

# A genetic map of human metabolism across the allele frequency spectrum

Martijn Zoodmsa<sup>1</sup>, Carl Beuchel<sup>1</sup>, Summaira Yasmeen<sup>1</sup>, Leonhard Kohleick<sup>1</sup>, Aakash Nepal<sup>1</sup>, Mine Koprulu<sup>2</sup>, Florian Kronenberg<sup>3</sup>, Manuel Mayr<sup>4</sup>, Alice Williamson<sup>1,2</sup>, Maik Pietzner<sup>1,2\*</sup> & Claudia Langenberg<sup>1,2\*</sup>

<sup>1</sup> Computational Medicine, Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Berlin, Germany

<sup>2</sup> Precision Healthcare Institute, Queen Mary University of London, London, UK

<sup>3</sup> Institute of Genetic Epidemiology, Medical University of Innsbruck, Innsbruck, Austria

<sup>4</sup> National Heart and Lung Institute, Imperial College London, London, United Kingdom

Corresponding authors:

Maik Pietzner (maik.pietzner [at] bih-charite.de)

Claudia Langenberg (claudia.langenberg [at] qmul.ac.uk)

\*These authors contributed equally

## Abstract

Genetic studies of human metabolism identified unknown disease processes and novel metabolic regulators, but have been limited in scale and allelic breadth. Here, we provide a data-driven map of the genetic regulation of circulating small molecules and lipoprotein characteristics (249 metabolic traits) measured using protein nuclear magnetic resonance spectroscopy (<sup>1</sup>H-NMR) across the allele frequency spectrum in ~450.000 individuals. In trans-ancestry analyses, we identify 29,824 locus–metabolite associations mapping to 753 regions with effects largely consistent between men and women and major ancestral groups represented in UK Biobank. We develop a framework for classifying the observed extreme genetic pleiotropy, enabling identification of upstream ‘master’ regulators of lipid metabolism (‘proportional pleiotropy’), such as *ANGPTL3*. We establish rare-to-common allelic series by integrating machine-learning guided effector gene assignments with rare exonic variant analyses providing high confidence gene assignments at >100 loci, including less established regulators of lipid metabolism like *SIDT2*. At 17 such loci we observed phenotypic heterogeneity among variants mapping to the same gene indicating differential metabolic roles of the altered

medRxiv preprint doi: <https://doi.org/10.1101/2025.01.30.25321073>; this version posted February 2, 2025. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/) .

36 gene product. We identify *VEGFA* as a potential modulator of HDL-mediated risk for  
37 coronary artery disease. Our results demonstrate how rare-to-common genetic variation  
38 combined with deep molecular profiling can identify unknown and inform on poorly  
39 understood regulators of human metabolism to guide prevention and treatment of  
40 diseases.

## 41 Introduction

42  
43 Our understanding of human metabolism is mostly based on dedicated hypothesis  
44 testing in experimental settings, informed by model organisms or observations in rare  
45 diseases patients. Only recently, high-throughput profiling of small molecules in large-  
46 scale studies has enabled systematic testing of genetic variation across the genome and  
47 provided an agnostic approach for the discovery of genes that encode key metabolic  
48 regulators<sup>1-11</sup>. These efforts have provided important new insights into how genetic  
49 variation shapes human chemical and metabolic individuality<sup>1</sup> and have corroborated a  
50 large body of biochemical knowledge<sup>1,2,10,12</sup>.

51  
52 The value of such genome-metabolome-wide association studies (mGWAS) extends  
53 beyond the mapping of biochemical pathways, sometimes demonstrating almost  
54 immediate clinical value. They provided examples how readily available  
55 supplementation strategies may prevent disease or delay onset in high risk individuals,  
56 such as serine for the rare eye disorder macular telangiectasia type 2<sup>2</sup>. Others further  
57 identified unknown variants affecting the absorption, distribution, metabolism, and  
58 excretion of exogenous compounds, most importantly drugs<sup>1,13</sup>, providing pathways to  
59 mitigate adverse drug effects. However, there are several challenges that currently limit  
60 the potential of mGWAS studies, in particular for causal inference. These include 1) the  
61 still rather small number of, at most, a dozen genetic variants linked to single molecules,  
62 2) the inability to distinguish whether pleiotropic variants act on different molecules or  
63 pathways independently (horizontal pleiotropy), or whether they serve as ‘root causes’ of  
64 successive downstream changes (vertical pleiotropy), 3) the difficulty in distinguishing  
65 between locus-specific and metabolite abundance effects when colocalization at

66 disease-risk loci is observed<sup>1</sup>, and 4) the challenge of confidently assigning effector  
67 genes at newly identified loci.

68

69 Here, we integrated rare (based on whole exome sequencing) and common genetic  
70 variation with measures of 249 metabolic phenotypes, including small molecules and  
71 detailed lipoprotein characteristics, among >450,000 UK Biobank participants  
72 representing three distinct ancestries. We demonstrate largely consistent genetic  
73 regulation across ancestries and sexes for almost 30,000 locus – metabolite associations  
74 and systematically categorise abundant genetic pleiotropy. By integrating machine-  
75 learning derived effector gene assignments with rare exonic variation, we identify  
76 previously unknown regulators of metabolism and observe heterogeneity in association  
77 profiles for variants mapping to the same gene. Finally, we demonstrate how systematic  
78 integration of statistical colocalization and Mendelian randomization can identify  
79 pathways with the potential to mitigate cardiovascular disease risk beyond current  
80 approaches focused primarily on LDL-cholesterol lowering.

## 81 Results

82

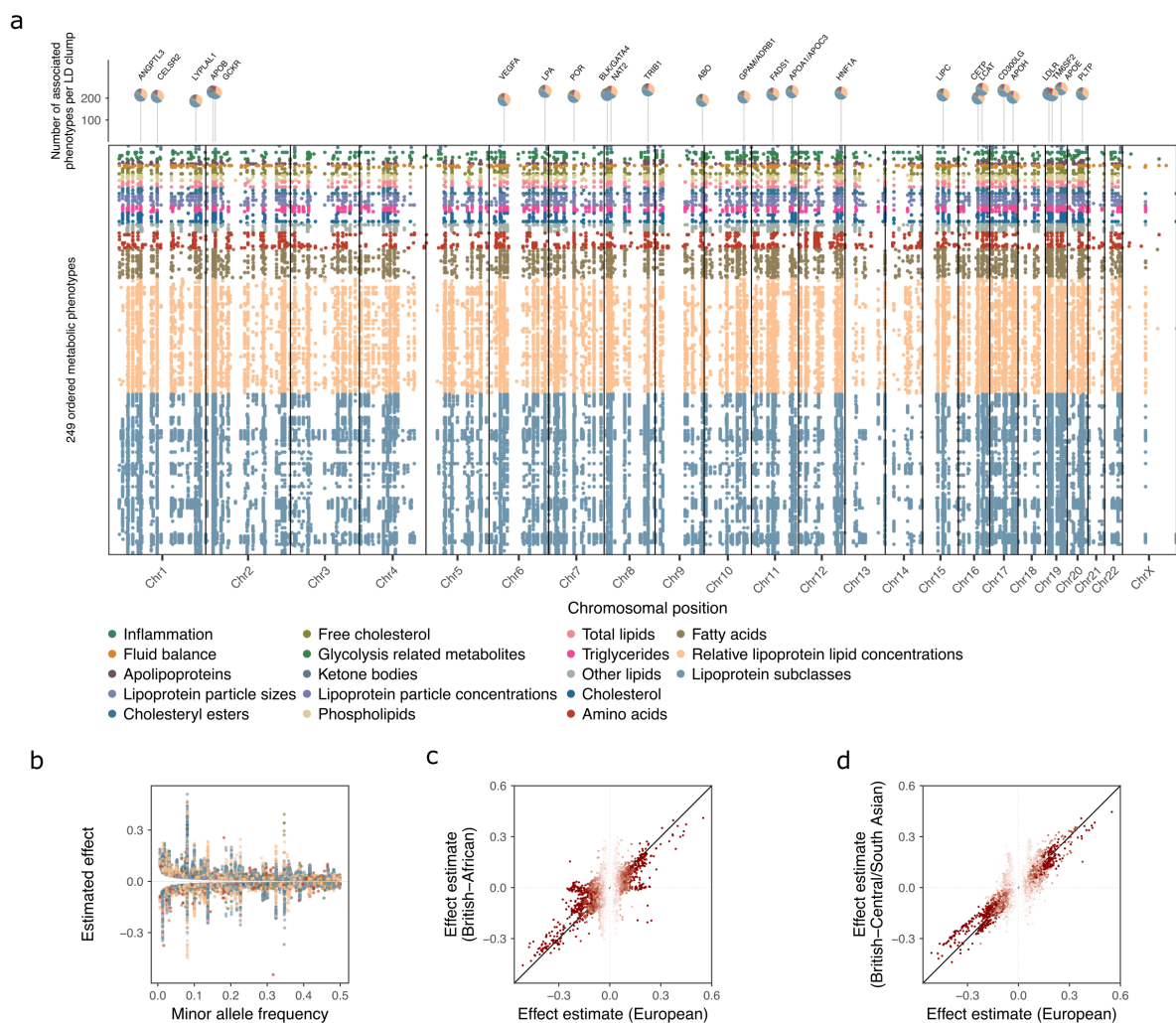
83 We integrated genome-wide association studies (GWAS; population-specific minor allele  
84 frequency (MAF) $\geq$ 0.5%) with rare exome-wide association studies (MAF $\leq$ 0.5%) on  
85 plasma concentrations of 249 metabolite phenotypes, quantified using <sup>1</sup>H-nuclear  
86 magnetic resonance spectroscopy (NMR). We included up to 450,000 UK Biobank (UKB)  
87 participants across three major ancestries (British White European – EUR (n=434,646);  
88 British African – BA (n=6,573); British Central South Asian – BSA (n=8,796);  
89 **(Supplementary Fig. 1)**. The NMR measures provided a detailed readout of lipoprotein  
90 particles along a range of lipoprotein sizes containing 14 subclasses (i.e., extra-large very-  
91 low density (VLDL) to small high-density (HDL) lipoprotein particles), along with small  
92 molecules such as amino acids and ketone bodies quantified in molar concentration  
93 units **(Supplementary Table 1)**.

94

95 ***Common genetic variation underlying circulating metabolites***

96 We identified 29,824 regional sentinel–NMR measure associations in trans-ancestral  
97 meta-analyses, representing 753 non-overlapping genomic regions (**Fig. 1a;**  
98 **Supplementary Table 2**). Nearly half of these regions (N=359, 47%) were associated with  
99 more than 10 NMR measures, demonstrating considerable pleiotropy cutting across  
100 metabolite classes for 350 regions. Characteristics of large HDL particles, such as  
101 concentration, particle size and (phospho)lipid, cholesterol, cholesteryl ester and  
102 triglyceride content, were associated with the largest number of regions (median: 166,  
103 IQR: 126-195), compared to median of 105 associated regions observed across all NMR  
104 measures (IQR: 68-142). Findings that considerably extended previous work<sup>3</sup> and  
105 replicated parallel efforts using UK Biobank<sup>9</sup> (**Supplementary Fig. 2**). Genes with well-  
106 characterised roles in human metabolism were significantly enriched among the closest  
107 genes to regional sentinels across different significance bins (adjusted p-values < 4.24 x  
108 10<sup>-9</sup>; **Supplementary Fig. 3**). This suggests that ever-larger studies of often considered  
109 omnigenic traits, such as metabolites, still yield biological plausible findings and not  
110 merely non-specific upstream regulators.

111  
112 Almost all regional sentinel associations (n=29,410, 98.6%) showed little evidence of  
113 heterogeneity (p>10<sup>-4</sup>) across ancestries. To rule out possible artefacts that might have  
114 masked ancestral-specific effects (e.g., variant coverage and statistical power in the  
115 smaller ancestry groups), we repeated the GWAS within each ancestry separately.  
116 Consistent with the trans-ancestral meta-analysis, we observed high correlations of  
117 effect estimates for regional sentinels identified in the largest subgroup of White  
118 European participants when compared to those of African and Central South Asian  
119 ancestry (**Fig. 1c-d; Supplementary Table 3, Supplementary Fig. 4**). Although we note  
120 that the limited sample size did not permit comprehensive replication, our ancestry-  
121 specific analyses also revealed one locus not seen in European participants. The  
122 previously reported<sup>14</sup> missense variant rs3211938 within *CD36* which is common in  
123 people of African ancestry (MAF<sub>BA</sub> = 0.12) but absent in European ancestry (MAF<sub>EUR</sub> = 0.0),  
124 was significantly associated (p-values < 1.49 x 10<sup>-10</sup>) with lower plasma concentrations of  
125 omega 3 fatty acids and 15 other NMR measures, including lipoprotein particle  
126 characteristics. This is in line with the role of *CD36* as a fatty acid translocase, facilitating  
127 the recognition and uptake of long-chain fatty acids.



128  
 129 **Figure 1: Common genetic regulation of circulating metabolites. A)** Top-down Manhattan plot  
 130 showing trans-ancestral sentinel variants for 249 metabolic phenotypes at a metabolome-adjusted  
 131 genome-wide significance threshold of  $p < 2.0 \times 10^{-10}$ . Each row represents an NMR measure, coloured  
 132 for biochemical class, chromosomal positions are shown on the x-axis. **B)** Weighted average allele  
 133 frequency compared to estimated effect size for trans-ancestral sentinel variants. Points are coloured  
 134 for biochemical classification. **C)** Comparison of effect sizes between White European samples (x-  
 135 axis) and British-African samples (y-axis). We considered variants that were significant in either  
 136 population. **D)** Similar to **C)** but comparing British-Central /South Asian samples. Dots are coloured  
 137 according to their absolute Z-score in White European samples.

### 139 **Refinement of regional associations through multi-ancestry fine-mapping**

140 We next employed a two-stage strategy to refine regional associations to a small number  
 141 of candidate causal variants. Firstly, we implemented fine-mapping in the largest group  
 142 of European-ancestry participants. We then further refined the subset of loci with at least  
 143 suggestive evidence across ancestries ( $p < 10^{-4}$ ) using trans-ancestral fine-mapping,  
 144 leveraging the differential blocks of linkage disequilibrium (LD) despite vastly different  
 145 sample sizes.

146

147 We first identified 3,007 statistically independent metabolite quantitative trait loci  
148 (mQTLs) associated with one or more NMR measure, representing a total of 43,322  
149 credible set – NMR measurement pairs (**Supplementary Table 4**). This successfully  
150 defined 16,170 credible sets with a high-confidence variant (posterior inclusion  
151 probability (PIP) > 0.5). Among these were low-to-common frequent variants with  
152 functional consequences in metabolic genes, such as rs78734745 (MAF=0.8%;  
153 PIP=67.9%), a splice donor variant for *ME1*, associated with plasma citrate levels (beta=-  
154 0.11; p-value<1.6x10<sup>-21</sup>). Lead fine-mapped mQTLs for a given NMR measure explained,  
155 on average, 6.9% (range: 0.57% - 13.42%) of the variance in plasma concentrations  
156 (**Supplementary Fig. 5**).

157

158 Secondly, we leveraged the different LD-block structure among participants of British  
159 African and British Central South Asian ancestry to further refine a total of 3,336 credible  
160 sets that still contained >1 variant and for which the locus had at least suggestive  
161 evidence for significance in either ancestry (P < 1.0 x 10<sup>-4</sup>). Trans-ethnic fine-mapping led  
162 to an increase in the number of credible sets containing high-confidence variants  
163 (Europeans: 997, multi-ancestral: 1,794) and decreased the median credible set size  
164 from 9 to 4 variants, while increasing the median posterior inclusion probability from 0.06  
165 to 0.16 (**Supplementary Fig. 6**). This included 1,107 (33.7%) credible sets with two or  
166 fewer variants, and 1,518 (45%) credible sets that were reduced in size by more than half.  
167 We note, however, that most eligible European credible sets were already comparatively  
168 small (median 9 variants), but sometimes still spanned multiple genes.

169

170 For example, a signal associated with mono-unsaturated fatty acids (MUFA)  
171 concentrations at 17q21.2 contained 76 genetic variants spread across several genes  
172 covering a 1Mb window in the European-only discovery. The signal was fine-mapped to  
173 as few as 4 variants (two intergenic, one <50kb distance to the gene body) after  
174 incorporating evidence from other ancestries (**Supplementary Fig. 6**). Three of these four  
175 variants mapped to the *PTRF/CAVIN1* gene, which plays a crucial role in the formation of  
176 caveolae that are particularly abundant in adipocytes. Thus, *PTRF/CAVIN1* has been

177 linked to generalized lipodystrophies<sup>15</sup>, providing a biologically plausible effector gene at  
178 this locus through trans-ancestral refinement of the credible set.

179

### 180 ***Sex-differential effects at loci encoding metabolic genes***

181 Many aspects of metabolism are known to vary by sex<sup>16,17</sup>, but only few genetic loci have  
182 been identified that may explain such differences<sup>18,19</sup>. While we observed highly  
183 correlated effect sizes across female and male participants (median  $R^2$ : 0.98, range: 0.90  
184 – 0.99), we also identified 360 putative sex-differential loci for 239 metabolic traits,  
185 representing 1,800 heterogeneous associations in sex-stratified meta-analyses  
186 (heterogeneity p-value  $< 5 \times 10^{-8}$ , see **Methods**). To rule out that sex-differential effects  
187 could be explained by other factors that differ between the sexes, we performed  
188 additional analyses identifying that sex-differential effects at one-third of loci ( $n=625$ ,  
189 34.7%) were attenuated when controlling for factors such as body mass index, tobacco  
190 use, alcohol intake, and the use of lipid-lowering or diabetes medication  
191 (**Supplementary Fig. 7, Supplementary Table 5**). For loci unaffected by such additional  
192 factors, effect estimates were generally directionally concordant between the sexes but  
193 showed differences in magnitude (**Fig. 3a**). This is consistent with results previously  
194 observed for proteomics<sup>20</sup> and suggests that the majority of significant sex interactions  
195 do not reflect sex-discordant effects. We observed pleiotropic sex-differential loci  
196 associated with 30 or more NMR measures near established lipoprotein genes (*APOE*,  
197 *APOC1*, *LPL*) but also less established genes (*SIDT2*), where sex was the most likely  
198 modifying factor. These finding may help to better understand sex-specific cut-offs in  
199 cardiovascular risk assessment in clinical guidelines to initiate treatment with lipid  
200 lowering medication<sup>21</sup>. We found *CPS1* on 2q34 to show the strongest sex differences, in  
201 line with previous reports<sup>18</sup>, with effect sizes for glycine being twice as large in females  
202 compared to males (rs1047891, beta females = 0.77, beta males = 0.34 s.d. units).

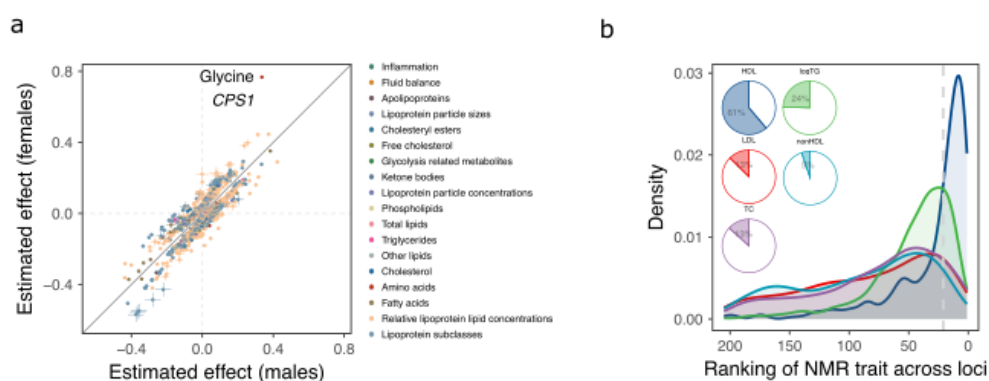
203

### 204 ***Biological reclassification of established 'lipid' loci***

205 To assess the value of metabogenomic studies involving lipoprotein profiling based on  
206 <sup>1</sup>H-NMR spectroscopy over standard clinical markers, we systematically classified the  
207 NMR metabolome association profiles for 1,657 genetic variants reported for commonly  
208 measured clinical markers (LDL-cholesterol, HDL-cholesterol, total cholesterol and

209 triglycerides) by the Global Lipids Genetics consortium (GLGC) in 1.6 million samples<sup>22</sup>.  
 210 Around 25% of associated variants had the corresponding NMR measure among the top  
 211 10% of the most strongly associated NMR measures, with 22.5% of genetic variants  
 212 showing significantly stronger associations with refined lipoprotein measures compared  
 213 to their matching measure on the NMR platform, an observation most pronounced for  
 214 non-HDL and LDL-cholesterol concentrations (**Fig. 3b**). While this indicated that relevant  
 215 loci for lipoprotein metabolism can be discovered using readily available clinical  
 216 measurements, it also demonstrates the necessity of refined lipoprotein profiles for  
 217 better understanding the relevant biological pathways, including any inference about  
 218 druggability or use for genetic causal inference methods. One such example was the  
 219 *PNPLA3* locus (tagged by rs3747207, associated with LDL-cholesterol by the GLGC;  $p =$   
 220  $2.3 \times 10^{-21}$ ,  $\beta = -0.014$ ), where we observed no evidence of association with LDL-  
 221 cholesterol ( $\beta = -0.001$ ,  $p = 0.49$ ) but LDL particle size ( $\beta = 0.045$ ,  $p\text{-value} = 1.04 \times 10^{-}$   
 222  $73$ ), and multiple characteristics of extra-large VLDL particles (**Supplementary Fig. 8**). The  
 223 intronic rs3747207 variant is in strong LD ( $r^2 = 0.98$ ) with the well-known missense variant  
 224 rs738409 (p.I148M) that has been demonstrated to confer hepatic lipid accumulation by  
 225 altering ubiquitination of patatin-like phospholipase domain-containing protein 3  
 226 (*PNPLA3*) encoded by *PNPLA3*<sup>23</sup>. Our results provide human genetic support for a  
 227 recently proposed role of *PNPLA3* in the secretion of large VLDL particles<sup>24</sup>. The  
 228 association with LDL-cholesterol in massive scale studies likely being a distant  
 229 downstream consequence.

230



231

232 **Figure 3: Putative sex-differential loci and reclassification of established lipid loci. A)**  
 233 Comparison of effect sizes of putatively sex-differential loci (defined as loci with heterogeneity  $p$ -  
 234 value  $< 5 \times 10^{-8}$  in a meta-analysis across the sexes). **B)** Rank distributions for each of the five matching  
 235 NMR traits compared to the Lipids Genetics traits across genetic loci. Per locus – trait combination,



236 205 lipid-related NMR traits were ranked based on their absolute effect size and compared to the NMR  
237 trait that corresponds the Lipids Genetics consortium trait. Pie charts show the percentage of loci  
238 where the corresponding NMR trait is ranked among the top 10% of associated traits.  
239

#### 240 **Machine-learning guided effector gene assignment**

241 Assigning effector genes to genetic variants remains one of the most important  
242 bottlenecks for translating GWAS results into tangible insights. We assigned effector  
243 genes for almost three-quarters of European fine-mapped mQTLs (73.6%; n=2,213) with  
244 at least moderate confidence (candidate gene score  $\geq 1.5$ , range 0 to 3), including about  
245 28.2% with high-confidence assignments (score  $\geq 2$ ; n=848), by training a machine  
246 learning model that integrates functional genomic resources with pathway information  
247 inspired by the ProGeM framework<sup>25</sup> (**Supplementary Table 6**). For example, we  
248 prioritised the fatty acid elongase gene *ELOVL6* for 16 different NMR measures (tagged by  
249 rs3813829), including the fraction of cholesterol and other fatty acids on very small VLDL  
250 and very large HDL particles in addition to the fraction of saturated fatty acids. The gene  
251 product, ELOVL fatty acid elongase 6, catalyses the rate-limiting step in long-chain fatty  
252 acid elongation, which are subsequently incorporated into lipoprotein particles. We also  
253 prioritized genes with upstream roles in metabolism, including a locus on 17q25.3 where  
254 we prioritized cytohesin-1 (*CYTH1*) as the candidate causal gene for five independent,  
255 genetic variants linked to 11 distinct NMR measures mostly comprising characteristics  
256 of VLDL particles. *CYTH1*, previously associated with type 2 diabetes<sup>26</sup>, promotes  
257 activation of ADP-ribosylation factors (ARF)1, ARF5 and ARF6, regulators of lipid vesicle  
258 transport, membrane lipid composition and modification<sup>27</sup>, demonstrating a relevant but  
259 indirect link to lipoprotein metabolism.

260  
261 We observed considerable overlap of machine-learning guided effector gene predictions  
262 (top three genes) with those reported based on manually curated biological plausibility  
263 (191 out of 283 loci)<sup>3</sup> or based on colocalization with protein quantitative trait loci that  
264 have not been used to train the algorithm<sup>28</sup> (81 out of 143; **Supplementary Table 6**). While  
265 missing overlap indicates room for improvement, 24 high-confidence assignments did  
266 strongly disagree with either external source (gene score  $> 2$  but no match among pQTL  
267 prioritised or manually curated ones). This included a locus on chromosome 19q13.11  
268 tagged by rs62102718 for which we prioritised *PEPD* with high-confidence (score=2.42)

269 as opposed to *CEBPA*<sup>3</sup>. *PEPD* encodes peptidase D, highly relevant for collagen turnover,  
270 that has been shown to promote adipose tissue fibrosis in mouse knock-out models and  
271 promoting insulin resistance<sup>29</sup>. Insulin resistance, in turn, being a very plausible  
272 explanation for the pleiotropic effect of the variant on diverse lipoprotein characteristics  
273 (n=31).

274

### 275 ***Tissue distribution of effector genes***

276 We next tested tissue-specific expression patterns of the identified effector genes to  
277 better understand organ sites contributing to (lipoprotein) metabolism. Strong clustering  
278 was observed at both the tissue and metabolite levels, reflecting both known and less  
279 established organ contributions (**Supplementary Fig. 9a, Supplementary Table 7**).  
280 Genes characteristic of the liver, adipose tissue, adrenal gland, but also female breast  
281 tissue (likely reflecting its high adipose tissue content) were significantly enriched among  
282 effector gene sets across the metabolic measures captured by NMR. This included  
283 significant enrichment of all amino acids in liver tissue (e.g., phenylalanine: odds ratio  
284 (OR): 14.8,  $p < 1.3 \times 10^{-8}$ , histidine: OR 7.9,  $p < 2.9 \times 10^{-11}$ ) but also for skeletal muscle in  
285 alanine metabolism (OR:3.82;  $p\text{-value} < 7.9 \times 10^{-9}$ ). Similar enrichments were observed  
286 when using the closest gene instead of our annotated effector genes for mQTLs  
287 (**Supplementary Fig. 9b**).

288

### 289 ***Modes of metabolic and systemic pleiotropy***

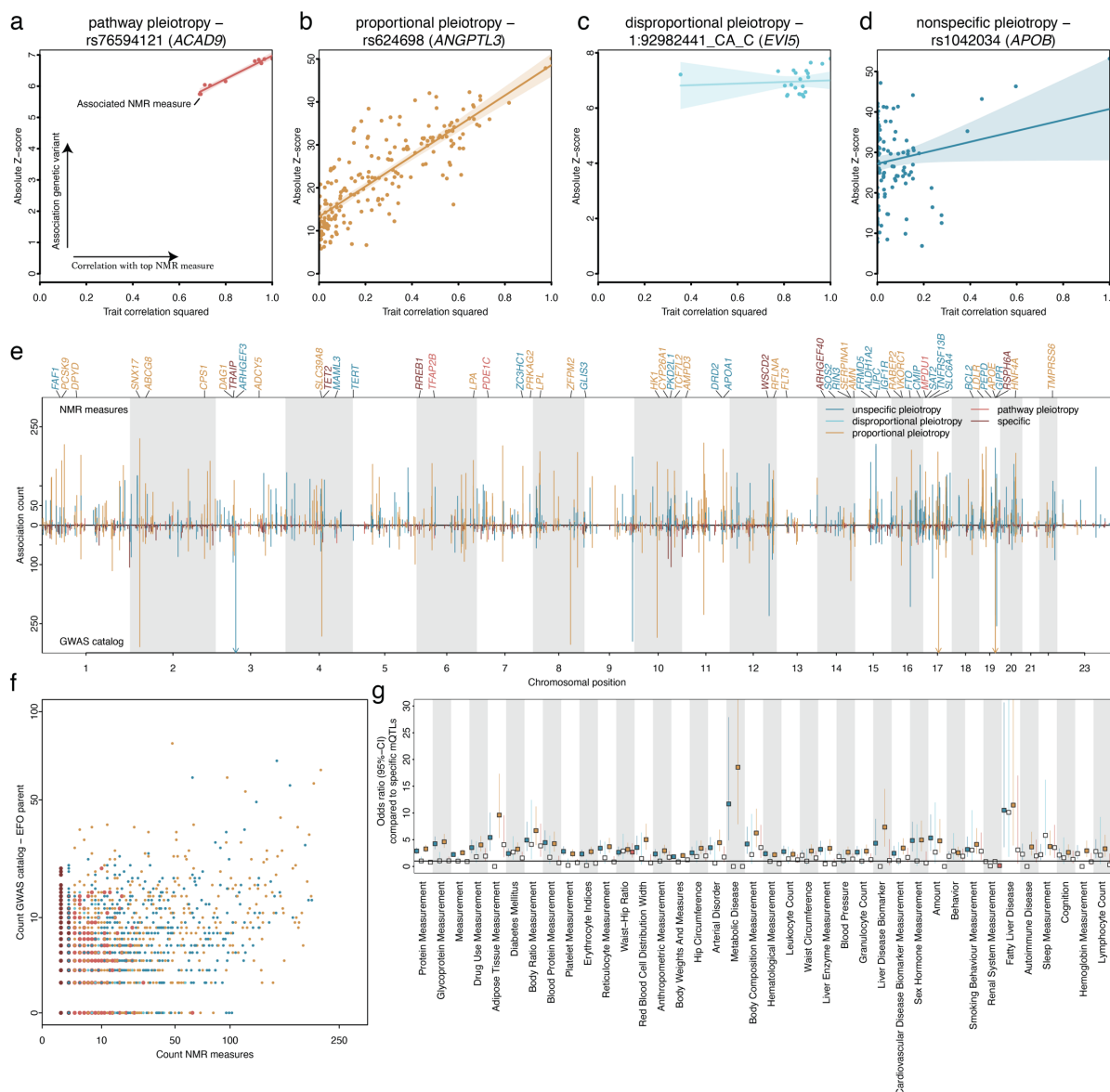
290 Pleiotropy is a widespread but poorly understood phenomenon and we developed a  
291 framework to characterise four different modes of metabolic pleiotropy for all fine-  
292 mapped mQTLs (**Fig. 4a-d; Supplementary Fig. 10 and Table 6; see Methods**). About  
293 half of the pleiotropic mQTLs (n=880;  $\geq 2$  NMR measures) showed evidence for two  
294 different modes of vertical pleiotropy. Firstly, within confined pathways (n=218; ‘pathway  
295 pleiotropy’) or, secondly, as a function of the correlation with the ‘lead’ NMR measure  
296 (n=662; ‘proportional pleiotropy’; **Fig. 4a**). For example, rs76594121 tagged an mQTL at  
297 3q21.3 associating with different characteristics of large HDL particles, for which we  
298 prioritized *ACAD9* as the most likely candidate gene (**Fig. 4a**). The gene product of *ACAD9*,  
299 acyl-CoA dehydrogenase family member 9, is part of complex I of the respiratory chain  
300 that catalyses the oxidation of fatty acids with a high affinity for long chain fatty acids that

301 are, amongst others, carried by HDL particles. A prototype example for ‘proportional  
302 pleiotropy’ was an mQTL tagged by rs624698 for which we prioritized *ANGPTL3* as the  
303 effector gene (**Fig. 4b**). Angiopoietin-like 3, encoded by *ANGPTL3*, inhibits lipoprotein  
304 lipase activity but also endothelial lipase, resulting in increased triglycerides, HDL-  
305 cholesterol, and phospholipid concentrations, consistent with HDL-particle  
306 characteristics being the most strongly associated NMR measures ( $p < 1.0 \times 10^{-546}$ ). Other  
307 associations being downstream effects on lipoprotein metabolism rather than acting on  
308 independent pathways (**Fig. 4b**), considerably expanding previous genetic  
309 observations<sup>30</sup>.

310  
311 The remaining half of pleiotropic mQTLs showed evidence for two modes of horizontal  
312 pleiotropy: those with evidence for ‘disproportional pleiotropy’ (n=68) and a larger group  
313 with evidence for ‘nonspecific pleiotropy’ (n=720). For example, a small deletion on  
314 chromosome 1 (chr1:92982441:CA>C) was associated with a highly correlated cluster of  
315 NMR measures, including characteristics of IDL, LDL, and VLDL particles (**Fig. 4c**), but  
316 for which we detected no correlation of association strengths according to the lead NMR  
317 measure, the concentration of esterified cholesterol in medium-sized VLDL particles  
318 ( $p < 6.8 \times 10^{-14}$ ). We prioritized *EVI5* as the most likely candidate gene, supported by  
319 previous studies on rare functional variants<sup>31</sup>. The gene product of *EVI5*, ecotropic viral  
320 integration site 5, has no apparent link to (lipoprotein) metabolism in line with most of the  
321 gene assignments for mQTLs with a similar nonspecific pleiotropy pattern. An example  
322 of ‘nonspecific pleiotropy’ was the *APOB* missense variant rs676210 (p.Pro2739Leu)  
323 associated with 126 NMR measures across the entire lipoprotein density range, but also  
324 creatinine and glycoprotein acetyl concentrations (**Fig. 4d**). The differential effects of the  
325 same genetic variation on distinct lipoprotein subgroups aligns with changes in lipid  
326 profiles seen with mipomersen, an antisense oligonucleotide against *APOB*, that  
327 demonstrated reductions in LDL-cholesterol but also subsequent increases in the  
328 triglyceride content of VLDL particles as hepatic adaption occurs<sup>32</sup>.

329  
330 Modes of molecular pleiotropy only partially translated into pleiotropy across the entire  
331 breath of phenotypes and diseases studied genetically (**Fig. 4e**). We observed a two-fold  
332 enrichment of ‘proportional pleiotropic’ (OR: 2.11;  $p < 2.0 \times 10^{-14}$ ) and to a lesser extent an

333 enrichment of ‘nonspecific pleiotropic’ (OR: 1.52;  $p < 1.1 \times 10^{-5}$ ) variants among variants  
334 reported in the GWAS catalog for  $\geq 5$  non-metabolomic trait categories (**see Methods**). In  
335 contrast, the set of pleiotropic GWAS catalog variants was significantly depleted for  
336 ‘specific’ mQTLs (odds ratio: 0.42;  $p < 1.6 \times 10^{-21}$ ). Some phenotypically specific variants  
337 thereby provided clues to understand non-specific molecular pleiotropy (**Fig. 4f**). For  
338 example, rs8101064, an intronic variant in *INSR*, encoding the insulin receptor, has been  
339 reported for type 2 diabetes among East Asians<sup>33</sup> and was associated with 40 NMR  
340 measures in a nonspecific manner, likely reflecting the broad effects of insulin resistance  
341 on whole body lipid metabolism. Systemic mechanisms explaining effects of  
342 ‘proportional’ and ‘nonspecific’ pleiotropic mQTLs were further evidenced by a more than  
343 20-fold significant enrichment of associated trait categories such as ‘metabolic disease’,  
344 ‘fatty liver disease’, and ‘arterial disorders’ (**Fig. 4g**).



345  
 346 **Figure 4: Modes of pleiotropy. a-d)** Exemplary scatterplots opposing the squared trait correlation of  
 347 the lead NMR measure for the listed variant against the absolute Z-score from linear regression  
 348 models for all associated NMR measures. The colours indicate different modes of pleiotropy and  
 349 correspond to the legend in e). For each plot, a linear regression fit with 95%-confidence interval is  
 350 given. **e)** Number of associated NMR measures for each of 3007 mQTL groups opposed to  
 351 associations reported in the GWAS catalog after pruning the GWAS catalog for metabolic phenotypes  
 352 (see Methods). Colouring is according to modes of pleiotropy. **f)** Scatterplot opposing the number of  
 353 associated NMR measures (x-axis) of each mQTL group with the number of reported EFO parent  
 354 categories in the GWAS catalog. **g)** Odds ratios and 95%-confidence intervals from logistic regression  
 355 models testing whether EFO categories (x-axis) are more frequently reported for pleiotropic mQTL  
 356 groups compared to specific ones. Darker colours indicated estimates passing corrected statistical  
 357 significance.

358  
 359 **Convergence of common and rare genetic variation shaping metabolism**

360 Previous investigations focussed on either large-scale common variant<sup>1-4,8</sup> or  
 361 comparatively small-scale rare exonic variant discovery efforts<sup>7,34</sup>, but these approaches

362 have not been able to integrate such information at scale to establish allelic series that  
363 confidently link genes to metabolites, including previously unknown regulators of  
364 metabolism. We identified rare variation ( $MAF \leq 0.05\%$ ) in a total of 209 genes to be  
365 significantly ( $p < 1.1 \times 10^{-8}$ ) linked to one or more of 249 NMR measures combining ultra-  
366 rare gene burden analysis (3,709 significant associations; **Supplementary Table 8**) and  
367 rare exonic variant analysis (4,131 significant associations; **Supplementary Table 9**).  
368 Effect sizes were significantly larger compared to more frequent variant effects (**Fig. 5a**).  
369 For example, people carrying rare predicted loss-of-function variants in *SLC13A5* had  
370 more than 1.4 s.d. units higher plasma citrate concentrations per copy of the possibly  
371 damaging allele ( $\beta: 1.41$ ;  $p\text{-value} < 2.6 \times 10^{-20}$ ).

372  
373 We also observed considerable pleiotropy, including 47 genes associated with 20 or more  
374 NMR measures. Many of these genes have well-known roles in metabolism or small  
375 molecule transport, such as half ( $n=23/51$ ) of the genes being involved in (peripheral)  
376 cholesterol metabolism (**Supplementary Fig. 11**). On the other hand, rare pleiotropic  
377 variants with large effect sizes ( $MAF < 0.02\%$  and  $\beta > 0.6$  s.d. units) pointed towards  
378 less-established regulators of metabolism including *SIDT2* (chr11:117186662:C>T,  
379  $n=124$  carriers), *JAK2* (chr9:5073770:G>T,  $n=73$  carriers) or *CEP164*  
380 (chr11:117356670:C>G,  $n=49$  carriers). Experimental work already suggested a role for  
381 the gene product of *SIDT2* (SID1 transmembrane family member 2) in hepatic lipid  
382 metabolism and apolipoprotein A1 (ApoA1) secretion, the main protein component of  
383 HDL particles which constituted the majority of associated NMR measures (**Fig. 5b**)<sup>35,36</sup>.  
384 In contrast, associations with *JAK2* variants indicate a link to clonal haematopoiesis of  
385 indeterminate potential (CHIP)<sup>37</sup> with uncertain causality.

386  
387 We observed strong overlap between our gene burden and common variant findings, with  
388 85.4% of rare variant ( $n=3528$ ) and 75.5% of gene burden ( $n=2802$ ) associations being no  
389 more than 100kb away from the nearest statistically independent lead credible set  
390 variant (**Fig. 5c**). In contrast, most common variant findings (92.3%) were not within  
391 500kb of matching rare variant/burden evidence . Notably, 12.1% of gene burden results  
392 were more than 1Mb away from the next common credible set variant for the respective

393 NMR measure, aligning with recent observations that both approaches partly prioritise  
394 different genes<sup>38</sup>.

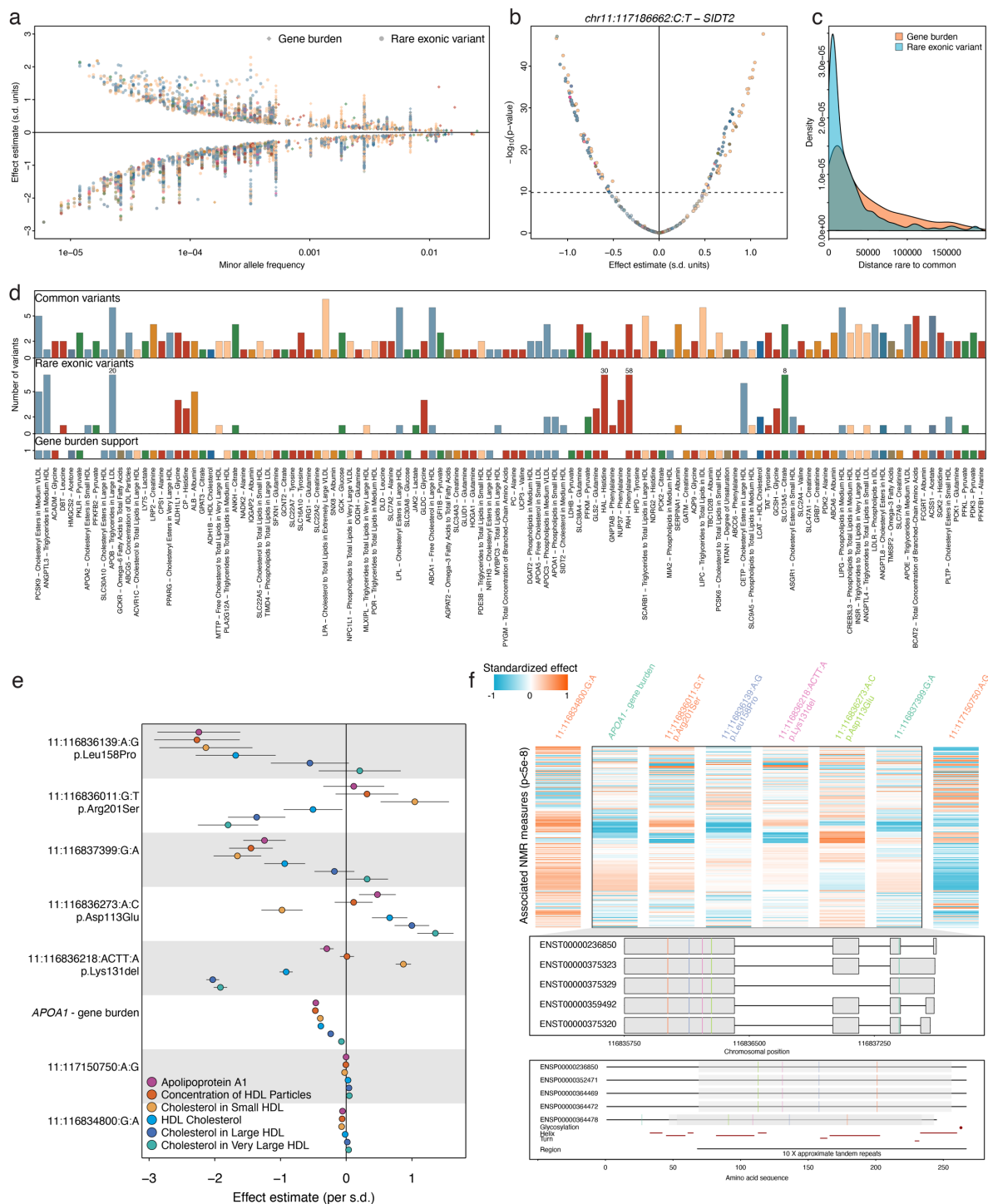
395

396 At 116 genes (55.5%), rare variant and/or burden evidence overlapped with effector gene  
397 predictions for closeby common credible set variants ( $\leq 200\text{kb}$ ) for one or more  
398 associated NMR measure (**Fig. 5d**), providing independent support for allelic series (**Fig.**  
399 **5d; Supplementary Table 10**). For example, we identified an allelic series composed of  
400 7 rare loss-of-function (LoF), 1 gain-of-function (GoF), and 4 common variants for serum  
401 citrate levels at *SLC13A5* encoding a sodium-dependent citrate co-transporter. Another  
402 allelic series at *ANKH* comprised four common variants (rs185448606 – MAF=1.3%;  
403 rs17250977 – MAF=4.0%; rs826351 – MAF=44.3%; rs2921604 – MAF=45.9%) and a rare  
404 missense variant chr5:14745916:T>C (MAF=0.0069%) being also associated with lower  
405 serum concentrations of citrate (beta=-2.18 s.d. units,  $p < 5.2 \times 10^{-11}$ ) (**Fig. 5d**). *ANKH*  
406 encodes for a multipass transporter, recently shown to transport citrate<sup>39</sup>, with an  
407 important role in bone health<sup>39</sup>.

408

409 We observed evidence that genetic variants even within an allelic series had differential  
410 metabolic consequences, covering a total of 17 genes associated with  $\geq 10$  NMR  
411 measures (**Supplementary Table 10**). The most outstanding example included 7 variants  
412 (5 rare; 2 common) and a cumulative burden of rare predicted LoF variants mapping to  
413 *APOA1*. They distinctively associated with one or more of 87 NMR measures, most  
414 strongly with diverse characteristics of HDL particles of which the gene product,  
415 Apolipoprotein A1 (ApoA1), is the major component (**Fig. 5e-f**). This included four rare  
416 missense variants (MAF $\leq 0.03\%$ ) encoded in exon 4 that each had partly differential  
417 effects on the number, size, and cholesterol content of HDL particles (**Fig. 5e**). Only one  
418 of which (p.Leu158Pro) primarily associated with serum ApoA1 concentrations and HDL  
419 particle number, micking the association with the cumulative burden of high-confidence  
420 predicted LoF variants in *APOA1*, suggesting a potentially dysfunctional protein that lacks  
421 interaction with lecithin cholesterol acyl transferase to facilitate cholesterol uptake<sup>40</sup>. In  
422 contrast, p.Lys131del and p.Arg201Ser seemed to rather predispose to a shift in  
423 cholesterol content from large towards small HDL particles, a pattern opposed by  
424 p.Asp113Glu (**Fig. 5e**). An observation consistent with amyloid formation by ApoA1 that

425 has been observed in early case reports of p.Lys131del (historically the ApoA-I<sub>Helsinki</sub><sup>41</sup>) in  
 426 which HDL-cholesterol or ApoA1 concentrations are only mildly changed but aggregation  
 427 of misfolded ApoA1 protein can confer organ damage later in life<sup>42</sup>). Since p.Asp113Glu  
 428 and p.Arg201Ser have not yet been identified to cause amyloidosis, we cannot rule out  
 429 the possibility that each variant maps to distinctive parts of ApoA1 with subsequently  
 430 different consequences on function and/or stability (**Supplementary Fig. 12**).  
 431



432



433 **Figure 5 Rare coding variation associated with NMR measures and convergence with common**  
434 **variant associations. a)** Effect estimates against minor allele frequency (MAF) of significantly  
435 associated gene burden (diamonds;  $p < 1.2 \times 10^{-8}$  and rare exonic variants (MAF < 0.05%; circles;  $p < 2.0$   
436  $\times 10^{-10}$ ). **b)** Effect estimates and  $-\log_{10}(p\text{-values})$  for associations of the rare intronic variant  
437 chr11:117186662:C>T within *SIDT2* across all 249 NMR measures. The dotted horizontal line indicates  
438 the multiple testing threshold ( $p < 2.0 \times 10^{-10}$ ). **c)** Genomic distance between gene burden (blue) or rare  
439 exonic variants (orange) towards the next common credible set variant. **d)** Evidence for allelic series  
440 based on i) gene burden analysis (bottom panel), ii) rare exonic variants (middle panel), and iii)  
441 common variants with prioritized effector gene matching to the evidence from exonic analysis. For  
442 each gene, only the NMR measure most significantly associated with the strongest common variant  
443 is shown in case multiple NMR measures were associated. Some bars for the number of associated  
444 rare exonic variants have been capped to fit into plotting margin but the number is given in the plot. **e)**  
445 Effect estimates (dots) and 95%-confidence intervals for 7 variants mapping to *APOA1* as well as a  
446 cumulative burden of high-confidence pLOF variants within *APOA1* and bespoke circulating measures  
447 of ApoA1 and HDL particles (colour gradient). **f)** The top displays a heatmap of standardized effect  
448 estimates (per variant) across 87 NMR measures for each associated variant and a cumulative burden  
449 within *APOA1*. Variants mapping into the region encoding the protein are surrounded by a rectangle.  
450 Variant effects have been aligned to the minor allele. The middle panel maps the corresponding  
451 variants to their respective transcripts encoding different forms of *APOA1*, while the lowest panels  
452 maps missense variants onto the amino acid sequence of the protein. Variant names coloured  
453 similarly had highly correlated association profiles.

454

#### 455 ***Phenotypic consequences of rare variation in metabolic genes***

456 Rare inborn errors of metabolism are among the few disorders screened for at birth by  
457 most healthcare systems globally, as early intervention – such as appropriate  
458 substitution or dietary regimens – can prevent developmental issues and diseases later  
459 in life. We observed a more than 3-fold enrichment of genes previously linked to  
460 Mendelian diseases<sup>43</sup> ('OMIM genes') among those associated with NMR measures in  
461 gene burden and rare exonic variant analyses (odds ratio: 3.30;  $p\text{-value} < 6.5 \times 10^{-17}$ ;  
462 **Supplementary Table 11**), in line with results reported from previous mGWAS<sup>1,2,7,8</sup>. For  
463 15 out of 106 genes, we found evidence of significantly associated disease risk ( $p < 7.5 \times 10^{-$   
464  $7$ ), largely replicating signs and symptoms of corresponding rare disorders  
465 (**Supplementary Table 12**). Associations with NMR measures thereby represented  
466 different modes of action. For cardiovascular diseases, most prominently familial  
467 hypercholesterolemia (e.g., via *APOB*), they likely acted as mediators, whereas  
468 associations converging on *PKD1* for cystic kidney disease likely indicated disease  
469 consequences. We further observed less understood pleiotropic roles of OMIM genes.  
470 For example, rare predicted loss-of-function variants within *SMAD6* are known to cause,  
471 amongst others, malformations of bones, e.g., Craniosynostosis 7, characterised by

472 malformations of the skull and subsequent brain damage, and we observed a strongly  
473 increased risk for other disorders of the cervical region (OR:28.8; 95%-CI: 10.3 – 80.5; p-  
474 value<1.4x10<sup>-10</sup>), as well as significantly smaller VLDL particles (beta:-0.13; 95%-CI: -0.16  
475 - -0.09; p-value< 1.5x10<sup>-9</sup>) among rare variant carriers in UKB. The gene product, SMAD  
476 Family Member 6, suppresses TGF-beta signalling, which has known effects on bone  
477 morphogenetic proteins<sup>44</sup>. Independent evidence suggests that *SMAD6* downregulation  
478 reduces the expression of core genes involved in lipoprotein metabolism, such as *LDLR*  
479 <sup>45</sup>, that may explain the disease-unrelated association.

480

481 When we tested more generally whether a rare variant burden in metabolic genes was  
482 associated with disease susceptibility, we observed a significant enrichment among  
483 susceptibility genes for endocrine and metabolic disorders, such as type 2 diabetes and  
484 different lipidemias but not among other disease categories (**Supplementary Fig. 13**).

485

#### 486 ***Risk mitigation of atherosclerotic cardiovascular disease beyond LDL-cholesterol***

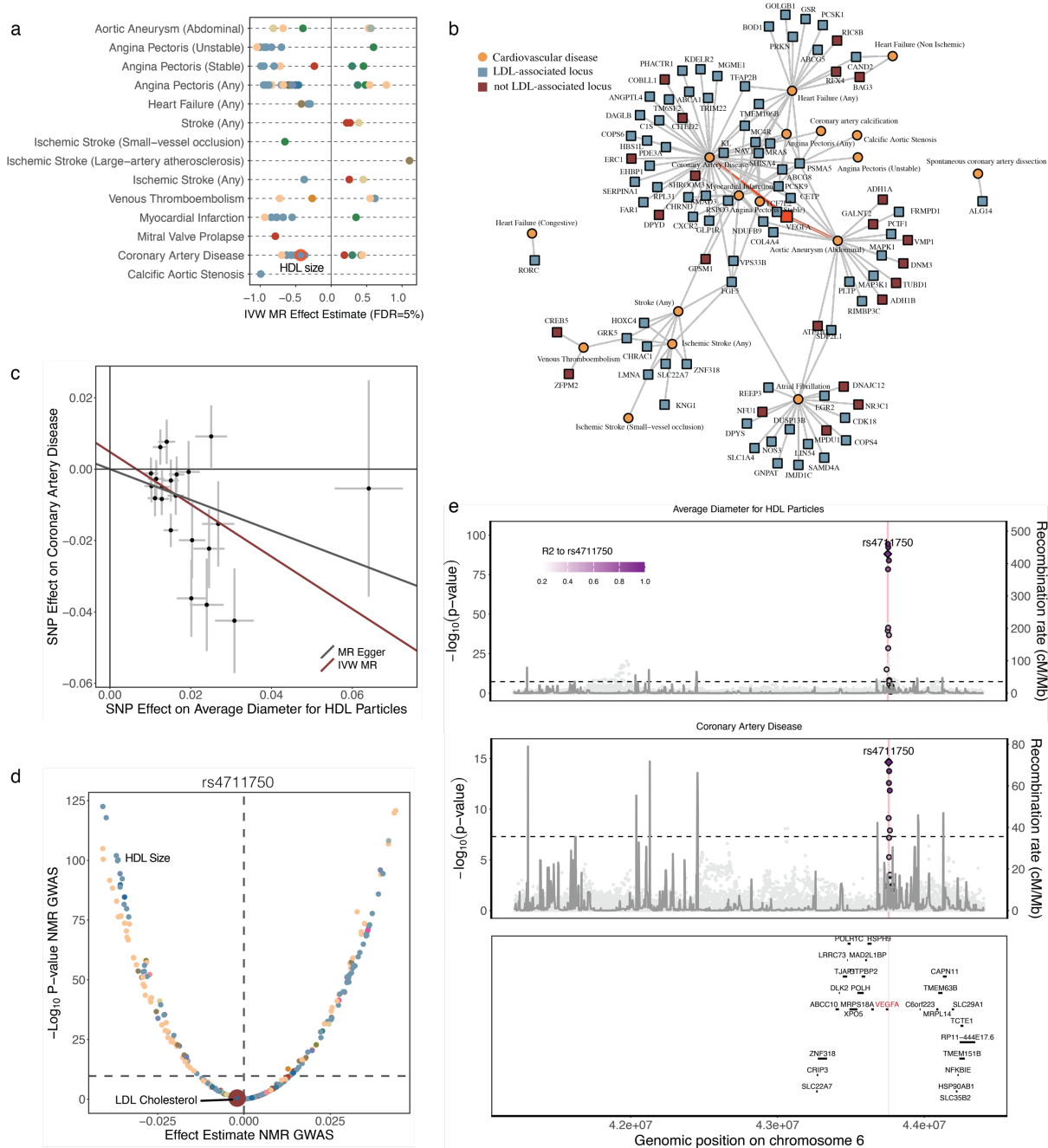
487 The success of LDL-cholesterol-lowering drugs for the prevention of atherosclerotic  
488 cardiovascular disease (ACVD) can be effectively recapitulated by genetic evidence.  
489 Genetic predisposition to high LDL-cholesterol is strongly associated with an increased  
490 risk of ACVD ('level effect'), and genetic variations that mimic potent drug targets, such  
491 as at *PCSK9*, show strong evidence of shared effects on both LDL-cholesterol and ACVD  
492 ('locus effect')<sup>46</sup>. To identify potential pathways to mitigate the residual risk not addressed  
493 by lowering of LDL-cholesterol<sup>47</sup>, we systematically integrated outcome data across 25  
494 CVD phenotypes<sup>48-62</sup>, including non-atherosclerotic diseases to test specificity of  
495 disease associations, with NMR phenotypes and assessed the convergence of locus and  
496 level effects (**Supplementary Table 13**).

497

498 We identified significant evidence (false-discovery rate (FDR)<5%) for 1,146 'level effects'  
499 across 218 NMR measures with one or more of 22 CVD phenotypes using pleiotropy  
500 curated genetic instruments in Mendelian randomization (**Fig. 6a; Supplementary Table**  
501 **14**). Independently, we observed evidence for 5,527 'locus effects', establishing a shared  
502 genetic architecture (posterior probability (PP)>80%) between 87 mQTL associated with  
503 247 NMR measures and 17 CVD phenotypes (**Fig. 6b; Supplementary Table 15**). For a

504 total of 46 NMR measure – CVD combinations we found converging evidence for level-  
505 and locus-effects, including 23 not associated in our study with parameter of LDL-  
506 metabolism (see Methods; **Fig. 6b**), providing potential alternatives for addressing  
507 residual risk (**Supplementary Table 16**).

508  
509 For example, we observed robust evidence that, among other measures related to HDL  
510 size and composition, genetic susceptibility to larger HDL particle size was associated  
511 with a 35% reduced risk of coronary artery disease (CAD; odds ratio=0.65; 95%-CI: 0.50  
512 – 0.83;  $p_{adj}<0.007$ , **Fig. 6c**) along with robust evidence of a shared and directionally  
513 concordant genetic signal tagged by rs4711750 at the *VEGFA* locus (PP = 99%, **Fig. 6e**).  
514 The locus has previously been implicated in CAD risk<sup>48</sup>, and our results now suggest that  
515 one likely pathway to modulate CAD risk might be via HDL particle size or characteristics  
516 of large HDL particles not captured by HDL-cholesterol. Vascular endothelial growth  
517 factor A (VEGFA), encoded at *VEGFA*, is primarily known for its role in angiogenesis<sup>63</sup>, but  
518 it has also been described as a regulatory factor of transendothelial transport of  
519 esterified cholesterol from HDL but not LDL particles via activation of scavenger receptor  
520 BI (SR-BI) during reverse cholesterol transport<sup>64</sup>. Inhibition of VEGFA is a major  
521 pharmaceutical target to suppress vascularisation of malignant tumours<sup>63</sup>, and agents  
522 targeting VEGF signalling are well-known for adverse cardiovascular effects<sup>65</sup>, suggesting  
523 that activation of VEGFA, rather than inhibition, might be necessary to potentially reduce  
524 CAD risk. Our observations contribute to a growing body of evidence that more tailored  
525 approaches - rather than increasing HDL cholesterol content – will likely be needed for  
526 potential cardiovascular benefits, given the discouraging trials for most agents increasing  
527 HDL-cholesterol<sup>66</sup>. We note, however, that HDL-particle size might still only be a  
528 ‘measurable’ surrogate, rather than being the true underlying mechanism. For example,  
529 inhibition of reverse cholesterol transport via dysfunctional SR-BI increased HDL particle  
530 size as well as CAD risk<sup>67</sup>.



531

532

533

534

535

536

537

538

539

540

541

542

**Figure 6 Genetic prioritisation to target residual cardiovascular risk. a)** Summary of two-sample Mendelian randomization analysis testing for putatively causal effects of NMR measures on the risk on diverse cardiovascular diseases (CVD). Shown are effect estimates for NMR – disease pairs passing multiple testing. Metabolites are coloured according to the scheme from Figure 1. **b)** Locus – disease network highlighting loci for which at least one NMR measure showed evidence of colocalization with one or more CVDs (PP≥80%). Only loci without evidence for unspecific pleiotropy are depicted. Loci were annotated with the most likely causal gene. Loci coloured in blue showed evidence for being associated with LDL-Cholesterol whereas red did not. **c)** Dose-response plot for SNPs associated with HDL particle size (after filtering for pleiotropic SNPs) against the risk for coronary artery disease. Effect estimates (dots) and 95%-Cis are given and MR-regression lines added. **d)** Effect of rs4711750 across the NMR metabolome. **e)** Locuszoom plot centred around

543 *VEGFA* demonstrating colocalization for the genetic signal for HDL particle size and coronary artery  
544 disease.

545

### 546 ***Disease-wide Mendelian randomization screen for non-lipoprotein measures***

547 Having established categories of pleiotropy for mQTLs beyond simple association  
548 counting, we finally aimed to demonstrate its application in a disease-wide screen using  
549 1394 disease outcomes from the FinnGen study<sup>68</sup> (release 11) for non-lipoprotein  
550 measures. We observed a strong decline, 29 to 13 metabolite – disease association with  
551 significant evidence (adjusted p-value < 0.05) from two-sample MR (‘level effect’) once  
552 subsetting to metabolite-specific instruments, indicating false-positive results due to  
553 pleiotropy (**Supplementary Table 17**).

554

555 We observed evidence for convergence of locus and level convergence for a risk-  
556 increasing effect of genetically predicted plasma glycoprotein acetyl concentrations on  
557 type 2 diabetes risk (odds ratio per 1 s.d. increase: 1.67; p-value<3.9x10<sup>-7</sup>). The  
558 association persisted even after additional exclusion of variants with evidence for  
559 pleiotropy in the GWAS catalog (odds ratio: 1.69; p-value<9.1x10<sup>-5</sup>). Notably, ‘locus’  
560 convergence was based on the consistent effect of the rare loss-of-function variant  
561 chr20:44413714:C>T (MAF = 0.02%) within *HNF4A* on plasma glycoprotein acetyl  
562 concentrations (beta: 0.60; p-value<8.3x10<sup>-15</sup>) and the cumulative effect of ultra-rare  
563 loss-of-function variants on type 2 diabetes risk (odds ratio: 2.68; p-value: 6.5x10<sup>-10</sup>).  
564 However, we note that plasma glycoprotein acetyl concentrations proxy a complex  
565 chronic inflammatory state<sup>69</sup> warrants further follow-up analysis to establish  
566 mechanistic links to type 2 diabetes. In contrast, previously reported associations  
567 between genetically predicted levels of branched-chain amino acids and type 2 diabetes  
568 reached at best nominal significance with a smaller effect size than previously  
569 estimated<sup>70</sup> (e.g., plasma leucine concentrations: odds ratio per s.d. unit: 1.19; p-  
570 value<0.02).

## 571 Discussion

572

573 The genetic basis of circulating metabolites provides insights into the complexion of  
574 human metabolomic regulation and its subsequent influence on health and disease. By

575 integrating common and rare genetic variation with circulating metabolite  
576 concentrations in 450,000 individuals from three different ancestries, we provide here a  
577 data-driven map of the circulating metabolome across the allele frequency spectrum.  
578 This map identifies previously unrecognized modulators of metabolism with potential  
579 health implications.

580  
581 By combining ML-guided common variant-to-gene annotation with rare exonic variation,  
582 we provided high-confidence effector gene assignments at >100 loci, including some  
583 with less established roles in (lipoprotein) metabolism, such as *SIDT2*. These findings  
584 present compelling candidates for further functional studies, with a strong incentive that  
585 they are likely relevant to human biology, in contrast to species differences frequently  
586 encountered in animal models<sup>71,72</sup>. Large-scale studies similar to ours, but with a broader  
587 coverage of the plasma metabolome, will likely uncover many more genes with yet  
588 undefined roles in metabolism, complementing hypothesis-driven research in  
589 experimental models.

590  
591 After more than two decades of GWAS, it has become clear that pleiotropic effects of  
592 genetic variants are ubiquitous (see, e.g.,<sup>73</sup>). Little distinction has been possible beyond  
593 the generic concepts of ‘vertical’ and ‘horizontal’ pleiotropy or measures of simple  
594 counting. We refine these concepts by observing variants associated with dozens of NMR  
595 measures but consistent with the concept of effects diluting/propagating along pathways  
596 (‘proportional pleiotropy’). Conversely, we also observe variants associated with  
597 comparatively few NMR measures in an inconsistent pattern (‘disproportional  
598 pleiotropy’) suggesting distinct effects on otherwise highly correlated traits. Our data-  
599 driven approach thereby augments previous concepts focussed around biochemical  
600 pathways reporting directionally discordant pleiotropy to discover metabolic bottlenecks  
601 <sup>74</sup>.

602  
603 Disturbance in metabolism or rearrangements thereof are a hallmark of many diseases,  
604 including those not classically considered as ‘metabolic’, such as eye disorders<sup>2</sup>, but  
605 whether these are pathways for prevention or intervention rather than a consequence of  
606 the disease remains often elusive in humans. We demonstrated considerable overlap

607 between mQTLs with disease risk loci, including rare-to-common allelic series that can  
608 reveal unknown effector genes. However, many such ‘locus effects’ were characterised  
609 by nonspecific pleiotropy, implicating the plasma metabolite as a bystander rather than  
610 cause of the disease. This observation aligns with the relatively few notable exceptions,  
611 such as HDL particle characteristics and CAD, from two-sample MR analyses that  
612 contrasted the broad spectrum of observed disease-associations described for the same  
613 NMR platform<sup>75</sup>. These observations might be best explained by the concept of metabolic  
614 flexibility, which includes built-in redundance in key pathways to combat various intrinsic  
615 and extrinsic perturbations.

616

617 An important distinction of our study compared to most previous efforts was the  
618 availability of highly standardized measurements in a well-designed single large cohort,  
619 mitigating influences of preanalytical variables and enabling analyses of even ultra-rare  
620 variants. However, this also meant that we had little opportunity to investigate the  
621 influence of different states of metabolism on our genetic results (such as an overnight  
622 fast) or investigate robustness of findings in different environments or at scale in other  
623 ancestries. For example, UKB participants were not asked to fast overnight prior to their  
624 baseline visit, which has been shown to impact genetic findings<sup>3</sup>. Consequently, sentinel  
625 variants capture relatively little variance (0.57% - 1.07%) in circulating ketone body  
626 concentrations that are highly dependent on the time since last food consumption. Other  
627 limitations included the sensitivity and coverage of the <sup>1</sup>H-NMR platform, and future  
628 efforts are likely to reveal more diverse phenotypic consequences of genetically  
629 constrained flexibility of human metabolism. Another technical aspect to consider in the  
630 interpretation of our results is the indirect nature of <sup>1</sup>H-NMR derived measurements of  
631 certain analytes, including apolipoproteins, that may no longer be reliable in the  
632 presence of rare damaging variants that change the properties of apolipoproteins as  
633 observed for ApoA1.

## 634 Acknowledgements

635 The authors acknowledge the Scientific Computing of the IT Division at the Charité -  
636 Universitätsmedizin Berlin for providing computational resources that have contributed  
637 to the research results reported in this paper  
638 ([https://www.charite.de/en/research/research\\_support\\_services/research\\_infrastructur](https://www.charite.de/en/research/research_support_services/research_infrastructur)

639 e/science\_it/#c30646061). We acknowledge Nightingale Health Plc for access to the UK  
640 Biobank NMR biomarker data. We are deeply grateful to the participants, investigators  
641 and teams of the UKB and FinnGen studies. We thank Benjamin Wild for assistance in  
642 data processing and helpful discussions. This work was supported by DZHK (German  
643 Centre for Cardiovascular Research) and BMBF (German Ministry of Education and  
644 Research) grants to C.L. and co-funded by a European Union grant to M.P (ERC, GenDrug,  
645 101116072). Views and opinions expressed are however those of the author(s) only and  
646 do not necessarily reflect those of the European Union or the European Research  
647 Council. Neither the European Union nor the granting authority can be held responsible  
648 for them. The funders had no role in study design, data collection and analysis, decision  
649 to publish or preparation of the manuscript.

## 650 Author Contributions Statement

651 Conceptualization: MZ, MP, CL  
652 Data curation/Software: MZ, CB, YS, LK, MP  
653 Formal Analysis: MZ, CB, YS, LK, AW, MP  
654 Methodology: MZ, CB, YS, LK, MK, AW, MP, CL  
655 Visualization: MZ, CB, LK, AW, MP, CL  
656 Funding acquisition: CL, MP  
657 Project administration: CL  
658 Supervision: MP, CL  
659 Writing – original draft: MZ, CB, YS, MP, CL  
660 Writing – review & editing: MZ, CB, YS, LK, MK, FK, MM, AW, MP, CL

## 661 Competing Interests Statement

662 None of the authors declare a conflict of interest.

## 663 Data availability

664 All individual-level data is publicly available to bona fide researchers from the UK Biobank  
665 (<https://www.ukbiobank.ac.uk/>). Full summary statistics for all analyses are publicly  
666 available through the NHGRI-EBI GWAS Catalogue (see Github repository for GWAS  
667 Catalogue identifiers).  
668

## 669 Code availability

670 Code for the main analyses is freely available on Github ([https://github.com/comp-](https://github.com/comp-med/ukb-mgwas)  
671 [med/ukb-mgwas](https://github.com/comp-med/ukb-mgwas)) and permanently archived on Zenodo  
672 ([doi.org/10.5281/zenodo.14716599](https://doi.org/10.5281/zenodo.14716599)).

## 673 Methods

### 674 Study design

675



676 UK Biobank is a prospective cohort study from the UK that contains more than 500,000  
677 volunteers between 40 and 69 years of age at inclusion. The study design, sample  
678 characteristics, and genotype data have been described elsewhere<sup>76,77</sup>. The UKBB was  
679 approved by the National Research Ethics Service Committee North West Multi-Centre  
680 Haydock and all study procedures were performed in accordance with the World Medical  
681 Association Declaration of Helsinki ethical principles for medical research. We included  
682 460,036 individuals across the three major ancestries in the UK Biobank in our analyses  
683 for whom inclusion criteria (given consent to further usage of the data, availability of  
684 genetic data, and passed quality control of genetic data) applied. Data from the UKBB  
685 were linked to death registries and hospital episode statistics (HES). We used the  
686 ancestry assignments as defined by the pan-UKB<sup>78</sup>, and further made an effort to assign  
687 unclassified individuals to their respective ancestries based on a k-nearest neighbour  
688 approach using genetic principal components. All analyses were conducted under UKBB  
689 application 44448 and 30418.  
690

## 691 Metabolomic measurements

692  
693 Up to 249 targeted metabolomic measurements were quantified using the Nightingale  
694 NMR platform in human EDTA plasma samples. Detailed experimental procedures for the  
695 NMR platform are described elsewhere<sup>75,79</sup>. The NMR platform covers a wide range of  
696 metabolic biomarkers, including lipoprotein lipids, fatty acids as well as small molecules  
697 such as amino acids, ketone bodies and glycolysis metabolites. All metabolites are  
698 quantified in molar concentration units. We combine here three data releases that cover  
699 the full breadth of the UKBB. Metabolomics data was available for 482276 individuals,  
700 including 19699 samples with data from both the baseline and repeat visit.  
701

702 Metabolites were reliably detected, with only one biomarker over 2.5% missingness in  
703 releases 1/2 (creatinine) and release 3 (3-Hydroxybutyrate). 98% of the samples had < 5%  
704 missingness over all biomarkers in releases 1/2 and release 3. We used the *ukbnmr*<sup>80</sup> R  
705 package (v2.2, R v4.3.2) for quality control and removal of technical variation in the NMR  
706 data. This includes technical confounders such as sample preparation time, shipping  
707 plate well, spectrometer effects, time drift within spectrometers and outlier plates.  
708

709 We removed samples that were flagged by Nightingale for poor quality and used the MICE  
710 (Multivariate Imputation by Chained Equations)<sup>81</sup> R package to impute the remaining  
711 dataset. In total, we imputed 0.16% and 0.17% of data in releases 1/2 and release 3,  
712 respectively.  
713

714 We observed overall good consistency with the overlapping routine blood biomarkers  
715 previously measured in the same cohort (median  $R^2$ : 0.9, range: 0.62 – 0.94)  
716 (**Supplementary Fig. 14**).  
717

## 718 Adjustment of metabolomic data for medication use

719

720 We sought to adjust the NMR data for medication use, especially cholesterol-lowering  
721 medication to avoid false positive results driven by medication use in downstream  
722 genetic analyses. For male and female participants separately, we fitted linear models to  
723 quantify the impact of 6 drug categories on each NMR phenotype: cholesterol-lowering  
724 medicine, blood pressure medication, diabetic medication including Metformin usage,  
725 oral contraceptive pill or minipill (female only), hormone replacement therapy (female  
726 only) (UKB fields 6177 and 6153) (**Supplementary Fig. 15, Supplementary Table 18**).

727  
728 We used data from individuals with both baseline and repeat assessment metabolic data  
729 available and estimated the effect of medication in individuals that did not take any drugs  
730 at the time of the baseline visit (N = 6,312 male, N = 6,713 female participants). We fitted  
731 a linear model predicting the follow-up metabolic data from baseline metabolic data,  
732 sex, age and medication use:

$$733 \quad \text{NMR}_{\text{baseline}} \sim \text{NMR}_{\text{follow-up}} + \text{age} + \text{bmi} + \text{med}_{\text{cholesterol}} + \text{med}_{\text{diabetic}} + \text{med}_{\text{contraception}} + \text{med}_{\text{hormone}} + \text{error}$$

735  
736 We note that the sample sizes for diabetic medication ( $N_{\text{male}} = 45$ ,  $N_{\text{female}} = 29$ ), oral  
737 contraceptive medication (N = 27) and hormone replacement therapy (N=148) were too  
738 small to reliably estimate any effects. Effect estimates for diabetic medication were  
739 correlated to estimates for cholesterol-lowering medicine. The effect estimates for blood  
740 pressure medication were minimal across the phenotypes. We considered thus only the  
741 impact of cholesterol-lowering medicine and corrected the metabolic data in a sex-  
742 specific manner.

## 744 Genotyping and GWAS analyses

745  
746 GWAS was performed on 249 metabolic traits measured by the NMR platform on  
747 European (n = 434,646), British-Asian (n = 8,796) and British-African participants  
748 (n=6,573) that had complete phenotypic, covariate and genetic information available. We  
749 performed GWAS under the additive model using REGENIE (v3.2.5)<sup>82</sup> that employs a two-  
750 step procedure to account for population structure. We derived a set of high-quality  
751 genotyped variants per population by applying following filters: (MAF > 1%, MAC > 100,  
752 missingness rate < 10%, p<sub>HWE</sub> > 10<sup>-15</sup>). Further, linkage disequilibrium pruning was  
753 performed using a 1000 kb window, shifting by 100 variants and removing variants with  
754 LD(r<sup>2</sup>) > 0.8. We used these variants as input for the first step of REGENIE, to generate  
755 individual trait predictions using the leave-one-chromosome-out scheme. These  
756 predictions are used in the second step where individual variants are tested. Models were  
757 adjusted for age, sex and the first ten genetic principal components. We tested variants  
758 with a minor allele frequency > 0.5%, amounting to 11.5M variants in European  
759 individuals, 11.5M variants in British-Asian individuals and 19.3M variants in British-  
760 African individuals.

761  
762 For initial discovery, we performed a meta-analysis across the three ancestral groups  
763 using METAL<sup>83</sup>. We required variants to be present in at least two ancestral groups. To  
764 declare significance, we considered a stringent p-value threshold (2.0 x 10<sup>-10</sup>) by dividing

765 the standard genome-wide threshold by the number of metabolic phenotypes ( $5.0 \times 10^{-8}$   
766 / 249).

767

768 We tested our results for genomic inflation and calculated the SNP-based heritability  
769 using LD-score regression (LDSC)<sup>84</sup> (**Supplementary Table 19**).

770

## 771 Regional clumping and fine-mapping

772

773 We used regional clumping ( $\pm 500\text{kb}$ ) around sentinel variants from the analyses  
774 including White European samples to select independent genomic regions associated  
775 with a metabolic phenotype and collapsed neighbouring regions using BEDtools  
776 (v2.30.0). We treated the extended MHC region (chr6:25.5-34.0Mb) as one region.

777

778 Within each region of interest, excluding the MHC region, we performed statistical fine-  
779 mapping for all phenotypes associated with that region using the ‘Sum of single effects’  
780 model (SuSiE) implemented in the *susieR* (v0.12.35) R package<sup>85</sup>. Briefly, SuSiE employs  
781 a Bayesian framework for variable selection in a multiple regression problem with the aim  
782 to identify sets of independent variants each of which likely contains the true causally  
783 underlying genetic variant. We implemented the workflow using default prior and  
784 parameter settings, apart from the minimum absolute correlation, which we set to 0.1.  
785 Since SuSiE is implemented in a linear regression framework, we used the GWAS  
786 summary statistics with a matching correlation matrix of dosage genotypes instead of  
787 individual level data to implement fine-mapping (*susie\_rss()*) as recommended by the  
788 authors<sup>85</sup>.

789

790 Within a given region, a phenotype can be associated to multiple, independent genetic  
791 loci. To determine the appropriate number of credible sets, we iterated over the  
792 maximum credible sets parameter in *susieR* from two to ten, thus generating fine-  
793 mapped results constrained to a range of maximum number of credible sets. For each  
794 collection of credible sets, we pruned sets where the lead variant was correlated to the  
795 lead variant of other credible sets ( $R^2 > 0.25$ ). After pruning, we considered the fine-  
796 mapped results constrained to the largest number of credible sets that still contained  
797 one or more credible sets.

798

799 We performed several sensitivity analyses by computing joint models per locus –  
800 phenotype combination. We obtained all lead variants across the credible sets provided  
801 by SuSiE based on posterior inclusion probability and fitted a single linear model, jointly  
802 modelling the effect of all distinct credible sets in the locus for a given phenotype.  
803 Subsequently, we retained only credible sets where the lead variant reached genome-  
804 wide significance ( $p = 5.0 \times 10^{-8}$ ) in both marginal and joint statistics. Furthermore, we  
805 ensured the estimated coefficients were directionally concordant and of similar  
806 magnitude between joint and marginal models ( $\pm 25\%$ ). Linear models were  
807 implemented in R using the *glm()* function and used only unrelated white-European  
808 participants and the same set of covariates as described above.

809

810 Finally, we used LD-clumping ( $r^2 > 0.6$ ) to identify credible sets shared across metabolic  
811 phenotypes.

812  
813 We computed the correlation matrix with LDscore v2.0 using genetic data from 50,000  
814 randomly selected, unrelated White European UKB participants. In situations where  
815 SuSiE did not deliver a credible set, we used the Wakefield approximation<sup>86</sup> to compute  
816 95%-credible sets.  
817

## 818 Multi-ancestry finemapping

819  
820 We aimed to use ancestrally diverse genetic data to refine the credible sets identified in  
821 the White European analyses. For 35577 credible sets that contained two or more  
822 variants in the European analyses, we checked for evidence of a genetic signal in the  
823 British-African and British-Central/South Asian ancestries ( $P < 1.0 \times 10^{-4}$ ) at the same  
824 locus ( $\pm 25\text{kb}$  on either side of the credible set). After these filters, we considered 6979  
825 credible sets for finemapping across ancestries using MultiSuSiE<sup>87</sup>. We considered only  
826 quality-controlled variants that were prevalent in all populations ( $\text{MAF} > 0.5\%$ ) and used  
827 the posterior inclusions probabilities from the European analyses as priors. LD matrices  
828 were calculated from a random subset of 50,000 White European participants for  
829 Europeans, and using all available individuals for the British-African and British-  
830 Central/South Asian ancestries.  
831

## 832 Replication of genetic associations

833  
834 We replicated our trans-ancestral genetic signals using two independent studies: i) the  
835 so-far largest published mGWAS<sup>3</sup>, and ii) a parallel effort using overlapping UK Biobank  
836 data<sup>9</sup>, both using the same NMR platform. We considered a set of metabolic traits that  
837 were directly measured by the NMR platform and not inferred from other traits to avoid  
838 multiplicative errors in these more sensitive phenotypes. In total, we were able to match  
839 144 (Karjalainen et al) and 169 (Tambets et al) metabolic traits, for which we compared  
840 sentinel variants that passed metabolome-adjusted, genome-wide significance in our  
841 trans-ancestral meta-analysis and that overlapped between the studies.  
842

## 843 Sex-specific genetic analyses

844  
845 To assess whether our genetic analyses were driven by sex differences and whether our  
846 results were transferrable to both sexes, we performed sex-stratified GWAS within the  
847 largest ancestry (EUR). We defined ‘female’ and ‘male’ sex including participants where  
848 the recorded sex and sex chromosomes aligned (XX for females and XY for males). The  
849 recorded sex was self-reported, and it was not possible to distinguish sex from gender.  
850 We acknowledge the importance of distinguishing between sex and gender in research  
851 and that chromosomal make-up does not always align with self-identified gender. In  
852 total, phenotypic, covariate and genetic data was fully available for 198,796 males and  
853 235,850 females. GWASs were performed using REGENIE as described above. Per

854 metabolic trait, we meta-analysed the sex-stratified results using the inverse-variance  
855 weighted model in METAL (v2020-05-05)<sup>83</sup> and finally clumped results on the  
856 heterogeneity p-value using plink (v2.00)<sup>88</sup> (--clump-p1 5 x 10<sup>-8</sup>, --clump-r2 0, --clump-kb  
857 2500). We considered loci putatively sex-differential if the meta-analysis heterogeneity  
858 p-values were genome-wide significant (p < 5 x 10<sup>-8</sup>). We performed additional sensitivity  
859 analyses on the putatively sex-differential loci by assessing the influence of covariates  
860 confounded by sex (BMI, tobacco usage, alcohol consumption, lipid-lowering  
861 medication and diabetic medication). To properly model all gene by environment  
862 interactions<sup>89</sup>, we fitted the following model per clump lead variant and associated  
863 metabolic phenotype, including the same set of covariates used in the original GWAS:

864  
865 NMR phenotype ~ SNP + confounder + sex + age + fasting duration + PC1-10 + SNP\*sex + SNP\*confounder  
866 + sex\*confounder + error  
867

## 868 Causal gene assignment

869  
870 To assign candidate genes for all metabolite-QTLs residing outside the MHC region, we  
871 first collected annotations for each genetic variant or proxies thereof (r<sup>2</sup> > 0.6), including  
872 1) distance to the gene body and 2) putative functional consequences based on the  
873 Variant Effect Predictor (VEP) tool offered by Ensembl. We further collated up to 10  
874 closest genes within a 2 Mb window and subsequent gene features such as: 1) eQTL  
875 evidence for a given variant-gene pair for each tissue available in the eQTL Catalogue  
876 release 7<sup>90</sup>, 2) evidence of being annotated as metabolic in the MGI or Orphanet  
877 databases as defined in ProGem<sup>25</sup>, 3) evidence of being listed in the OMIM database<sup>43</sup> 4)  
878 and evidence of being an already assigned drug target in Open Targets<sup>91</sup> clinical stage III  
879 and IV.

880  
881 With no universally accepted standard for variant-to-gene assignments, we relied on  
882 prior biological and genomic information to create three sets of “putative true positive”  
883 (PTP) set: 1) genes annotated as part of a cholesterol pathway in the KEGG<sup>92</sup> or  
884 REACTOME<sup>93</sup> database (n=6791 , 722 unique SNPs), 2) genes annotated as part of a lipid  
885 pathway (n=5670 , 603 unique SNPs) and 3) genes annotated as part of an amino acid-  
886 related pathway (n=8349, 895 unique SNPs). We used all fine-mapped SNPs associated  
887 with metabolites classified in the respective NMR metabolite class (Cholesterol:  
888 Cholesterol, Cholesteryl esters, Free cholesterol; Lipid: Total lipids, other lipids, Relative  
889 lipid concentration, Phospholipids); Amino Acid: Amino acid) in the PTP set and used  
890 overlapping SNPs in only one PTP set. The dataset was split in a 7:3 ratio to obtain training  
891 and test sets without overlapping variants. We trained a Random Forest classifier using  
892 5-fold cross-validation with implemented subsampling to account for the unbalanced  
893 datasets. The implementation was carried out using python scikit-learn v1.4.1. We used  
894 the balanced accuracy score to choose the best-performing forest from each training set.  
895 Subsequently, we used the best-performing Random Forest classifiers from each PTP set  
896 to assign candidate scores for all putative effector genes across the entire set of  
897 metabolite-QTLs. We then calculated the median score of these classifiers and selected  
898 the highest-scoring gene as the assigned gene for the variant. Within each PTP set, we  
899 omitted features used to define true positive sets. Each of the three classifiers exhibited

900 consistent performance with a mean ROC AUC of 0.80 and a mean balanced accuracy  
901 score of 0.69 (**Supplementary Fig. 16**).

902  
903 To provide another layer of evidence for assignment of causal genes at metabolic loci, we  
904 performed cis-colocalisation with protein targets measured in the independent Fenland  
905 study<sup>28</sup>. Cis (e.g. gene body  $\pm$  500kb) summary statistics were preprocessed using  
906 MungeSumStats<sup>94</sup>. To relax the single causal variant assumption, we employed a  
907 colocalization approach where we fine-mapped all traits with SuSiE and then performed  
908 colocalization among all credible sets using functionality of the ``coloc`` (v5.2.3)<sup>95,96</sup> and  
909 ``susier`` (v0.12.35)<sup>85</sup> R packages. For this, we set the prior probability that a SNP is  
910 associated with both traits to  $5 \times 10^{-6}$  and restricted the maximum number of credible sets  
911 for the outcome data to 5<sup>95</sup>.  
912

### 913 Tissue enrichment of metabolic loci

914  
915 We tested whether genes proximal to metabolic loci and assigned effector genes were  
916 enriched in tissue compartments by leveraging data from the Human Protein Atlas<sup>97</sup>.  
917 Specifically, we used a two-sided Fisher's test whether metabolic genes were enriched  
918 among tissue-specific genes (tissue-enriched or tissue-enhanced as defined by the  
919 Protein Atlas) against all protein-coding genes as background.  
920

### 921 Pleiotropy assignment and overlap with the GWAS catalog

922  
923 To assign modes of pleiotropy for each mQTL, we first clumped lead credible set variants  
924 across NMR measures by LD, collating variants with  $r^2 \geq 0.6$  as a single signal, referred to  
925 hereafter as 'mQTL. This was done based on dosage files of all unrelated White European  
926 UKB participants and implemented with the *igraph* (v.2.0.1.1) package in R. For each  
927 mQTL, we then computed all possible Pearson correlation coefficients among  
928 associated NMR measures. To classify each mQTL-group, generated two metrics: 1) the  
929 25<sup>th</sup> percentile of all correlations among associated NMR measures, and 2) the Pearson  
930 correlation coefficient between the association strengths for each measure ( $-\log_{10}(p$ -  
931 value) and its correlation coefficient with the most strongly associated measure within  
932 the mQTL. The latter is a measure to what extent the association between NMR  
933 measures at a given locus ('pleiotropy') can be explained by being correlated with the  
934 most proximal associated measure. Based on opposing those two measures for all  
935 mQTLs we opted to threshold each at 0.6 to define the following five groups: 1) 'specific'  
936 mQTLs associated with only  $\leq 3$  highly correlated NMR measures ( $\rho \geq 0.6$ ), 2) 'pathway  
937 pleiotropic' mQTLs associated with highly correlated NMR measures ( $\rho \geq 0.6$ ) that also  
938 followed the described association pattern ( $\rho \geq 0.6$ ), 3) 'proportional pleiotropic' mQTL  
939 groups associated with, in part, uncorrelated NMR measures but highly correlated  
940 association statistics ( $\rho \geq 0.6$ ), 4) 'disproportional pleiotropic' mQTLs associated with  
941 highly correlated NMR measures ( $\rho \geq 0.6$ ), but without evidence that this translated into  
942 a correlation of association statistics ( $\rho < 0.6$ ), and 5) all remaining mQTLs as 'unspecific  
943 pleiotropic' groups.  
944

945 To quantify the extent to which our pleiotropy assignment extends beyond the NMR  
946 measured analysed here, we intersected mQTLs and proxies thereof with results reported  
947 in the GWAS catalog (download: 20/05/2024). We first pruned GWAS catalog entries for  
948 those with mapped traits (to minimize double counting), results that met genome-wide  
949 significance ( $p < 5 \times 10^{-8}$ ) and had location information available. We further dropped  
950 results similar to NMR measures based on broad EFO terms (e.g., EFO:0005105 and child  
951 terms indicating 'lipid or lipoprotein measurement'). To further account for traits mapping  
952 to similar categories, we iteratively traced back mapped EFO terms to broader parent  
953 terms. We finally classified mQTLs to be 'specific' in the GWAS catalog, if they associated  
954 with less than five parent EFO-terms and 'unspecific' otherwise. This information was  
955 primarily used to define instruments for Mendelian randomization analysis.  
956

## 957 Integration of metabolomic measurements with cardiovascular 958 endpoints

959  
960 We next aimed to utilize the mQTLs to investigate the shared genetic basis of the 249 NMR  
961 and 25 selected cardiovascular disease (CVD) traits. We utilized public databases (GWAS  
962 Catalog, openGWAS, CVD-KP) to collect CVD data comprising the largest currently  
963 publicly available GWASs on coronary artery disease and myocardial infarction, angina  
964 pectoris, aortic aneurysm, heart failure and stroke, peripheral arterial disease including  
965 2-5 subtypes for each phenotype. An additional 10 CVD traits had no subtype data  
966 available (**Supplementary Table 13**) Data was harmonized and if necessary, lifted over  
967 to GRCh37 using the `MungeSumstats` (v1.13.2) R package<sup>94</sup>. We queried mQTL lead  
968 variants and proxies in strong LD ( $r^2 > 0.8$ ; LD backbone based on UK Biobank, as described  
969 above) of each NMR trait in each region and corresponding summary statistics for each CVD  
970 trait.

971  
972 To investigate `variant` effects on NMR metabolite concentrations and CVD outcomes,  
973 we performed statistical colocalization screens for all combinations of the NMR traits in  
974 regions with at least one credible set and CVD traits with matching summary statistics<sup>98</sup>.  
975 We applied statistical colocalization as described before (see 'Causal gene assignment').  
976

977 To estimate `level` effects of NMR metabolite concentrations on CVD outcomes, we  
978 performed Mendelian Randomization analysis using the `TwoSampleMR` package  
979 (v0.5.1), implementing the inverse-variance weighted and the MR-Egger methods. We  
980 used all 249 NMR metabolites as exposure variables, the 25 CVDs as outcome variables  
981 and assessed separately four sets of instruments: 1) sentinel variants, 2) lead credible  
982 set variants, 3) lead credible set variants restricted for molecular pleiotropy (e.g.  
983 'pathway pleiotropy') and 4) lead credible set variants restricted for both molecular and  
984 phenotypic pleiotropy. We used the Wald ratio method to estimate the effect of NMR  
985 concentrations on CVD outcomes using only single genetic variants<sup>99</sup>. We used MR Egger  
986 to test for evidence of a pleiotropic association, an intercept p-value of  $p > 0.0001$   
987 indicating evidence of no pleiotropy and checked for concordance between the effect  
988 estimates of IVW-MR, MR-Egger and single genetic variant MR. We controlled the false  
989 discovery rate (FDR) at  $FDR = 5\%$ <sup>100</sup>. To further limit the possible extend of pleiotropic

990 associations, we only reported `level effects` passing these filters in the variant sets 2-  
991 4, prioritizing the association in the more stringent variant set.

992  
993 The overlap of `locus effects` showing no `disproportional pleiotropy` according to the  
994 section **Pleiotropy assignment and overlap with the GWAS catalog** as well as a  
995 significant single variant MR (FDR=5%) and `level effects` calculated from metabolite-  
996 specific or metabolite- and phenome-specific variants was used to identify gene-  
997 metabolite pairs associated with cardiovascular disease risk independent of LDL-  
998 metabolism. We considered loci as independent from LDL-metabolism if they did not  
999 associate with clinical LDL-cholesterol at the locus with  $p < 2.0 \times 10^{-10}$  and the effect  
1000 estimate of any variant on clinical LDL-C ranked upwards the 80<sup>th</sup> percentile of all effect  
1001 estimates at the locus.

## 1002 Rare Variant Analyses with whole exome sequencing data

1003

### 1004 Whole Exome Sequencing data QC

1005 An in-depth description of whole exome sequencing, including experimental details,  
1006 variant calling, and standard quality control measures for the UK Biobank, has been  
1007 extensively reported by Backman et al.<sup>101</sup>. We performed additional quality control (QC)  
1008 steps at the UKB Research Analysis Platform (RAP; <https://ukbiobank.dnanexus.com/>).

1009

1010 We employed bcftools (v1.15.1) to process population-level Variant Call Format (pVCF)  
1011 files. Initially, we normalised the data using the reference sequence GRCh38 build,  
1012 followed by splitting multi-allelic variants. Subsequently, we conducted QC on these  
1013 variants using a set of parameters outlined below to filter high-quality variants for  
1014 downstream genetic analyses. Genotypes for SNPs were set to missing if the read depth  
1015 was less than 7 (or less than 10 for INDELS) or if the genotype quality was below 20.  
1016 Furthermore, we excluded variants if the allele balance (AB) was less than 0.25 or greater  
1017 than 0.8 in heterozygous carriers.

1018

1019 Finally, we computed the missingness rate for each variant and excluded those with  
1020 missing values in over 50% of the participants.

1021

### 1022 Variant Annotation and Gene burden Masks

1023 Variants were annotated using ENSEMBL Variant Effect Predictor (VEP)<sup>102</sup> (v106.1) with  
1024 the most severe consequence for each variant chosen across all protein-coding  
1025 transcripts. We further utilized additional plugins REVEL<sup>103</sup>, CADD v1.6<sup>104</sup>, and LOFTEE<sup>105</sup>,  
1026 for variant annotation. Based on these scores we defined six partially overlapping variant  
1027 masks: 1) high-confidence predicted loss-of-function (pLOF, based on LOFTEE and  
1028 includes stop-gained, splice site disrupting, and frameshift variants), 2) any pLOF  
1029 assigned high impact by VEP, 3) pLOF and high-impact missense variants (CADD score >  
1030 20 or REVEL score > 0.5), 4) pLOF and any missense variants, 5) only high-impact  
1031 variants, and 6) any missense variants but not pLOF. We tested synonymous variants



1032 separately as a negative control. We tested each mask in different minor allele frequency  
1033 bins, using 0.5% and 0.005% as thresholds.

1034  
1035 We performed rare variant association testing (RVAT) using WES data across 249  
1036 quantitative NMR phenotypes using REGENIE (v3.1.1) via the DNAnexus Swiss Army Knife  
1037 tool (v4.9.1). Similar to common variant GWASs, we used a two-step approach by  
1038 REGENIE. However, we additionally generated step1 LOCO files with and without  
1039 adjusting for common signals via a polygenic score (PGS) in the RVAT models per  
1040 phenotype. In practice, we computed a PGS for each phenotype using effect sizes of lead  
1041 variants from the GWAS summary statistics and corresponding dosages of variants from  
1042 imputed data. All RVAT models were then adjusted for PGS in addition to age, biological  
1043 sex, fasting duration and the first ten genetic PCs. We first performed aggregated gene  
1044 burden testing across for 19,026 genes using a set of masks as defined above. For the  
1045 gene burden testing we used aggregated Cauchy association test (ACAT) to estimate a p-  
1046 value for each gene across all combinations of masks and allele frequency bins. ACAT  
1047 first computes p-values for all sets defined by various masks within a gene and then takes  
1048 these p-values as input to compute one p-value for the respective gene via a well  
1049 approximated Cauchy distribution.

1050  
1051 We have also performed single variant association testing for exonic variants commonly  
1052 referred to as exome wide association study (ExWAS). For the ExWAS, we only tested  
1053 variants with MAC >5 and reported results for variants with a MAF < 0.0005. We have  
1054 performed these analyses in individuals of White Europeans, British African and British  
1055 South Asian ancestry.

1056  
1057 We considered findings as robust, if they passed multiple testing corrected statistical  
1058 significance (gene burden:  $p < 1.2 \times 10^{-8}$  [corrected for the number of genes x number of  
1059 traits]; ExWAS:  $p < 2.0 \times 10^{-10}$  [same as for common variant GWAS, conventional genome-  
1060 wide significance corrected for the number of traits]) in both the model with and without  
1061 adjusting for the common variant PGS and effect sizes did not differ by more than 20%  
1062 between these models, since this might otherwise indicate that rare variant findings  
1063 cannot clearly distinguished from common variant effects.

1064

## 1065 Phenotype definition

1066 To systematically test for phenotypic consequences of genes identified through rare  
1067 variant analysis, we collated 626 disease entities following previous work<sup>1</sup> by aggregating  
1068 information from self-report, hospital episode statistics, death certificates, and primary  
1069 care data (45% of the UKB population). Each of the disease entities had at least one  
1070 common variant finding passing statistical significance, and we employed a similar  
1071 analysis workflow using REGENIE as described for NMR measures but using logistic  
1072 regression with saddle point approximation.

1073

## 1074 Integration of OMIM (Online Mendelian Inheritance in Man)

1075 We downloaded the OMIM gene – disease list (09/11/2023) and kept 7,327 unique entries  
1076 after filtering for gene entries with high confidence (level 3). We computed the  
1077 enrichment of genes associated with any NMR measure from rare variant or gene burden  
1078 analysis against a background of 19,989 protein coding genes using Fisher’s exact test.  
1079

## 1080 Mendelian Randomisation analyses

1081  
1082 We performed a phenome-wide Mendelian Randomisation screen using outcome  
1083 summary statistics from the independent FinnGenn<sup>68</sup> cohort, release 11 (June 2024). We  
1084 assessed only outcomes with genome-wide significant signals ( $5 \times 10^{-8}$ ), yielding 1394  
1085 phenotypic outcomes. We selected 21 non-lipid NMR biomarkers as exposure variables  
1086 and assessed separately four sets of instruments as described previously for the  
1087 cardiovascular Mendelian Randomisation analyses. We included two well-characterized  
1088 lipid biomarkers (LDL-C and ApoB) as positive controls in the MR analyses.  
1089

1090 We performed MR using the TwoSampleMR package (v0.5.1), implementing the inverse-  
1091 variance weighted and the MR-Egger methods. We discarded results with MR Egger  $p <$   
1092  $0.001$ , Cochran’s Q  $p$ -value  $< 1.0 \times 10^{-6}$  and results where the estimated effect was  
1093 directionally discordant between the IVW and Egger methods.

## 1094 References

- 1095  
1096 1. Surendran, P. *et al.* Rare and common genetic determinants of metabolic  
1097 individuality and their effects on human health. *Nat. Med.* **28**, 2321–2332 (2022).  
1098 2. Lotta, L. A. *et al.* A cross-platform approach identifies genetic regulators of human  
1099 metabolism and health. *Nat. Genet.* **53**, 54–64 (2021).  
1100 3. Karjalainen, M. K. *et al.* Genome-wide characterization of circulating metabolic  
1101 biomarkers. *Nature* **628**, 130–138 (2024).  
1102 4. Kettunen, J. *et al.* Genome-wide association study identifies multiple loci influencing  
1103 human serum metabolite levels. *Nat. Genet.* **44**, 269–276 (2012).  
1104 5. Chen, Y. *et al.* Genomic atlas of the plasma metabolome prioritizes metabolites  
1105 implicated in human diseases. *Nat. Genet.* **55**, 44–53 (2023).

- 1106 6. Nag, A. *et al.* Effects of protein-coding variants on blood metabolite measurements  
1107 and clinical biomarkers in the UK Biobank. *Am. J. Hum. Genet.* **110**, 487–498 (2023).
- 1108 7. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants  
1109 associated with human blood metabolites. *Nat. Genet.* **49**, 568–578 (2017).
- 1110 8. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat.*  
1111 *Genet.* **46**, 543–550 (2014).
- 1112 9. Tambets, R. *et al.* Genome-wide association study for circulating metabolites in  
1113 619,372 individuals. *medRxiv* (2024) doi:10.1101/2024.10.15.24315557.
- 1114 10. Yin, X. *et al.* Genome-wide association studies of metabolites in Finnish men  
1115 identify disease-relevant loci. *Nat. Commun.* **13**, 1644 (2022).
- 1116 11. Van Der Meer, D. *et al.* Pleiotropic and sex-specific genetic architecture of  
1117 circulating metabolic markers. *medRxiv* (2024) doi:10.1101/2024.07.30.24311254.
- 1118 12. Khan, A. *et al.* Metabolic gene function discovery platform GeneMAP identifies  
1119 SLC25A48 as necessary for mitochondrial choline import. *Nat. Genet.* (2024)  
1120 doi:10.1038/s41588-024-01827-2.
- 1121 13. Schlosser, P. *et al.* Genetic studies of paired metabolomes reveal enzymatic and  
1122 transport processes at the interface of plasma and urine. *Nat. Genet.* **55**, 995–1008  
1123 (2023).
- 1124 14. Love-Gregory, L. *et al.* Variants in the CD36 gene associate with the metabolic  
1125 syndrome and high-density lipoprotein cholesterol. *Hum. Mol. Genet.* **17**, 1695–1704  
1126 (2008).
- 1127 15. Adiyaman, S. C. *et al.* Congenital generalized lipodystrophy type 4 due to a novel  
1128 PTRF/CAVIN1 pathogenic variant in a child: effects of metreleptin substitution. *J.*  
1129 *Pediatr. Endocrinol. Metab. JPEM* **35**, 946–952 (2022).

- 1130 16. Mauvais-Jarvis, F. Sex differences in energy metabolism: natural selection,  
1131 mechanisms and consequences. *Nat. Rev. Nephrol.* **20**, 56–69 (2024).
- 1132 17. Gerdtts, E. & Regitz-Zagrosek, V. Sex differences in cardiometabolic disorders.  
1133 *Nat. Med.* **25**, 1657–1666 (2019).
- 1134 18. Wittemans, L. B. L. *et al.* Assessing the causal association of glycine with risk of  
1135 cardio-metabolic diseases. *Nat. Commun.* **10**, 1060 (2019).
- 1136 19. Mittelstrass, K. *et al.* Discovery of sexual dimorphisms in metabolic and genetic  
1137 biomarkers. *PLoS Genet.* **7**, e1002215 (2011).
- 1138 20. Koprulu, M. *et al.* Similar and different: systematic investigation of  
1139 proteogenomic variation between sexes and its relevance for human diseases.  
1140 *medRxiv* (2024) doi:10.1101/2024.02.16.24302936.
- 1141 21. Michos, E. D., McEvoy, J. W. & Blumenthal, R. S. Lipid Management for the  
1142 Prevention of Atherosclerotic Cardiovascular Disease. *N. Engl. J. Med.* **381**, 1557–  
1143 1567 (2019).
- 1144 22. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association  
1145 studies of lipids. *Nature* **600**, 675–679 (2021).
- 1146 23. BasuRay, S., Wang, Y., Smagris, E., Cohen, J. C. & Hobbs, H. H. Accumulation of  
1147 PNPLA3 on lipid droplets is the basis of associated hepatic steatosis. *Proc. Natl.*  
1148 *Acad. Sci. U. S. A.* **116**, 9521–9526 (2019).
- 1149 24. Johnson, S. M. *et al.* PNPLA3 is a triglyceride lipase that mobilizes  
1150 polyunsaturated fatty acids to facilitate hepatic secretion of large-sized very low-  
1151 density lipoprotein. *Nat. Commun.* **15**, 4847 (2024).
- 1152 25. Stacey, D. *et al.* ProGeM: a framework for the prioritization of candidate causal  
1153 genes at molecular quantitative trait loci. *Nucleic Acids Res.* **47**, e3 (2019).

- 1154 26. Vujkovic, M. *et al.* Discovery of 318 new risk loci for type 2 diabetes and related  
1155 vascular outcomes among 1.4 million participants in a multi-ancestry meta-analysis.  
1156 *Nat. Genet.* **52**, 680–691 (2020).
- 1157 27. Donaldson, J. G. & Jackson, C. L. ARF family G proteins and their regulators: roles  
1158 in membrane transport, development and disease. *Nat. Rev. Mol. Cell Biol.* **12**, 362–  
1159 375 (2011).
- 1160 28. Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases.  
1161 *Science* **374**, eabj1541 (2021).
- 1162 29. Pellegrinelli, V. *et al.* Dysregulation of macrophage PEPD in obesity determines  
1163 adipose tissue fibro-inflammation and insulin resistance. *Nat. Metab.* **4**, 476–494  
1164 (2022).
- 1165 30. Wang, Q. *et al.* Metabolic profiling of angiotensin-like protein 3 and 4 inhibition:  
1166 a drug-target Mendelian randomization analysis. *Eur. Heart J.* **42**, 1160–1169 (2021).
- 1167 31. Hindy, G. *et al.* Rare coding variants in 35 genes associate with circulating lipid  
1168 levels-A multi-ancestry analysis of 170,000 exomes. *Am. J. Hum. Genet.* **109**, 81–96  
1169 (2022).
- 1170 32. Sjouke, B., Balak, D. M. W., Beuers, U., Ratziu, V. & Stroes, E. S. G. Is mipomersen  
1171 ready for clinical implementation? A transatlantic dilemma. *Curr. Opin. Lipidol.* **24**,  
1172 301–306 (2013).
- 1173 33. Spracklen, C. N. *et al.* Identification of type 2 diabetes loci in 433,540 East Asian  
1174 individuals. *Nature* **582**, 240–245 (2020).
- 1175 34. Rhee, E. P. *et al.* An exome array study of the plasma metabolome. *Nat.*  
1176 *Commun.* **7**, 12360 (2016).

- 1177 35. Chen, X., Gu, X. & Zhang, H. Sidt2 regulates hepatocellular lipid metabolism  
1178 through autophagy. *J. Lipid Res.* **59**, 404–415 (2018).
- 1179 36. Sampieri, A., Asanov, A., Méndez-Acevedo, K. M. & Vaca, L. SIDT2 Associates  
1180 with Apolipoprotein A1 (ApoA1) and Facilitates ApoA1 Secretion in Hepatocytes.  
1181 *Cells* **12**, 2353 (2023).
- 1182 37. Jaiswal, S. *et al.* Clonal Hematopoiesis and Risk of Atherosclerotic  
1183 Cardiovascular Disease. *N. Engl. J. Med.* **377**, 111–121 (2017).
- 1184 38. Weiner, D. J. *et al.* Polygenic architecture of rare coding variation across 394,783  
1185 exomes. *Nature* **614**, 492–499 (2023).
- 1186 39. Szeri, F. *et al.* The membrane protein ANKH is crucial for bone mechanical  
1187 performance by mediating cellular export of citrate and ATP. *PLoS Genet.* **16**,  
1188 e1008884 (2020).
- 1189 40. Chroni, A. & Kardassis, D. HDL Dysfunction Caused by Mutations in apoA-I and  
1190 Other Genes that are Critical for HDL Biogenesis and Remodeling. *Curr. Med. Chem.*  
1191 **26**, 1544–1575 (2019).
- 1192 41. Tilly-Kiesi, M. *et al.* ApoA-I<sub>Helsinki</sub> (Lys<sub>107</sub> →0) Associated With Reduced HDL  
1193 Cholesterol and LpA-I:A-II Deficiency. *Arterioscler. Thromb. Vasc. Biol.* **15**, 1294–1306  
1194 (1995).
- 1195 42. Zanoni, P. & Von Eckardstein, A. Inborn errors of apolipoprotein A-I metabolism:  
1196 implications for disease, research and development. *Curr. Opin. Lipidol.* **31**, 62–70  
1197 (2020).
- 1198 43. Amberger, J. S., Bocchini, C. A., Schiettecatte, F., Scott, A. F. & Hamosh, A.  
1199 OMIM.org: Online Mendelian Inheritance in Man (OMIM®), an online catalog of human  
1200 genes and genetic disorders. *Nucleic Acids Res.* **43**, D789–D798 (2015).

- 1201 44. Horiki, M. *et al.* Smad6/Smurf1 overexpression in cartilage delays chondrocyte  
1202 hypertrophy and causes dwarfism with osteopenia. *J. Cell Biol.* **165**, 433–445 (2004).
- 1203 45. Zhang, F., Sodroski, C., Cha, H., Li, Q. & Liang, T. J. Infection of Hepatocytes With  
1204 HCV Increases Cell Surface Levels of Heparan Sulfate Proteoglycans, Uptake of  
1205 Cholesterol and Lipoprotein, and Virus Entry by Up-regulating SMAD6 and SMAD7.  
1206 *Gastroenterology* **152**, 257-270.e7 (2017).
- 1207 46. Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets  
1208 through human genetics. *Nat. Rev. Drug Discov.* **12**, 581–594 (2013).
- 1209 47. Hoogeveen, R. C. & Ballantyne, C. M. Residual Cardiovascular Risk at Low LDL:  
1210 Remnants, Lipoprotein(a), and Inflammation. *Clin. Chem.* **67**, 143–153 (2021).
- 1211 48. Aragam, K. G. *et al.* Discovery and systematic characterization of risk variants  
1212 and genes for coronary artery disease in over a million participants. *Nat. Genet.* **54**,  
1213 1803–1815 (2022).
- 1214 49. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human  
1215 phenotypes. *Nat. Genet.* **53**, 1415–1424 (2021).
- 1216 50. Roychowdhury, T. *et al.* Genome-wide association meta-analysis identifies risk  
1217 loci for abdominal aortic aneurysm and highlights PCSK9 as a therapeutic target. *Nat.*  
1218 *Genet.* **55**, 1831–1842 (2023).
- 1219 51. Roychowdhury, T. *et al.* Regulatory variants in TCF7L2 are associated with  
1220 thoracic aortic aneurysm. *Am. J. Hum. Genet.* **108**, 1578–1589 (2021).
- 1221 52. Miyazawa, K. *et al.* Cross-ancestry genome-wide analysis of atrial fibrillation  
1222 unveils disease biology and enables cardioembolic risk prediction. *Nat. Genet.* **55**,  
1223 187–197 (2023).

- 1224 53. Yu Chen, H. *et al.* Dyslipidemia, inflammation, calcification, and adiposity in  
1225 aortic stenosis: a genome-wide study. *Eur. Heart J.* **44**, 1927–1939 (2023).
- 1226 54. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic  
1227 discovery across human disease. *Cell Genomics* **2**, 100192 (2022).
- 1228 55. Kavousi, M. *et al.* Multi-ancestry genome-wide study identifies effector genes  
1229 and druggable pathways for coronary artery calcification. *Nat. Genet.* **55**, 1651–1664  
1230 (2023).
- 1231 56. Henry, A. *et al.* Mapping the aetiological foundations of the heart failure  
1232 spectrum using human genetics. *medRxiv* (2023) doi:10.1101/2023.10.01.23296379.
- 1233 57. Ishigaki, K. *et al.* Large-scale genome-wide association study in a Japanese  
1234 population identifies novel susceptibility loci across different diseases. *Nat. Genet.*  
1235 **52**, 669–679 (2020).
- 1236 58. Mishra, A. *et al.* Stroke genetics informs drug discovery and risk prediction  
1237 across ancestries. *Nature* **611**, 115–123 (2022).
- 1238 59. Roselli, C. *et al.* Genome-wide association study reveals novel genetic loci: a  
1239 new polygenic risk score for mitral valve prolapse. *Eur. Heart J.* **43**, 1668–1680 (2022).
- 1240 60. Hartiala, J. A. *et al.* Genome-wide analysis identifies novel susceptibility loci for  
1241 myocardial infarction. *Eur. Heart J.* **42**, 919–933 (2021).
- 1242 61. van Zuydam, N. R. *et al.* Genome-Wide Association Study of Peripheral Artery  
1243 Disease. *Circ. Genomic Precis. Med.* **14**, e002862 (2021).
- 1244 62. Adlam, D. *et al.* Genome-wide association meta-analysis of spontaneous  
1245 coronary artery dissection identifies risk variants and genes related to artery integrity  
1246 and tissue-mediated coagulation. *Nat. Genet.* **55**, 964–972 (2023).



- 1247 63. Pérez-Gutiérrez, L. & Ferrara, N. Biology and therapeutic targeting of vascular  
1248 endothelial growth factor A. *Nat. Rev. Mol. Cell Biol.* **24**, 816–834 (2023).
- 1249 64. Velagapudi, S. *et al.* VEGF-A Regulates Cellular Localization of SR-BI as Well as  
1250 Transendothelial Transport of HDL but Not LDL. *Arterioscler. Thromb. Vasc. Biol.* **37**,  
1251 794–803 (2017).
- 1252 65. Chen, H. X. & Cleck, J. N. Adverse effects of anticancer agents that target the  
1253 VEGF pathway. *Nat. Rev. Clin. Oncol.* **6**, 465–477 (2009).
- 1254 66. Tall, A. R., Thomas, D. G., Gonzalez-Cabodevilla, A. G. & Goldberg, I. J.  
1255 Addressing dyslipidemic risk beyond LDL-cholesterol. *J. Clin. Invest.* **132**, e148559  
1256 (2022).
- 1257 67. Zanoni, P. *et al.* Rare variant in scavenger receptor BI raises HDL cholesterol and  
1258 increases risk of coronary heart disease. *Science* **351**, 1166–1171 (2016).
- 1259 68. Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped  
1260 isolated population. *Nature* **613**, 508–518 (2023).
- 1261 69. Ritchie, S. C. *et al.* The Biomarker GlycA Is Associated with Chronic Inflammation  
1262 and Predicts Long-Term Risk of Severe Infection. *Cell Syst.* **1**, 293–301 (2015).
- 1263 70. Lotta, L. A. *et al.* Genetic Predisposition to an Impaired Metabolism of the  
1264 Branched-Chain Amino Acids and Risk of Type 2 Diabetes: A Mendelian  
1265 Randomisation Analysis. *PLoS Med.* **13**, e1002179 (2016).
- 1266 71. Zhao, Y. *et al.* Small rodent models of atherosclerosis. *Biomed. Pharmacother.*  
1267 *Biomedecine Pharmacother.* **129**, 110426 (2020).
- 1268 72. Shim, J., Al-Mashhadi, R. H., Sørensen, C. B. & Bentzon, J. F. Large animal models  
1269 of atherosclerosis--new tools for persistent problems in cardiovascular medicine. *J.*  
1270 *Pathol.* **238**, 257–266 (2016).

- 1271 73. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in  
1272 complex traits. *Nat. Genet.* **51**, 1339–1348 (2019).
- 1273 74. Smith, C. J. *et al.* Integrative analysis of metabolite GWAS illuminates the  
1274 molecular basis of pleiotropy and genetic correlation. *eLife* **11**, e79348 (2022).
- 1275 75. Julkunen, H. *et al.* Atlas of plasma NMR biomarkers for health and disease in  
1276 118,461 individuals from the UK Biobank. *Nat. Commun.* **14**, 604 (2023).
- 1277 76. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic  
1278 data. *Nature* **562**, 203–209 (2018).
- 1279 77. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes  
1280 of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779  
1281 (2015).
- 1282 78. Karczewski, K. J. *et al.* Pan-UK Biobank GWAS improves discovery, analysis of  
1283 genetic architecture, and resolution into ancestry-enriched effects. *medRxiv* (2024)  
1284 doi:10.1101/2024.03.13.24303864.
- 1285 79. Würtz, P. *et al.* Quantitative Serum Nuclear Magnetic Resonance Metabolomics  
1286 in Large-Scale Epidemiology: A Primer on -Omic Technologies. *Am. J. Epidemiol.* **186**,  
1287 1084–1096 (2017).
- 1288 80. Ritchie, S. C. *et al.* Quality control and removal of technical variation of NMR  
1289 metabolic biomarker data in ~120,000 UK Biobank participants. *Sci. Data* **10**, 64  
1290 (2023).
- 1291 81. Buuren, S. van & Groothuis-Oudshoorn, K. mice: Multivariate Imputation by  
1292 Chained Equations in R. *J. Stat. Softw.* **45**, 1–67 (2011).
- 1293 82. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for  
1294 quantitative and binary traits. *Nat. Genet.* **53**, 1097–1103 (2021).

- 1295 83. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of  
1296 genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
- 1297 84. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from  
1298 polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
- 1299 85. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to  
1300 variable selection in regression, with application to genetic fine mapping. *J. R. Stat.*  
1301 *Soc. Ser. B Stat. Methodol.* **82**, 1273–1300 (2020).
- 1302 86. Wakefield, J. Bayes factors for genome-wide association studies: comparison  
1303 with P-values. *Genet. Epidemiol.* **33**, 79–86 (2009).
- 1304 87. Rossen, J. *et al.* MultiSuSiE improves multi-ancestry fine-mapping in All of Us  
1305 whole-genome sequencing data. *medRxiv* (2024) doi:10.1101/2024.05.13.24307291.
- 1306 88. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and  
1307 richer datasets. *GigaScience* **4**, 7 (2015).
- 1308 89. Keller, M. C. Gene × environment interaction studies have not properly controlled  
1309 for potential confounders: the problem and the (simple) solution. *Biol. Psychiatry* **75**,  
1310 18–24 (2014).
- 1311 90. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression  
1312 and splicing quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
- 1313 91. Ochoa, D. *et al.* The next-generation Open Targets Platform: reimaged,  
1314 redesigned, rebuilt. *Nucleic Acids Res.* **51**, D1353–D1359 (2023).
- 1315 92. Kanehisa, M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*  
1316 *Res.* **28**, 27–30 (2000).
- 1317 93. Milacic, M. *et al.* The Reactome Pathway Knowledgebase 2024. *Nucleic Acids*  
1318 *Res.* **52**, D672–D678 (2024).

- 1319 94. Murphy, A. E., Schilder, B. M. & Skene, N. G. MungeSumstats: a Bioconductor  
1320 package for the standardization and quality control of many GWAS summary  
1321 statistics. *Bioinformatics* **37**, 4593–4596 (2021).
- 1322 95. Wallace, C. Eliciting priors and relaxing the single causal variant assumption in  
1323 colocalisation analyses. *PLOS Genet.* **16**, e1008720 (2020).
- 1324 96. Wallace, C. A more accurate method for colocalisation analysis allowing for  
1325 multiple causal variants. *PLOS Genet.* **17**, e1009440 (2021).
- 1326 97. Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science* **347**,  
1327 1260419 (2015).
- 1328 98. Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of  
1329 Genetic Association Studies Using Summary Statistics. *PLoS Genet.* **10**, e1004383  
1330 (2014).
- 1331 99. Burgess, S., Small, D. S. & Thompson, S. G. A review of instrumental variable  
1332 estimators for Mendelian randomization. *Stat. Methods Med. Res.* **26**, 2333–2355  
1333 (2017).
- 1334 100. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical  
1335 and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **57**,  
1336 289–300 (1995).
- 1337 101. Backman, J. D. *et al.* Exome sequencing and analysis of 454,787 UK Biobank  
1338 participants. *Nature* **599**, 628–634 (2021).
- 1339 102. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122  
1340 (2016).
- 1341 103. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the  
1342 Pathogenicity of Rare Missense Variants. *Am. J. Hum. Genet.* **99**, 877–885 (2016).

- 1343 104. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of  
1344 human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
- 1345 105. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from  
1346 variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- 1347