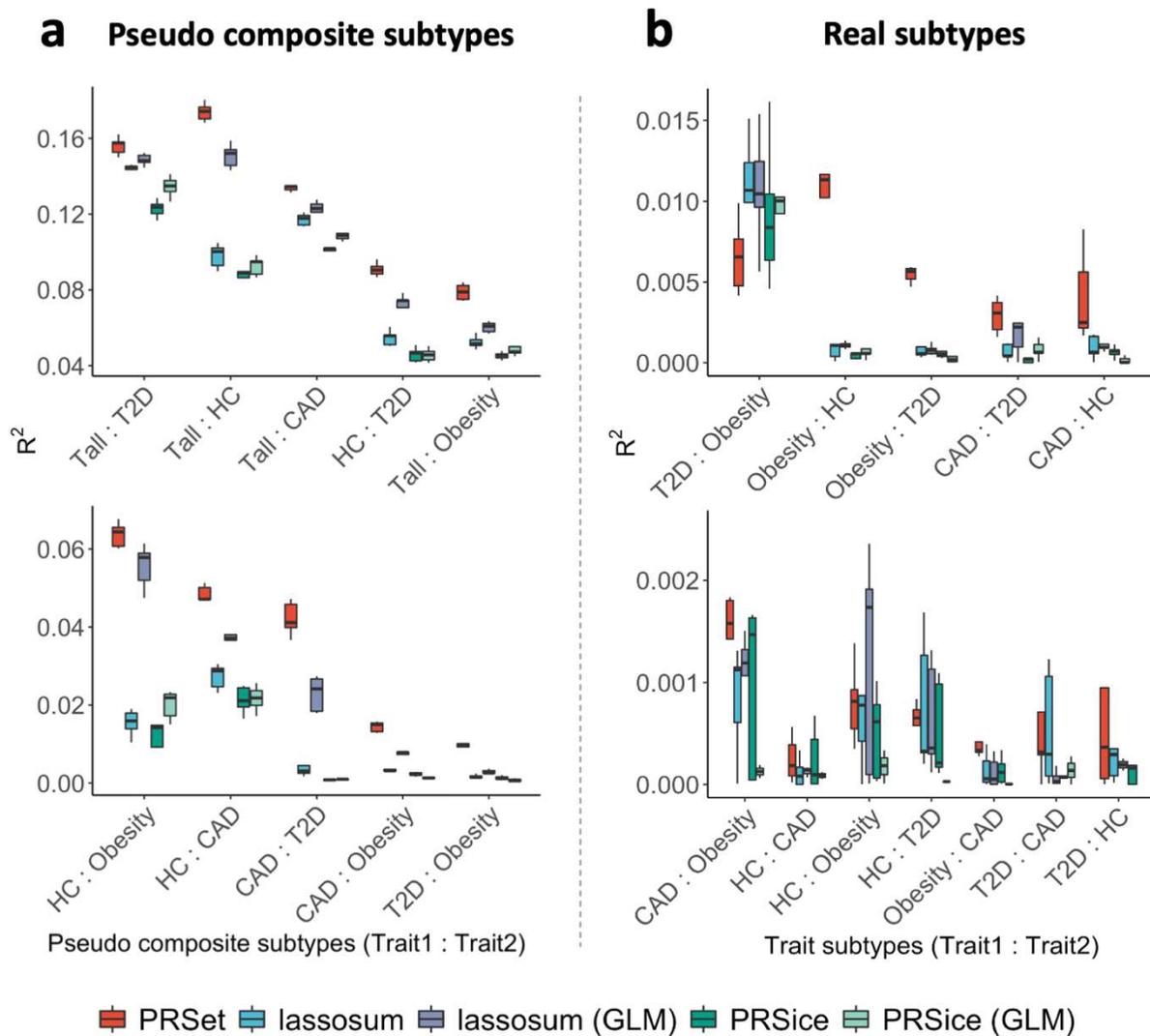
1003 **Supplementary Note 4. Performance of PRSet vs genome-wide PRS methods**
1004 **with lasso regression**

1005 *Rationale:* The inclusion of multiple pathway specific PRSs into a lasso regression may contribute
1006 to the improvement on PRSet classification performance. To test whether the regularization
1007 method also improves classification in genome wide PRSs, we calculated genome wide PRSs
1008 with lassosum and PRSice, and included the PRSs into a lasso regression.

1009 *Methods:* Genome-wide PRS analyses were performed using lassosum and PRSice-2. For
1010 PRSice-2, ~500 PRSs were calculated at different *P*-value thresholds using high resolution
1011 scoring (step size of the threshold=0.001). Whereas in standard PRS calculations PRS at each
1012 *P*-value threshold includes all the SNPs with *P*-value below the threshold (i.e. each PRS is a
1013 subset of each other), in this analysis we calculated PRSs non-cumulatively, so that each PRS
1014 includes SNPs below the current *P*-value threshold, but it does not include SNPs at *P*-values
1015 higher than the previous threshold. Using non-cumulative PRSs followed by lasso regression is a
1016 better benchmarking strategy than using standard PRSs, because for PRSet the overlap of SNPs
1017 across pathways is only partial. For lassosum, PRSs were calculated at different values of penalty
1018 factor λ and soft-thresholding parameter *s*. PRSs for each software were then included in a
1019 generalized linear model with lasso regularization using the `cv.glmnet` function from the glmnet
1020 package (v4.0-2).

1021 *Results:* For the pseudo composite subtype classification analyses, the application of lasso
1022 regularization improved the performance of both lassosum and PRSice (**Supplementary Fig. 7a**).
1023 Improvement in performance was stronger for lassosum (GLM) than for PRSice, but none of the
1024 two methods outperformed PRSet, giving supportive evidence that the use of pathways improves
1025 classification.

1026 For the real subtype classification analyses, there was not a clear improvement in performance
 1027 when GLM and regularization was applied (**Supplementary Fig. 7b**). In some cases, the best
 1028 model for PRSice with GLM and regularization was the intercept, meaning that the R2 could not
 1029 be obtained. For T2D with presence/absence of hypercholesterolemia (T2D : HC), intercept was
 1030 the best model for the five-fold cross validations, therefore R2 estimates are omitted in the figure.



1031

1032 **Supplementary Fig 7.** Performance of PRSet vs genome-wide PRS methods with a generalized
 1033 linear model and lasso regularization. **a**, pseudo composite subtype classification results. **b**, real
 1034 subtype classification results. For the lassosum (GLM) and PRSice (GLM) models, whole-genome

1035 PRSs calculated with different parameters were subsequently used in a generalized linear model
1036 with lasso regularization. HC; hypercholesterolemia.

1037 **Supplementary References:**

- 1038 1. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass
1039 index in ~700000 individuals of European ancestry. *Hum. Mol. Genet.* **27**, 3641–3649 (2018).
- 1040 2. Sanchez-Roige, S. *et al.* Genome-Wide Association Study Meta-Analysis of the Alcohol Use
1041 Disorders Identification Test (AUDIT) in Two Population-Based Cohorts. *AJP* **176**, 107–118
1042 (2018).
- 1043 3. Sinnott-Armstrong, N. *et al.* Genetics of 35 blood and urine biomarkers in the UK Biobank.
1044 *Nature Genetics* **53**, 185–194 (2021).
- 1045 4. Global Lipids Genetics Consortium *et al.* Discovery and refinement of loci associated with
1046 lipid levels. *Nature Genetics* **45**, 1274–1283 (2013).
- 1047 5. Inouye, M. *et al.* Genomic Risk Prediction of Coronary Artery Disease in 480,000 Adults:
1048 Implications for Primary Prevention. *Journal of the American College of Cardiology* **72**, 1883–
1049 1893 (2018).
- 1050 6. Marioni, R. E. *et al.* GWAS on family history of Alzheimer’s disease. *Transl Psychiatry* **8**, 1–7
1051 (2018).

1052

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [ExtendedDataTables.xlsx](#)