

Assessing the contribution of rare-to-common protein-coding variants to circulating metabolic biomarker levels via 412,394 UK Biobank exome sequences

Abhishek Nag¹, Lawrence Middleton¹, Ryan S. Dhindsa^{2,3}, Dimitrios Vitsios¹, Eleanor Wigmore¹, Erik L. Allman¹, Anna Reznichenko⁴, Keren Carss¹, Katherine R. Smith¹, Quanli Wang², Benjamin Challis⁴, Dirk S. Paul¹, Andrew R. Harper¹, Slavé Petrovski¹

¹Centre for Genomics Research, Discovery Sciences, BioPharmaceuticals R&D, AstraZeneca, Cambridge, UK

²Centre for Genomics Research, Discovery Sciences, BioPharmaceuticals R&D, AstraZeneca, Waltham, USA

³Department of Molecular and Human Genetics, Baylor College of Medicine and Jan and Dan Duncan Neurological Research Institute at Texas Children's Hospital, Houston, TX 77030, USA

⁴Translational Science and Experimental Medicine, Research and Early Development, Cardiovascular, Renal and Metabolism (CVRM), BioPharmaceuticals R&D, AstraZeneca, Gothenburg, Sweden

Corresponding author:

Slavé Petrovski

Vice-President, Centre for Genomics Research,

Discovery Sciences, BioPharmaceuticals R&D

AstraZeneca

Cambridge

United Kingdom

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.
Email: slav.petrovski@astrazeneca.com

1 **Abstract**

2
3 Genome-wide association studies have established the contribution of common and low
4 frequency variants to metabolic biomarkers in the UK Biobank (UKB); however, the role of
5 rare variants remains to be assessed systematically. We evaluated rare coding variants for
6 198 metabolic biomarkers, including metabolites assayed by Nightingale Health, using
7 exome sequencing in participants from four genetically diverse ancestries in the UKB
8 (N=412,394). Gene-level collapsing analysis – that evaluated a range of genetic
9 architectures – identified a total of 1,303 significant relationships between genes and
10 metabolic biomarkers ($p < 1 \times 10^{-8}$), encompassing 207 distinct genes. These include
11 associations between rare non-synonymous variants in *GIGYF1* and glucose and lipid
12 biomarkers, *SYT7* and creatinine, and others, which may provide insights into novel disease
13 biology. Comparing to a previous microarray-based genotyping study in the same cohort, we
14 observed that 40% of gene-biomarker relationships identified in the collapsing analysis were
15 novel. Finally, we applied Gene-SCOUT, a novel tool that utilises the gene-biomarker
16 association statistics from the collapsing analysis to identify genes having similar biomarker
17 fingerprints and thus expand our understanding of gene networks.

18 **Introduction**

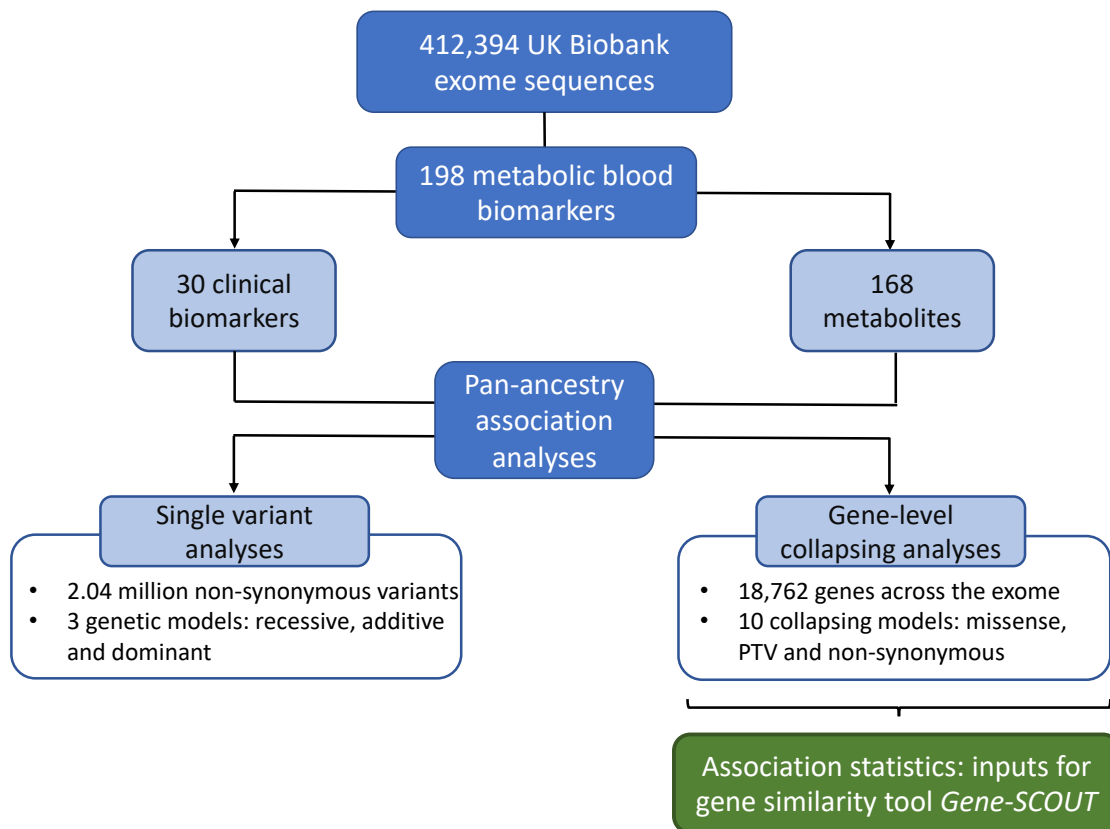
19 Metabolic blood biomarkers represent intermediate or end products of biochemical pathways
20 that can be used to diagnose and monitor human disease. The application of metabolic
21 biomarkers as intermediate traits to dissect the genetic basis of complex human diseases is
22 well-established. Investigating the genetic underpinnings of blood biomarkers can offer novel
23 insights into human disease mechanisms and, in turn, provide potential therapeutic targets.
24 Large-scale genome-wide association studies (GWAS) have so far identified hundreds of
25 genetic loci that regulate blood biomarker and metabolite levels¹⁻¹¹; however, difficulty in
26 mapping these loci to causal genes and interpreting functional effects of non-coding variants
27 have stymied the clinical impact for many of these associations¹².

28
29 The UK Biobank (UKB)¹³ is a large population-based resource of ~500,000 participants with
30 genetic data linked to a diverse set of phenotypic measurements. Genotype data from
31 microarrays and large population-based imputation panels have helped establish the
32 contribution of common and low frequency variants towards blood biomarkers in the UKB¹⁴.
33 The availability of exome sequences in the same population now allows for the exploration of
34 rare coding variants regulating metabolic blood biomarkers. Associations for rare coding
35 variants have demonstrably greater translational potential given their larger effect sizes¹⁵
36 and our ability to more directly interpret their functional impact¹⁶.

37
38 Using exome sequences from 412,394 unrelated participants across multiple genetic
39 ancestries in the UKB, we present findings of variant-level and gene-level (collapsing)
40 association tests for 198 metabolic blood biomarkers. We then introduce a novel tool, Gene-
41 SCOUT, that utilises this rich catalogue of gene-biomarker association statistics to identify
42 genes with similar biomarker fingerprints as a given (target) gene of interest and expand our
43 understanding of gene networks.

44 Results

45 In this study, we analysed 198 metabolic blood biomarkers, including 30 clinical blood
46 biomarkers related to glucose and lipid metabolism, renal and liver function (**Table S1A**),
47 and an additional 168 Nightingale assay blood metabolite measurements related to
48 lipoprotein lipids, fatty acids and their compositions, and various other low-molecular weight
49 metabolites¹⁷ (**Table S1B**). Most of the metabolic biomarkers pertain to lipid metabolism
50 (77%) and correlate highly with each other (**Figure 2a**). Many metabolic biomarkers also
51 demonstrate strong associations with clinical traits documented in the UKB (**Figure 2b**).

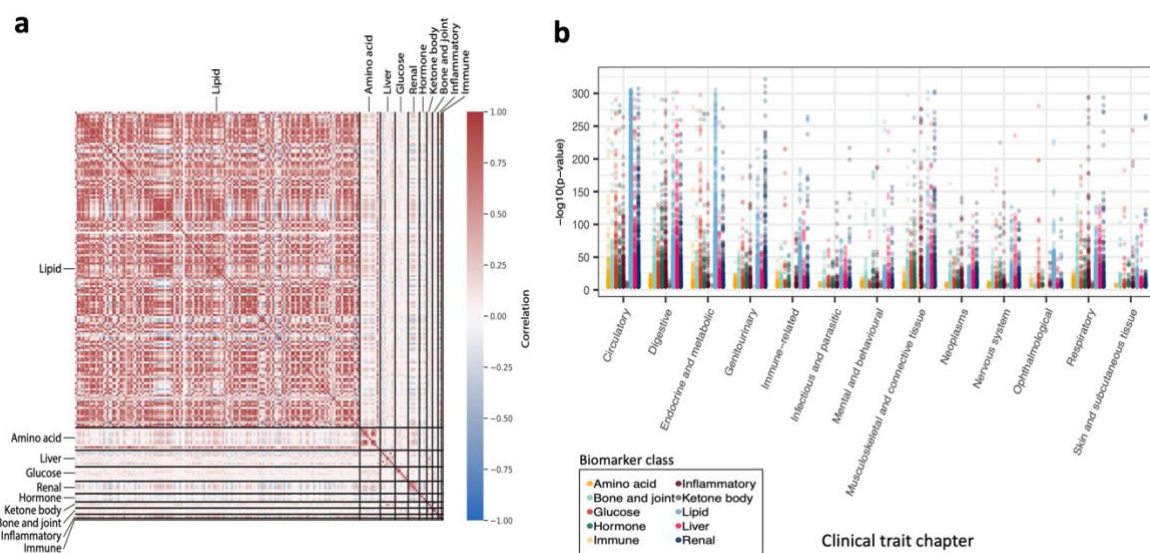


52
53 **Figure 1: A schematic of the association analyses that were conducted for the**
54 **metabolic blood biomarkers using the UK Biobank exome sequences**

55 The UK Biobank exome sequences were used to conduct single variant (under 3 genetic models) and
56 gene-level (under 10 collapsing models) association analyses for the clinical blood biomarkers (N=30)
57 and blood metabolite measurements (N=168). The gene-level association statistics for these
58 metabolic biomarkers were used as inputs for the gene similarity tool *Gene-SCOUT*.

59
60 We first conducted a single variant analysis between the non-synonymous coding
61 variants (N=2,043,019 for the European ancestry subset) and the 198 metabolic biomarkers
62 (**Figure 1**). Excluding the MHC region, 19,351 significant variant-biomarker associations

63 ($p < 1 \times 10^{-8}$) were identified in the European subset of UK Biobank, which mapped to 12,217
 64 significant relationships between genes and biomarkers (**Tables S2A, S2B**). Pruning
 65 variants in linkage disequilibrium (r^2 threshold of 0.5) resulted in 9,738 significant gene-
 66 biomarker relationships. Notably, 243 distinct PTVs accounted for 1,366 significant
 67 associations, of which 602 (44%) were attributable to rare PTVs (MAF < 0.1%) with large
 68 effect sizes (>0.5 SD) (**Figure 4a**) (**Tables S3A, S3B**). We identified 28 PTVs with MAF as
 69 low as 0.001% that achieved significance in the ExWAS: these include associations relating
 70 to several well-established and biologically plausible relationships such as *GOT1* and
 71 aspartate aminotransferase, *CST3* and cystatin C, *APOB* and cholesterol biomarkers, *ALPL*
 72 and alkaline phosphatase (**Table S3A**). Among other PTV findings that may provide new
 73 insights into important biology, a rare frameshift variant (MAF=0.03%) in *PLIN1* – a gene
 74 known to cause familial lipodystrophy¹⁸ – was associated with HDL-cholesterol (beta=0.40
 75 [0.27,0.53], $p = 1.6 \times 10^{-9}$), and a rare splice variant (MAF=0.09%) in *TNFRSF10B* – loss of
 76 which has been reported to promote survival of virus-infected liver cells¹⁹ – was associated
 77 with gamma glutamyltransferase (beta=0.21 [0.14,0.28], $p = 3.8 \times 10^{-9}$) (**Table S3A**).



78
 79 **Figure 2: Characteristics of metabolic blood biomarkers analysed in this study**
 80 The 198 metabolic blood biomarkers analysed in this study were grouped into the following 10
 81 biological classes: lipid, amino acid, liver, glucose, renal, hormone, ketone body, bone and joint,
 82 inflammatory, and immune. (a) The plot demonstrates that metabolic biomarkers belonging to the
 83 same biological class are correlated with each other. (b) Strong associations (plotted on the Y-axis)
 84 were observed between the metabolic biomarkers and 15,719 clinical traits (grouped by chapter)
 85 documented in the UKB²⁰.

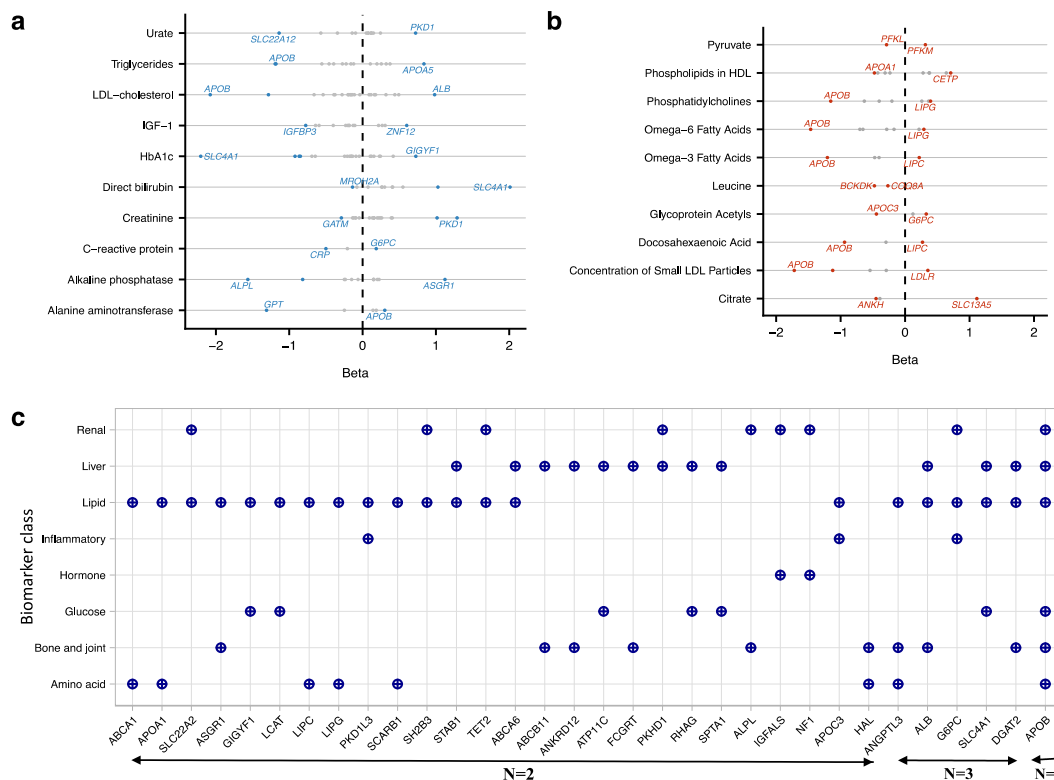
86 Next, we performed a gene-level collapsing analysis that tests the aggregate effect of
87 rare functional variants in each gene. We employed 10 different models to capture a diverse
88 range of genetic architectures (**Methods**). In the analysis involving individuals of European
89 ancestry alone, we identified 1,303 significant relationships between genes and metabolic
90 biomarkers ($p < 1 \times 10^{-8}$) (**Tables S4A, S4B; Figures 3a, 3b**). Most (68%, 880/1,303) gene-
91 biomarker relationships detected via the collapsing analysis were captured through models
92 that exclusively focused on PTV classes (“ptv” and “ptv5pct”), while the remaining 32%
93 were attributable to models that incorporated missense variants. We detected more
94 significant associations using our “ptv” and “ptv5pct” models than a prior study²¹ that also
95 performed gene-level collapsing analysis using the UKB exome sequence data, albeit with a
96 different analytical framework. For instance, associations between PTVs in 12 genes and
97 HbA1c that we detected were not reported in the other study: this includes the glucose
98 metabolism genes *HK1* and *G6PC2* (**Figure S1**). We also extended our gene-level
99 collapsing analysis to include all ancestral groups in the UKB (**Methods**). This detected an
100 additional 51 significant gene-biomarker relationships (**Tables S5A, S5B**). For the gene-
101 biomarker relationships that were significant only in the pan-ancestry analysis, we did not
102 observe a significant difference in the estimated effect size between the European-only and
103 the pan-ancestry analyses ($p=0.83$), suggesting that increased statistical power rather than
104 ancestry-specific effects is the more likely reason why these associations were identified in
105 the pan-ancestry analysis. One such association detected exclusively in the pan-ancestry
106 analysis was between recessive carriers of nonsynonymous variants in the membrane
107 transport gene *SYT7* and blood creatinine levels (number of QV carriers=5, $\beta=2.17$
108 [1.46,2.87], $p=1.6 \times 10^{-9}$). With 3 of the 5 carriers observed in the South Asian and African
109 ancestry participants, the pan-ancestry analysis facilitated detection of this association,
110 which was not study-wide significant in the European subset (number of QV carriers=2,
111 $\beta=1.17$ [0.06,2.28], $p=0.04$). Remarkably consistent with the biomarker findings,
112 recessive carriers of *SYT7* PTVs demonstrate a increased risk of glomerular disease in the

113 pan-ancestry analysis (OR=92.1 [12.1,713.2], $p=2.6 \times 10^{-5}$), but the clinical association on its
114 own is not yet study-wide significant.

115 The significant gene-level relationships from the collapsing analyses encompassed
116 207 distinct genes, of which 32 were associated with biomarkers across different biological
117 classes (**Figure 3c**). This includes *GIGYF1*, a tyrosine kinase receptor signalling protein, in
118 which rare PTVs were associated with biomarkers of glucose [glucose (beta=0.59
119 [0.42,0.76], $p=7.9 \times 10^{-12}$) and HbA1c (beta=0.73 [0.57,0.88], $p=4.5 \times 10^{-20}$)] and cholesterol
120 metabolism [total cholesterol (beta=-0.66 [-0.82,-0.50], $p=2.0 \times 10^{-15}$), LDL-cholesterol (beta=-
121 0.61 [-0.78,-0.45], $p=3.4 \times 10^{-13}$) and apolipoprotein B (beta=-0.60 [-0.77,-0.44], $p=1.3 \times 10^{-12}$)].
122 Additionally, among clinical traits documented in the UKB²⁰, significant associations were
123 observed for rare PTVs in *GIGYF1* with the risk of hypothyroidism (OR=4.2 [2.7,6.6],
124 $p=7.1 \times 10^{-9}$) and type 2 diabetes (OR=4.0 [2.7,5.8], $p=1.0 \times 10^{-10}$). Since hypothyroidism is
125 known to raise LDL-cholesterol levels, we subsequently tested the *GIGYF1*-LDL-cholesterol
126 association adjusted for a diagnosis of hypothyroidism. The signal between *GIGYF1* PTVs
127 and LDL-cholesterol (adjusted for the effect of statins) remained significant upon adjusting
128 for hypothyroidism (beta=-0.55 [-0.71,-0.38]; $p=6.2 \times 10^{-11}$), suggesting that the *GIGYF1* locus
129 likely influences cholesterol levels independent of solely thyroid hormone-mediated
130 pathways. Thus, by leveraging information from over 400,000 UKB exomes, our study
131 provides a more comprehensive picture regarding *GIGYF1*'s biomarker fingerprint and
132 associated clinical traits, expanding on previously reported common⁷ and rare variant
133 associations^{21,22} at this locus.

134 We observed that adjusting biomarkers for medications that influence their levels can
135 also improve detection of associations: 31/84 (37%) significant gene-biomarker relationships
136 for apolipoprotein B, LDL-cholesterol, total cholesterol, and urate from the collapsing
137 analysis were detected only after we adjusted their values for commonly prescribed
138 medications (**Table S4A**). This includes association between putatively damaging missense
139 variants and PTVs in *HMGCR* ("flexdmg" model) and LDL-cholesterol (medication-adjusted:
140 beta=-0.19 and $p=1.7 \times 10^{-11}$; medication-unadjusted: beta=-0.15 and $p=6.1 \times 10^{-8}$), which

141 validates the value of medication adjustment to untangle the effects of therapeutic
 142 intervention vs natural aberration of *HMGCR*. Moreover, for gene-biomarker relationships
 143 that were significantly associated in both the medication-unadjusted and the medication-
 144 adjusted analyses (N=52), the absolute effect sizes were observably higher in the latter
 145 (**Figure S2**), but the difference was not statistically significant in the current sample (Mann
 146 Whitney $p=0.28$).



147 **Figure 3: Significant relationships between genes and metabolic blood biomarkers**
 148 **identified in the collapsing analysis**
 149 (a, b) Significant gene relationships ($p < 1 \times 10^{-8}$) identified for select clinical biomarkers and metabolites
 150 in the collapsing analysis have been shown. The genes with the highest absolute effect sizes for each
 151 have been labelled.
 152 (c) The plot lists the 32 genes that were significantly associated ($p < 1 \times 10^{-8}$) with metabolic biomarkers
 153 across two or more biological classes in the collapsing analysis. For each such gene, the
 154 corresponding biological classes have been indicated.

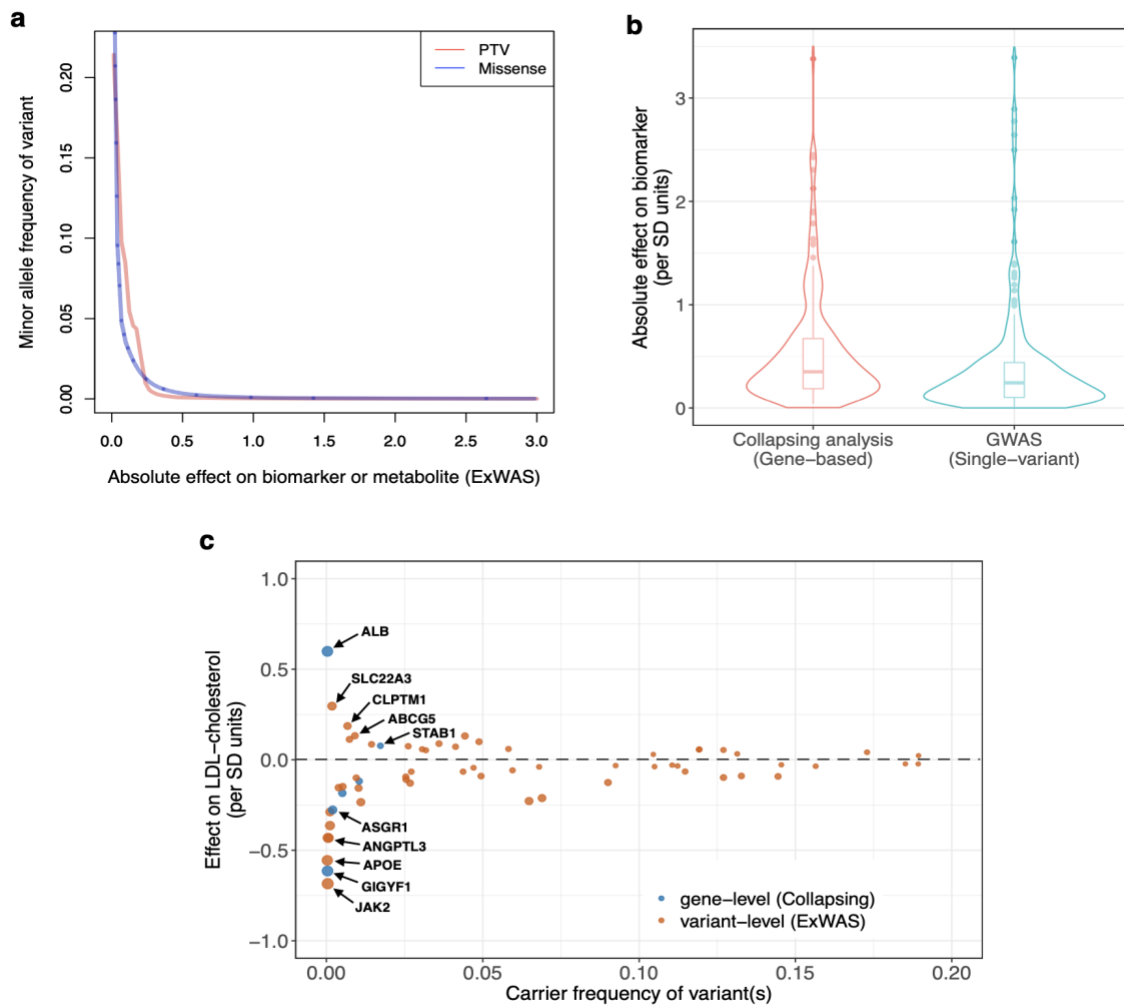
156
 157 **Gene-level collapsing analysis: capturing allelic series**

158 We observed that 17% (215/1,303) of significant relationships between genes and
 159 metabolic biomarkers from the collapsing analysis did not achieve significance in the
 160 respective variant-level ExWAS (**Table S6**). Next, we also compared the gene-biomarker
 161 relationships that achieved significance in the collapsing analysis (**Tables S4A, S7**) and the

162 microarray-based GWAS¹⁴ (as per a less stringent significance threshold: $p < 1 \times 10^{-7}$) for the
163 32 biomarkers (28 blood and 4 urinary biomarkers) analysed in both studies. Of the
164 significant gene-biomarker relationships identified in the collapsing analysis, 40% (142/357)
165 were not detected in the microarray-based GWAS (**Table S8**). These include associations
166 for well-known drug target genes such as *HMGCR* (with LDL-cholesterol) and *PPARG* (with
167 HDL-cholesterol). Furthermore, the effect size estimates were significantly higher in the
168 collapsing analysis than in microarray-based GWAS for 215 gene-biomarker relationships
169 detected via both approaches (Mann-Whitney $p = 8.0 \times 10^{-6}$) (**Table S9; Figure 4b**). One likely
170 explanation for this is that by testing aggregate effects of rare putative functional variants in
171 a gene, associations arising from collapsing analysis are enriched for larger effects (**Figure**
172 **4c**). Collectively, these results highlight that application of a gene-based rare variant
173 collapsing analysis to large-scale exome sequencing can increase power to capture
174 associations that are driven by an allelic series, and thus expand our understanding of the
175 genetic architecture of traits, especially where a lot of success has already been achieved
176 through traditional microarray-based GWAS.

177 *SLC4A1*, which encodes a chloride/bicarbonate anion exchange protein in the red cell
178 membrane, represents one such gene for which multiple signals were detected in the gene-
179 level collapsing analysis but not in the ExWAS. We observed 32 carriers for 28 distinct *SLC4A1*
180 PTVs, of which 25 (89%) were private (i.e., observed in a single carrier) (**Figure S3**). Overall,
181 *SLC4A1* PTVs were significantly associated with a strong reduction in HbA1c (beta=-2.2 [-2.6,-
182 1.8], $p = 1.4 \times 10^{-25}$) and LDL-cholesterol (beta=-1.0 [-1.4,-0.7], $p = 8.0 \times 10^{-9}$), while also showing
183 strong increases in total bilirubin (beta=1.7 [1.3,2.0], $p = 1.1 \times 10^{-22}$) and direct bilirubin
184 (beta=2.0 [1.7,2.4], $p = 1.8 \times 10^{-28}$). Among clinical phenotypes, *SLC4A1* PTVs are significantly
185 associated with disorders of reduced red cell membrane stability such as hereditary
186 spherocytosis and hereditary haemolytic anaemia, but not with any phenotype related to
187 glucose or lipid metabolism ($p < 1 \times 10^{-5}$). Similarly, in ClinVar, several missense and loss-of-

188 function mutations in this gene are reported as pathogenic for hereditary spherocytosis.
189 Therefore, we further tested the *SLC4A1*–biomarker associations after adjusting for the
190 diagnosis of hereditary spherocytosis or hereditary haemolytic anaemia and relevant blood
191 cell indices, including red cell distribution width (RDW) and mean corpuscular haemoglobin



192
193
194
195
196
197
198
199
200
201
202
203
204
205
206

Figure 4: Effects of coding variants on metabolic blood biomarkers

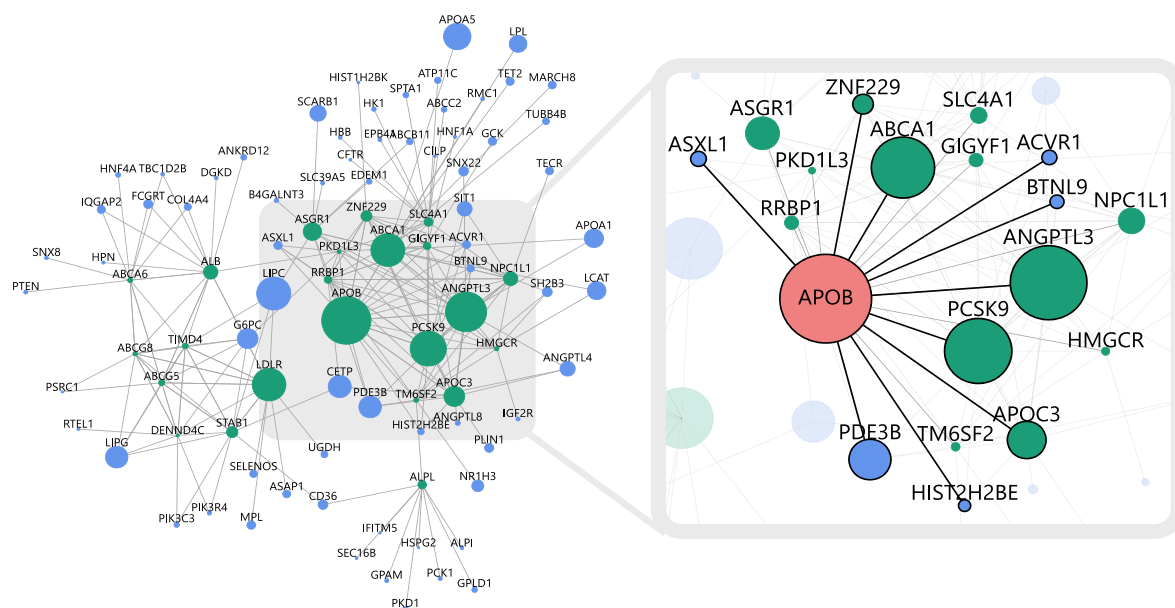
(a) Absolute effect sizes for missense variants and PTVs significantly associated ($p < 1 \times 10^{-8}$) with metabolic biomarkers in the single variant analysis (ExWAS) as a function of their minor allele frequency (in cases where a missense variant or a PTV was significantly associated with more than one biomarker, the association with the highest absolute effect size was selected). (b) The effect sizes estimated in the gene-based collapsing analysis and the Sinnott-Armstrong *et al* microarray-based GWAS¹⁴ were compared for the gene-biomarker relationships that were significantly associated in both ($N=215$). For each significant gene-biomarker relationship, the collapsing model (from the collapsing analysis) and the individual variant (from the microarray-based GWAS) with the highest absolute effect sizes were selected. The effect sizes estimated in the collapsing analysis were significantly higher than that in the GWAS (Mann-Whitney $p=8.0 \times 10^{-6}$). (c) Comparing effect sizes for individual variants and aggregate of rare variants (in a gene) that were significantly associated ($p < 1 \times 10^{-8}$) with LDL-cholesterol. Some examples of genes significantly associated with LDL-cholesterol have been highlighted. The Y-axis has been capped at 1 SD units for visual clarity.

207
208 concentration (MCHC). The gene-based *SLC4A1* PTV signals remained significant in the
209 adjusted analyses (**Table S10**). Although PTVs in this gene may be independently associated
210 with biomarkers of glucose, lipid and bilirubin metabolism, we cannot rule out the
211 possibility of under-reporting of hereditary spherocytosis and hereditary haemolytic
212 anaemia in the UKB that explains these observations. The *SLC4A1* enigma is consistent with
213 previous reports of other red blood cell loci that have also been significantly associated with
214 HbA1c²³.

215 **Gene-SCOUT: estimating gene similarity based on cohort statistics from collapsing analysis**

216 We considered the opportunity to leverage this new and rich catalogue of gene-level
217 association statistics from the collapsing analysis to determine genes with similar biomarker
218 fingerprints. To achieve this, we developed a gene similarity tool 'Gene-SCOUT'²⁴, that
219 solely uses the gene-level collapsing analysis statistics across the studied biomarkers to
220 identify genes with the most comparable biomarker genetic associations as a given gene of
221 interest. No other information is used in constructing the gene similarity scores. Since this
222 tool estimates gene similarity for an index gene by selecting features based on the
223 significance cut-off of $p < 1 \times 10^{-5}$, gene neighbours could not be determined for genes that did
224 not achieve association $p < 1 \times 10^{-5}$ with any biomarker feature. Accordingly, for our feature set
225 comprising of 198 biomarkers, we were able to determine gene similarity for 3%
226 (536/18,762) of human protein-coding genes. To illustrate Gene-SCOUT's application, we
227 selected the 24 genes that were significantly associated ($p < 1 \times 10^{-8}$) with LDL-cholesterol in
228 the collapsing analysis. We used each gene in this set as a seed gene to construct a
229 network figure that demonstrates their respective gene neighbours (**Figure 5**). Using *APOB*
230 as an example, we observe that genes with the most comparable biomarker fingerprint as
231 *APOB* include: *ABCA1*, *ACVR1*, *APOC3*, *ANGPTL3*, *ASGR1*, *ASXL1*, *BTNL9*, *GIGYF1*,
232 *HIST2H2BE*, *HMGCR*, *NPC1L1*, *PCSK9*, *PDE3B*, *PKD1L3*, *RRBP1*, *SLC4A1*, *TM6SF2* and
233 *ZNF229*. For some of these genes (e.g., *ZNF229*, *ACVR1*), the links with lipid metabolism

234 appear to be novel, in addition to the recently described relationships for *GIGYF1*^{21,22}.
235 Inhibition of *APOB*, such as through mipomersen, is known to be clinically effective in
236 reducing blood cholesterol levels. Remarkably, 5 (namely, *APOC3*, *ANGPTL3*, *HMGCR*,
237 *NPC1L1* and *PCSK9*) of the 18 genes (28%) determined to have similar cohort-level genetic
238 associations with biomarkers as *APOB* are also targets of lipid-lowering drugs that are already
239 approved or in various stages of development (<https://www.fda.gov/drugs>).



240
241 **Figure 5: Network figure demonstrating the gene neighbours i.e., genes with most**
242 **similar biomarker genetic signals, as the set of genes that were significantly**
243 **associated with LDL-cholesterol in the collapsing analysis**
244 The 24 genes that were significantly associated ($p < 1 \times 10^{-8}$) with LDL-cholesterol in the collapsing
245 analysis were used as seed genes (green nodes) to construct a network figure demonstrating
246 respective gene neighbours (edges). Non-seed genes are represented using blue nodes. The size of
247 a gene node corresponds to the number of features (of total 198) that the gene is associated with
248 at $p < 1 \times 10^{-5}$. The inset demonstrates the genes with most similar biomarker signature as *APOB* – these
249 include the ten closest genes for *APOB* as the seed gene (black edges) and other seed genes that
250 have *APOB* among their ten closest genes (grey edges).

251 Discussion

252 We used the 454,796 UK Biobank exome sequences to explore the contribution of private-to-
253 rare-to-common protein-coding variation for 30 clinical biomarkers and 168 metabolite
254 measurements. By adopting variant- and gene-level analysis frameworks and assessing the
255 full allelic frequency spectrum, we have expanded our understanding of the genetic
256 architecture of metabolic biomarkers that have previously been studied through microarray
257 data. The finding that 17% of gene-biomarker relationships detected in the gene-level
258 collapsing analysis were not identified in the single variant analysis demonstrates the power
259 of testing an aggregate effect of rare variants in a gene encompassing a range of genetic
260 architectures. We also illustrated how adjusting biomarker values for commonly prescribed
261 medications can improve signal detection.

262 There are several strengths of our study that might have implications for identifying or
263 validating drug targets. First, by virtue of focusing on coding variants, the observed
264 associations could provide a more causal link between a gene and a blood biomarker^{25–28}.
265 Moreover, association signals emerging from collapsing analysis are driven by an aggregate
266 effect of multiple rare variants (allelic series) that tend to be less impacted by local LD
267 structure. This contrasts with associations identified in microarray-based GWAS that often
268 map to non-coding regions of the genome or to regions of extensive LD, making it more
269 challenging to pinpoint the underlying causal variants.

270 Associations involving putative functional variants can also indicate the desired
271 modulation of the target gene e.g., upregulation or downregulation of the target gene
272 product, required to mitigate the risk of the disease related to the associated biomarker. For
273 instance, we observed a total of 182 associations for rare (MAF<0.1%) PTVs with the 30
274 blood biomarkers, which is >3-fold more than the 53 conditionally independent PTV
275 associations (for the same set of blood biomarkers) reported in the microarray-based
276 analysis¹⁴.

277 We also introduce a novel tool (Gene-SCOUT) that utilises all the gene-level
278 collapsing analysis statistics across the 198 studied biomarkers to estimate a ‘similarity’
279 metric between genes. With the aid of specific examples, we were able to demonstrate that
280 this approach can successfully identify genes with similar biomarker fingerprints.

281 While there are certain advantages of using blood biomarkers to dissect the genetics
282 of complex human diseases, including greater statistical power offered by quantitative traits
283 and better insights into biological pathways underlying associations, further work is
284 necessary to establish the causal relationship between genetic loci identified using
285 biomarkers or metabolites and the related disease(s). For instance, we observed
286 associations between certain biomarkers and variants in genes that encode them (e.g., *ALB*
287 with albumin, and *CST3* with cystatin C) – although such associations serve as excellent
288 positive controls that demonstrate the robustness of our analysis framework, they may not
289 offer novel insights into disease pathophysiology.

290 Using the largest collection of exome sequences linked to a diverse set of circulating
291 metabolic biomarkers, we demonstrate the value of this resource to enhance our
292 understanding of human diseases, and potentially, provide novel therapeutic targets focused
293 on mimicking natural human genetic discoveries. Our study also strongly supports the use of
294 a gene-based collapsing framework to uncover gene-biomarker relationships that are driven
295 by an aggregate effect of multiple rare, non-synonymous variants.

296

297 **Methods**

298

299 **UK Biobank (UKB) Resource**

300 The UKB resource¹³ is a prospective cohort study of ~500,000 individuals from across the
301 United Kingdom, aged between 40 and 69 years. The average age at recruitment for the
302 sequenced participants was 56.5 years and 54% of the sequenced cohort are females.
303 Participant data, obtained through questionnaires and assessment visits, include health
304 records that are periodically updated by the UKB, self-report survey information, linkage to
305 death and cancer registries, urine and blood biomarkers, imaging data, accelerometer data
306 and various other phenotypic endpoints¹³. All study participants provided informed consent.
307 For this study, data from the UKB resource was accessed under the application number
308 26041.

309

310 **Metabolic blood biomarkers**

311 Routine clinical blood biomarkers related to glucose and lipid metabolism, renal and liver
312 function, among others (N=30), were measured in the majority of the ~500,000 UKB
313 participants (**Table S1A**). Additionally, 168 blood metabolites, including lipoprotein lipids,
314 fatty acids and their compositions, and various low-molecular weight metabolites, were
315 profiled in a subset of ~120,000 UKB participants by Nightingale Health using nuclear
316 magnetic resonance spectroscopy¹⁷ (**Table S1B**). Samples with a 'quality control (QC) flag'
317 for the blood metabolites were excluded. In total, we analysed 198 metabolic blood
318 measures: 30 clinical biomarkers and 168 metabolites. We applied rank-based inverse-
319 normal transformation to the measurements prior to performing association analyses.

320 For four blood biomarkers (LDL-cholesterol, total cholesterol, apolipoprotein B and
321 urate) we adjusted for the effect of commonly prescribed medications known to influence
322 their levels. For LDL-cholesterol, total cholesterol and apolipoprotein B, we adjusted for the
323 effect of statins based on their 'statin adjustment factors', previously estimated in the UKB as
324 0.684, 0.749 and 0.719, respectively¹⁴. Similarly, we adjusted urate for the effect of

325 allopurinol based on an ‘allopurinol adjustment factor (0.810)’, calculated using an approach
326 identical to that described for statins¹⁴.

327

328 **Whole-exome sequencing and bioinformatics pipeline**

329 Whole-exome sequences for 454,988 UKB participants were generated at the Regeneron
330 Genetics Center as part of a pre-competitive data generation collaboration between AbbVie,
331 Alnylam Pharmaceuticals, AstraZeneca, Biogen, Bristol-Myers Squibb, Pfizer, Regeneron
332 and Takeda²⁹. The exome sequencing procedure and the relevant QC steps have been
333 detailed previously in Szustakowski *et al* (2021)²⁹ and Wang *et al* (2021)²⁰. The FASTQ
334 sequences that were made available were first aligned, following which, single nucleotide
335 variants (SNVs) and small indels were called using Illumina’s DRAGEN Bio-IT Platform
336 Germline Pipeline v3.0.7 on the Amazon Web Services cloud compute platform available at
337 AstraZeneca’s Centre for Genomics Research. SNPEff v4.3³⁰ was used to annotate the
338 ‘most damaging effect’ predicted for each protein coding variant. In addition, we used certain
339 other bioinformatic tools such as missense tolerance ratio (MTR) scores³¹ to identify regions
340 of protein coding genes under constraint for missense variants, and REVEL³² to prioritise
341 coding variants based on their predicted deleteriousness. Further details on how these tools
342 were applied to the UKB exome sequencing dataset have been previously described²⁰.

343

344 **Selection of UKB samples for the association analyses**

345 Prior to performing the association analyses, we excluded samples from the available UKB
346 exome sequencing dataset (N=454,796) based on the following QC measures²⁰ (**Figure S4**):
347 (i) *DNA contamination*: VerifyBAMID freemix (measure of DNA contamination) >4%.
348 (ii) *Coverage depth*: $\geq 10x$ for <94.5% of the consensus coding sequence (CCDS release
349 22).
350 (iii) *Relatedness*: 2nd-degree relatives or closer (equivalent to kinship coefficient >0.0884), as
351 estimated using the --kinship function in KING v2.2.2³³.

352 Additionally, to perform analyses accounting for differing genetic ancestry, we
353 assigned samples to one of the four major ancestral groups (minimum 1,000 participants):

354 European (N=394,695), South Asian (N=8,078), East Asian (N=2,209) and African
355 (N=7,412). This was done by excluding participants: (i) with predicted genetic ancestry
356 <0.99 (for European ancestry) or <0.95 (for the remaining ancestries), as estimated using
357 PEDDY v0.4.2; or (ii) lying outside four standard deviations for the top four principal
358 components for each of the genetic ancestry collections.

359 **Association analysis for metabolic blood biomarkers**

361 A number of stringent variant-level QC steps, detailed previously²⁰, were applied to select
362 variant calls with highest confidence for association testing. Briefly, the variant-level QC
363 criteria included coverage depth, genotype and mapping quality scores, DRAGEN variant
364 status, read position rank sum score (RPRS), mapping quality rank sum score (MQRS),
365 alternate allele read proportion for heterozygous calls, proportion of samples failing any of
366 these QC criteria, and gnomAD-related filters.

367 Association testing between the metabolic blood biomarkers and the variants in the
368 exome sequencing dataset was conducted using two complementary analytical approaches
369 **(Figure 1)**:

- 370 (i) Single variant exome-wide association study (ExWAS)
- 371 (ii) Gene-level collapsing analysis

372 We conducted the association analyses separately in the European ancestry
373 participants as this comprised the single largest ancestral group in this resource and for all
374 four ancestries combined ('pan-ancestry' analysis).

375 Single variant exome-wide association study (ExWAS)

377 In the single-variant analysis (hereafter referred to as 'ExWAS'), variants that passed the QC
378 steps were filtered further to include those that had a minimum of six carriers (equivalent to
379 MAF>0.0008% in the European ancestry subset). We additionally excluded variants that had
380 one of the following annotations as their most damaging effect as per SNPEff:

381 *3_prime_UTR, 5_prime_UTR, initiator_codon_variant, non_coding_transcript_exon_variant,*

382 and *synonymous_variant*. The remaining non-synonymous coding variants (N=2,043,019 in
383 the European ancestry subset) were used to perform the ExWAS.

384 The ExWAS was conducted by fitting a linear regression model adjusted for age, sex
385 and BMI (for blood metabolites only), using the tool PEACOK that was developed as a
386 modification of the R package PHESANT³⁴. For the pan-ancestry analysis, we additionally
387 included the categorical ancestral group and top five ancestry principal components as
388 covariates. For each of the 198 biomarkers, three different genetic models were evaluated in
389 the ExWAS: (i) genotypic (AA vs AB vs BB), (ii) dominant (AA+AB vs BB), and (iii) recessive
390 (AA vs AB+BB), where A and B denote the reference and alternative alleles, respectively. A
391 significance cut-off of $p < 1 \times 10^{-8}$ was adopted for the ExWAS³⁵.

392 Gene-level collapsing analysis 393

394 In order to boost power to detect associations for rare variants (including private mutations)
395 having the same direction of effect, we adopted a collapsing framework to test the aggregate
396 effect of rare functional variants in a gene. Overall, 10 different collapsing models (9
397 dominant and one recessive) were implemented per gene to evaluate a range of genetic
398 architectures. Additionally, a synonymous collapsing model was used for the purpose of
399 establishing an empirical negative control²⁰.

400 As outlined in **Table S11**, the criteria for qualifying variants (QVs)³⁶ for the collapsing
401 models were based on the following parameters: type of variant (missense, non-
402 synonymous or PTV), minor allele frequency, *in silico* deleteriousness predictors (REVEL
403 and MTR), and type of genetic model (dominant or recessive). The following variant
404 annotations were used to define PTVs: *exon_loss_variant*, *frameshift_variant*, *start_lost*,
405 *stop_gained*, *stop_lost*, *splice_acceptor_variant*, *splice_donor_variant*, *gene_fusion*,
406 *bidirectional_gene_fusion*, *rare_amino_acid_variant* and *transcript_ablation*. Hemizygous
407 genotypes for the X chromosome also qualified for the recessive model.

408 For a given collapsing model, the effect of QVs in each gene (N=18,762) was
409 calculated as the difference in the mean of a blood biomarker between carriers and non-

410 carriers of the QVs, using a linear regression model in PEACOCK. Covariates used in the
411 linear regression model were identical to that described for the ExWAS.

412 A significance cut-off of $p < 1 \times 10^{-8}$ was set for the collapsing analysis based on the
413 observed p-value distribution for the synonymous model and an n-of-1 permutation, as
414 described previously²⁰.

415 **Association analysis of clinical phenotypes documented in the UKB**

417 We harmonized and union mapped the clinical phenotypes available in the UKB, as previously
418 described²⁰. Phenome-wide collapsing analysis for 15,719 clinical phenotypes was performed
419 for the 11 collapsing models, as described in our previously published study²⁰. We queried the
420 results of this analysis for genes of interest that emerged from the analysis of the metabolic
421 biomarkers.

422 Additionally, we also performed an association analysis between the each of the 198 metabolic
423 biomarkers and the clinical phenotypes using a linear regression model adjusted for age and
424 sex.

425 **Comparison of results from collapsing analyses to microarray-based genome-wide** 426 **association study**

427 We explored the hypothesis that the application of a collapsing framework – that tests the
428 aggregate effect of rare functional variants in a gene identified using exome sequencing –
429 detected gene-biomarker relationships that were previously not identified in microarray-
430 based studies. In order to do that, we compared our findings with the results from a recent
431 study¹⁴ that conducted single variant association analysis (GWAS) for clinical biomarkers in
432 the UKB using microarray data, including directly genotyped coding variants. Besides the
433 28/30 clinical blood biomarkers that we studied, seven other biomarkers (mainly, urine-
434 related) were analysed in the GWAS. These seven biomarkers comprised of four urinary
435 biomarkers that were directly measured in the UKB and an additional three derived
436 measurements. For the purpose of comparing findings, we additionally performed gene-level
437 collapsing analysis for the four urinary biomarkers for which data were directly available in
438

439 the UKB (i.e. ‘sodium in urine’, ‘potassium in urine’, ‘microalbumin in urine’, and ‘creatinine
440 (enzymatic) in urine’). To be consistent with the microarray-based GWAS, we used the
441 statin-adjusted values for LDL-cholesterol, total cholesterol, and apolipoprotein B, and the
442 medication-unadjusted values for the remaining biomarkers. Thereafter, for the set of 32
443 biomarkers (28 blood and 4 urinary biomarkers) common to both studies, we compared
444 gene-biomarker relationships that achieved significance ($p < 1 \times 10^{-8}$) in the collapsing analysis
445 with gene-biomarker relationships corresponding to the significant coding variant
446 associations reported in the GWAS. We considered a comparatively relaxed significance
447 threshold of $p = 1 \times 10^{-7}$ for the GWAS results in order to be stringent when attributing a gene-
448 biomarker relationship as being specific to the collapsing analysis.

449 We also hypothesised that the various variant-level “purifying” filters implemented for
450 QV selection in the collapsing analysis can enable a more direct estimate for the effect of
451 gene aberrations (e.g., PTVs) on biomarker levels. To investigate this hypothesis, we
452 compared the effect sizes for gene-biomarker relationships that achieved significance in both
453 the gene-level collapsing analysis and the microarray-based GWAS. For each such gene-
454 biomarker relationship, we selected: (i) the *model* with the highest absolute beta in the
455 collapsing analysis, and (ii) the individual *variant* with the highest absolute beta as reported
456 in the Sinnott-Armstrong *et al* GWAS¹⁴. For the latter, we adopted the absolute beta
457 estimated in the genotypic model in our ExWAS (for the corresponding gene-biomarker
458 relationship) as a substitute, to account for possible differences in trait transformation,
459 association model or covariates between our study and the Sinnott-Armstrong *et al* GWAS.
460 Nonetheless, the absolute betas were highly correlated between the Sinnott-Armstrong *et al*
461 GWAS and our ExWAS (Spearman’s $\rho = 0.99$) (**Figure S5**). We then compared the
462 absolute beta of the collapsing model [step (i)] with that of the individual variant [step (ii)].
463 This approach provides a means to compare the effect size of aberrations in genes on
464 biomarker levels estimated from individual coding variants captured by microarrays with that
465 estimated from an aggregate of rare coding variants identified using exome sequencing.
466

467 **Estimating gene similarity based on association signatures from collapsing analysis**

468 We aimed to leverage the rich catalogue of gene-level association statistics from the
469 collapsing analysis – ascertained for the set of studied metabolic biomarkers and under
470 different QV models – to identify genes that possess similar metabolic biomarker fingerprint
471 as a (target) gene of interest. Such a ‘gene similarity’ metric can provide opportunities to not
472 only expand our understanding of gene networks, but also offer alternative candidates in
473 cases of difficult-to-drug targets. Gene-SCOUT (Gene Similarity from Continuous Traits)²⁴,
474 the tool that we developed for this purpose, can also estimate “similarity” between genes
475 based on any set of quantitative traits of interest.

476 Rather than calculating similarities between genes directly, Gene-SCOUT estimates
477 distances between genes, which it then uses as a proxy for their similarity. Based on that,
478 the set of genes having the smallest distance from a given seed gene represent those that
479 are most ‘similar’ to it. We applied the cosine distance method – which is commonly used in
480 natural language processing³⁷ – to calculate distances between genes³⁸ based on their
481 effects on the metabolic biomarkers (referred to as ‘features’) estimated in the collapsing
482 analysis. In order to minimise the impact of stochastic effects on the gene similarity
483 estimations, for a given seed gene of interest, only those features that the genes is
484 associated with at $p < 1 \times 10^{-5}$ are selected (‘feature selection’ step), guided by sensitivity
485 analyses performed for a range of p-value thresholds²⁴. Thus, distances from genes having
486 $p > 1 \times 10^{-5}$ for all features in common with the seed gene are not considered.

487 The feature set used to generate the Gene-SCOUT results comprised of the 198
488 metabolic blood biomarkers. Though there is a degree of correlation in our feature set
489 (**Figure 2a**), we have demonstrated through simulations that correlation between features
490 has minimal impact on gene similarity estimations²⁴.
491 To illustrate the tool’s utility, we generated a network figure showing the genes that were
492 most similar to each of the 24 genes that were significantly associated with LDL-cholesterol
493 in the collapsing analysis.

494 **Ethics Reporting**

495 The protocols for UKB are overseen by The UK Biobank Ethics Advisory Committee (EAC);
496 for more information see: <https://www.ukbiobank.ac.uk/ethics/> and
497 <https://www.ukbiobank.ac.uk/wp-content/uploads/2011/05/EGF20082.pdf>.

498
499 **Acknowledgements**

500 We thank the participants and investigators in the UKB study who made this work possible
501 (Resource Application Number 26041); the UKB Exome Sequencing Consortium (UKB-ESC)
502 members AbbVie, Alnylam Pharmaceuticals, AstraZeneca, Biogen, Bristol-Myers Squibb,
503 Pfizer, Regeneron and Takeda for funding the generation of the exome sequence data; the
504 Regeneron Genetics Center for completing the sequencing and initial quality control of the
505 exome sequencing data; and the AstraZeneca Centre for Genomics Research Analytics and
506 Informatics team for processing and analysis of sequencing data.

507
508 **Author Contributions**

509 S.P. designed the study. A.N., L.M., R.S.D., D.V., E.W., Q.W. and S.P. performed the
510 analyses and statistical interpretation. A.N., R.S.D., A.R.H. and S.P. drafted the manuscript.
511 All authors contributed to the review and critical revision of the manuscript.

512
513 **Competing interests**

514 A.N., L.M., R.S.D., D.V., E.W., E.L.A., A.R., K.C., K.R.S., Q.W., B.C., D.S.P., A.R.H. and
515 S.P. are current employees and/or stockholders of AstraZeneca.

References

- 516 1. Willer, C. J. *et al.* Discovery and refinement of loci associated with lipid levels. *Nature*
517 *genetics* **45**, (2013).
- 518 2. Wuttke, M. *et al.* A catalog of genetic loci associated with kidney function from
519 analyses of a million individuals. *Nature genetics* **51**, (2019).
- 520 3. Kettunen, J. *et al.* Genome-wide association study identifies multiple loci influencing
521 human serum metabolite levels. *Nature genetics* **44**, (2012).
- 522 4. Yet, I. *et al.* Genetic Influences on Metabolite Levels: A Comparison across
523 Metabolomic Platforms. *PloS one* **11**, (2016).
- 524 5. Suhre, K. *et al.* A genome-wide association study of metabolic traits in human urine.
525 *Nature genetics* **43**, (2011).
- 526 6. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nature*
527 *genetics* **46**, (2014).
- 528 7. Klarin, D. *et al.* Genetics of blood lipids among ~300,000 multi-ethnic participants of
529 the Million Veteran Program. *Nature genetics* **50**, (2018).
- 530 8. Chambers, J. C. *et al.* Genome-wide association study identifies loci influencing
531 concentrations of liver enzymes in plasma. *Nature genetics* **43**, (2011).
- 532 9. Prins, B. P. *et al.* Genome-wide analysis of health-related biomarkers in the UK
533 Household Longitudinal Study reveals novel associations. *Scientific reports* **7**, (2017).
- 534 10. Wheeler, E. *et al.* Impact of common genetic determinants of Hemoglobin A1c on type
535 2 diabetes risk and diagnosis in ancestrally diverse populations: A transethnic
536 genome-wide meta-analysis. *PLoS medicine* **14**, (2017).
- 537 11. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants
538 associated with human blood metabolites. *Nature genetics* **49**, (2017).
- 539 12. Gallagher, M. D. & Chen-Plotkin, A. S. The Post-GWAS Era: From Association to
540 Function. *American journal of human genetics* **102**, (2018).
- 541 13. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
542 *Nature* **562**, 203–209 (2018).
- 543 14. Sinnott-Armstrong, N. *et al.* Genetics of 35 blood and urine biomarkers in the UK
544 Biobank. *Nature genetics* **53**, (2021).
- 545 15. UK10K Consortium *et al.* The UK10K project identifies rare variants in health and
546 disease. *Nature* **526**, (2015).
- 547 16. MacArthur, D. G. *et al.* A systematic survey of loss-of-function variants in human
548 protein-coding genes. *Science (New York, N.Y.)* **335**, (2012).
- 549 17. Ritchie, S. C. *et al.* Quality control and removal of technical variation of NMR
550 metabolic biomarker data in ~120,000 UK Biobank participants. *medRxiv*
551 2021.09.24.21264079 (2021) doi:10.1101/2021.09.24.21264079.
- 552 18. Gandotra, S. *et al.* Perilipin deficiency and autosomal dominant partial lipodystrophy.
553 *The New England journal of medicine* **364**, (2011).
- 554 19. Shin, G.-C., Kang, H. S., Lee, A. R. & Kim, K.-H. Hepatitis B virus-triggered
555 autophagy targets TNFRSF10B/death receptor 5 for degradation to limit
556 TNFSF10/TRAIL response. *Autophagy* **12**, (2016).
- 557 20. Wang, Q. *et al.* Rare variant contribution to human disease in 281,104 UK Biobank
558 exomes. *Nature* (2021) doi:10.1038/s41586-021-03855-y.
- 559 21. Aimee M. Deaton *et al.* Gene-level analysis of rare variants in 379,066 whole exome
560 sequences identifies an association of GIGYF1 loss of function with type 2 diabetes.
561 *Scientific Reports* **11**, (2021).
- 562 22. Jurgens, S. J. *et al.* Rare Genetic Variation Underlying Human Diseases and Traits:
563 Results from 200,000 Individuals in the UK Biobank. *bioRxiv* 2020.11.29.402495
564 (2020) doi:10.1101/2020.11.29.402495.
- 565 23. Chen, J. *et al.* The trans-ancestral genomic architecture of glycemic traits. *Nature*
566 *genetics* **53**, (2021).

- 567 24. Lawrence Middleton *et al.* Gene-SCOUT: identifying genes with similar continuous
568 trait fingerprints from phenome-wide association analyses. *Nucleic Acids Res (in*
569 *submission)* (2021).
- 570 25. Cohen, J. C., Boerwinkle, E., Mosley, T. H. & Hobbs, H. H. Sequence variations in
571 PCSK9, low LDL, and protection against coronary heart disease. *The New England*
572 *journal of medicine* **354**, (2006).
- 573 26. Abul-Husn, N. S. *et al.* A Protein-Truncating HSD17B13 Variant and Protection from
574 Chronic Liver Disease. *The New England journal of medicine* **378**, (2018).
- 575 27. Akbari, P. *et al.* Sequencing of 640,000 exomes identifies GPR75 variants associated
576 with protection from obesity. *Science (New York, N. Y.)* **373**, (2021).
- 577 28. Nag, A. *et al.* Human genetic evidence supports MAP3K15 inhibition as a therapeutic
578 strategy for diabetes. *medRxiv* 2021.11.14.21266328 (2021)
579 doi:10.1101/2021.11.14.21266328.
- 580 29. Szustakowski, J. D. *et al.* Advancing human genetics research and drug discovery
581 through exome sequencing of the UK Biobank. *Nature genetics* **53**, (2021).
- 582 30. Cingolani, P. *et al.* A program for annotating and predicting the effects of single
583 nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster*
584 strain w1118; iso-2; iso-3. *Fly* **6**, 80–92 (2012).
- 585 31. Traynelis, J. *et al.* Optimizing genomic medicine in epilepsy through a gene-
586 customized approach to missense variant interpretation. *Genome research* **27**, 1715–
587 1729 (2017).
- 588 32. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity
589 of Rare Missense Variants. *American journal of human genetics* **99**, 877–885 (2016).
- 590 33. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association
591 studies. *Bioinformatics* **26**, (2010).
- 592 34. Millard, L. A. C., Davies, N. M., Gaunt, T. R., Davey Smith, G. & Tilling, K. Software
593 Application Profile: PHEASANT: a tool for performing automated phenome scans in UK
594 Biobank. *International journal of epidemiology* **47**, (2018).
- 595 35. Fadista, J., Manning, A. K., Florez, J. C. & Groop, L. The (in)famous GWAS P-value
596 threshold revisited and updated for low-frequency variants. *European journal of*
597 *human genetics : EJHG* **24**, (2016).
- 598 36. Petrovski, S. *et al.* An Exome Sequencing Study to Assess the Role of Rare Genetic
599 Variation in Pulmonary Fibrosis. *American journal of respiratory and critical care*
600 *medicine* **196**, (2017).
- 601 37. Huang A. Similarity Measures for Text Document Clustering. *NZCSRSC* (2008).
- 602 38. Kittipong Chomboon, Pasapitch Chujai, Pongsakorn Teerarassamee, Kittisak
603 Kerdprasop & Nittaya Kerdprasop. An Empirical Study of Distance Metrics for k-
604 Nearest Neighbor Algorithm. *Proceedings of the 3rd International Conference on*
605 *Industrial Application Engineering* (2015).
- 606