

1 **A cross-platform approach identifies genetic regulators of human metabolism and health**

2
3 Luca A. Lotta^{1#}, Maik Pietzner^{1#}, Isobel D. Stewart¹, Laura B.L. Wittemans^{1,2}, Chen Li¹, Roberto
4 Bonelli^{3,4}, Johannes Raffler⁵, Emma K. Biggs⁶, Clare Oliver-Williams^{7,8}, Victoria P.W. Auyeung¹, Jian'an
5 Luan¹, Eleanor Wheeler¹, Ellie Paige⁹, Praveen Surendran^{7,10,11,12}, Gregory A. Michelotti¹³, Robert A.
6 Scott¹, Stephen Burgess^{14,15}, Verena Zuber^{14,16}, Eleanor Sanderson¹⁷, Albert Koulman^{1,5,18}, Fumiaki
7 Imamura¹, Nita G. Forouhi¹, Kay-Tee Khaw¹⁵, MacTel Consortium, Julian L. Griffin¹⁹, Angela M.
8 Wood^{7,10,11,20,21}, Gabi Kastenmüller⁵, John Danesh^{7,10,11,20,22,23}, Adam S. Butterworth^{7,10,11,20,22,23}, Fiona
9 M. Gribble⁶, Frank Reimann⁶, Melanie Bahlo^{3,4}, Eric Fauman²⁴, Nicholas J. Wareham¹, Claudia
10 Langenberg^{1,11*}

- 11
12
13
14 1) MRC Epidemiology Unit, University of Cambridge, Cambridge, UK
15 2) The Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford
16 3) Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical Research,
17 Parkville, Australia
18 4) Department of Medical Biology, The University of Melbourne, Parkville, Australia
19 5) Institute of Computational Biology, Helmholtz Zentrum München – German Research Center for
20 Environmental Health, Neuherberg, Germany
21 6) Metabolic Research Laboratories, University of Cambridge, Cambridge, United Kingdom
22 7) British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and Primary
23 Care, University of Cambridge, Cambridge, UK
24 8) Homerton College, University of Cambridge, Cambridge, UK
25 9) National Centre for Epidemiology and Population Health, The Australian National University, Canberra,
26 Australia
27 10) British Heart Foundation Centre of Research Excellence, University of Cambridge, Cambridge, UK
28 11) Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge,
29 Cambridge, UK
30 12) Rutherford Fund Fellow, Department of Public Health and Primary Care, University of Cambridge, UK
31 13) Metabolon Inc, Durham, North Carolina USA
32 14) MRC Biostatistics Unit, University of Cambridge, Cambridge, United Kingdom
33 15) Department of Public Health and Primary Care, University of Cambridge, Cambridge, United Kingdom
34 16) Department of Epidemiology and Biostatistics, Imperial College London, UK
35 17) MRC Integrative Epidemiology Unit, Bristol Medical School, University of Bristol, UK
36 18) NIHR BRC Nutritional Biomarker Laboratory, University of Cambridge, UK
37 19) Biomolecular Medicine, Department of Metabolism, Digestion and Reproduction, Imperial College
38 London, UK
39 20) National Institute for Health Research Blood and Transplant Research Unit in Donor Health and Genomics,
40 University of Cambridge, Cambridge, UK
41 21) The Alan Turing Institute, London, UK
42 22) National Institute for Health Research Cambridge Biomedical Research Centre, University of Cambridge
43 and Cambridge University Hospitals, Cambridge, UK
44 23) Department of Human Genetics, Wellcome Sanger Institute, Hinxton, UK
45 24) Internal Medicine Research Unit, Pfizer Worldwide Research, Cambridge, MA 02142, USA

46
47 *# these authors contributed equally*

48
49 *Corresponding author:
50 Claudia Langenberg
51 MRC Epidemiology Unit
52 University of Cambridge School of Clinical Medicine
53 Institute of Metabolic Science
54 Cambridge, UK
55 claudia.langenberg@mrc-epid.cam.ac.uk

56 **Abstract**

57 In cross-platform analyses of 174 metabolites we identify 499 associations ($p < 4.9 \times 10^{-10}$)
58 characterized by pleiotropy, allelic heterogeneity, large and non-linear effects, and enrichment for
59 nonsynonymous variation. We identify a signal at *GLP2R* (p.Asp470Asn) shared between higher
60 citrulline levels, body mass index, fasting glucose-dependent insulinotropic peptide and type 2
61 diabetes, with beta-arrestin signalling as the underlying mechanism. Genetically-higher serine levels
62 are shown to reduce the likelihood (by 95%) and predict development of macular telangiectasia type
63 2, a rare degenerative retinal disease. Integration of genomic and small molecule data across
64 platforms enables discovery of regulators of human metabolism and translation into clinical insights.
65

66 Introduction

67 Metabolites are small molecules that reflect biological processes and are widely measured in
68 clinical medicine as diagnostic, prognostic or treatment response biomarkers¹. Blood levels of
69 metabolites are highly heritable with twin studies reporting a median explained variance in plasma
70 levels of 6.9% and a maximum of 50% depending on the metabolite^{2,3}. Several earlier studies have
71 started to characterise the genetic architecture of metabolite variation in the general population²⁻¹⁰,
72 but been limited in size and scope by focussing on metabolites assessed using a single method.
73 Integration of genetic association results for metabolites measured on different platforms can help
74 maximise the power for a given metabolite and provide a more refined understanding of genetic
75 influences on blood metabolite levels and human physiology.

76 To identify genomic regions regulating metabolite levels and systematically study their
77 relevance for disease, we conducted a cross-platform meta-analysis of genetic effects on levels of
78 174 blood metabolites measured in large-scale population-based studies. We included metabolites
79 covered by the targeted Biocrates AbsoluteIDQ™ p180 platform and measured in the Fenland Study.
80 We integrated unpublished data for any of these metabolites that were covered by the Nightingale
81 (¹H-NMR, INTERVAL Study) or Metabolon (Discovery HD4™, EPIC-Norfolk and Interval Studies)
82 platforms, or had previously been reported^{2,4,5}. The focus on this targeted set of 'platform-specific'
83 metabolites enabled us to clearly map metabolites across platforms and maximise the sample size
84 for each of the 174 metabolites for this proof of concept cross-platform genome-wide association
85 study (mGWAS). To facilitate rapid sharing of our results, we developed a webserver
86 (<https://omicscience.org/apps/crossplatform/>) that allows flexible interrogation of our results.

87 Results

88 *Associations with blood metabolites at 144 genomic regions*

89 Genome-wide meta-analyses were conducted for 174 metabolites from 7 biochemical classes
90 (i.e. amino acids, biogenic amines, acylcarnitines, lyso-phosphatidylcholines, phosphatidylcholines,
91 sphingomyelins and hexose) commonly measured using the Biocrates p180 kit in up to 86,507
92 individuals, contributing over 3.7 million individual-metabolite data points (70% from unpublished
93 studies; **Fig. 1**). For each of the 174 metabolites, this was the largest GWAS to date, with at least a
94 doubling of sample size (**Fig. 1A**). Sample sizes ranged from 8,569 to 86,507 individuals for
95 metabolites depending on the platform used in each contributing study. Using GWAS analyses we
96 estimated the association of up to 10.2 million single nucleotide variants with a minor allele
97 frequency (MAF) >0.5%, including 6.1 million with MAF ≥ 5%.

98 We identified 499 variant-metabolite associations (362 unreported) from 144 loci (94
99 unreported) at a metabolome-adjusted genome-wide significance threshold of $p < 4.9 \times 10^{-10}$
100 (correcting the usual GWAS-threshold, $p < 5 \times 10^{-8}$, for 102 principal components explaining 95% of the
101 variance in metabolite levels using principal component analysis; **Fig. 1**). The vast majority of these
102 associations were consistent across studies and measurement platforms [median I^2 : 26.8
103 (interquartile range: 0 – 70.1) for 465 associations with at least two contributing studies]
104 (**Supplementary Tab. S1-2**). To identify possible sources of heterogeneity, we investigated the
105 influence of differences by cohort, measurement platform, metabolite class, and association
106 strength in a joint meta-regression model (**Supplementary Tab. S3**). This showed that heterogeneity
107 was mainly due to the overall strength of the signal, i.e. associations with higher z-scores showed
108 greater heterogeneity ($p < 1.05 \times 10^{-9}$). However, the majority of these statistically heterogeneous
109 associations were directionally consistent and nominally significant across and within each stratum
110 for 146 of 170 associations with a z-score > 10 , demonstrating the feasibility of pooling association
111 estimates across metabolomics platforms for the purpose of genetic discovery. Genetic variants at
112 the *NLRP12* locus, e.g. rs4632248, were a notable exception with large estimates of heterogeneity
113 ($I^2 > 90\%$). The *NLRP12* locus is known to affect the monocyte count¹¹ and has been shown to have
114 pleiotropic effects on the plasma proteome in the INTERVAL study¹². Monocytes, or at least a
115 subpopulation subsumed under this cell count measure, release a wide variety of biomolecules upon
116 activation or may die during the sample handling process and hence releasing intracellular
117 biomolecules, such as taurine¹³, into the plasma. In brief, one specific source of heterogeneity in
118 mGWAS associations might relate to sample handling differences across studies.

119 This highlights the utility of our genetic cross-platform approach to maximise power for a given
120 metabolite, substantially extending previous efforts for any given metabolite¹⁴. Previously reported
121 associations from platform-specific studies were also found to generally be consistent in our cross-
122 platform meta-analysis (**Supplementary Tab. S2**; <https://omicscience.org/apps/crossplatform/>).

123 *Insights in the genetic architecture of metabolite levels*

124 We identified a median of 2 (range: 1-67, **Fig. 2A**) associated metabolites for each locus and a
125 median of 3 (range: 1-20, **Fig. 2B**) locus associations for each metabolite, reflecting pleiotropy and
126 the extensive contribution of genetic loci to circulating metabolite levels. The number of associations
127 was proportional to the estimated heritability and the sample size of the meta-analysis for a given
128 metabolite (**Fig. 2C**).

129 We applied a multi-trait statistical colocalisation method¹⁵ and identified between 1-30
130 (median: 2) metabolites that did not meet the discovery p-value threshold, but showed high

131 posterior probability (>75%) of a shared genetic signal for 49 out of the 144 loci (**Supplemental Fig.**
132 **S1**). Two distinct variants (rs2414577 and rs261334) nearby *LIPC* showed the largest gain in
133 additionally associated metabolites, in line with previous reports of extensive pleiotropy and allelic
134 heterogeneity at this locus⁹. We note that a low posterior probability for the alignment of multiple
135 metabolites at other loci might be explained by the presence of multiple causal variants shared
136 across multiple metabolites.

137 To systematically classify pleiotropic variants taking into account the correlation structure
138 among metabolites we derived a data-driven metabolic network and performed community
139 detection (see **Methods** and **Supplemental Fig. S2**). A total of 129 (60.5%) of 214 variants
140 (associated with at least two metabolites at $p < 5 \times 10^{-8}$) were associated with metabolites from at
141 least two of the 14 communities (range: 2 – 11; **Supplemental Fig. S2**), i.e. showed evidence for
142 ‘horizontal’ or broad pleiotropy. The most extreme variants included those near *FADS1* (e.g.
143 rs17455) associated with 61 metabolites across 11 communities at $p < 5 \times 10^{-8}$. In contrast, rs2638315
144 (likely tagging a missense variant rs2657879 at *GLS2*) was associated with nine metabolites within a
145 single community and would therefore be considered as ‘vertical pleiotropic’ for a well-defined
146 group of correlated metabolites (**Supplemental Fig. S2**).

147 Similar to what is routinely observed in GWAS literature, effect size estimates increased with
148 decreasing minor allele frequency (MAF) (**Fig. 3A**). However, there were 26 associations (**Tab. 1**) for
149 common lead variants with per-allele differences in metabolites levels greater than 0.25 standard
150 deviations (SD), a per-allele effect size that is >3-fold larger than the strongest common variants
151 associated with SDs of body mass index at the *FTO* locus.

152 Variants identified in this study explained up to 23% of the variance (median: 1.4%; interquartile
153 range: 0.5% - 2.8%) and up to 99.8% of the chip-based heritability (median 9.2%; interquartile range:
154 4.7% - 17.1%) for the 141 metabolites with at least one genetic association (**Fig. 2D**). The 26 common
155 variants with large effect sizes (>0.25 SD per allele) were identified for metabolites with higher
156 heritability (**Fig. 2D**) and accounted for up to 74% of the heritability explained in those metabolites.

157 GWAS analyses generally assume a linear relationship between genotypes and phenotypes, i.e.
158 an additive dose-response model. The identification of several metabolite-associated variants with
159 large effect sizes and availability of individual-level data in the Fenland cohort allowed us to test
160 whether the metabolite-associated variants showed evidence of deviation from a linear model. Of
161 499 associations tested, 9 showed evidence of departure from a linear association (**Fig. 2E-M**).
162 Modelling actual genotypes rather than assuming ‘additive’ linear associations in these instances
163 explained a median of 7.4% more (range: 1.4-15.2%) of the heritability in metabolite levels (**Fig. 2N**).

164 Associations better described by an autosomal recessive or dominant model of inheritance might be
165 the most likely explanation for this. Variant rs3916, for example, which showed a more than additive
166 positive effect on butyrylcarnitine, is in perfect LD with a missense variant within *ACADS* (rs1799958,
167 MAF=26%), which encodes for short-chain acyl-CoA dehydrogenase (SCAD). SCAD deficiency is an
168 autosomal recessive disease diagnosed by elevated butyrylcarnitine concentrations in blood and
169 homozygous carrier status for established pathogenic variants¹⁶.

170 In 61 of the 499 associations the lead association signal was a nonsynonymous variant, a 40-fold
171 enrichment compared to what would be expected by chance given the annotation of ascertained
172 genetic variants (two-tailed binomial test, $p=5\times 10^{-30}$, **Fig. 3D**). For a further 59 associations, the lead
173 variant was in high LD with a nonsynonymous variant ($r^2>0.8$). Lead variants that were
174 nonsynonymous, or variants in high LD with a nonsynonymous variant, generally had lower MAF,
175 larger effect sizes, and smaller 99%-credible sets (**Supplemental Tab. S4**) than variants that were not
176 in these categories (**Fig 3B-D**).

177 We identified 22 loci harbouring two ($n=21$) or three ($n=1$) independent signals, i.e. different
178 plasma metabolites were associated with distinct genetic variants within the same genomic region
179 (**Supplementary Tab. S2**). For six regions, our two different annotations approaches assigned only
180 one causal gene (see below and **Methods**), including *ACADM*, *GLDC*, *ARG1*, *MARCH8*, *SLC7A2*, and
181 *LIPC* (**Supplementary Tab. S2**). We found evidence that allelic heterogeneity, i.e. conditionally
182 independent variants at a locus for a specific metabolite, explains the association pattern at 3 of
183 those loci (*ACADM*, *ARG1*, and *LIPC*; **Supplementary Tab. S5**). We identified another 16 loci
184 harbouring at least one (range: 2–6) additional conditionally independent variant(s) in exact
185 conditional analyses (see **Methods**, **Supplementary Tab. S5**).

186 *Effector genes, tissues, pathways*

187 We used two complementary strategies to prioritize likely causal genes for the observed
188 associations: (1) a hypothesis-free genetic approach based on physical distance, genomic annotation
189 and integration of expression quantitative trait loci (eQTLs) to prioritize genes in a systematic and
190 standardised way (see **Methods**), and (2) a biological knowledge-based approach integrating existing
191 knowledge about specific metabolites or related pathways to identify biologically plausible
192 candidate genes from the 20 genes closest to the lead variant (**Fig. 4A**). Using the hypothesis-free
193 genetic approach, we identified 249 unique likely causal genes for the 499 associations, with at least
194 one gene per association and some genes prioritized as likely causal for multiple metabolite
195 associations. The knowledge-based approach identified 130 biologically plausible genes for 349 out
196 of 499 associations. We asked whether the hypothesis-free genetic approach identified biologically

197 plausible genes (prioritized by strategy 2) more often than expected by chance. Amongst 9,980
198 possible gene-metabolite pairs (20 genes x 499 associations), 420 (4.2%) were biologically plausible,
199 condensed to 350 gene(s)-metabolite assignments after accounting for overlapping annotations. Of
200 the latter, 126 pairs (36%) were identical to genetically-prioritized gene-metabolite pairs,
201 representing a significant enrichment of biologically plausible genes among those prioritised by the
202 hypothesis-free algorithm (~8-fold more than expected by chance; two-tailed binomial test,
203 $p=2.3\times 10^{-80}$; **Fig. 4B**). Among the consistently assigned genes between both approaches, assignment
204 of the nearest gene (124 times out of 126, X^2 -test, $p<2.5\times 10^{-45}$) was the strongest shared factor, as
205 might be expected, followed by being (or in LD with) a missense variant ($R^2>0.8$, 30 times out of 126,
206 X^2 -test, $p<1.3\times 10^{-07}$) and only a minor contribution of eQTL data (20 times out of 126, X^2 -test,
207 $p<0.001$). Over 70% of genetically prioritized genes were enzymes or transporters (**Fig. 4C**).
208 Inconsistencies between the approaches might be explained by non-consideration of information on
209 biological pathways in the hypothesis-free genetic approach, as well as variants acting more distal to
210 the biological determinants of plasma metabolite levels not being considered in the knowledge-
211 based approach. The missense variant rs1260326 within *GCKR*, for example, colocalised with 49
212 metabolites across diverse biochemical classes (**Supplemental Fig. S1**) and likely confers its effects
213 on glucose metabolism through impaired inhibition of glucokinase by glucokinase regulatory protein
214 and might hence be considered as putative causal candidate by the knowledge-driven approach for
215 plasma glucose only. However, impairments in glucose metabolism result in numerous downstream
216 consequences including more distal metabolic branches such as amino acid and lipid metabolism.

217 In addition to being enriched in genes previously implicated in the biology of these metabolites,
218 the genetically prioritized genes were also enriched in genes known for mutations to cause rare
219 inborn errors of metabolism (IEMs), i.e. monogenic defects in the metabolism of small molecules
220 with very specific metabolite changes (**Fig. 4B**).

221 Integrating GWAS statistics across cohorts and platforms allowed us to identify three genes that
222 have never been associated with any metabolite level so far. At the *CERS6* locus, rs4143279
223 associates with levels of sphingomyelin (d18:1/16:0) ($p = 4.2\times 10^{-10}$). *CERS6* encodes a ceramide
224 synthase facilitating formation of ceramide, a precursor of sphingomyelins¹⁷. At the *ASNS* locus,
225 rs17345286 associates with levels of asparagine ($p = 4.7\times 10^{-20}$). The lead variant is in high LD ($R^2=1$)
226 with a missense mutation in *ASNS* (rs1049674, p.Val210Glu). *ASNS* encodes an asparagine
227 synthase¹⁸. Finally, at the *SLC43A1* locus, rs2649667 associates with levels of phenylalanine ($p =$
228 3.6×10^{-13}). *SLC43A1* encodes a liver-enriched transporter of large neutral amino acids, including
229 phenylalanine¹⁹.

230 *A shared functional variant in GLP2R is linked to type 2 diabetes*

231 Because several of the metabolites captured in this GWAS have been associated with incident
232 type 2 diabetes (T2D), we sought to investigate whether the association between metabolite-
233 associated loci and diabetes could provide insights into underlying pathophysiologic mechanisms.
234 We observed a significant enrichment of T2D-associations among our metabolite variants (p-
235 value= 2.8×10^{-7} , **Fig. 5**) using a meta-analysis of 80,983 T2D cases and 842,909 controls (see
236 **Methods**).

237 Amongst the T2D- and metabolite-associated loci was a missense p.Asp470Asn (rs17681684)
238 variant in the *GLP2R* gene encoding the receptor for glucagon-like peptide 2, a 33 amino acid
239 peptide hormone encoded by the proglucagon gene (*GCG*) that stimulates the growth of intestinal
240 tissue. Common variants at *GLP2R* are associated with an increased risk of T2D²⁰. The previously
241 reported lead variant for T2D (rs78761021) is in high LD ($r^2 > 0.87$) with our lead citrulline association
242 signal at *GLP2R* (rs17681684), which was associated with a 4% higher T2D risk (per-allele odds ratio,
243 1.04; 95%-confidence interval, 1.02, 1.05; $p = 1.1 \times 10^{-8}$), comparable to previous reports²⁰.
244 Considering eleven phenotypes related to glucose homeostasis and metabolic health²¹⁻²³, the A-
245 allele of rs17681684 was significantly associated with insulin disposition index (beta=-0.067,
246 $p < 0.002$)²², corrected insulin response (beta=-0.061, $p < 0.004$)²², glycated haemoglobin 1c (HbA1c)
247 (beta=0.006, $p < 0.0003$)²¹, and body mass index (beta=0.010, $p < 5.3 \times 10^{-9}$), in addition to the
248 previously reported positive association with fasting glucose-dependent insulinotropic peptide (GIP)
249 and the suggestive inverse association with post-glucose load GLP-1 (beta=-0.035, $p < 4.6 \times 10^{-4}$)²⁴.
250 While sample sizes and hence significance levels for insulin traits were not sufficient to support
251 formal colocalisation analysis, we still obtained a high posterior probability (PP>75%) for a shared
252 genetic signal across plasma citrulline, T2D risk, body mass index, and fasting levels of GIP (**Fig. 5B**).
253 The *GLP2R* p.Asp470Asn variant was the only of 6 independent genome-wide significant citrulline-
254 raising loci that was associated with a higher risk of T2D, which indicates that the association does
255 not reflect a general effect of blood citrulline levels on T2D risk but rather a locus-specific association
256 at *GLP2R* (**Fig. 5C**). Plasma citrulline levels have been shown to reflect the volume of intestinal cells
257 and are a marker of GLP2R target engagement in the treatment of short-bowel syndrome with
258 glucagon-like peptide 2 analogues²⁵. Taken together, this suggests that genetically higher GLP2R-
259 signalling, indicated by the higher citrulline levels among *GLP2R* 470Asn carriers, may lead to
260 chronically elevated GIP (though increased enteroendocrine mass and number of GIP-secreting K-
261 cells), which has been shown to downregulate GIP receptors on pancreatic beta cells²⁶, thereby
262 contributing to the observed reduction in the insulin secretory response and increase in T2D risk.

263 G-protein coupled receptors like GLP2R may signal via G-protein-dependent cyclic adenosine
264 monophosphate (cAMP) production or via G-protein-independent beta-arrestin mediated

265 signalling²⁷. To investigate if the *GLP2R* p.Asp470Asn variant affects signalling via either of these
266 pathways, we expressed the *GLP2R* p.Asp470Asn variant in different *in vitro* models (see **Methods**).
267 We show that the variant allele is significantly associated with reduced recruitment of beta-arrestin
268 to GLP2R upon glucagon-like peptide 2 stimulation, but not with cAMP signalling, which suggests a
269 potential role for impaired beta-arrestin recruitment to GLP2R in the pathophysiology of T2D (**Fig.**
270 **5E-G**).

271 *Serine levels are causally related to a rare eye disease*

272 A recent GWAS of macular telangiectasia type 2 (MacTel), a rare neurovascular degenerative
273 retinal disease, identified three genome-wide susceptibility loci (*PHGDH*, *CPS1*, and *TMEM161B*–
274 *LINC00461*) of which the same variants at *PHGDH* and *CPS1* were associated with levels of the amino
275 acids serine and glycine in this GWAS²⁸. More recently, it was shown that low serine availability is
276 linked to both MacTel as well as hereditary sensory and autonomic neuropathy type 1 through
277 elevated levels of atypical deoxyshingolipids²⁹. Whether genetic predisposition to low serine and
278 glycine levels affects MacTel more generally or has predictive utility has not been investigated. To
279 test this and to explore the specificity of associations between genetic influences on metabolite
280 levels and the risk of MacTel, we generated genetic scores using the sentinel variants for each of the
281 141 metabolites with at least one significantly associated locus identified in this GWAS and tested
282 their associations with the risk of MacTel. Genetic scores for serine and glycine were the only scores
283 associated with risk for MacTel after removal of the known highly pleiotropic *GCKR* variant (**Fig. 6A**).
284 A one SD increase in the genetic score for serine was associated with a 95% lower risk of MacTel
285 (odds ratio (95%-confidence interval), 0.05 (0.03-0.08); $p=9.5\times 10^{-30}$; **Fig. 6A**). Each of five serine
286 associated variants was individually associated with lower MacTel risk, with a clear dose-response
287 relationship and no evidence of heterogeneity (**Fig. 6B**). The association was unchanged when
288 removing the *GCKR* locus. To disentangle the effect of these two highly correlated metabolites on
289 MacTel risk, we used multivariable Mendelian randomization analysis, which allowed us to test for a
290 causal effect of both measures simultaneously. In this analysis, the effect of serine (odds ratio: 0.10,
291 $p<2.9\times 10^{-9}$) remained strong, while the effect of glycine (odds ratio: 0.50, $p<0.01$) was attenuated
292 (**Extended Data 2**). These results provide genetic evidence that the link between glycine and MacTel
293 is *via* serine through glycine conversion. This hypothesis is supported by the evidence of a log-linear
294 relationship between associations with serine and risk of MacTel among glycine-associated variants
295 (**Fig. 6B**). These findings provide strong evidence that pathways indexed by genetically higher serine
296 levels are strongly and causally associated with protection against MacTel.

297 Given the large observed effect size, we estimated whether using serine and glycine-associated
298 loci might improve the prediction of this rare disease. Adding genetically predicted glycine and

299 serine levels, based on newly discovered metabolite instruments from the present study and
300 previous MacTel variants linked to glycine and serine metabolism, substantially improved prediction
301 of MacTel based on an area under the receiver operating characteristic curve from 0.65 (CI 95%:
302 0.626-0.682) to 0.73 (0.702-0.753) (**Fig. 6**).

303 *From rare to common, the role of inborn errors of metabolism*

304 In his seminal 1902 work on alkaptonuria³⁰, also known as dark or black urine disease, Archibald
305 Garrod was the first to hypothesise that inborn errors of metabolism (IEMs) are “extreme examples
306 of variations of chemical behaviour which are probably everywhere present in minor degrees”.
307 Previous studies have shown enrichment of metabolite quantitative trait loci (mQTLs) in genes
308 known to cause IEMs³¹. Whether or not common variants at IEM-causing loci translate into clinically
309 manifest disease remains unknown. The identification of several metabolite-associated variants at
310 IEM-linked genes in this GWAS meta-analysis allows an investigation of the health consequences of
311 genetically determined differences in metabolism for more frequently occurring variants,
312 representing potentially milder forms of the metabolic and other clinical symptoms of IEMs, and
313 providing new candidate genes for rare extreme metabolic disorders that currently lack a genetic
314 basis (**Fig. 7A**). We identified 153 locus-metabolite associations for which 53 unique IEM-associated
315 genes were prioritized as likely causal using either the hypothesis-free genetic approach or the
316 knowledge-based approach on the basis of the Orphanet database³². In 89% of these associations
317 (136 of 153) the metabolite associated with a given GWAS locus perfectly matched, or was closely
318 related to, the metabolite affected in patients with the corresponding IEM (**Fig. 7B**).

319 To test whether IEM-mirroring lead variants from our metabolite GWAS may increase the risk of
320 common manifestations of diseases seen in patients with the corresponding IEM (**Fig. 7A**) we
321 obtained a list of electronic health record diagnosis codes (International Statistical Classification of
322 Diseases and Related Health Problems 10th Revision [ICD-10]) and mapped those based on
323 symptoms seen in both, IEM patients and patients with common, complex disease manifestations
324 (see **Methods**). We identified 93 ICD-10 codes with at least 500 cases within the UK Biobank study
325 that aligned with the symptoms or presentations seen in patients with IEMs mapping to mQTLs in
326 the present study. We obtained the association statistics of 85 unique metabolite-associated lead
327 variants at the 136 locus-metabolite associations with these 93 clinical diagnoses and observed 36
328 associations that met statistical significance (false discovery rate < 5%, **Supplemental Table S6 and**
329 **Fig. 7B**). For 15 out of those we obtained strong evidence of a shared genetic metabolite-phenotype
330 signal using colocalisation analyses (posterior probability of a shared signal >80%; **Fig. 7D and**
331 **Supplemental Fig. S3**). These instances linked common genetic variants in or near *APOE*, *PCSK9*, *LPL*,
332 and *LDLR* associated with sphingomyelins (SM 16:0, SM 18:0, and SM-OH 24:1) with atherosclerotic

333 heart disease diagnosis codes (I21, I25), mirroring what is observed in rare familial forms of
334 dyslipidaemia in which these sphingomyelins are elevated and the risk of ischemic heart disease is
335 greatly increased^{33,34}. These results provide further evidence that common variation at IEM genes
336 can lead to clinical phenotypes and diseases that correspond to those that patients with rare
337 mutations in those same genes are severely affected by. Further studies with detailed follow-up for
338 specific outcomes may provide greater power and help clarify the medical consequences of genetic
339 differences in metabolism caused by metabolite altering variants in the general population.

340 **Discussion**

341 This large-scale genome-wide meta-analysis has integrated genetic associations for 174
342 metabolites across different measurement platforms, an approach that has resulted in a three-fold
343 increase in our knowledge of genetic loci regulating levels of these metabolites. We assign likely
344 causal genes for many of the identified associations using a dual approach that combined automated
345 database mining with manual curation.

346 Previous platform-specific genetic studies of blood metabolites have been substantially smaller
347 in size due to being restricted to a single platform and/ or study²⁻¹⁰. We build on these earlier studies
348 to identify and demonstrate enrichment of rare and low-frequency coding variants in enzyme and
349 transporter genes with large effects and reveal the importance of non-linear associations at several
350 loci.

351 Our results not only provide detailed insight into the genetic determinants of human
352 metabolism but consider their relevance for disease aetiology and prediction. We explore both
353 locus-specific and polygenic score effects and provide tangible examples with clear translational
354 potential. We discovered a strong link between GLP2R, citrulline metabolism and T2D, and
355 demonstrate that the p.Asp470Asn variant underlying the citrulline and T2D associations leads to
356 significantly reduced recruitment of beta-arrestin to GLP2R in various cellular models, providing an
357 explanation for a possible pathological mechanism of a variant previously predicted to be benign²⁴.

358 The finding that a standard deviation increase in serine levels via a genetic score is associated
359 with 95% lower risk of MacTel shows that genetic differences resulting in very specific metabolic
360 consequences can have profound effects on health. Our results suggest that inclusion of genetic
361 scores for metabolite levels can improve identification of high risk individuals. Serine and glycine
362 supplementation and/ or pharmacologic modulation of serine metabolism may help to reduce
363 development or alter the prognosis of this rare, severe eye disease, specifically if targeted to people
364 genetically with a genetic susceptibility to low serine levels. It is important to note, that randomized

365 control trials are needed testing this hypothesis before any recommendations on supplementations
366 could be made.

367 We finally show specific examples where common genetic variation in IEM-related genes is
368 associated with phenotypes that are also caused by rare highly penetrant mutations. These results
369 suggest that rare variants in metabolite regulating genes newly identified in our study may be
370 valuable candidate genes in patients without a genetic diagnosis but severe alterations in the
371 corresponding or related metabolites. Hence these results provide a new starting point for further
372 investigations into the relationships between human metabolism and common and rare disorders.

373

374 **Acknowledgement/Funding**

375 M.P. was supported by a fellowship from the German Research Foundation (DFG PI 1446/2-1). C.O.
376 was founded by an early career fellowship at Homerton College, University of Cambridge. L. B. L. W.
377 acknowledges funding by the Wellcome Trust (WT083442AIA). J.G. was supported by grants from
378 the Medical Research Council (MC_UP_A090_1006, MC_PC_13030, MR/P011705/1 and
379 MR/P01836X/1). Work in the Reimann/Gribble laboratories was supported by the Wellcome Trust
380 (106262/Z/14/Z and 106263/Z/14/Z), UK Medical Research Council (MRC_MC_UU_12012/3) and
381 PhD funding for EKB from MedImmune/AstraZeneca. Praveen Surendran is supported by a
382 Rutherford Fund Fellowship from the Medical Research Council (MR/S003746/1). A. W. is supported
383 by a BHF-Turing Cardiovascular Data Science Award and by the EC-Innovative Medicines Initiative
384 (BigData@Heart). J.D. is funded by the National Institute for Health Research [Senior Investigator
385 Award] [*]. The EPIC-Norfolk study (<https://doi.org/10.22025/2019.10.105.00004>) has received
386 funding from the Medical Research Council (MR/N003284/1 and MC-UU_12015/1) and Cancer
387 Research UK (C864/A14136). The genetics work in the EPIC-Norfolk study was funded by the Medical
388 Research Council (MC_PC_13048). Metabolite measurements in the EPIC-Norfolk study were
389 supported by the MRC Cambridge Initiative in Metabolic Science (MR/L00002/1) and the Innovative
390 Medicines Initiative Joint Undertaking under EMIF grant agreement no. 115372. The Fenland Study is
391 supported by the UK Medical Research Council (MC_UU_12015/1 and MC_PC_13046). Nightingale
392 Health NMR assays were funded by the European Commission Framework Programme 7 (HEALTH-
393 F2-2012-279233). Metabolon Metabolomics assays and The academic coordinating centre,
394 Metabolon metabolomics, and DNA extraction and genotyping for INTERVAL was supported by core
395 funding from: National Institute for Health Research (NIHR) Blood and Transplant Research Unit in
396 Donor Health and Genomics (NIHR BTRU-2014-10024), UK Medical Research Council
397 (MR/L003120/1), British Heart Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946), the
398 NIHR [Cambridge Biomedical Research Centre at the Cambridge University Hospitals NHS Foundation
399 Trust], and NIHR BioResource (<http://bioresource.nihr.ac.uk>) [*]. This work was supported by Health
400 Data Research UK, which is funded by the UK Medical Research Council, Engineering and Physical
401 Sciences Research Council, Economic and Social Research Council, Department of Health and Social
402 Care (England), Chief Scientist Office of the Scottish Government Health and Social Care
403 Directorates, Health and Social Care Research and Development Division (Welsh Government),
404 Public Health Agency (Northern Ireland), British Heart Foundation and Wellcome. We are grateful to
405 all the participants who have been part of the project and to the many members of the study teams
406 at the University of Cambridge who have enabled this research.

407 *The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or
408 the Department of Health and Social Care.

409 UK Biobank: This research has been conducted using the UK Biobank resource under
410 Application Number 44448.

411 **Author Contribution**

412 L.A.L. and C.L. designed the study. L.A.L., M.P., and CL drafted the manuscript. L.A.L., M.P., I.D.S.,
413 L.B.L.W., R.B., C.O., V.P.W.A., J.L., E.W., E.P., P.S., S.B., V.Z., and E.S. analysed the data. J.R. and G.K.
414 designed and implemented the webserver. K.K. and N.J.W. are PIs of the EPIC-Norfolk cohort. G.A.M.
415 advised on metabolite mapping across platforms. A.L. and F.I. provided metabolite measurements
416 and quality control in Fenland. E.K.B., F.M.G., and F.R. performed all experimental work on GLP2R.
417 M.B. contributed data on MacTel. E.F. performed knowledge-based annotation of genes to variants.
418 J.D. and A.S.B. were responsible for the INTERVAL study. All authors contributed to the
419 interpretation of results and critically reviewed the manuscript.

420 **Competing Interests statement**

421 A.S.B. has received grants from AstraZeneca, Biogen, Bioverativ, Merck, Novartis, and Sanofi. J. D.
422 sits on the International Cardiovascular and Metabolic Advisory Board for Novartis (since 2010), the
423 Steering Committee of UK Biobank (since 2011), the MRC International Advisory Group (ING)
424 member, London (since 2013), the MRC High Throughput Science 'Omics Panel Member, London
425 (since 2013), the Scientific Advisory Committee for Sanofi (since 2013), the International
426 Cardiovascular and Metabolism Research and Development Portfolio Committee for Novartis and
427 the Astra Zeneca Genomics Advisory Board (2018). E.B.F. is an employee and stock holder of Pfizer.
428 L.A.L. is presently an employee and shareholder of Regeneron Pharmaceuticals Inc. The remaining
429 authors declare no competing interests.

430

431 **REFERENCES (Main text)**

- 432 1. Wishart, D. S. Metabolomics for investigating physiological and pathophysiological processes.
433 *Physiol. Rev.* **99**, 1819–1875 (2019).
- 434 2. Shin, S.-Y. Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.* **46**,
435 543–550 (2014).
- 436 3. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants associated with
437 human blood metabolites. *Nat. Genet.* **49**, 568–578 (2017).
- 438 4. Draisma, H. H. M. *et al.* Genome-wide association study identifies novel genetic variants
439 contributing to variation in blood metabolite levels. *Nat. Commun.* **6**, 7208 (2015).
- 440 5. Kettunen, J. *et al.* Genome-wide study for circulating metabolites identifies 62 loci and
441 reveals novel systemic effects of LPA. *Nat. Commun.* **7**, 11122 (2016).

- 442 6. Illig, T. *et al.* A genome-wide perspective of genetic variation in ... [Nat Genet. 2010] -
443 PubMed result. *Nat. Genet.* **42**, 137–41 (2010).
- 444 7. Suhre, K. *et al.* Human metabolic individuality in biomedical and pharmaceutical research.
445 *Nature* **477**, 54–60 (2011).
- 446 8. Rhee, E. P. P. *et al.* A genome-wide association study of the human metabolome in a
447 community-based cohort. *Cell Metab.* **18**, 130–43 (2013).
- 448 9. Gallois, A. *et al.* A comprehensive study of metabolite genetics reveals strong pleiotropy and
449 heterogeneity across time and context. *Nat. Commun.* **10**, 1–13 (2019).
- 450 10. Rhee, E. P. *et al.* An exome array study of the plasma metabolome. *Nat. Commun.* **7**, 12360
451 (2016).
- 452 11. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to
453 Common Complex Disease. *Cell* **167**, 1415–1429.e19 (2016).
- 454 12. Bansal, N. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
- 455 13. Learn, D. B., Fried, V. A. & Thomas, E. L. Taurine and hypotaurine content of human
456 leukocytes. *J. Leukoc. Biol.* **48**, 174–182 (1990).
- 457 14. Yet, I. *et al.* Genetic Influences on Metabolite Levels: A Comparison across Metabolomic
458 Platforms. *PLoS One* **11**, e0153672 (2016).
- 459 15. Foley, C. N. *et al.* A fast and efficient colocalization algorithm for identifying shared genetic
460 risk factors across multiple traits. **44**, 1–47 (2019).
- 461 16. Pedersen, C. B. *et al.* The ACADS gene variation spectrum in 114 patients with short-chain
462 acyl-CoA dehydrogenase (SCAD) deficiency is dominated by missense variations leading to
463 protein misfolding at the cellular level. *Hum Genet* **124**, 43–56 (2008).
- 464 17. Lahiri, S. *et al.* Kinetic characterization of mammalian ceramide synthases: Determination of
465 Km values towards sphinganine. *FEBS Lett.* **581**, 5289–5294 (2007).
- 466 18. Horowitz, B. *et al.* Asparagine synthetase activity of mouse leukemias. *Science (80-)*. **160**,
467 533–535 (1968).
- 468 19. Babu, E. *et al.* Identification of a Novel System L Amino Acid Transporter Structurally Distinct
469 from Heterodimeric Amino Acid Transporters. *J. Biol. Chem.* **278**, 43838–43845 (2003).
- 470 20. Scott, R. A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in
471 Europeans. *Diabetes* **66**, 2888–2902 (2017).
- 472 21. Wheeler, E. *et al.* Impact of common genetic determinants of Hemoglobin A1c on type 2
473 diabetes risk and diagnosis in ancestrally diverse populations: A transethnic genome-wide
474 meta-analysis. *PLoS Med.* **14**, (2017).
- 475 22. Prokopenko, I., Poon, W., Mägi, R., Prasad, B. R. & Salehi, S. A. A Central Role for GRB10 in
476 Regulation of Islet Function in Man. *Claire Levy-Marchal* **17**,.
- 477 23. Manning, A. K. *et al.* A genome-wide approach accounting for body mass index identifies
478 genetic variants influencing fasting glycemic traits and insulin resistance. *Nat. Genet.* **44**, 659–
479 669 (2012).
- 480 24. Almgren, P. *et al.* Genetic determinants of circulating GIP and GLP-1 concentrations. (2017)
481 doi:10.1172/jci.insight.93306.
- 482 25. Fragkos, K. C. & Forbes, A. *Citrulline as a marker of intestinal function and absorption in*
483 *clinical settings: A systematic review and meta-analysis. United European Gastroenterology*
484 *Journal* vol. 6 181–191 (SAGE Publications Ltd, 2018).
- 485 26. Tseng, C. C. & Zhang, X. Y. The cysteine of the cytoplasmic tail of glucose-dependent
486 insulinotropic peptide receptor mediates its chronic desensitization and down-regulation.
487 *Mol. Cell. Endocrinol.* **139**, 179–186 (1998).
- 488 27. Estall, J. L., Koehler, J. A., Yusta, B. & Drucker, D. J. The glucagon-like peptide-2 receptor C
489 terminus modulates β -arrestin-2 association but is dispensable for ligand-induced
490 desensitization, endocytosis, and G-protein-dependent effector activation. *J. Biol. Chem.* **280**,
491 22124–22134 (2005).
- 492 28. Scerri, T. S. *et al.* Genome-wide analyses identify common variants associated with macular

- 493 telangiectasia type 2. *Nat. Genet.* **49**, 559–567 (2017).
494 29. Gantner, M. L. *et al.* Serine and lipid metabolism in macular disease and peripheral
495 neuropathy. *N. Engl. J. Med.* **381**, 1422–1433 (2019).
496 30. Garrod, A. E. The incidence of alkaptonuria: a study in chemical individuality. *Lancet* **160**,
497 1616–1620 (1902).
498 31. Shin, S. Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.* **46**,
499 543–550 (2014).
500 32. Rath, A. *et al.* Representation of rare diseases in health information systems: The orphanet
501 approach to serve a wide range of end users. *Hum. Mutat.* **33**, 803–808 (2012).
502 33. Stübiger, G. *et al.* Targeted profiling of atherogenic phospholipids in human plasma and
503 lipoproteins of hyperlipidemic patients using MALDI-QIT-TOF-MS/MS. *Atherosclerosis* **224**,
504 177–186 (2012).
505 34. Van Der Graaf, A., Kastelein, J. J. P. P. & Wiegman, A. *Heterozygous familial*
506 *hypercholesterolaemia in childhood: Cardiovascular risk prevention. Journal of Inherited*
507 *Metabolic Disease* vol. 32 (2009).
508
509

510 **FIGURE LEGENDS**

511

512 **Figure 1A** Sample size by contributing study and technique. Metabolites with similar numbers have
513 been collapsed and exact numbers are given in the Extended Data Table 1. **B** A three-dimensional
514 Manhattan plot displaying chromosomal position (x-axis) of significant associations ($p < 4.9 \times 10^{-10}$)
515 accounting for multiple testing, z-axis) across all metabolites (y-axis). Colours indicate metabolite
516 groups. P-values were obtained from a meta-analysis of genome-wide summary statistics from linear
517 regression models using genetic variants as exposures and metabolite levels as outcome run within
518 each contributing study. **C** A top view of the 3D-Manhattan plot. Dots indicate significantly
519 associated loci. Colours indicate novelty of metabolite – locus associations. Loci with indication for
520 pleiotropy have been annotated.

521

522 **Figure 2A** Distribution of pleiotropy, i.e. number of associated metabolites, among loci identified in
523 the present study. **B** Distribution of polygenicity of metabolites, i.e. number of identified loci for
524 each metabolite under investigation. **C** Scatterplot comparing the estimated heritability of each
525 metabolite against the number of associated loci. Size of the dots indicates samples sizes. **D**
526 Heritability estimates for single metabolites. Colours indicate the proportion of heritability
527 attributed to single nucleotide polymorphisms (SNPs) with large effect sizes ($\beta > 0.25$ per allele). **E –**
528 **M** SNP – metabolite association with indication of non-additive effects. Beta is an estimate from the
529 departure of linearity. Point estimates and 95%-confidence intervals were obtained from linear
530 regression models run among up to 8,714 participants in the Fenland cohort. **N** Barplot showing the
531 increase in heritability and explained variance for each SNP – metabolite pair when including non-
532 additive effects.

533

534 **Figure 3A** Scatterplot comparing the minor allele frequencies (MAF) of associated variants with
535 effect estimates from linear regression models (N loci=499). Colours indicate possible functional
536 consequences of each variant: maroon – nonsynonymous variant; blue – in strong LD ($r^2 > 0.8$) with a
537 nonsynonymous variant and grey otherwise. **B-D** Distribution of effect sizes (B), allele frequencies
538 (C), and width of credible sets (D) based on the type of single nucleotide polymorphism (SNP) (0 –
539 non-coding or synonymous, 1 – in strong LD with nonsynonymous, 2 - nonsynonymous) identified as
540 metabolite quantitative trait loci (N=499). Data are represented as boxplots where the middle line is
541 the median, the lower and upper hinges correspond to the first and third quartiles, the upper
542 whisker extends from the hinge to the largest value no further than $1.5 \times$ IQR from the hinge (where
543 IQR is the inter-quartile range) and the lower whisker extends from the hinge to the smallest value
544 at most $1.5 \times$ IQR of the hinge, while data beyond the end of the whiskers are outlying points that
545 are plotted individually. **E** Distribution of functional annotations of metabolite associated variants
546 (red), trait-associated variants (blue – continuous, purple – diseases) obtained from the GWAS
547 catalogue, and all SNPs included in the present genome-wide association studies. The inlet for exonic
548 variants distinguishes between synonymous (syn) and nonsynonymous variants (nsyn).

549

550 **Figure 4A** Comparison between the hypothesis-free genetically prioritized versus biologically
551 plausible approaches used in the present study to assign candidate genes to metabolite associated
552 single nucleotide polymorphisms. The Venn-diagram displays the overlap between both approaches.
553 **B** Enrichment of genetically prioritized genes among biologically plausible or genes linked to inborn
554 errors of metabolism (IEM). Enrichment was tested using a two-sided binomial test. **C** Proportion of
555 genetically prioritized genes encoding for either enzymes or transporters.

556

557 **Figure 5A** Enrichment of associations with type 2 diabetes (T2D: 80,983 cases, 842,909 controls)
558 among metabolite-associated SNPs. Blue dots indicate metabolite-SNPs and grey dots indicate a

559 random selection of matched control SNPs. **B** Regional association plots for plasma citrulline, type 2
560 diabetes, body mass index, and fasting levels of glucose-dependent insulinotropic peptide (GIP)
561 focussing on the *GLPR2* gene. Variants are coloured based on linkage disequilibrium with the lead
562 variant (rs17681684) for plasma citrulline. *Summary statistics for GIP were obtained from the more
563 densely genotyped study included in Almgren et al.²⁴ (to increase coverage of genetic variants for
564 multi-trait colocalisation). **C** Individual association summary statistics for all citrulline associated
565 SNPs (coded by the citrulline increasing allele) for T2D and an inverse-variance weighted (IVW)
566 estimate pooling all effects. **D** Schematic sketch for the location of the missense variant induces
567 amino acid substitution in the glucagon-like peptide-2 receptor (GLP2R). **E** GLP-2 dose response
568 curves in cAMP assay for GLP2R wild-type and mutant receptors. The dose response curves of cAMP
569 stimulation by GLP-2 in CHO K1 cells transiently transfected with either GLP2R wild-type or mutant
570 constructs. Data were normalised to the wild-type maximal and minimal response, with 100% being
571 GLP-2 maximal stimulation of the wild-type GLP2R, and 0% being wild-type GLP2R cells with buffer
572 only. Mean \pm standard errors are presented (n=4). **F-G** Summary of wild-type and mutant GLP2R
573 beta-arrestin 1 and beta-arrestin 2 responses. Area under the curve (AUC) summary data of
574 individuals wells and geometric mean/SD from n=3-4 independent experiments (1 to 3 wells per
575 experiment) displayed for beta-arrestin 1 recruitment (E) and beta-arrestin 2 recruitment (F). AUCs
576 were calculated using the 5 minutes prior to ligand addition as the baseline value. Mean \pm standard
577 errors are presented. Normal distribution of log₁₀-transformed data was determined by the
578 D'Agostino & Pearson normality test. For statistical analysis, experiments were summarized by
579 means of log₁₀-transformed results across technical replicates and analysed using one-way ANOVA
580 with Bonferroni correction for post-hoc tests.

581

582 **Figure 6A** Results from genetic scores for each metabolite on risk for macular telangiectasia type 2
583 (MacTel). The dotted line indicates the level of significance after correction for multiple testing. The
584 inset shows the same results but after dropping the pleiotropic variants in *GCKR* and *FADS1-2*. **B**
585 Effect estimates of serine-associated genetic variants on the risk for MacTel. **C** Comparison of effect
586 sizes for lead variants associated with plasma serine levels and the risk for MacTel. **D** Receiver
587 operating characteristic curves (ROC) comparing the discriminative performance for MacTel using a)
588 sex, the first genetic principal component, and two MacTel variants (rs73171800 and rs9820286) not
589 associated with metabolite levels, and b) additionally including genetically predicted serine and
590 glycine at individual levels as described in the methods. The area under the curve (AUC) is given in
591 the legend.

592

593 **Figure 7A** Scheme of the workflow to link common variation in genes causing inborn errors of
594 metabolism (IEM) to complex diseases. **7B** Flowchart for the systematic identification of metabolite-
595 associated variants to genes and diseases related to inborn errors of metabolism (IEM). **C** P-values
596 from phenome-wide association studies among UK Biobank using variants mapping to genes
597 knowing to cause IEMs and binary outcomes classified with the ICD-10 code. Colours indicate
598 disease classes. The dotted line indicates the significance threshold controlling the false discovery
599 rate at 5%. **D** Posterior probabilities (PPs) from statistical colocalisation analysis for each significant
600 triplet consisting of a metabolite, a variant, and a ICD-10 code among UK Biobank. The dotted line
601 indicates high likelihood (>80%) for one of the four hypothesis tested: H0 – no signal; H1 – signal
602 unique to the metabolite; H2 – signal unique to the trait; H3 – two distinct causal variants in the
603 same locus and H4 – presence of a shared causal variant between a metabolite and a given trait.

604

605 TABLES

606

607 **Table 1** Genomic loci with effect sizes larger than 0.25 units in standard deviation of metabolite
608 levels per allele.

rsID	Position*	Metabolite	EA/OA	EAF	N	MA p-value	Beta**	SE	Candidate genes	Expl. var. (%)
rs13538	2:73868328	Acetylorntithine	A/G	0.78	30692	1.99E-1984	0.85	0.010	<i>NAT8, ACTG2</i>	18.4
rs3916	12:121177272	Butyrylcarnitine	C/G	0.26	30694	1.67E-2010	0.81	0.010	<i>ACADS,</i>	16.9
rs12587599	14:104575130	Asparagine	T/C	0.14	23606	8.98E-294	0.49	0.013	<i>ASPG, ADSSL1</i>	8.2
rs3970551	22:18906839	Proline	G/A	0.11	23618	1.10E-224	0.48	0.015	<i>PRODH</i>	5.0
rs174547	11:61570783	lysoPC a C20:4	T/C	0.67	16829	4.42E-398	0.47	0.015	<i>FADS1, DAGLA</i>	9.9
rs174545	11:61569306	PC aa C38:4	C/G	0.67	16828	1.37E-361	0.45	0.015	<i>FADS1,</i>	9.2
rs715	2:211543055	Glycine	C/T	0.31	80000	3.00E-1632	0.44	0.006	<i>CPS1, IDH1</i>	12.9
rs174564	11:61588305	PC ae C42:3	A/G	0.66	9363	5.72E-183	0.44	0.015	<i>FADS1, DAGLA</i>	8.9
rs174547	11:61570783	PC aa C36:4	T/C	0.67	16830	3.25e-313	0.43	0.015	<i>FADS1, DAGLA</i>	8.6
rs1171617	10:61467182	Carnitine	T/G	0.77	31001	2.06E-444	0.43	0.011	<i>SLC16A9,</i>	7.0
rs102275	11:61557803	PC ae C40:5	T/C	0.67	16839	8.23E-202	0.43	0.015	<i>C11orf10, DAGLA</i>	8.7
rs7157785	14:64235556	PC aa C28:1	T/G	0.16	16833	4.60E-136	0.35	0.019	<i>SGPP1,SYNE2</i>	3.3
rs174547	11:61570783	PC ae C36:5	T/C	0.67	16828	2.48E-185	0.33	0.015	<i>FADS1, DAGLA</i>	5.1
rs102275	11:61557803	PC aa C38:5	T/C	0.67	16836	8.31E-198	0.33	0.015	<i>C11orf10, DAGLA</i>	5.0
rs174564	11:61588305	PC ae C42:2	A/G	0.66	9363	7.04E-99	0.32	0.015	<i>FADS1, DAGLA</i>	4.8
rs174564	11:61588305	lysoPC a C26:1	A/G	0.66	9363	1.38E-91	0.32	0.016	<i>FADS1, DAGLA</i>	4.6
rs7157785	14:64235556	SM (OH) C14:1	T/G	0.16	16833	1.65E-96	0.29	0.019	<i>SGPP1</i>	2.2
rs174546	11:61569830	PC aa C24:0	C/T	0.67	13184	4.16E-89	0.29	0.016	<i>FADS1, DAGLA</i>	3.6
rs174546	11:61569830	PC ae C38:5	C/T	0.67	16839	8.98E-146	0.29	0.015	<i>FADS1, DAGLA</i>	3.9
rs7552404	1:76135946	Octanoylcarnitine	A/G	0.69	31969	2.30E-260	0.28	0.010	<i>ACADM</i>	2.8
rs1171615	10:61469090	Propionylcarnitine	T/C	0.77	32590	7.09E-185	0.27	0.011	<i>SLC16A9</i>	3.1
rs1171617	10:61467182	Acetylcarnitine	T/G	0.77	31008	1.92E-156	0.27	0.011	<i>SLC16A9</i>	3.3
rs2286963	2:211060050	Nonaylcarnitine	G/T	0.36	13925	5.46E-159	0.26	0.016	<i>ACADL</i>	3.2
rs12210538	6:110760008	Octadecandienylcarnitine	A/G	0.77	30227	1.69E-144	0.26	0.011	<i>SLC22A16</i>	1.0
rs102275	11:61557803	PC aa C36:5	T/C	0.66	16835	2.09E-120	0.25	0.015	<i>C11orf10, DAGLA</i>	3.0
rs174550	11:61571478	PC ae C36:3	C/T	0.33	16830	2.05E-105	0.25	0.015	<i>FADS1, DAGLA</i>	2.7

609

610

611

612

613

EA = effect allele; OA = other allele; MA = meta-analysis; SE = standard error; *Chromosome:Position based on Genome Reference Consortium Human Build 37; **based on a meta-analysis of linear regression models with genetic variants as exposure and metabolite levels as outcome run across cohorts for which individual-level data was available (more information is provided in Supplementary Tab. S2).

615 **Methods**

616 **Study design and participating cohorts**

617 We performed genome-wide meta-analyses of the levels of 174 metabolites from 7 biochemical
618 categories (amino acids, biogenic amines, acylcarnitines, phosphatidylcholines,
619 lysophosphatidylcholines, sphingomyelins, and sum of hexoses) captured by the Biocrates p180 kit
620 measured using mass spectrometry (MS). As described in more detail below, a total of 174
621 metabolites were successfully measured in up to 9,363 plasma samples from genotyped participants
622 of the Fenland study³⁵.

623 To maximise sample size and power, we meta-analysed genome-wide association (GWAS)
624 results from the Fenland cohort with those run in the EPIC-Norfolk³⁶ and INTERVAL³⁷ studies, in
625 which metabolites were profiled using MS (Metabolon Discovery HD4 platform) or protein nuclear
626 magnetic resonance (¹H-NMR) spectrometry^{38,39} (**Supplementary Tab. 1**). Ten of the 174 Biocrates
627 metabolites were covered across all platforms, while 38 were available on the Biocrates and
628 Metabolon platforms and 126 were unique to Biocrates (**Fig. 1**). We integrated publicly available
629 summary statistics from genome-wide meta-analyses of the same metabolites measured using MS
630 (with Biocrates or Metabolon platforms) or ¹H-NMR spectrometry (**Supplementary Tab. 1**).
631 Metabolites were matched across platforms by comparing metabolite names and biochemical
632 formulas. Mapping across different Metabolon platforms was done based on retention time/index
633 (RI), mass to charge ratio (m/z), and chromatographic data (including MS/MS spectral data).
634 Scientists at Metabolon Inc. independently reviewed and confirmed metabolite matches.

635 A summary of the characteristics of participating cohorts is given in **Supplemental Table S1** and
636 in the **Supplemental Methods**.

637 **Metabolomics measurements**

638 The levels of 174 metabolites were measured in the Fenland study by the AbsoluteIDQ®
639 Biocrates p180 Kit (Biocrates Life Sciences AG, Innsbruck, Austria) as reported elsewhere in
640 detail^{39,40}.

641 The levels of up to 38 metabolites were measured in EPIC-Norfolk and INTERVAL using the
642 Metabolon HD4 Discovery platform. Measurements were carried out using MS/MS instruments and
643 more details can be found in the Supplemental Methods.

644 The serum levels of 230 metabolites were measured in the INTERVAL study using ¹H-NMR
645 spectroscopy^{38,41}. Among those, 10 metabolites (creatinine, alanine, glutamine, glycine, histidine,
646 isoleucine, leucine, valine, phenylalanine, and tyrosine) overlapped with what is captured by the

647 Biocrates p180 Kit and were used in the present study. Further details of the ¹H-NMR spectroscopy,
648 quantification data analysis and identification of the metabolites have been described previously^{38,42}.
649 Participants with >30% of metabolite measures missing and duplicated individuals were removed.
650 Metabolite data more than 10 SD from the mean was also removed.

651 **GWAS and meta-analysis**

652 In Fenland and EPIC-Norfolk, metabolite levels were natural log-transformed, winsorised to
653 five standard deviations and then standardised to a mean of 0 and a standard deviation of 1.
654 Genotypes were measured using Affymetrix Axiom or Affymetrix SNP5.0 genotyping arrays and
655 further genotypes were imputed using 1000 Genomes Phase 3 as a reference (**Supplemental Tab. S1
656 and Methods**). GWAS were run in BOLT-LMM v2.2 or SNPTTEST v2.4.1 (Fenland, N=9,736, EPIC-
657 Norfolk, N=5,841) adjusting for age, sex, and the top four genetic principal components when
658 SNPTTEST was used. A similar workflow was used for metabolite data from the INTERVAL cohort (¹H-
659 NMR: N=40,818, Metabolon HD4: N=8,455).

660 For each metabolite, we performed a meta-analysis of z-scores (betas divided by standard
661 errors) as a measure of association, signals and loci (see below), using METAL software.
662 Heterogeneity between studies for each association was estimated by Cochran's Q-test. For each
663 metabolite, we also performed a meta-analysis of beta and standard errors for the subset of studies
664 (Fenland and, when available, EPIC-Norfolk and/or INTERVAL) where we had access to individual
665 level data and standardised phenotype preparation to estimate effect sizes. Quality filters
666 implemented after meta-analysis included exclusion of SNPs not captured by at least 50% of the
667 participating studies and 50% of the maximum sample size for that metabolite and variants with a
668 minor allele frequency below 0.5%. As a result, meta-analyses assessed the associations of up to
669 13.1 million common or low-frequency autosomal SNPs. Chromosome and base pair positions are
670 determined referring to GRCh37 annotation. To define associations between genetic variants and
671 metabolites, we corrected the conventional threshold of genome wide significance for 102 tests (i.e.
672 $p < 4.9 \times 10^{-10}$), corresponding to the number of principal components explaining 95% of the variance
673 of the 174 metabolites in the Fenland cohort, as previously described⁴³.

674 **Signal selection**

675 For each metabolite, we ranked associated SNPs ($p < 4.9 \times 10^{-10}$) by z-score to select trait-sentinel
676 SNPs and defined an "association" region as the region extending 1 Mb to each side of the trait-
677 sentinel SNP. During forward selection of trait-sentinel SNPs and loci for each trait, adjacent and
678 partially overlapping association regions were merged by extending region boundaries to a further 1

679 Mb. We defined overall lead-sentinel SNP and loci for any metabolite using a similar approach. Trait-
680 sentinel SNPs were sorted by z-score for the forward selection of lead-sentinel SNPs and a “locus”
681 was defined as the region extending 1 Mb each side of the lead-sentinel SNP. Regions larger than 2
682 Mb defined in the trait-sentinel association region definition were carried over in the definition of
683 lead-sentinel SNP loci. As a result, all lead-sentinel SNPs were >1Mb apart from each other and had
684 very low or no linkage disequilibrium ($R^2 < 0.05$).

685 For a given locus, independent signals across metabolites were determined based on linkage
686 disequilibrium (LD)-clumping of SNPs that reached the Bonferroni corrected p-value. SNPs with the
687 smallest p-values and an R^2 less than 0.05 were identified as independent signals. LD patterns were
688 estimated with SNP genotype data imputed using the haplotype reference consortium (HRC)
689 reference panel, with additional variants from the combined UK10K plus 1000 Genomes Phase 3
690 reference panel in the EPIC-Norfolk study ($n = 19,254$ after removing ancestry outliers and related
691 individuals).

692 Throughout the manuscript, the term “locus” indicates a genomic region (≥ 1 Mb each side) of a
693 lead-sentinel SNP harbouring one or more trait-sentinel SNPs; “signal” indicates a group of trait-
694 sentinel SNPs in LD with each other but not with other trait-sentinel SNPs in the locus ($R^2 < 0.05$);
695 “association” indicates trait-sentinel SNP to metabolite associations defined by a trait-lead SNP and
696 its surrounding region (≥ 1 Mb each side).

697 We tested at each locus for conditional independent variants using exact stepwise conditional
698 analysis in the largest Fenland sample ($n = 8,714$) using SNPTTEST v2.5 restricting to loci with evidence
699 for a genome-wide signal in this data set ($p < 5 \times 10^{-8}$, see **Supplemental Methods**).

700 **Investigation of heterogeneity**

701 We used a meta-regression model to identify factors associated with larger I^2 values across all
702 499 identified SNP-metabolite associations. To this end, a vector of heterogeneity estimates, I^2 , from
703 the meta-analysis was obtained as outcome and the following explanatory variables were
704 considered: strength of effect (absolute Z-score of the SNP – metabolite association), biochemical
705 class, dummy variables indicating the study of origin (related to the measurement platform), and the
706 number of contributing studies as an estimate of sample size. A significant effect of any of those
707 terms in a linear regression model was taken to indicate a source of heterogeneity across SNP-
708 metabolite associations and hence identified systematic factors contributing to any observed cross-
709 platform heterogeneity.

710 **Statistical fine-mapping**

711 We used statistical fine mapping to determine 99%-credible intervals for all independently
712 associated SNPs using the R package ‘corrcoverage’ using the subset of GWAS results for which we
713 had access to individual level data and incorporating results from conditional analysis to account for
714 multiple independent signals at a locus.

715 **Multi-trait colocalisation across metabolites**

716 We used hypothesis prioritisation in multi-trait colocalisation (HyPrColoc)¹⁵ at each of the
717 identified 144 loci 1) to identify metabolites sharing a common causal variant over and above what
718 could be identified in the meta-analysis to increase statistical power, and 2) to identify loci with
719 evidence of multiple causal variants with distinct associated metabolite clusters. HyPrColoc provides
720 for each cluster three different types of output: 1) a posterior probability (PP) that all traits in the
721 cluster share a common genetic signal, 2) a regional association probability, i.e. that all the
722 metabolites share an association with one or more variants in the region, and 3) the proportion of
723 the PP explained by the candidate variant. We considered a highly likely alignment of a genetic signal
724 across various traits if the PP > 75% or the regional association probability > 80% and the PP > 50%.
725 The second criterion takes into account that metabolites may share multiple causal variants at the
726 same locus. We used the same set of summary statistics as described for statistical fine-mapping.
727 We further filtered metabolites with no evidence of a likely genetic signal ($p > 10^{-5}$) in a region before
728 performing HyPrColoc, which improved clustering across traits by minimizing noise. We used the
729 same workflow to test for the alignment of a genetic signal at the *GLPR2* locus using summary
730 statistics from T2D (see below), a meta-analysis for body mass index across GIANT and UK Biobank,
731 plasma GIP, and plasma citrulline.

732 **Testing for non-linear effects**

733 We tested each of the 499 identified SNP (j) – metabolite (i) pairs for the deviation from an
734 additive linear model by introducing a dummy variable encoding heterozygous carriers (D), i.e. D = 1
735 if heterozygous and 0 otherwise, in the following regression model:

$$736 \text{Metabolite}_i \sim \beta_1 + \beta_2 * \text{SNP}_j + \beta_3 * D + \dots \text{Confounder} \dots + \epsilon$$

737 A significant estimate β_3 indicates departure from linearity. In a more formal framework this test
738 allows to test for either a dominant negative or positive model of inheritance depending on the
739 coding of the effect allele. We implemented this test in STATA version 14 using individual level data
740 from the Fenland cohort.

741 **Metabolic network and community detection**

742 We used Gaussian graphical modelling (GGMs) to construct a metabolic network across all 174
743 metabolites in a data-driven manner² as implemented in the R package *GeneNet*. The final network
744 comprised 167 metabolites and 554 significant ($p < 3.3 \times 10^{-6}$) edges. We performed community
745 detection using the Girvan-Newman algorithm as implemented in the R package *igraph* and
746 obtained 14 distinct communities including those covering metabolites of distinct biochemical
747 species as well as subdividing larger metabolite classes (**Supplemental Fig. S2**).

748 **Hypothesis-free (genetic) assignment of causal genes**

749 To assign likely causal genes to lead SNPs at each locus we generated a scoring system. We
750 identified the nearest gene for each variant by querying HaploReg⁴⁴. Next we integrated expression
751 quantitative trait loci (eQTL) studies (GTEx v6p) to identify genes whose expression levels are
752 associated with metabolite levels using TWAS/FUSION (Transcriptome-wide association study /
753 Functional summary-based imputation)⁴⁵. In doing so, we assigned to each variant-metabolite
754 association one or more associated genes using the variant as common anchor. We further assigned
755 higher impact for a causal gene if either the metabolite variant itself or a proxy in high linkage
756 disequilibrium ($R^2 > 0.8$) was a missense variant for a known gene again using the HaploReg database
757 to obtain relevant information. Based on those three criteria we ranked all possible candidate genes
758 and kept those with the highest score as putative causal gene.

759 **Knowledge-based (biological) assignment of causal genes**

760 Metabolite traits are unique among genetically evaluated phenotypes in that the functional
761 characterization of the relevant genes has often already been carried out using classic biochemical
762 techniques. The objective for the knowledge-based assignment strategy was to find the
763 experimental evidence that has previously linked one of the genes proximal to the GWAS lead
764 variant to the relevant metabolite. For many loci and metabolites this 'retrospective' analysis has
765 already been carried out^{31,46}. For these cases, previous causal gene assignments were generally
766 adopted. For novel loci, we employed a dual strategy that combined automated database mining
767 with manual curation. In the automated phase, seven approaches were employed to identify
768 potential causal genes among the 20 protein-coding genes closest to each lead variant, as described
769 in detail below, using the shortest distance determined from the lead SNP to each gene's
770 transcription start site (TSS) or transcription end site (TES), with a distance value of 0 assigned if the
771 SNP fell between the TSS and TES.

772 These 7 approaches were as follows:

773 1) HMDB metabolite names⁴⁷ were compared to each entrez gene name;

- 774 2) Metabolite names were compared to the name and synonyms of the protein encoded by each
775 gene⁴⁸
- 776 3) HMDB metabolite names and their parent terms (class) were compared to the names for the
777 protein encoded by each gene (UniProt).
- 778 4) Metabolite names were compared to rare diseases linked to each gene in OMIM³² after
779 removing the following non-specific substrings from disease names: uria, emia, deficiency, disease,
780 transient, neonatal, hyper, hypo, defect, syndrome, familial, autosomal, dominant, recessive, benign,
781 infantile, hereditary, congenital, early-onset, idiopathic;
- 782 5) HMDB metabolite names and their parent terms were compared to all GO biological processes
783 associated with each gene after removing the following non-specific substrings from the name of the
784 biological process: metabolic process, metabolism, catabolic process, response to, positive
785 regulation of, negative regulation of, regulation of. For this analysis only gene sets containing fewer
786 than 500 gene annotations were retained.
- 787 6) KEGG maps⁴⁹ containing the metabolite as defined in HMDB were compared to KEGG maps
788 containing each gene, as defined in KEGG. For this analysis the large “metabolic process” map was
789 omitted.
- 790 7) Each proximal gene was compared to the list of known interacting genes as defined in HMDB.
791 For each text-matching based approach, a fuzzy text similarity metric (pair coefficient) as encoded in
792 the ruby gem “fuzzy_match” was used with a score greater than 0.5 considered as a match.

793 In the next step, all automated hits at each locus were manually reviewed for plausibility. In
794 addition, other genes at each locus were reviewed if the Entrez gene or UniProt description of the
795 gene suggested it could potentially be related to the metabolite. If existing experimental evidence
796 could be found linking one of the 20 closest genes to the metabolite, that gene was selected as the
797 biologically most likely causal gene. If no clear experimental evidence existed for any of the 20
798 closest protein coding genes, no causal gene was manually selected. In a few cases multiple genes at
799 a locus had existing experimental evidence. This frequently occurs in the case of paralogs with
800 similar molecule functions. In these cases, all such genes were flagged as likely causal genes.

801 **Enrichment of type 2 diabetes associations among metabolite associated lead variants**

802 We examined whether the set of independent lead metabolite associated variants (N=168)
803 were enriched for associations with T2D. We plotted observed versus expected $-\log_{10}(p\text{-values})$ for
804 the 168 lead variants in a QQ-plot, using association statistics from a T2D meta-analysis including
805 80,983 cases and 842,909 non-cases from the DIAMANTE study⁵⁰ (55,005 T2D cases, 400,308 non-

806 cases), UK Biobank⁵¹ (24,758 T2D cases, 424575 non-cases, application number 44448) and the EPIC-
807 Norfolk study (additional T2D cases not included in DIAMANTE study: 1,220 T2D cases and 18,026
808 non-cases). This QQ-plot was compared to those for 1000 sets of variants, where variants in each set
809 were matched to the index metabolite variants in terms of MAF, the number of variants in LD
810 ($R^2 > 0.5$), gene density and distance to nearest gene (for all parameters +/- 50% of the index variant
811 value), but otherwise randomly sampled from across the autosome excluding the HLA region. MAF
812 and LD parameters for individual variants were determined from the EPIC-Norfolk study (using the
813 combined HRC, UK10K and 1000G imputation as previously described) and gene information was
814 derived from GENCODE v19 annotation⁵². A one-tailed Wilcoxon rank sum test was used to compare
815 the distribution of association $-\log_{10}$ p-values for the metabolite associated variants with that for
816 the randomly sampled, matched, variants.

817 **Functional characterisation of D470N mutant GLP2R**

818 To investigate the functional differences between wild-type (WT) GLP2R and the D470N
819 mutant GLP2R we generated D470N GLP2R mutant constructs using site-directed mutagenesis and
820 characterised canonical GLP2R signalling pathways via cAMP as well as alternative signalling
821 pathways via β -arrestin and P-ERK.

822 Human GLP2R cDNA within the pcDNA3.1+ vector was purchased, and Gibson cloning was
823 completed to insert an internal ribosome entry site (IRES) and venus gene downstream of the GLP2R
824 sequence. Following this, QuikChange Lightning site directed mutagenesis was used to perform a
825 single base change from GAC (encoding aspartic acid) to AAC (encoding asparagine) at amino acid
826 position 470 (**Supplemental Fig. 4A-B**). Successful mutagenesis was confirmed by DNA Sanger
827 sequencing (**Supplemental Fig. 4C**), and the successful products were scaled up for use in functional
828 assays. The WT and mutant GLP2R constructs within the pcDNA3.1+ vector were used to assess
829 signalling by cAMP and P-ERK. To determine β -arrestin recruitment using NanoBiT[®] technology, an
830 alternative vector was required for lower expression of GLP2R, and fusion of GLP2R to the Large BiT
831 subunit of NanoBiT[®]. For this, GLP2R was cloned into the pBiT1.1_C[TK/LgBiT] vector using
832 restriction cloning and ligation. DNA Sanger sequencing was then used for confirmation of successful
833 cloning.

834 After generation of WT and D470N GLP2R containing constructs, these were used to assess
835 differences in WT and mutant GLP2R signalling. The initial signalling pathway to be assessed was G α s
836 signalling via cAMP. CHO K1 cells were transiently transfected with WT or mutant GLP2R constructs,
837 then after 16-24 hours were treated with a dose response of GLP-2. cAMP levels were measured
838 following 30 minutes of GLP-2 treatment, in an end-point lysis HitHunter[®] cAMP assay. The presence

839 of IRES-Venus within the GLP2R expressing vectors allowed transfection efficiency to be determined
840 for each construct. Transfection efficiency was approximately 60-70%, with no differences between
841 the WT and mutant constructs. Comparison of the GLP-2 dose-response in WT and mutant GLP2R
842 expressing cells revealed no significant differences in signalling, with an almost overlapping dose
843 response curve (**Fig. 5E**).

844 Both β -arrestin 1 and β -arrestin 2 recruitment were assessed using a Nano-Glo[®] live cell
845 assay in transiently transfected HEK293 cells. Briefly, the recruitment of β -arrestin to GLP2R brings
846 the large and small BiT subunit of NanoBiT[®] together, resulting in increased luciferase activity. The
847 top concentrations from the GLP-2 dose response in the cAMP assay (1–100 nmol/l GLP-2) were
848 chosen for stimulation of the GLP2R and observation of β -arrestin recruitment. Both β -arrestin 1 and
849 β -arrestin 2 were recruited to the WT GLP2R upon GLP-2 stimulation, in a dose-dependent manner
850 (**Supplemental Fig. 5a, c**). The maximal luciferase activity for both β -arrestin 1 and β -arrestin 2
851 recruitment to the mutant GLP2R was significantly decreased when compared to the WT GLP2R,
852 indicating the extent of β -arrestin recruitment was markedly decreased (**Supplemental Fig. 5b, d**).
853 The example traces indicate that neither β -arrestin 1 or β -arrestin 2 were recruited to the mutant
854 GLP2R upon stimulation with 1 nmol/l GLP-2, however the same concentration of GLP-2 induced β -
855 arrestin recruitment to the WT GLP2R. Overall there was a significant decrease in β -arrestin 1 and β -
856 arrestin 2 recruitment to the D470N GLP2R mutant (**Figure 5F-G**).

857 **Genetic score and Mendelian randomization analysis for MacTel**

858 For each metabolite a genetic score was calculated using all variants meeting genome-wide
859 significance and their beta-estimates as weights obtained from the meta-analysis of studies for
860 which individual level data was available. We used fixed-effect meta-analysis to test for the effect of
861 the genetic score on MacTel risk using the summary statistics from the most recent GWAS. A
862 conservative Bonferroni-correction for the number of tested scores was used to declare significance
863 ($p < 3.5 \times 10^{-4}$). Sensitivity analyses were performed where the pleiotropic *GCKR* variant was removed.

864 To test for causality between circulating levels of glycine and serine for MacTel we
865 performed two types of Mendelian randomization (MR) analysis. In a two-sample univariable MR⁵³
866 we tested for an individual effect of serine (n=4 SNPs) or glycine (n=15 SNPs) on the risk of MacTel
867 using independent non-pleiotropic (i.e. the variant in *GCKR*) genome-wide SNPs as instruments. To
868 this end, we used the inverse variance weighted method to pool SNP ratio estimates using random
869 effects as implemented in the R package *MendelianRandomization*. SNP effects on the risk for
870 MacTel were obtained from²⁸. To disentangle the individual effect of those two highly correlated
871 metabolites at the same time we used a multivariable MR model⁵⁴ including all SNPs related to

872 serine or glycine (n=15 SNPs). Beta estimates and standard errors for both metabolites and all SNPs
873 were obtained from the summary statistics and mutually used as exposure variables in multivariable
874 MR. Effect estimates were again pooled using a random effect model as implemented in the R
875 package *MendelianRandomization*. This procedure allowed us to obtain causal estimates for both
876 metabolites while accounting for the effect on each other. Estimates can be interpreted as increase
877 in risk for MacTel per 1 SD increase in metabolite levels while holding the other metabolite constant.

878 To estimate a potential clinical usefulness of the identified variants we constructed two
879 genetic risk scores for MacTel using a) sex, the first genetic principal component, and the SNPs
880 rs73171800 and rs9820286 which were identified by the MacTel GWAS study²⁸ but not found to be
881 related to either glycine or serine in our study and b) all the previous but additionally including
882 genetically predicted serine and glycine at individual levels, via genetic scores, to the model. An
883 interaction between serine and sex was included²⁸. To assess the predictive ability of both models,
884 receiver operating characteristic curves were computed based on prediction values in 1,733 controls
885 and 476 MacTel cases.

886 **Identification of genes related to inborn errors of metabolism**

887 Biologically or genetically assigned candidate genes were annotated for IEM association
888 using the Orphanet database³². Using a binomial two-tailed test, enrichment of metabolic loci was
889 assessed by comparing the annotated list with the full list of 784 IEM genes in Orphanet against a
890 backdrop of 19,817 protein-coding genes⁵⁵. IEM-annotated loci for which the associated metabolite
891 matched or was closely biochemically related to the IEM corresponding metabolite(s) based on
892 IEMBase⁵⁶ were considered further for analysis.

893 We hypothesised that IEM-annotated loci with metabolite-specific consequences could also
894 have phenotypic consequences similar to the IEM. To test this, we first obtained terms describing
895 each IEM and translated them into IEM-related ICD-10 codes using the Human Phenotype Ontology
896 and previously-generated mappings^{57,58}. We obtained association statistics from the 85 IEM SNPs for
897 phenotypic associations with corresponding ICD-codes among UK Biobank restricting to diseases
898 with at least 500 cases (N=93, **Fig. 7B**, <http://www.nealelab.is/uk-biobank>). We tested locus-disease
899 pairs meeting statistical significance (controlling the false discovery rate at 5% to account for
900 multiple testing) for a common genetic signal with the corresponding locus-metabolite association
901 using statistical colocalisation. We report only those examples with strong evidence for a shared
902 genetic signal (see below).

903 **Colocalisation analyses**

904 We used statistical colocalisation⁵⁹ to test for a shared genetic signal between a metabolite and
905 a disease of interest. We obtained posterior probabilities (PP) of: H0 – no signal; H1 – signal unique
906 to the metabolite; H2 – signal unique to the trait; H3 – two distinct causal variants in the same locus
907 and H4 – presence of a shared causal variant between a metabolite and a given trait. PPs above 80%
908 were considered highly likely. We used p-values and MAFs obtained from the summary statistics
909 with default priors to perform colocalisation using the R package *coloc*.

910 DATA AVAILABILITY

911 All genome-wide summary statistics will be made available through an interactive webserver upon
912 publication of the manuscript.

913 CODE AVAILABILITY

914 Each use of software programs has been clearly indicated and information on the options that were
915 used is provided in the Methods section. Source code to call programs is available upon request.

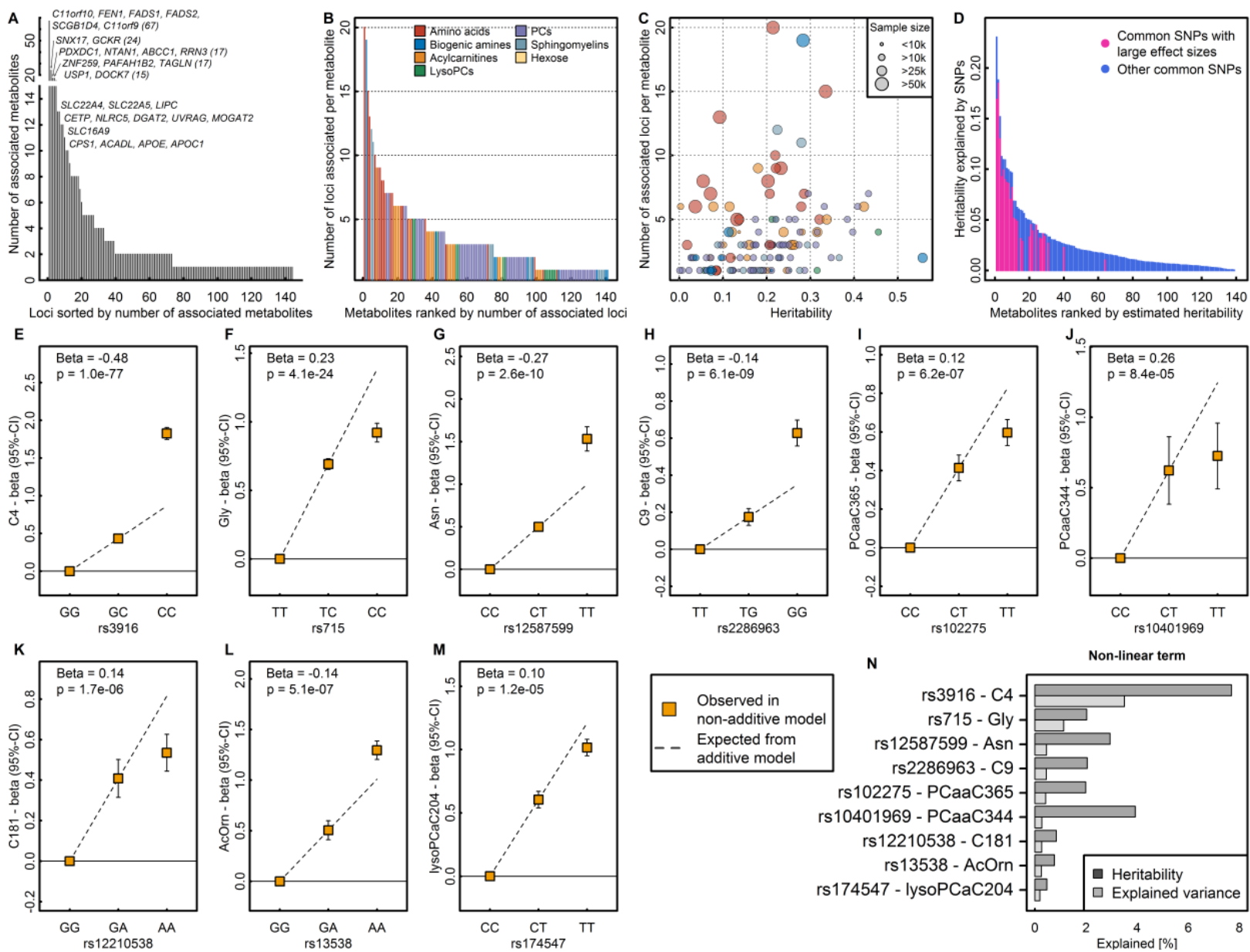
916

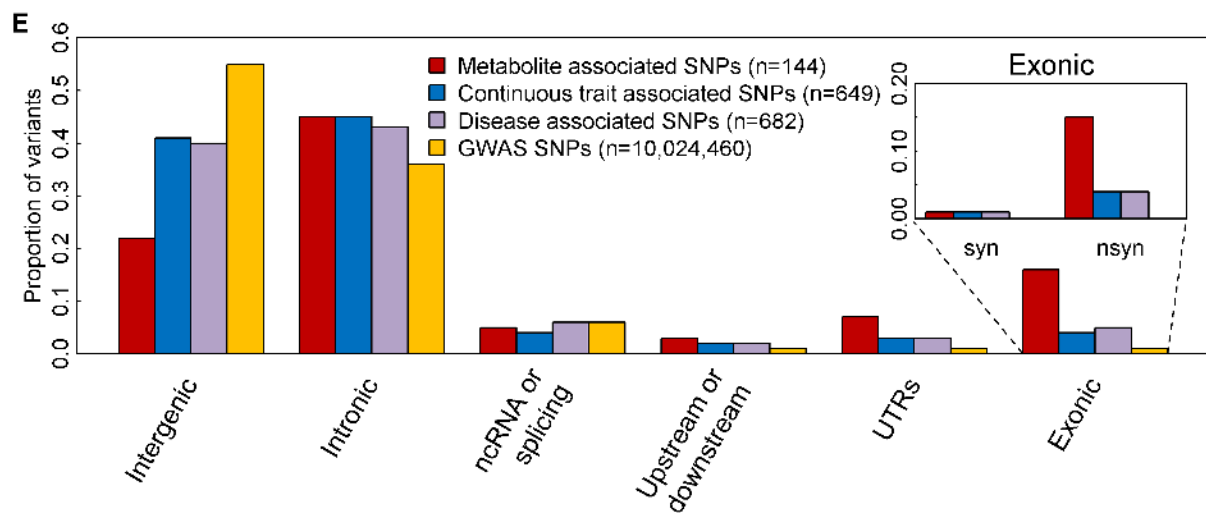
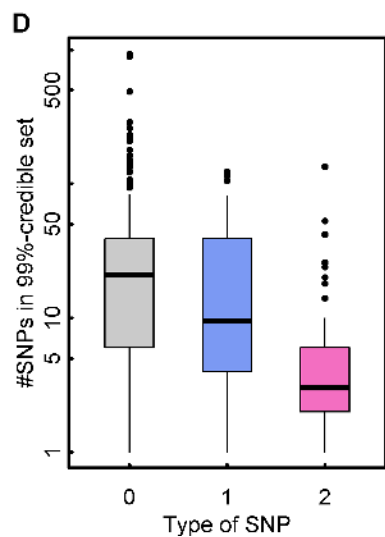
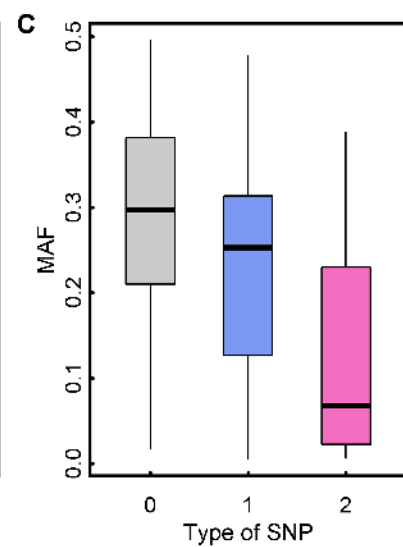
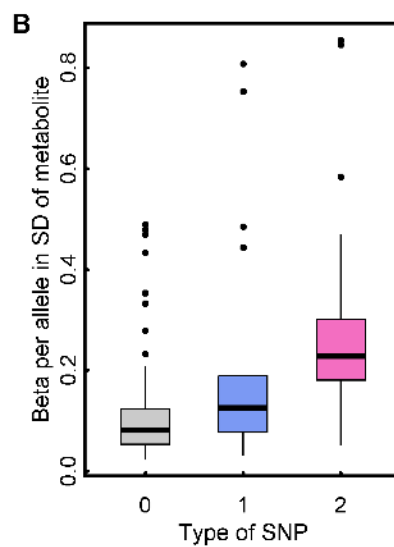
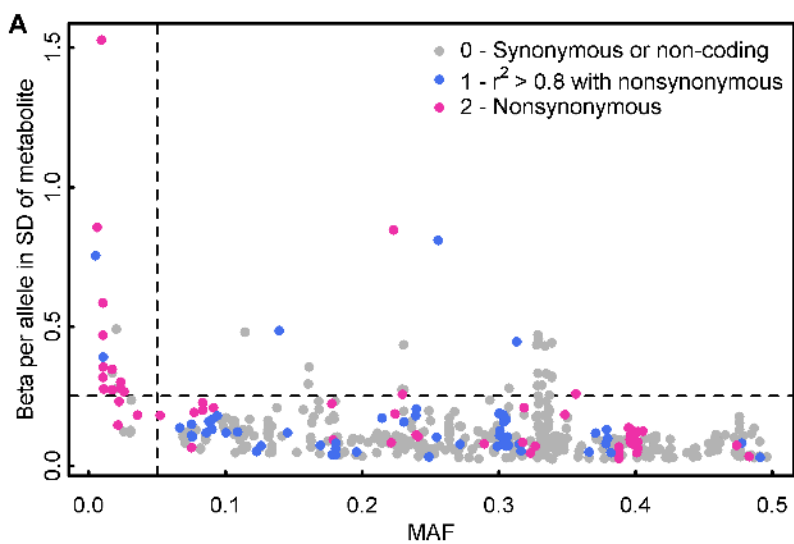
917 METHODS-ONLY REFERENCES

918

- 919 35. Lindsay, T. *et al.* Descriptive epidemiology of physical activity energy expenditure in UK adults
920 (The Fenland study). *Int. J. Behav. Nutr. Phys. Act.* **16**, 126 (2019).
- 921 36. Day, N. *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European
922 Prospective Investigation of Cancer. *Br. J. Cancer* **80 Suppl 1**, 95–103 (1999).
- 923 37. Moore, C. *et al.* The INTERVAL trial to determine whether intervals between blood donations
924 can be safely and acceptably decreased to optimise blood supply: Study protocol for a
925 randomised controlled trial. *Trials* **15**, (2014).
- 926 38. Soininen, P. *et al.* High-throughput serum NMR metabonomics for cost-effective holistic
927 studies on systemic metabolism. *Analyst* **134**, 1781–5 (2009).
- 928 39. Wittemans, L. B. L. *et al.* Assessing the causal association of glycine with risk of cardio-
929 metabolic diseases. *Nat. Commun.* **10**, (2019).
- 930 40. Lotta, L. A. *et al.* Genetic Predisposition to an Impaired Metabolism of the Branched-Chain
931 Amino Acids and Risk of Type 2 Diabetes: A Mendelian Randomisation Analysis. *PLoS Med.*
932 (2016) doi:10.1371/journal.pmed.1002179.
- 933 41. Di Angelantonio, E. *et al.* Efficiency and safety of varying the frequency of whole blood
934 donation (INTERVAL): a randomised trial of 45 000 donors. *Lancet* **390**, 2360–2371 (2017).
- 935 42. Inouye, M. *et al.* Metabonomic, transcriptomic, and genomic variation of a population cohort.
936 *Mol. Syst. Biol.* **6**, 441 (2010).
- 937 43. Li, J. & Ji, L. Adjusting multiple testing in multilocus analyses using the eigenvalues of a
938 correlation matrix. *Heredity (Edinb.)* **95**, 221–227 (2005).
- 939 44. Ward, L. D. & Kellis, M. HaploReg: A resource for exploring chromatin states, conservation,
940 and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.*
941 **40**, (2012).
- 942 45. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies.
943 *Nat. Genet.* **48**, 245–252 (2016).
- 944 46. Stacey, D. *et al.* ProGeM: A framework for the prioritization of candidate causal genes at
945 molecular quantitative trait loci. *Nucleic Acids Res.* **47**, (2019).

- 946 47. Wishart, D. S. *et al.* HMDB 4.0: The human metabolome database for 2018. *Nucleic Acids Res.*
947 **46**, D608–D617 (2018).
- 948 48. Bateman, A. *et al.* UniProt: The universal protein knowledgebase. *Nucleic Acids Res.* **45**,
949 D158–D169 (2017).
- 950 49. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: New perspectives
951 on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361 (2017).
- 952 50. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-
953 density imputation and islet-specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
- 954 51. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide
955 Range of Complex Diseases of Middle and Old Age. *PLoS Med.* **12**, (2015).
- 956 52. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.
957 *Nucleic Acids Res.* **47**, D766–D773 (2019).
- 958 53. Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis with
959 multiple genetic variants using summarized data. *Genet. Epidemiol.* **37**, 658–665 (2013).
- 960 54. Burgess, S. & Thompson, S. G. Multivariable Mendelian randomization: The use of pleiotropic
961 genetic variants to estimate causal effects. *Am. J. Epidemiol.* **181**, 251–260 (2015).
- 962 55. Harrow, J. *et al.* GENCODE: The reference human genome annotation for the ENCODE
963 project. *Genome Res.* **22**, 1760–1774 (2012).
- 964 56. Lee, J. J. Y., Wasserman, W. W., Hoffmann, G. F., Van Karnebeek, C. D. M. & Blau, N.
965 Knowledge base and mini-expert platform for the diagnosis of inborn errors of metabolism.
966 *Genet. Med.* **20**, 151–158 (2018).
- 967 57. Köhler, S. *et al.* The human phenotype ontology in 2017. *Nucleic Acids Res.* **45**, D865–D876
968 (2017).
- 969 58. Wu, P. *et al.* Mapping ICD-10 and ICD-10-CM codes to phecodes: Workflow development and
970 initial evaluation. *J. Med. Internet Res.* **21**, (2019).
- 971 59. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association
972 studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- 973
- 974





A

Genetic prioritization approach

- Physically proximity
- Genetically predicted gene expression associated with metabolite
- Nonsynonymous variant is the lead or in LD with the lead variant

Biological knowledge-based approach

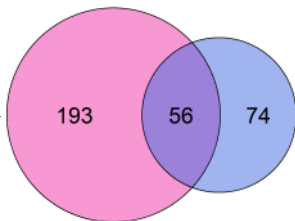
- Metabolite mapped to HMDB
 - Genes in locus considered for relevance to metabolite
-
- Entrez genes
 - Uniprot
 - OMIM
 - Go-terms
 - KEGG

Number of potential causal genes

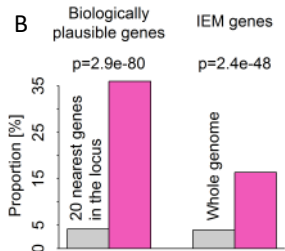
249

Number of potential causal genes

130



B



C

