

1 **Cross-platform genetic discovery of small molecule products of metabolism and application to**  
2 **clinical outcomes**

3  
4 Luca A. Lotta<sup>1#</sup>, Maik Pietzner<sup>1#</sup>, Isobel D. Stewart<sup>1</sup>, Laura B.L. Wittemans<sup>1,2</sup>, Chen Li<sup>1</sup>, Roberto  
5 Bonelli<sup>3,4</sup>, Johannes Raffler<sup>5</sup>, Emma K. Biggs<sup>6</sup>, Clare Oliver-Williams<sup>7,8</sup>, Victoria P.W. Auyeung<sup>1</sup>,  
6 Praveen Surendran<sup>7,9,10,11</sup>, Gregory A. Michelotti<sup>12</sup>, Robert A. Scott<sup>1</sup>, Stephen Burgess<sup>13,14</sup>, Verena  
7 Zuber<sup>13,15</sup>, Eleanor Sanderson<sup>16</sup>, Albert Koulman<sup>1,5,17</sup>, Fumiaki Imamura<sup>1</sup>, Nita G. Forouhi<sup>1</sup>, Kay-Tee  
8 Khaw<sup>14</sup>, MacTel Consortium, Julian L. Griffin<sup>18</sup>, Angela M. Wood<sup>7,9,10,19,20</sup>, Gabi Kastenmüller<sup>5</sup>, John  
9 Danesh<sup>7,9,10,19,21,22</sup>, Adam S. Butterworth<sup>7,9,10,19,21,22</sup>, Fiona M. Gribble<sup>6</sup>, Frank Reimann<sup>6</sup>, Melanie  
10 Bahlo<sup>3,4</sup>, Eric Fauman<sup>23</sup>, Nicholas J. Wareham<sup>1</sup>, Claudia Langenberg<sup>1,10\*</sup>

11  
12 *# these authors contributed equally*

- 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) British Heart Foundation Centre of Research Excellence, University of Cambridge, Cambridge, UK  
26 10) Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge,  
27 Cambridge, UK  
28 11) Rutherford Fund Fellow, Department of Public Health and Primary Care, University of Cambridge, UK  
29 12) Metabolon Inc, Durham, North Carolina USA  
30 13) MRC Biostatistics Unit, University of Cambridge, Cambridge, United Kingdom  
31 14) Department of Public Health and Primary Care, University of Cambridge, Cambridge, United Kingdom  
32 15) Department of Epidemiology and Biostatistics, Imperial College London, UK  
33 16) MRC Integrative Epidemiology Unit, Bristol Medical School, University of Bristol, UK  
34 17) NIHR BRC Nutritional Biomarker Laboratory, University of Cambridge, UK  
35 18) Biomolecular Medicine, Department of Metabolism, Digestion and Reproduction, Imperial College  
36 London, UK  
37 19) National Institute for Health Research Blood and Transplant Research Unit in Donor Health and Genomics,  
38 University of Cambridge, Cambridge, UK  
39 20) The Alan Turing Institute, London, UK  
40 21) National Institute for Health Research Cambridge Biomedical Research Centre, University of Cambridge  
41 and Cambridge University Hospitals, Cambridge, UK  
42 22) Department of Human Genetics, Wellcome Sanger Institute, Hinxton, UK  
43 23) Internal Medicine Research Unit, Pfizer Worldwide Research, Cambridge, MA 02142, USA

44  
45 \*Corresponding author:

46 Claudia Langenberg  
47 MRC Epidemiology Unit  
48 University of Cambridge School of Clinical Medicine  
49 Institute of Metabolic Science  
50 Addenbrooke's Treatment Centre  
51 Cambridge, UK  
52 [claudia.langenberg@mrc-epid.cam.ac.uk](mailto:claudia.langenberg@mrc-epid.cam.ac.uk)

53 **Abstract**

54       Circulating levels of small molecules or metabolites are highly heritable, but the impact of  
55 genetic differences in metabolism on human health is not well understood. In this cross-platform,  
56 genome-wide meta-analysis of 174 metabolite levels across six cohorts including up to 86,507  
57 participants (70% unpublished data), we identify 499 (362 novel) genome-wide significant  
58 associations ( $p < 4.9 \times 10^{-10}$ ) at 144 (94 novel) genomic regions. We show that inheritance of blood  
59 metabolite levels in the general population is characterized by pleiotropy, allelic heterogeneity, rare  
60 and common variants with large effects, non-linear associations, and enrichment for  
61 nonsynonymous variation in transporter and enzyme encoding genes. The majority of identified  
62 genes are known to be involved in biochemical processes regulating metabolite levels and to cause  
63 monogenic inborn errors of metabolism linked to specific metabolites, such as *ASNS* (rs17345286,  
64 MAF=0.27) and asparagine levels. We illustrate the influence of metabolite-associated variants on  
65 human health including a functional variant (rs17681684) in *GLP2R* associated with citrulline levels,  
66 impaired insulin secretion and type 2 diabetes risk. We link genetically-higher serine levels to a 95%  
67 reduction in the likelihood of developing macular telangiectasia type 2 [odds ratio (95% confidence  
68 interval) per standard deviation higher levels 0.05 (0.03-0.08;  $p = 9.5 \times 10^{-30}$ )]. We further demonstrate  
69 the predictive value of genetic variants identified for serine or glycine levels for this rare and difficult  
70 to diagnose degenerative retinal disease [area under the receiver operating characteristic curve:  
71 0.73 (95% confidence interval: 0.70-0.75)], for which low serine availability, through generation of  
72 deoxysphingolipids, has recently been shown to be causally relevant. These results show that  
73 integration of human genomic variation with circulating small molecule data obtained across  
74 different measurement platforms enables efficient discovery of genetic regulators of human  
75 metabolism and translation into clinical insights.

76

## 77 Introduction

78 Metabolites are small molecules that reflect biological processes and are widely measured in  
79 clinical medicine as diagnostic, prognostic or treatment response biomarkers<sup>1</sup>. Blood levels of  
80 metabolites are highly heritable as shown by previous studies which attempted to characterise the  
81 genetic architecture of metabolite variation in the general population<sup>2,3,4,5,6,7,8</sup>. Previous studies have  
82 been limited in scope by a focus on metabolites assessed using a single method. Integration of  
83 genetic association results for metabolites measured on different platforms can help maximise the  
84 power for a given metabolite and provide a more refined understanding of genetic influences on  
85 blood metabolite levels and human physiology.

86 To identify genomic regions regulating metabolite levels and systematically study their  
87 relevance for disease, we conducted a cross-platform meta-analysis of genetic effects on levels of  
88 174 blood metabolites measured in large-scale population-based studies on the Biocrates  
89 (AbsoluteIDQ™ p180, Fenland Study), Nightingale (<sup>1</sup>H-NMR, Interval Study) or Metabolon (Discovery  
90 HD4™, EPIC-Norfolk and Interval Studies) platforms combined with previously reported results<sup>9104</sup>.  
91 We launch with this publication a webserver to easily query our results for the purpose of targeted  
92 genetic studies, such as Mendelian randomization.

93

## 94 Results

### 95 *Associations with blood metabolites at 144 genomic regions*

96 Genome-wide meta-analyses were conducted for 174 metabolites from 7 biochemical classes  
97 (i.e. amino acids, biogenic amines, acylcarnitines, lyso-phosphatidylcholines, phosphatidylcholines,  
98 sphingomyelins and hexose) commonly measured using the Biocrates p180 kit in up to 86,507  
99 individuals, contributing over 3.7 million individual-metabolite data points (70% from unpublished  
100 studies; **Fig. 1**). For each of the 174 metabolites, this was the largest genome-wide association  
101 analyses (GWAS) to date, with at least a doubling of sample size (**Fig. 1C**). Sample sizes ranged from  
102 8,569 to 86,507 individuals for metabolites depending on the platform used in each contributing  
103 study. Using GWAS analyses we estimated the association of up to 10.2 million single nucleotide  
104 variants with a minor allele frequency (MAF) >0.5%, including 6.1 million with MAF ≥ 5%.

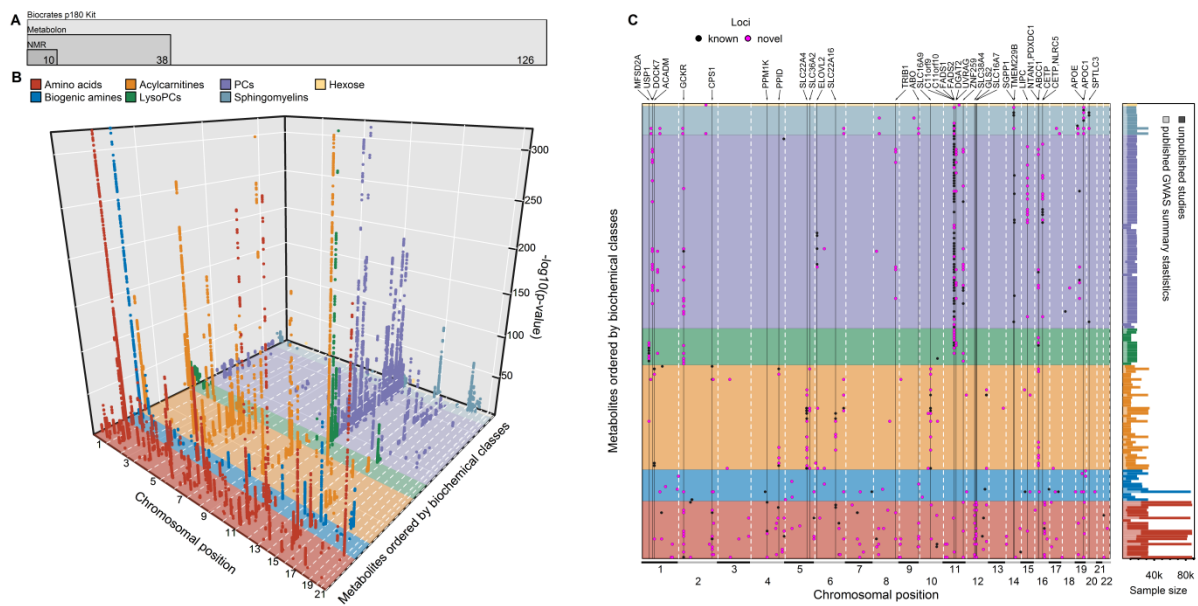
105 We identified 499 variant-metabolite associations (362 novel) from 144 loci (94 novel) at a  
106 metabolome-adjusted genome-wide significance threshold of  $p < 4.9 \times 10^{-10}$  (correcting the usual  
107 GWAS-threshold,  $p < 5 \times 10^{-8}$ , for 102 principle components explaining 95% of the variance in  
108 metabolite levels using principle component analysis; **Fig. 1**). The vast majority of these associations

109 were consistent across studies and measurement platforms [median  $I^2$ : 26.8 (interquartile range: 0 –  
110 70.1) for 465 associations with at least two contributing studies] (**Supplementary Tables 1-2**). To  
111 identify possible sources of heterogeneity, we investigated the influence of differences by cohort,  
112 measurement platform, metabolite class, and association strength in a joint meta-regression model  
113 (**Supplementary Table 3**). This showed that heterogeneity was mainly due to the overall strength of  
114 the signal, i.e. associations with higher z-scores showed greater heterogeneity ( $p < 1.05 \times 10^{-9}$ ).  
115 However, the majority of these statistically heterogeneous associations were directional consistent  
116 and nominally significant across and within each strata for 146 of 170 associations with a z-score >  
117 10, demonstrating the feasibility of pooling association estimates across metabolomics platforms for  
118 the purpose of genetic discovery. Genetic variants at the *NLRP12* locus, e.g. rs4632248, were a  
119 notable exception with large estimates of heterogeneity ( $I^2 > 90\%$ ). The *NLRP12* locus is known to  
120 affect the monocyte count<sup>11</sup> and has been shown to have pleiotropic effects on the plasma  
121 proteome in the INTERVAL study<sup>12</sup>. Monocytes, or at least a subpopulation subsumed under this cell  
122 count measure, release a wide variety of biomolecules upon activation or may die during the sample  
123 handling process and hence releasing intracellular biomolecules, such as taurine<sup>13</sup>, into the plasma.  
124 In brief, one specific source of heterogeneity in mGWAS associations might relate to sample  
125 handling differences across studies.

126 This highlights the utility of our genetic cross-platform approach to maximise power for a given  
127 metabolite, substantially extending previous efforts for any given metabolite<sup>14</sup>. Previously reported  
128 associations from platform-specific studies were also found to generally be consistent in our cross-  
129 platform meta-analysis (**Supplementary table 2; [webserver link to be inserted](#)**).

130

131



132

133 **Figure 1A** Overlap among the 174 plasma metabolites investigated in the present study across three different techniques:  
 134 Biocrates p180 Kit, Metabolon HD4, and proton nuclear magnetic resonance spectroscopy (NMR). **B** A three-dimensional  
 135 Manhattan plot displaying chromosomal position (x-axis) of significant associations ( $p < 4.9 \times 10^{-10}$ , z-axis) across all  
 136 metabolites (y-axis). Colours indicate metabolite groups. **C** A top view of the 3D-Manhattan plot. Dots indicate significantly  
 137 associated loci. Colours indicate novelty of metabolite – locus associations. Loci with indication for pleiotropy have been  
 138 annotated. Barplot on the right-hand side indicates sample sizes for each metabolite comparing the present samples sizes  
 139 with the previous largest genome-wide association study for a specific metabolite. PCs = phosphatidylcholines, LysoPCs =  
 140 lysophosphatidylcholines

141

### 142 *Insights in the genetic architecture of metabolite levels*

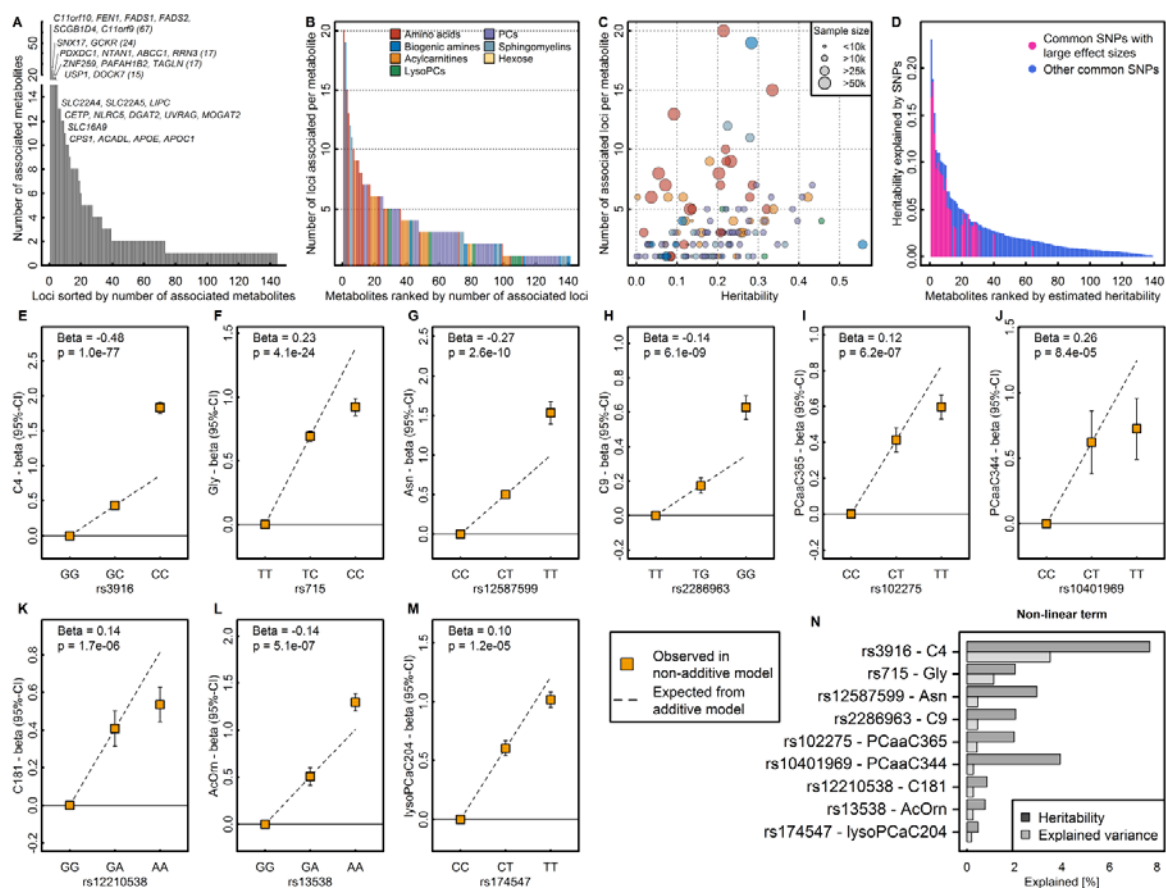
143 We found a median of 2 (range: 1-67, **Fig. 2A**) associated metabolites for each locus and a  
 144 median of 3 (range: 1-20, **Fig. 2B**) locus associations for each metabolite, reflecting pleiotropy and  
 145 the extensive contribution of genetic loci to circulating metabolite levels. The number of associations  
 146 was proportional to the estimated heritability and the sample size of the meta-analysis for a given  
 147 trait (**Fig. 2C**).

148 Similar to what is routinely observed in GWAS literature, effect size estimates increased with  
 149 decreasing minor allele frequency (MAF) (**Fig. 3A**). However, there were 26 associations  
 150 (**Supplemental Table 2**) for common lead variants with per-allele differences in metabolites levels  
 151 greater than 0.25 standard deviations (SD), a per-allele effect size that is >3-fold larger than the  
 152 strongest common variants associated with SDs of body mass index at the *FTO* locus.

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

158 GWAS analyses generally assume a linear relationship between genotypes and phenotypes, i.e.  
 159 an additive dose-response model. The identification of several metabolite-associated variants with  
 160 large effect sizes and availability of individual-level data in the Fenland cohort allowed us to test  
 161 whether the metabolite-associated variants showed evidence of deviation from a linear model. Of  
 162 499 associations tested, 9 showed evidence of departure from a linear association ( $p < 0.0001$ ; 180-  
 163 fold more than expected by chance; two-tailed binomial  $p < 2.2 \times 10^{-16}$ ; **Fig. 2E-M**). Modelling actual  
 164 genotypes rather than assuming linear associations in these instances explained a median of 7.4%  
 165 more (range: 1.4-15.2%) of the heritability in metabolite levels (**Fig. 2N**).

166



167

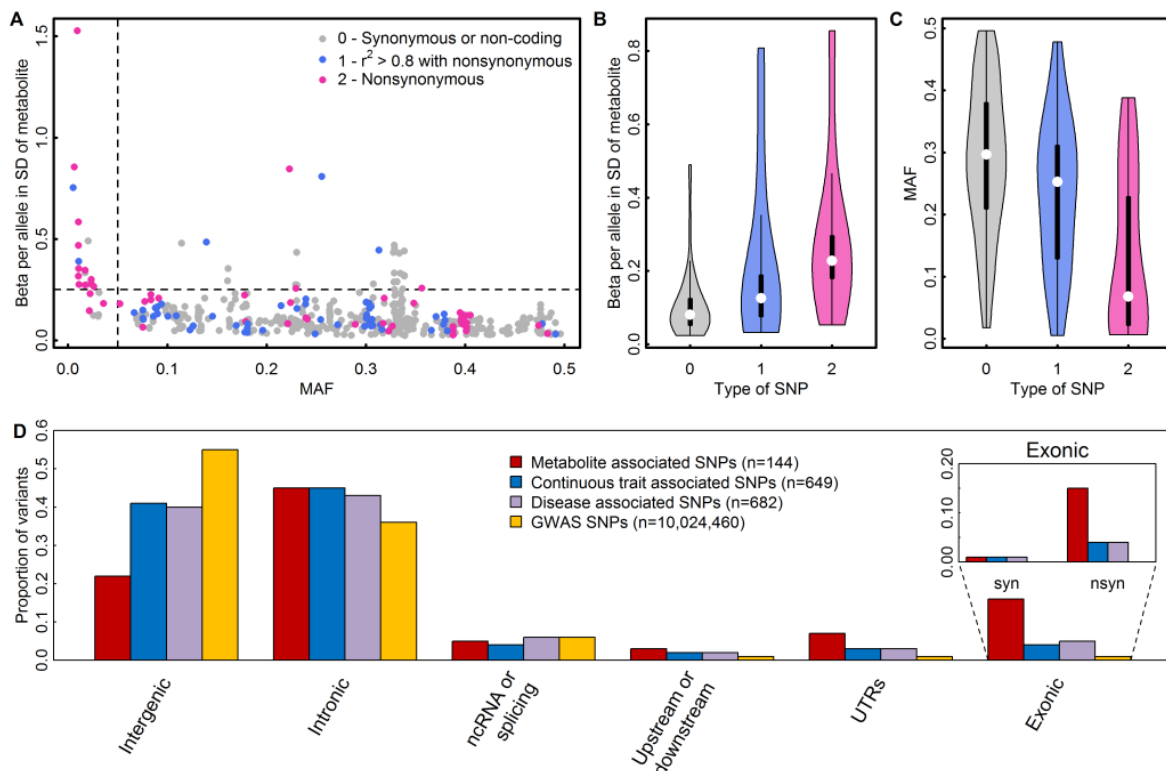
168 **Figure 2A** Distribution of pleiotropy, i.e. number of associated metabolites, among loci identified in the present study. **B**  
 169 Distribution of polygenicity of metabolites, i.e. number of identified loci for each metabolite under investigation. **C**  
 170 Scatterplot comparing the estimated heritability of each metabolite against the number of associated loci. Size of the dots  
 171 indicates sample sizes. **D** Heritability estimates for single metabolites. Colours indicate the proportion of heritability  
 172 attributed to single nucleotide polymorphisms (SNPs) with large effect sizes ( $\beta > 0.25$  per allele). **E – M** SNP – metabolite  
 173 association with indication of non-additive effects. Beta is an estimate from the departure of linearity. **N** Barplot showing  
 174 the increase in heritability and explained variance for each SNP – metabolite pair when including non-additive effects.

175

176

177

178 In 61 of the 499 associations the lead association signal was a nonsynonymous variant, a 40-fold  
179 enrichment compared to what would be expected by chance given the annotation of ascertained  
180 genetic variants (two-tailed binomial test,  $p=5 \times 10^{-30}$ , **Fig. 3D**). For a further 59 associations, the lead  
181 variant was in high LD with a nonsynonymous variant ( $R^2 > 0.8$ ). Lead variants that were  
182 nonsynonymous, or variants in high LD with a nonsynonymous variant, generally had lower MAF and  
183 larger effect sizes than variants that were not in these categories (**Fig 3B-C**).



184

185 **Figure 3A** Scatterplot comparing the minor allele frequencies (MAF) of associated variants with effect estimates from linear  
186 regression models (N loci=499). Colours indicate possible functional consequences of each variant: maroon –  
187 nonsynonymous variant; blue – in strong LD ( $r^2 > 0.8$ ) with a nonsynonymous variant and grey otherwise. **B/C** Distribution of  
188 effect sizes (B) and allele frequencies (C) based on the type of single nucleotide polymorphism (SNP) (0 – non-coding or  
189 synonymous, 1 – in strong LD with nonsynonymous, 2 - nonsynonymous). **D** Distribution of functional annotations of  
190 metabolite associated variants (red), trait-associated variants (blue – continuous, purple – diseases) obtained from the  
191 GWAS catalogue, and all SNPs included in the present genome-wide association studies. The inlet for exonic variants  
192 distinguishes between synonymous (syn) and nonsynonymous variants (nsyn).

193

194 We identified 22 loci harbouring two (n=21) or three (n=1) independent signals, i.e. different  
195 plasma metabolites were associated with distinct genetic variants within the same genomic region  
196 (**Supplementary Table S2**). For six regions our two different annotations approaches assigned only  
197 one causal gene (see below and **Methods**), including *ACADM*, *GLDC*, *ARG1*, *MARCH8*, *SLC7A2*, and  
198 *LIPC* (**Supplementary Table S2**). We found evidence that allelic heterogeneity, i.e. conditionally  
199 independent variants at a locus for a specific metabolite, explains the association pattern at 3 of  
200 those loci (*ACADM*, *ARG1*, and *LIPC*; **Supplementary Table S4**). We identified another 16 loci

201 harbouring at least one (range: 2–6) additional conditionally independent variant(s) in exact  
202 conditional analyses (see **Methods, Supplementary Table S4**).

203

#### 204 *Effector genes, tissues, pathways*

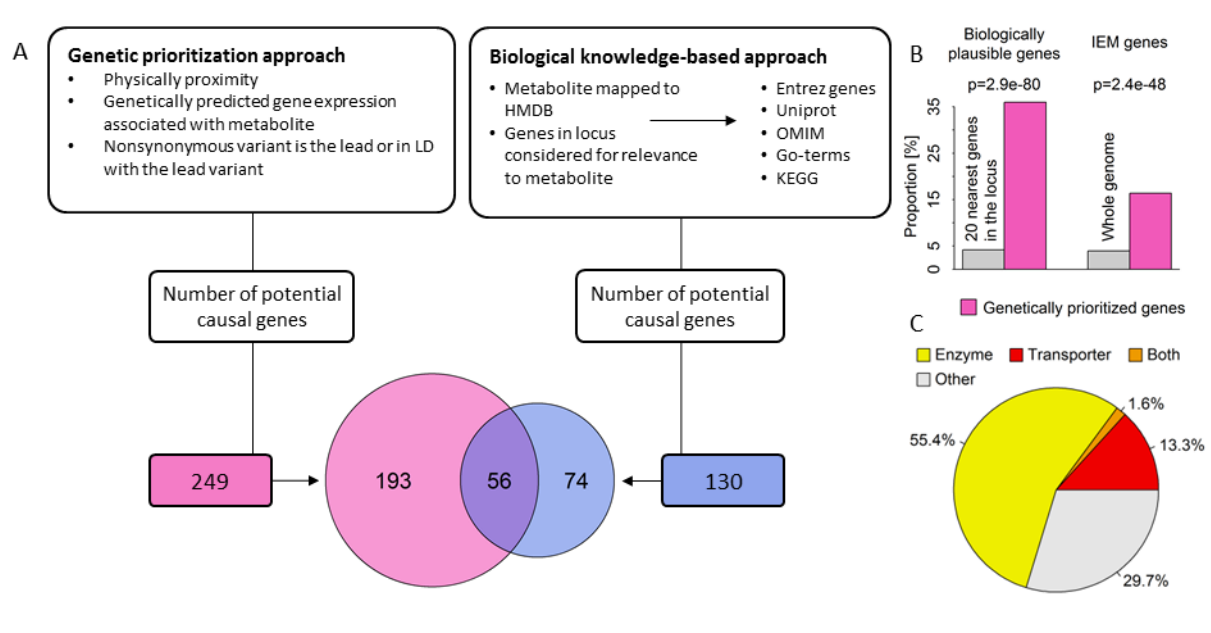
205 We used two complementary strategies to prioritize likely causal genes for the observed  
206 associations: (1) a hypothesis-free genetic approach based on physical distance, genomic annotation  
207 and integration of expression quantitative trait loci (eQTLs) to prioritize genes in a systematic and  
208 standardised way (see **Methods**), and (2) a biological knowledge-based approach integrating existing  
209 knowledge about specific metabolites or related pathways to identify biologically plausible  
210 candidate genes from the 20 genes closest to the lead variant (**Fig. 4A**). Using the hypothesis-free  
211 genetic approach, we identified 249 unique likely causal genes for the 499 associations, with at least  
212 one gene per association and some genes prioritized as likely causal for multiple metabolite  
213 associations. The knowledge-based approach identified 130 biologically plausible genes for 349 out  
214 of 499 associations. We asked whether the hypothesis-free genetic approach identified biologically  
215 plausible genes (prioritized by strategy 2) more often than expected by chance. An excess of  
216 biologically plausible genes amongst those prioritized by the hypothesis-free genetic algorithm  
217 would suggest that the approach is able to prioritize true positive associations that map to known  
218 pathways or underlying biology. Amongst 9,980 genes screened at the 499 associations, 420 (4.2%)  
219 were biologically plausible. A total of 350 gene-metabolite assignments from the first approach were  
220 also annotated in the knowledge-based approach with 126 pairs (36%) termed biologically plausible  
221 (~8-fold more than expected by chance; two-tailed binomial test,  $p=2.3\times 10^{-80}$ ; **Fig. 4B**). Among the  
222 consistently assigned genes between both approaches assignment of the nearest gene (124 times  
223 out of 126,  $X^2$ -test,  $p<2.5\times 10^{-45}$ ) was the strongest shared factor, as might be expected, followed by  
224 presence of a missense variant at least in LD ( $R^2>0.8$ , 30 times out of 126,  $X^2$ -test,  $p<1.3\times 10^{-07}$ ) and  
225 only a minor contribution of eQTL data (20 times out of 126,  $X^2$ -test,  $p<0.001$ ). Over 70% of  
226 genetically prioritized genes were enzymes or transporters (**Fig. 4C**).

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

231

232





233

234 **Figure 4A** Comparison between the hypothesis-free genetically prioritized versus biologically plausible approaches used in  
 235 the present study to assign candidate genes to metabolite associated single nucleotide polymorphisms. The Venn-diagram  
 236 displays the overlap between both approaches. **B** Enrichment of genetically prioritized genes among biologically plausible  
 237 or genes linked to inborn errors of metabolism (IEM). **C** Proportion of genetically prioritized genes encoding for either  
 238 enzymes or transporters.

239

240 Integrating GWAS statistics across cohorts and platforms allowed us to identify three genes that  
 241 have never been associated with any metabolite level so far. At the *CERS6* locus, rs4143279  
 242 associates with levels of sphingomyelin (d18:1/16:0) ( $p = 4.2 \times 10^{-10}$ ). *CERS6* encodes a ceramide  
 243 synthase facilitating formation of ceramide, a precursor of sphingomyelins<sup>15</sup>. At the *ASNS* locus,  
 244 rs17345286 associates with levels of asparagine ( $p = 4.7 \times 10^{-20}$ ). The lead variant is in high LD ( $R^2=1$ )  
 245 with a missense mutation in *ASNS* (rs1049674, p.Val210Glu). *ASNS* encodes an asparagine  
 246 synthase<sup>16</sup>. Finally, at the *SLC43A1* locus, rs2649667 associates with levels of phenylalanine ( $p =$   
 247  $3.6 \times 10^{-13}$ ). *SLC43A1* encodes a liver-enriched transporter of large neutral amino acids, including  
 248 phenylalanine<sup>17</sup>.

#### 249 *Insights into the causes of common and rare diseases from metabolite-associated loci*

250 The phenotypic consequences of metabolite-associated variants are currently not well  
 251 characterized. Here, we systematically investigate the contribution of individual loci and polygenic  
 252 predisposition associated with differences in metabolite levels to the risk of common and rare  
 253 diseases.

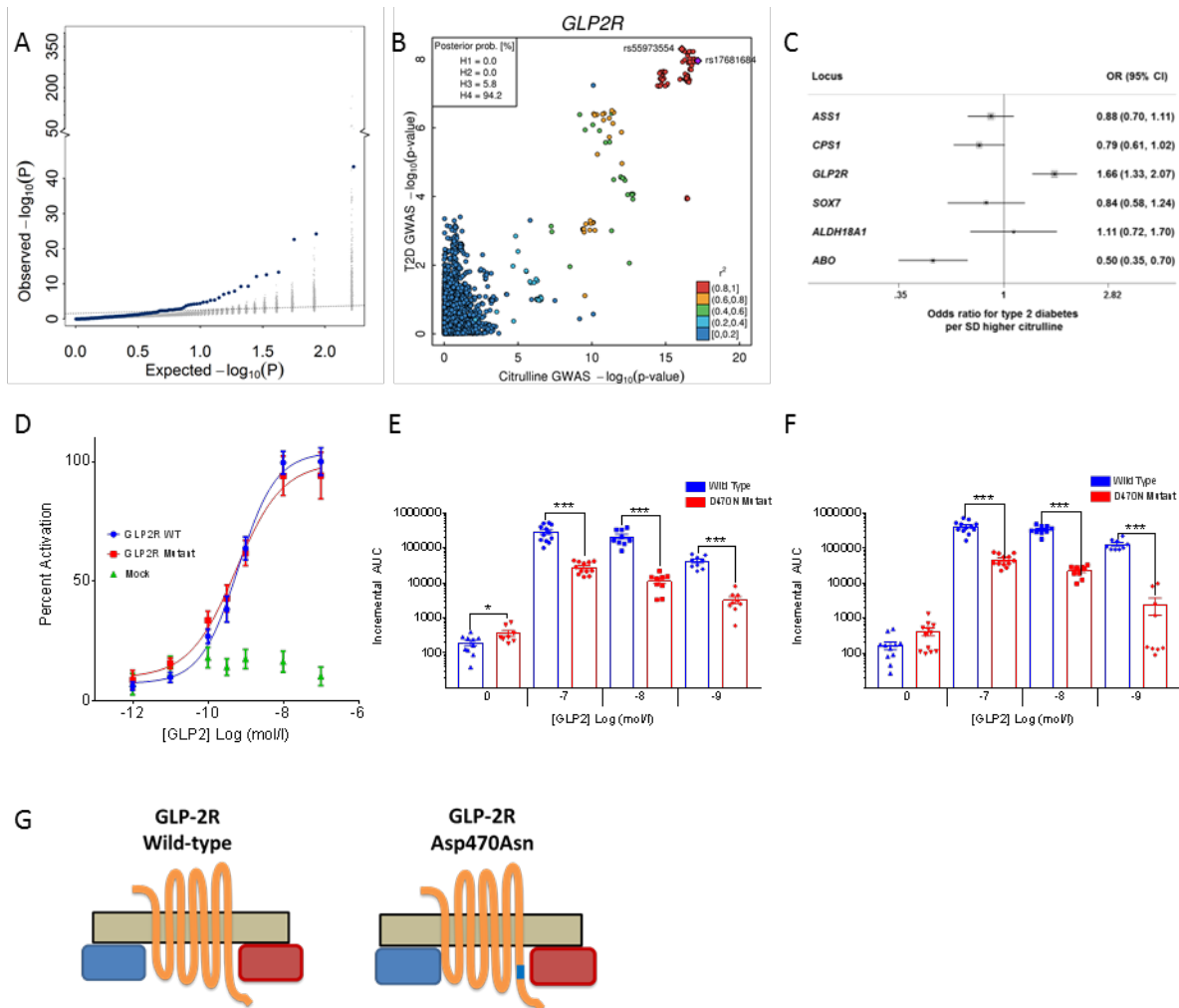
#### 254 *A citrulline-raising functional variant in GLP2R increases type 2 diabetes risk*

255 Because several of the metabolites captured in this GWAS have been associated with incident  
 256 type 2 diabetes, we sought to investigate whether the association between metabolite-associated

257 loci and diabetes could provide insights into underlying pathophysiologic mechanisms. Using  
258 estimates of effect for association with type 2 diabetes based on a meta-analysis of 80,983 cases and  
259 842,909 controls (see **Methods**), we observed a significant enrichment for associations with type 2  
260 diabetes ( $p\text{-value}=2.8\times 10^{-7}$ ) of metabolite-associated variants compared to a matched control set of  
261 variants (**Fig. 5A**).

262 Amongst the diabetes- and metabolite-associated loci was a missense p.Asp470Asn  
263 (rs17681684) variant in the *GLP2R* gene encoding the receptor for glucagon-like peptide 2, a 33  
264 amino acid peptide hormone encoded by the proglucagon gene (*GCG*) that stimulates the growth of  
265 intestinal tissue. Common variants at *GLP2R* are associated with an increased risk of type 2 diabetes  
266 (T2D) and lower levels of insulin secretion<sup>18</sup>. The previously reported lead variant for T2D  
267 (rs78761021) is in high LD (0.87) with our lead citrulline association signal at *GLP2R* (rs17681684),  
268 which was associated with a 4% higher type 2 diabetes risk (per-allele odds ratio, 1.04; 95%  
269 confidence interval, 1.02, 1.05;  $p=1.1\times 10^{-08}$ ), comparable to previous reports<sup>18</sup>. Our results show that  
270 the genetic signal for association with T2D statistically colocalises (posterior probability=0.94) with  
271 that for citrulline (**Fig. 5B**), levels of which reflect the volume of intestinal cells and are a marker of  
272 *GLP2R* target engagement in the treatment of short-bowel syndrome with glucagon-like peptide 2  
273 analogues<sup>19</sup>. The *GLP2R* p.Asp470Asn variant was the only of 6 independent genome-wide significant  
274 citrulline-raising loci that was associated with a higher risk of type 2 diabetes, which indicates that  
275 the association does not reflect a general impact of citrulline levels on diabetes risk but rather a  
276 locus-specific association at *GLP2R* (**Fig. 5C**). Taken together, this suggests that genetically higher  
277 *GLP2R* signalling, indicated by the higher citrulline levels among *GLP2R* 470Asn carriers, is associated  
278 with an increased risk of diabetes via lower insulin secretion.

279 G-protein coupled receptors like *GLP2R* may signal via G-protein-dependent cyclic adenosine  
280 monophosphate (cAMP) production or via G-protein-independent beta-arrestin mediated  
281 signalling<sup>20</sup>. To investigate if the *GLP2R* p.Asp470Asn variant affects signalling via either of these  
282 pathways, we expressed the *GLP2R* p.Asp470Asn variant in different *in vitro* models (see **Methods**).  
283 We show that the variant allele is significantly associated with reduced recruitment of beta-arrestin  
284 to *GLP2R* upon glucagon-like peptide 2 stimulation, but not with cAMP signalling, which suggests a  
285 potential role for impaired beta-arrestin recruitment to *GLP2R* in the pathophysiology of type 2  
286 diabetes (**Fig. 5D-F**).



287

288 **Figure 5A** Enrichment of associations with type 2 diabetes (T2D: 80,983 cases, 842,909 controls) among metabolite-  
 289 associated SNPs. Blue dots indicate metabolite-SNPs and grey dots indicate a random selection of matched control SNPs.  
 290 **B** Opposing  $-\log_{10}(p\text{-values})$  from the genome-wide association study of plasma citrulline with those from the T2D-GWAS  
 291 for SNPs located around *GLP2R*. The legend in the upper left gives the posterior probabilities (PP) from statistical  
 292 colocalisation analysis. H4 = PP for the hypothesis of a shared causal variant. **C** Individual association summary statistics for  
 293 all citrulline associated SNPs (coded by the citrulline increasing allele) for T2D. **D** GLP-2 dose response curves in cAMP assay  
 294 for GLP2R wild-type and mutant receptors. The dose response curves of cAMP stimulation by GLP-2 in CHO K1 cells  
 295 transiently transfected with either GLP2R wild-type or mutant constructs. Data were normalised to the wild-type maximal  
 296 and minimal response, with 100% being GLP-2 maximal stimulation of the wild-type GLP2R, and 0% being wild-type GLP2R  
 297 cells with buffer only. Mean  $\pm$  standard errors are presented ( $n=4$ ). **E-F** Summary of wild-type and mutant GLP2R beta-  
 298 arrestin 1 and beta-arrestin 2 responses. Area under the curve (AUC) summary data ( $n=3-4$ ) displayed for beta-arrestin 1  
 299 recruitment (**E**) and beta-arrestin 2 recruitment (**F**). AUCs were calculated using the 5 minutes prior to ligand addition as  
 300 the baseline value. Mean  $\pm$  standard errors are presented. Normal distribution of log10 transformed data was determined  
 301 by the D'Agostino & Pearson normality test. Following this statistical significance was assessed by one-way ANOVA with  
 302 post hoc Bonferroni test. \*\*\* $p<0.001$ , \* $p<0.05$ . **G** Schematic sketch for the location of the missense variant induces amino  
 303 acid substitution in the glucagon-like peptide-2 receptor (GLP2R).

304

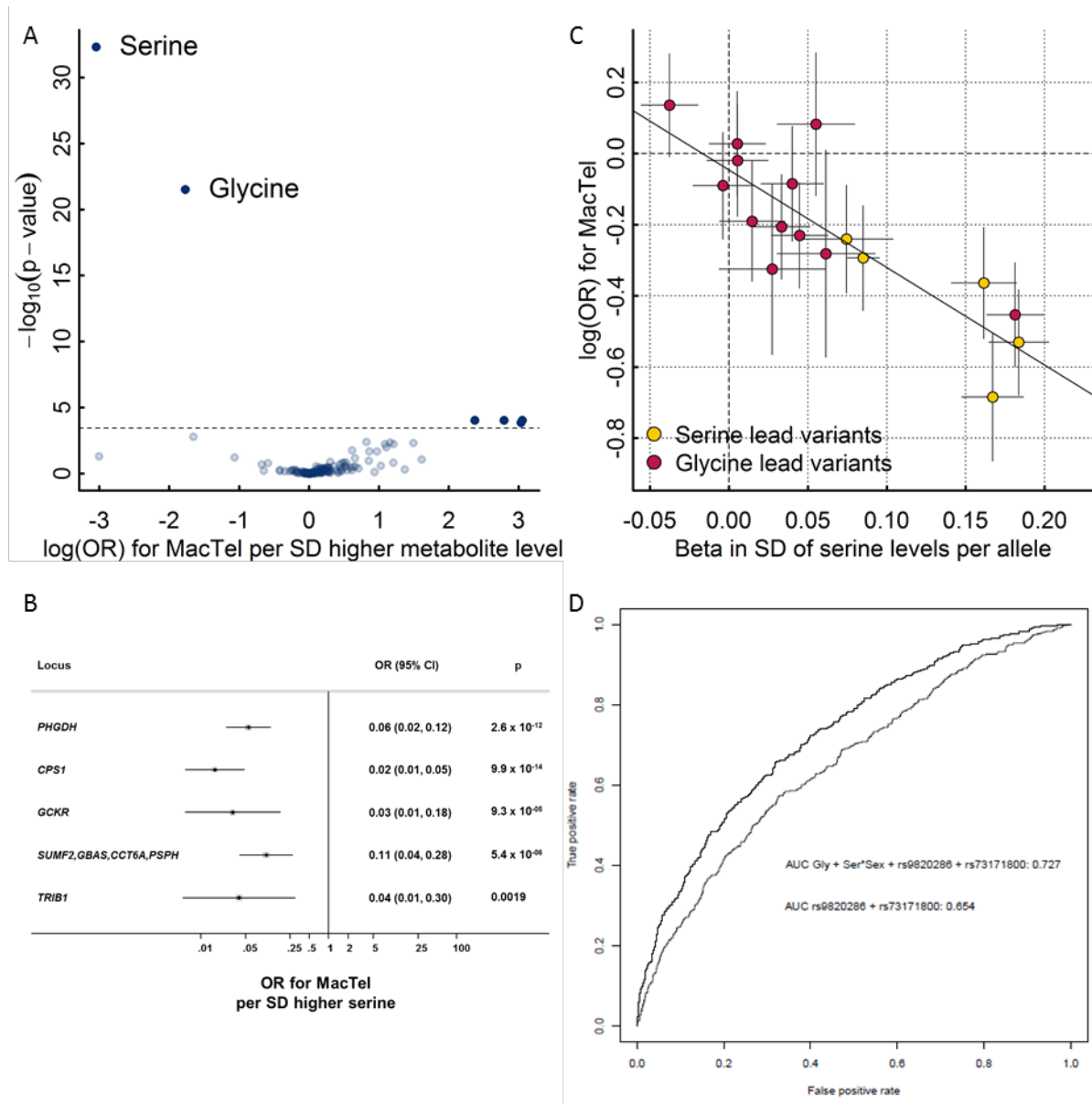
305 *Serine and glycine levels play a critical role in the aetiology of a rare eye disease*

306 A recent GWAS of macular telangiectasia type 2 (MacTel), a rare neurovascular degenerative  
 307 retinal disease, identified three genome-wide susceptibility loci (*PHGDH*, *CPS1*, and *TMEM161B-LINC00461*) of which the same variants at *PHGDH* and *CPS1* were associated with levels of the amino  
 308 acids serine and glycine in this GWAS<sup>21</sup>. More recently, it was shown that low serine availability is  
 309

310 linked to both MacTel as well as hereditary sensory and autonomic neuropathy type 1 through  
311 elevated levels of atypical deoxyshingolipids<sup>22</sup>. Whether genetic predisposition to low serine and  
312 glycine levels affects MacTel more generally or has predictive utility has not been investigated. To  
313 test this and to explore the specificity of associations between genetic influences on metabolite  
314 levels and the risk of MacTel, we generated genetic risk scores (GRS) using the sentinel variants for  
315 each of the 141 metabolites with at least one significantly associated locus identified in this GWAS  
316 and tested their associations with the risk of MacTel. GRS's for serine and glycine were the only  
317 scores associated with risk for MacTel after removal of the known highly pleiotropic *GCKR* variant  
318 (**Fig. 6A**). Each standard deviation higher serine levels via the serine GRS was associated with a 95%  
319 lower risk of MacTel (odds ratio (95% confidence interval), 0.05 (0.03-0.08);  $p=9.5\times 10^{-30}$ ; **Fig. 6A**).  
320 Each of five serine associated variants was individually associated with lower MacTel risk, with a  
321 clear dose-response relationship and no evidence of heterogeneity (**Fig. 6B**). The association was  
322 unchanged when removing the *GCKR* locus. To disentangle the effect of these two highly correlated  
323 metabolites on MacTel risk, we used multivariable Mendelian randomization analysis, which allowed  
324 us to test for a causal effect of both measures simultaneously. In this analysis, the effect of serine  
325 remained strong, while the effect of glycine was attenuated (**Tab. 1**). Glycine and serine can be  
326 interconverted and these results provide genetic evidence that the link between glycine and MacTel  
327 is via serine levels through glycine conversion. This hypothesis is supported by the evidence of a log-  
328 linear relationship between associations with serine and risk of MacTel among glycine-associated  
329 variants (**Fig. 6B**). These findings provide strong evidence that pathways indexed by genetically  
330 higher serine levels are strongly and causally associated with protection against MacTel.

331         Given the large observed effect size, we estimated whether using serine and glycine-associated  
332 loci might improve the prediction of this rare disease. Adding genetically predicted glycine and  
333 serine levels substantially improved prediction of MacTel based on an area under the receiver  
334 operating characteristic curve from 0.65 (CI 95%: 0.626-0.682) to 0.73 (0.702-0.753) (**Fig. 6**).

335



336

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

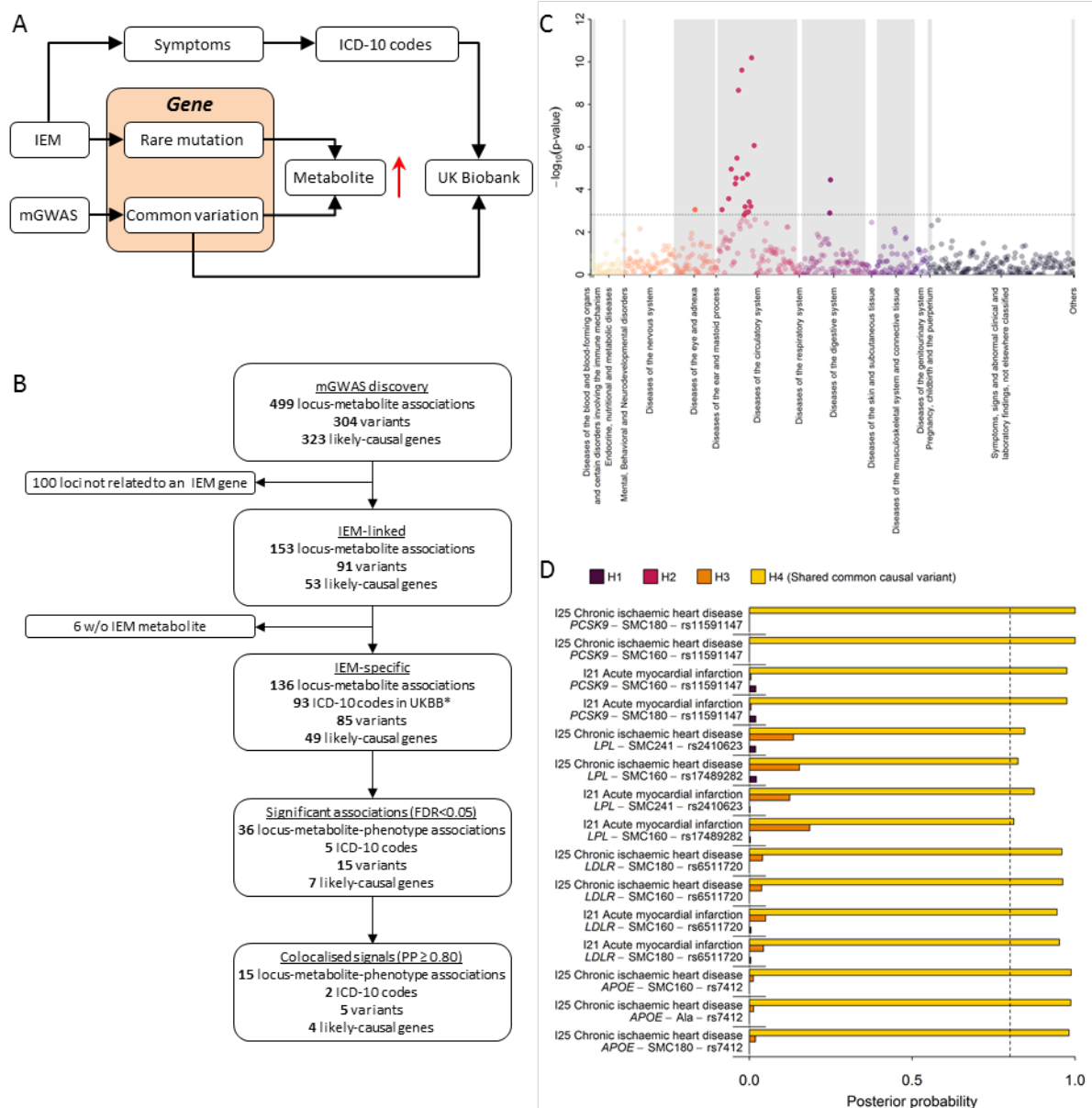
345

346 *Common variation at inborn error of metabolism (IEM) associated genes influences the risk of*  
 347 *common manifestations of diseases related to the phenotypic presentation of those IEMs*

348 In his seminal 1902 work on alkaptonuria<sup>23</sup>, also known as dark or black urine disease, Archibald  
 349 Garrod was the first to hypothesise that inborn errors of metabolism are “extreme examples of  
 350 variations of chemical behaviour which are probably everywhere present in minor degrees”.  
 351 Previous studies have shown enrichment of metabolite quantitative trait loci in genes known to

352 cause IEMs<sup>9</sup>. Whether or not these differences translate into clinically manifest disease remains  
353 unknown. The identification of several metabolite-associated variants at IEM-linked genes in this  
354 GWAS meta-analysis allows an investigation of the health consequences of genetically determined  
355 differences in metabolism for more frequently occurring variants, representing potentially milder  
356 forms of the metabolic and other clinical symptoms of IEMs, and providing new candidate genes for  
357 rare extreme metabolic disorders that currently lack a genetic basis (**Fig. 7A**). In this study, there  
358 were 153 locus-metabolite associations for which 53 unique IEM-associated genes were prioritized  
359 as likely causal using either the hypothesis-free genetic approach or the knowledge-based approach  
360 on the basis of the Orphanet database<sup>24</sup>. In 89% of these associations (136 of 153) the metabolite  
361 associated with a given GWAS locus perfectly matched, or was closely related to, the metabolite  
362 affected in patients with the corresponding IEM (**Fig. 7B**).

363 To test whether IEM-mirroring lead variants from our metabolite GWAS may increase the risk of  
364 common manifestations of diseases known to exist in patients with the corresponding IEM (**Fig. 7A**)  
365 we obtained a list of electronic health record diagnosis codes (International Statistical Classification  
366 of Diseases and Related Health Problems 10th Revision [ICD-10]) and mapped those based on  
367 symptoms seen in both, IEM patients and patients with common, complex disease manifestations  
368 (see **Methods**). We identified 93 ICD-10 codes with at least 500 cases within the UK Biobank study  
369 which aligned with the symptomatic seen in IEMs caused by mutations in genes associated with  
370 metabolite levels in the present study. We obtained the association statistics of 85 unique  
371 metabolite-associated lead variants at the 136 locus-metabolite associations with these 93 clinical  
372 diagnoses and observed 36 associations that met statistical significance (false discovery rate < 5%,  
373 **Supplemental Table S5 and Fig. 7B**). For 15 out of those we obtained strong evidence of a shared  
374 genetic metabolite-phenotype signal using colocalisation analyses (posterior probability of a shared  
375 signal >80%; **Fig. 7D and Supplemental Fig. S1**). These instances linked common genetic variants in  
376 or near *APOE*, *PCSK9*, *LPL*, and *LDLR* associated with sphingomyelins (SM 16:0, SM 18:0, and SM-OH  
377 24:1) with atherosclerotic heart disease diagnosis codes (I21, I25), mirroring what is observed in rare  
378 familial forms of dyslipidaemia in which these sphingomyelins are elevated and the risk of ischemic  
379 heart disease is greatly increased<sup>25,26</sup>. These results provide further evidence that common variation  
380 at IEM genes can lead to clinical phenotypes and diseases that correspond to those that patients  
381 with rare mutations in those same genes are severely affected by. Further studies with detailed  
382 follow-up for specific outcomes may provide greater power and help clarify the medical  
383 consequences of genetic differences in metabolism caused by metabolite altering variants in the  
384 general population.



385

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

395

## 396 Discussion

397 This large-scale genome-wide meta-analysis has integrated genetic associations for 174  
 398 metabolites across different measurement platforms, an approach that has resulted in a three-fold  
 399 increase in our knowledge of genetic loci regulating levels of these metabolites. We assign likely

400 causal genes for many of the identified associations using a dual approach that combined automated  
401 database mining with manual curation.

402 Previous platform-specific genetic studies of blood metabolites have been substantially smaller  
403 in size due to being restricted to a single platform and/ or study<sup>2,3,4,5,6,7,8</sup>. We build on these earlier  
404 studies to identify and demonstrate enrichment of rare and low-frequency coding variants in  
405 enzyme and transporter genes with large effects and reveal the importance of non-linear  
406 associations at several loci.

407 Our results not only provide detailed insight into the genetic determinants of human  
408 metabolism but consider their relevance for disease aetiology and prediction. We explore both  
409 locus-specific and polygenic score effects and provide tangible examples with clear translational  
410 potential. We discovered a strong link between GLP2R, citrulline metabolism and type 2 diabetes,  
411 and demonstrate that the p.Asp470Asn variant underlying the citrulline and diabetes associations  
412 leads to significantly reduced recruitment of beta-arrestin to GLP2R in various cellular models.  
413 GLP2R is related to GLP1R, a target of glucose-lowering therapeutic agonists approved for type 2  
414 diabetes that mimic the other proglucagon gene encoded peptide glucagon-like peptide 1.

415 The finding that a standard deviation increase in serine levels via a genetic risk score is  
416 associated with 95% lower risk of MacTel shows that genetic differences resulting in very specific  
417 metabolic consequences can have profound effects on health. Our results suggest that inclusion of  
418 genetic scores for metabolite levels can improve identification of high risk individuals. Serine and  
419 glycine supplementation and/ or pharmacologic modulation of serine metabolism may help to  
420 reduce development or alter the prognosis of this rare, severe eye disease, specifically if targeted to  
421 people genetically with a genetic susceptibility to low serine levels. It is important to note, that  
422 randomized control trials are needed testing this hypothesis before any recommendations on  
423 supplementations could be made.

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

430



## 431 **Methods**

### 432 **Study design and participating cohorts**

433 We performed genome-wide meta-analyses of the levels of 174 metabolites from 7 biochemical  
434 categories (amino acids, biogenic amines, acylcarnitines, phosphatidylcholines,  
435 lysophosphatidylcholines, sphingomyelins, and sum of hexoses) captured by the Biocrates p180 kit  
436 measured using mass spectrometry (MS). As described in more detail below, a total of 174  
437 metabolites were successfully measured in up to 9,363 plasma samples from genotyped participants  
438 of the Fenland study<sup>27</sup>.

439 To maximise sample size and power, we meta-analysed genome-wide association (GWAS)  
440 results from the Fenland cohort with those run in the EPIC-Norfolk<sup>28</sup> and INTERVAL<sup>29</sup> studies, in  
441 which metabolites were profiled using MS (Metabolon Discovery HD4 platform) or protein nuclear  
442 magnetic resonance (<sup>1</sup>H-NMR) spectrometry<sup>30,31</sup> (**Supplementary Table 1**). Ten of the 174 Biocrates  
443 metabolites were covered across all platforms, while 38 were available on the Biocrates and  
444 Metabolon platforms and 126 were unique to Biocrates (**Fig. 1**). We integrated publicly available  
445 summary statistics from genome-wide meta-analyses of the same metabolites measured using MS  
446 (with Biocrates or Metabolon platforms) or <sup>1</sup>H-NMR spectrometry (**Supplementary Table 1**).  
447 Metabolites were matched across platforms by comparing metabolite names and biochemical  
448 formulas. Mapping across different Metabolon platforms was done based on retention time/index  
449 (RI), mass to charge ratio (m/z), and chromatographic data (including MS/MS spectral data).  
450 Scientists at Metabolon Inc. independently reviewed and confirmed metabolite matches.

451 A summary of the characteristics of participating cohorts is given in **Supplemental Table S1**. The  
452 Fenland study is a population-based cohort study of 12,435 participants without diabetes born  
453 between 1950 and 1975<sup>27</sup>. Participants were recruited from general practice surgeries in Cambridge,  
454 Ely and Wisbech (United Kingdom) and underwent detailed metabolic phenotyping and genome-  
455 wide genotyping. The European Prospective Investigation of Cancer (EPIC)-Norfolk study is a  
456 prospective cohort of 25,639 individuals aged between 40 and 79 and living in the county of Norfolk  
457 in the United Kingdom at recruitment<sup>28</sup>. INTERVAL is a randomised trial of approximately 50,000  
458 whole blood donors enrolled from all 25 static centres of NHS Blood and Transplant, aiming to  
459 determine whether donation intervals can be safely and acceptably decreased to optimise blood  
460 supply whilst maintaining the health of donors<sup>29</sup>.

### 461 **Metabolomics measurements**

462 The levels of 174 metabolites were measured in the Fenland study by the AbsoluteIDQ®  
463 Biocrates p180 Kit (Biocrates Life Sciences AG, Innsbruck, Austria) as reported elsewhere in detail

464 <sup>3231</sup>. We used a Waters Acquity ultra-performance liquid chromatography (UPLC; Waters Ltd,  
465 Manchester, UK) system coupled to an ABSciex 5500 Qtrap mass spectrometer (Sciex Ltd,  
466 Warrington, UK). Samples were derivatised and extracted using a Hamilton STAR liquid handling  
467 station (Hamilton Robotics Ltd, Birmingham, UK). Flow injection analysis coupled with tandem mass  
468 spectrometry (FIA-MS/MS) using multiple reaction monitoring (MRM) in positive mode ionisation  
469 was performed to measure the relative levels of acylcarnitines, phosphatidylcholines,  
470 lysophosphatidylcholines and sphingolipids. The level of hexose was measured in negative ionisation  
471 mode. Ultra-performance liquid chromatography coupled with tandem mass spectrometry using  
472 MRM was performed to measure the concentration of amino acids and biogenic amines. The  
473 chromatography consisted of a 5-minute gradient starting at 100% aqueous (0.2% Formic acid)  
474 increasing to 95% acetonitrile (0.2% Formic acid) over a Waters Acquity UPLC BEH C18 column (2.1 x  
475 50 mm, 1.7  $\mu$ m, with guard column). Isotopically labelled internal standards are integrated within  
476 the Biocrates p180 Kit for quantification. Data was processed in the Biocrates Met/*IDQ* software. Raw  
477 metabolite readings underwent extensive quality control procedures. Firstly, we excluded from any  
478 further analysis metabolites for which the number of measurements below the limit of  
479 quantification (LOQ) exceeded 5% of measured samples. Excluded metabolites were carnosine,  
480 dopamine, putrescine, asymmetric dimethyl arginine, dihydroxyphenylalanine, nitrotyrosine,  
481 spermine, sphingomyelins SM(22:3), SM(26:0), SM(26:1), SM(24:1-OH), phosphatidylcholine acyl-  
482 alky 44:4, and phosphatidylcholine diacyl C30:2. Secondly, in samples with detectable but not  
483 quantifiable peaks, we assigned random values between 0 and the run-specific LOQ of a given  
484 metabolite. Finally, we corrected for batch-effects with a “location-scale” approach, i.e. with  
485 normalization for mean and standard deviation of batches.

486 The levels of up to 38 metabolites were measured in EPIC-Norfolk and INTERVAL using the  
487 Metabolon HD4 Discovery platform. Measurements were carried out using MS/MS instruments. For  
488 these measurements, instrument variability, determined by calculating the median relative standard  
489 deviation, was of 6%. Data Extraction and Compound Identification: raw data was extracted, peak-  
490 identified and quality control-processed using Metabolon’s hardware and software. Compounds  
491 were identified by comparison to library entries of purified standards or recurrent unknown entities.  
492 Metabolon maintains a library, based upon authenticated standards, that contains the retention  
493 time/index (RI), mass to charge ratio (m/z), and chromatographic data (including MS/MS spectral  
494 data) of all molecules present in the library. Identifications were based on three criteria: retention  
495 index, accurate mass match to the library +/- 10 ppm, and the MS/MS forward and reverse scores  
496 between the experimental data and authentic standards. Metabolite Quantification and Data  
497 Normalization: Peaks were quantified using area-under-the-curve. A data normalization step was

498 performed to correct variation resulting from instrument inter-day tuning differences. Essentially,  
499 each compound was corrected in run-day blocks by registering the medians to equal one (1.00) and  
500 normalizing each data point proportionately (termed the “block correction”).

501 The levels of 10 metabolites were measured in the INTERVAL study using  $^1\text{H-NMR}$   
502 spectroscopy<sup>33</sup>. All samples were analysed using a high-throughput serum  $^1\text{H-NMR}$  metabolomics  
503 platform<sup>30</sup>. This provided information on 230 metabolites, including creatinine and several amino  
504 acids (alanine, glutamine, glycine, histidine, isoleucine, leucine, valine, phenylalanine, and tyrosine),  
505 data on which were used in the present study. Further details of the  $^1\text{H-NMR}$  spectroscopy,  
506 quantification data analysis and identification of the metabolites have been described previously  
507 <sup>30,34</sup>. Participants with >30% of metabolite measures missing and duplicated individuals were  
508 removed. Metabolite data more than 10 SD from the mean was also removed.

### 509 **GWAS and meta-analysis**

510 In Fenland and EPIC-Norfolk, metabolite levels were natural log-transformed, winsorised to  
511 five standard deviations and then standardised to a mean of 0 and a standard deviation of 1.  
512 Genotypes were measured using Affymetrix Axiom or Affymetrix SNP5.0 genotyping arrays. In brief,  
513 genotyping in Fenland was done in two waves including 1,500 (Affymetrix SNP5.0) and 9,369  
514 (Affymetrix Axiom) participants (Supplemental table S1) and imputation was done using IMPUTE2 to  
515 1000 Genomes Phase 1v3 (Affymetrix SNP5.0) or phase 3 (Affymetrix Axiom) reference panels  
516 (**Supplemental table S1**). In EPIC-Norfolk, 21,044 samples were forwarded to imputation using 1000  
517 Genomes Phase 3 (Oct. 2014) reference panels (**Supplemental table S1**). Imputed SNPs with  
518 imputation quality score less than 0.3 or minor allele account less than 2 were removed from the  
519 imputed dataset. Genome-wide association analyses were carried out using BOLT-LMM v2.2  
520 adjusting for age, sex, and study-specific covariates in mixed linear models. Alternatively (when the  
521 BOLT-LMM algorithm failed) analyses were performed using SNPTTEST v2.4.1 in linear regression  
522 models, additionally adjusting for the top 4 genetic ancestry principal components and excluding  
523 related individuals (defined by proportion identity by descent calculated in Plink<sup>35</sup>) > 0.1875). GWAS  
524 analyses in Fenland were performed within genotyping chip, and associations meta-analysed.

525 In INTERVAL, genotyping was conducted using the Affymetrix Axiom genotyping array.  
526 Standard quality control procedures were conducted prior to imputation. The data were phased and  
527 imputed to a joint 1000 Genomes Phase 3 (May 2013)-UK10K reference imputation panel. After QC,  
528 a total of 40,905 participants remained with data obtained by  $^1\text{H-NMR}$  spectroscopy. For variants  
529 with a MAF of >1% and imputed variants with an info score of >0.4 a univariate GWAS for each of  
530 the ten metabolic measures was conducted, after adjustment for technical and seasonal effects,

531 including age, sex, and the first 10 principal components, and rank-based inverse normal  
532 transformation. The association analyses were performed using BOLT-LMM v2.2 and R. Data based  
533 on the Metabolon HD4 platform was available for 8,455 participants and SNPTEST v2.5.1 was used to  
534 test for associations with metabolite levels adjusting for age, sex, and the first five genetic principle  
535 components. In SNPTEST analyses, related individuals (proportion identity by descent > 0.1875) were  
536 excluded.

537 For all GWAS analysis within Fenland, EPIC-Norfolk and INTERVAL, variants with Hardy-  
538 Weinberg equilibrium  $p < 1 \times 10^{-6}$  and associations with absolute value of effect size >5 or standard  
539 error (SE) >10 or <0 were excluded; insertions and deletions were excluded.

540 For each metabolite, we performed a meta-analysis of z-scores (betas divided by standard  
541 errors) as a measure of association, signals and loci (see below), adjusting for genomic control using  
542 METAL software. Heterogeneity between studies for each association was estimated by Cochran's Q-  
543 test. For each metabolite, we also performed a meta-analysis of beta and standard errors for the  
544 subset of studies (Fenland and, when available, EPIC-Norfolk and/or INTERVAL) where we had access  
545 to individual level data and standardised phenotype preparation to estimate effect sizes. Quality  
546 filters implemented after meta-analysis included exclusion of SNPs not captured by at least 50% of  
547 the participating studies and 50% of the maximum sample size for that metabolite and variants with  
548 a minor allele frequency below 0.5% percent. As a result, meta-analyses assessed the associations of  
549 up to 13.1 million common or low-frequency autosomal SNPs. Chromosome and base pair positions  
550 are determined referring to GRCh37 annotation. To define associations between genetic variants  
551 and metabolites, we corrected the conventional threshold of genome wide significance for 102 tests  
552 (i.e.  $p < 4.9 \times 10^{-10}$ ), corresponding to the number of principal components explaining 95% of the  
553 variance of the 174 metabolites, as previously described<sup>36</sup>.

#### 554 **Signal selection**

555 For each metabolite, we ranked associated SNPs ( $p < 4.9 \times 10^{-10}$ ) by z-score to select trait-sentinel  
556 SNPs and defined an “association” region as the region extending 1 Mb to each side of the trait-  
557 sentinel SNP. During forward selection of trait-sentinel SNPs and loci for each trait, adjacent and  
558 partially overlapping association regions were merged by extending region boundaries to a further 1  
559 Mb. After defining trait-sentinel SNPs and association regions we defined overall lead-sentinel SNP  
560 and loci for any metabolite using a similar approach. Trait-sentinel SNPs were sorted by z-score for  
561 the forward selection of lead-sentinel SNPs and a “locus” was defined as the region extending 1 Mb  
562 each side of the lead-sentinel SNP. Regions larger than 2 Mb defined in the trait-sentinel association  
563 region definition were carried over in the definition of lead-sentinel SNP loci. As a result, all lead-

564 sentinel SNPs were >1Mb apart from each other and had very low or no linkage disequilibrium ( $R^2 <$   
565 0.05).

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

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

578 We tested at each locus for conditional independent variants using exact stepwise conditional  
579 analysis in the largest Fenland sample (n = 8,714) using SNPTTEST v2.5 with the same baseline  
580 adjustment as in the discovery approach. To refine signals at those loci we used a more recent  
581 imputation for this analysis based on the HRC v1 reference panel and additional SNPs imputed using  
582 UK10K and 1000G phase 3. We defined secondary signals as those with a conditional p-value  $< 5 \times 10^{-8}$   
583 <sup>8</sup>. To avoid problems with collinearity we tested after each round if inclusion of a new variant  
584 changed associations of all previous variants with the outcome using a joint model. If this model  
585 indicated that one or more of the previously selected variants dropped below the applied  
586 significance threshold we stopped the procedure, otherwise we repeated this procedure until no  
587 further variant met the significance threshold in conditional models. We considered only locus-  
588 metabolite associations meeting the GWAS-threshold for significance in the Fenland analysis  
589 (n=228).

## 590 **Colocalisation analyses**

591 For specified loci of interest statistical colocalisation<sup>37</sup> was performed to obtain posterior  
592 probabilities (PP) of: H0 – no signal; H1 – signal unique to the metabolite; H2 – signal unique to the  
593 trait; H3 – two distinct causal variants in the same locus and H4 – presence of a shared causal variant  
594 between a metabolite and a given trait. PP above 80% were considered highly likely. We used p-  
595 values and MAFs obtained from the summary statistics with default priors to perform colocalisation.

## 596 **Hypothesis-free (genetic) assignment of causal genes**

597 To assign likely causal genes to lead SNPs at each locus we generated a scoring system. We  
598 identified the nearest gene for each variant by querying HaploReg<sup>38</sup>. Next we integrated expression  
599 quantitative trait loci (eQTL) studies (GTEx v6p) to identify genes whose expression levels are  
600 associated with metabolite levels using TWAS/FUSION (Transcriptome-wide association study /  
601 Functional summary-based imputation)<sup>39</sup>. In doing so, we assigned to each variant-metabolite  
602 association one or more associated genes using the variant as common anchor. We further assigned  
603 higher impact for a causal gene if either the metabolite variant itself or a proxy in high linkage  
604 disequilibrium ( $R^2 > 0.8$ ) was a missense variant for a known gene again using the HaploReg database  
605 to obtain relevant information. Based on those three criteria we ranked all possible candidate genes  
606 and kept those with the highest score as putative causal gene.

### 607 **Knowledge-based (biological) assignment of causal genes**

608 Metabolite traits are unique among genetically evaluated phenotypes in that the functional  
609 characterization of the relevant genes has often already been carried out using classic biochemical  
610 techniques. The objective for the knowledge-based assignment strategy was to find the  
611 experimental evidence that has previously linked one of the genes proximal to the GWAS lead  
612 variant to the relevant metabolite. For many loci and metabolites this ‘retrospective’ analysis has  
613 already been carried out<sup>940</sup>. For these cases, previous causal gene assignments were generally  
614 adopted. For novel loci, we employed a dual strategy that combined automated database mining  
615 with manual curation. In the automated phase, seven approaches were employed to identify  
616 potential causal genes among the 20 protein-coding genes closest to each lead variant, as described  
617 in detail below, using the shortest distance determined from the lead SNP to each gene’s  
618 transcription start site (TSS) or transcription end site (TES), with a distance value of 0 assigned if the  
619 SNP fell between the TSS and TES.

620 These 7 approaches were as follows:

- 621 1) HMDB metabolite names<sup>41</sup> were compared to each entrez gene name;
- 622 2) Metabolite names were compared to the name and synonyms of the protein encoded by each  
623 gene<sup>42</sup>
- 624 3) HMDB metabolite names and their parent terms (class) were compared to the names for the  
625 protein encoded by each gene (UniProt).
- 626 4) Metabolite names were compared to rare diseases linked to each gene in OMIM<sup>43</sup> after  
627 removing the following non-specific substrings from disease names: uria, emia, deficiency, disease,

628 transient, neonatal, hyper, hypo, defect, syndrome, familial, autosomal, dominant, recessive, benign,  
629 infantile, hereditary, congenital, early-onset, idiopathic;

630 5) HMDB metabolite names and their parent terms were compared to all GO biological processes  
631 associated with each gene after removing the following non-specific substrings from the name of the  
632 biological process: metabolic process, metabolism, catabolic process, response to, positive  
633 regulation of, negative regulation of, regulation of. For this analysis only gene sets containing fewer  
634 than 500 gene annotations were retained.

635 6) KEGG maps<sup>44</sup> containing the metabolite as defined in HMDB were compared to KEGG maps  
636 containing each gene, as defined in KEGG. For this analysis the large “metabolic process” map was  
637 omitted.

638 7) Each proximal gene was compared to the list of known interacting genes as defined in HMDB.  
639 For each text-matching based approach, a fuzzy text similarity metric (pair coefficient) as encoded in  
640 the ruby gem “fuzzy\_match” was used with a score greater than 0.5 considered as a match.

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

649 For each manually selected causal gene, the earliest experimental evidence linking the gene  
650 (preferably the human gene) to the metabolite was identified. The median publication year for the  
651 identified experimental evidence was 2000.

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

653 We examined whether the set of independent lead metabolite associated variants (N=168)  
654 were enriched for associations with type 2 diabetes. We plotted observed versus expected  $-\log_{10}(p$   
655 values) for the 168 lead variants in a QQ-plot, using association statistics from a type 2 diabetes  
656 meta-analysis including 80,983 cases and 842,909 non-cases from the DIAMANTE study<sup>45</sup> (55,005  
657 T2D cases, 400,308 non-cases), UK Biobank<sup>46</sup> (24,758 T2D cases, 424,575 non-cases, application  
658 number 44448) and the EPIC-Norfolk study (additional T2D cases not included in DIAMANTE study:  
659 1,220 T2D cases and 18,026 non-cases). This QQ-plot was compared to those for 1000 sets of

660 variants, where variants in each set were matched to the index metabolite variants in terms of MAF,  
661 the number of variants in LD ( $R^2 > 0.5$ ), gene density and distance to nearest gene (for all parameters  
662 +/- 50% of the index variant value), but otherwise randomly sampled from across the autosome  
663 excluding the HLA region. MAF and LD parameters for individual variants were determined from the  
664 EPIC-Norfolk study (using the combined HRC, UK10K and 1000G imputation as previously described)  
665 and gene information was derived from GENCODE v19 annotation<sup>48</sup>. A one-tailed Wilcoxon rank  
666 sum test was used to compare the distribution of association  $-\log_{10}$  p-values for the metabolite  
667 associated variants with that for the randomly sampled, matched, variants.

### 668 **Functional characterisation of D470N mutant GLP2R**

669 To investigate the functional differences between wild-type (WT) GLP2R and the D470N  
670 mutant GLP2R we generated D470N GLP2R mutant constructs using site-directed mutagenesis and  
671 characterised canonical GLP2R signalling pathways via cAMP as well as alternative signalling  
672 pathways via  $\beta$ -arrestin and P-ERK.

#### 673 *Generation of D470N GLP2R mutant expressing constructs*

674 Human GLP2R cDNA within the pcDNA3.1+ vector was purchased, and Gibson cloning was  
675 completed to insert an internal ribosome entry site (IRES) and venus gene downstream of the GLP2R  
676 sequence. Following this, QuikChange Lightning site directed mutagenesis was used to perform a  
677 single base change from GAC (encoding aspartic acid) to AAC (encoding asparagine) at amino acid  
678 position 470 (**Supplemental Figure 2A-B**). Successful mutagenesis was confirmed by DNA Sanger  
679 sequencing (**Supplemental Figure 2C**), and the successful products were scaled up for use in  
680 functional assays. The WT and mutant GLP2R constructs within the pcDNA3.1+ vector were used to  
681 assess signalling by cAMP and P-ERK. To determine  $\beta$ -arrestin recruitment using NanoBiT®  
682 technology, an alternative vector was required for lower expression of GLP2R, and fusion of GLP2R  
683 to the Large BiT subunit of NanoBiT®. For this, GLP2R was cloned into the pBiT1.1\_C[TK/LgBiT] vector  
684 using restriction cloning and ligation. DNA Sanger sequencing was then used for confirmation of  
685 successful cloning.

#### 686 *Comparison of WT and D470N GLP2R signalling via cAMP*

687 After generation of WT and D470N GLP2R containing constructs, these were used to assess  
688 differences in WT and mutant GLP2R signalling. The initial signalling pathway to be assessed was G $\alpha$ s  
689 signalling via cAMP. CHO K1 cells were transiently transfected with WT or mutant GLP2R constructs,  
690 then after 16-24 hours were treated with a dose response of GLP-2. cAMP levels were measured  
691 following 30 minutes of GLP-2 treatment, in an end-point lysis HitHunter® cAMP assay. The presence  
692 of IRES-Venus within the GLP2R expressing vectors allowed transfection efficiency to be determined



693 for each construct. Transfection efficiency was approximately 60-70%, with no differences between  
694 the WT and mutant constructs. Comparison of the GLP-2 dose-response in WT and mutant GLP2R  
695 expressing cells revealed no significant differences in signalling, with an almost overlapping dose  
696 response curve (Figure 5d).

#### 697 *Comparison of $\beta$ -arrestin recruitment to the WT and D470N GLP2R*

698 Both  $\beta$ -arrestin 1 and  $\beta$ -arrestin 2 recruitment were assessed using a Nano-Glo<sup>®</sup> live cell  
699 assay in transiently transfected HEK293 cells. Briefly, the recruitment of  $\beta$ -arrestin to GLP2R brings  
700 the large and small BiT subunit of NanoBiT<sup>®</sup> together, resulting in increased luciferase activity. The  
701 top concentrations from the GLP-2 dose response in the cAMP assay (1–100 nmol/l GLP-2) were  
702 chosen for stimulation of the GLP2R and observation of  $\beta$ -arrestin recruitment. Both  $\beta$ -arrestin 1 and  
703  $\beta$ -arrestin 2 were recruited to the WT GLP2R upon GLP-2 stimulation, in a dose-dependent manner  
704 (Supplemental Figure 3a, c). The maximal luciferase activity for both  $\beta$ -arrestin 1 and  $\beta$ -arrestin 2  
705 recruitment to the mutant GLP2R was significantly decreased when compared to the WT GLP2R,  
706 indicating the extent of  $\beta$ -arrestin recruitment was markedly decreased (Supplemental Figure 3b, d).  
707 The example traces indicate that neither  $\beta$ -arrestin 1 or  $\beta$ -arrestin 2 were recruited to the mutant  
708 GLP2R upon stimulation with 1 nmol/l GLP-2, however the same concentration of GLP-2 induced  $\beta$ -  
709 arrestin recruitment to the WT GLP2R. Overall there was a significant decrease in  $\beta$ -arrestin 1 and  $\beta$ -  
710 arrestin 2 recruitment to the D470N GLP2R mutant (**Figure 5e-f**).

#### 711 **Genetic score and Mendelian randomization analysis for macular telangiectasia type 2**

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

718 To test for causality between circulating levels of glycine and serine for MacTel we  
719 performed two types of Mendelian randomization (MR) analysis. In a two-sample univariable MR<sup>49</sup>  
720 we tested for an individual effect of serine (n=4 SNPs) or glycine (n=15 SNPs) on the risk of MacTel  
721 using independent non-pleiotropic (i.e. the variant in *GCKR*) genome-wide SNPs as instruments. To  
722 this end, we used the inverse variance weighted method to pool SNP ratio estimates using random  
723 effects as implemented in the R package *MendelianRandomization*. SNP effects on the risk for  
724 MacTel were obtained from<sup>21</sup>. To disentangle the individual effect of those two highly correlated  
725 metabolites at the same time we used a multivariable MR model<sup>50</sup> including all SNPs related to

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

732 To estimate a potential clinical usefulness of the identified variants we constructed two  
733 GRS's for MacTel using a) sex, the first genetic principal component, and the SNPs rs73171800 and  
734 rs9820286 which were identified by the MacTel GWAS study<sup>21</sup> but not found to be related to either  
735 glycine or serine in our study and b) all the previous but additionally including genetically predicted  
736 serine and glycine at individual levels, via PGS, to the model. An interaction between serine and sex  
737 at birth was included to reflect the interaction between SNP rs715 and sex as previously identified<sup>21</sup>.  
738 To assess the predictive ability of both models, receiver operating characteristic curves were  
739 computed based on prediction values in 1,733 controls and 476 MacTel cases.

#### 740 **Identification of genes related to inborn errors of metabolism**

741 Biologically or genetically assigned candidate genes were annotated for IEM association  
742 using the Orphanet database<sup>43</sup>. Using a binomial two-tailed test, enrichment of metabolic loci was  
743 assessed by comparing the annotated list with the full list of 784 IEM genes in Orphanet against a  
744 backdrop of 19,817 protein-coding genes<sup>51</sup>. IEM-annotated loci for which the associated metabolite  
745 matched or was closely biochemically related to the IEM corresponding metabolite(s) based on  
746 IEMBase<sup>52</sup> were considered further for analysis.

747 We hypothesised that IEM-annotated loci with metabolite-specific consequences could also  
748 have phenotypic consequences similar to the IEM. To test this, we first obtained terms describing  
749 each IEM and translated them into IEM-related ICD-10 codes using the Human Phenotype Ontology  
750 and previously-generated mappings<sup>53,54</sup>. We obtained association statistics from the 85 IEM SNPs for  
751 phenotypic associations with corresponding ICD-codes among UK Biobank restricting to diseases  
752 with at least 500 cases (N=93, **Fig. 7B**, <http://www.nealelab.is/uk-biobank>). We tested locus-disease  
753 pairs meeting statistical significance (controlling the false discovery rate at 5% to account for  
754 multiple testing) for a common genetic signal with the corresponding locus-metabolite association  
755 using statistical colocalisation.

#### 756 **Acknowledgement/Funding**

757 M.P. was supported by a fellowship from the German Research Foundation (DFG PI 1446/2-1). C.O.  
758 was founded by an early career fellowship at Homerton College, University of Cambridge. L. B. L. W.

759 acknowledges funding by the Wellcome Trust (WT083442AIA). J.G. was supported by grants from  
760 the Medical Research Council (MC\_UP\_A090\_1006, MC\_PC\_13030, MR/P011705/1 and  
761 MR/P01836X/1). Work in the Reimann/Gribble laboratories was supported by the Wellcome Trust  
762 (106262/Z/14/Z and 106263/Z/14/Z), UK Medical Research Council (MRC\_MC\_UU\_12012/3) and  
763 PhD funding for EKB from MedImmune/AstraZeneca. Praveen Surendran is supported by a  
764 Rutherford Fund Fellowship from the Medical Research Council grant MR/S003746/1. A. W. is  
765 supported by a BHF-Turing Cardiovascular Data Science Award and by the EC-Innovative Medicines  
766 Initiative (BigData@Heart). J.D. is funded by the National Institute for Health Research [Senior  
767 Investigator Award] [\*]. The EPIC-Norfolk study (<https://doi.org/10.22025/2019.10.105.00004>) has  
768 received funding from the Medical Research Council (MR/N003284/1 and MC-UU\_12015/1) and  
769 Cancer Research UK (C864/A14136). The genetics work in the EPIC-Norfolk study was funded by the  
770 Medical Research Council (MC\_PC\_13048). Metabolite measurements in the EPIC-Norfolk study  
771 were supported by the MRC Cambridge Initiative in Metabolic Science (MR/L00002/1) and the  
772 Innovative Medicines Initiative Joint Undertaking under EMIF grant agreement no. 115372. We are  
773 grateful to all the participants who have been part of the project and to the many members of the  
774 study teams at the University of Cambridge who have enabled this research. The Fenland Study is  
775 supported by the UK Medical Research Council (MC\_UU\_12015/1 and MC\_PC\_13046). Participants  
776 in the INTERVAL randomised controlled trial were recruited with the active collaboration of NHS  
777 Blood and Transplant England ([www.nhsbt.nhs.uk](http://www.nhsbt.nhs.uk)), which has supported field work and other  
778 elements of the trial. DNA extraction and genotyping was co-funded by the National Institute for  
779 Health Research (NIHR), the NIHR BioResource (<http://bioresource.nihr.ac.uk>) and the NIHR  
780 [Cambridge Biomedical Research Centre at the Cambridge University Hospitals NHS Foundation  
781 Trust] [\*]. Nightingale Health NMR assays were funded by the European Commission Framework  
782 Programme 7 (HEALTH-F2-2012-279233). Metabolite Metabolomics assays were funded by the NIHR  
783 BioResource and the National Institute for Health Research [Cambridge Biomedical Research Centre  
784 at the Cambridge University Hospitals NHS Foundation Trust] [\*]. The academic coordinating centre  
785 for INTERVAL was supported by core funding from: NIHR Blood and Transplant Research Unit in  
786 Donor Health and Genomics (NIHR BTRU-2014-10024), UK Medical Research Council  
787 (MR/L003120/1), British Heart Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946) and the  
788 NIHR [Cambridge Biomedical Research Centre at the Cambridge University Hospitals NHS Foundation  
789 Trust] [\*]. The academic coordinating centre would like to thank blood donor centre staff and blood  
790 donors for participating in the INTERVAL trial.

791 This work was supported by Health Data Research UK, which is funded by the UK Medical Research  
792 Council, Engineering and Physical Sciences Research Council, Economic and Social Research Council,

793 Department of Health and Social Care (England), Chief Scientist Office of the Scottish Government  
794 Health and Social Care Directorates, Health and Social Care Research and Development Division  
795 (Welsh Government), Public Health Agency (Northern Ireland), British Heart Foundation and  
796 Wellcome.

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

799 UK Biobank: This research has been conducted using the UK Biobank resource under

800 Application Number 44448.

#### 801 **Author Contribution**

802 Concept and design: L.A.L. and C.L.

803 Generation, acquisition, analysis and/or interpretation of data: all authors.

804 Drafting of the manuscript: L.A.L., M.P., and C.L.

805 Critical review of the manuscript for important intellectual content and approval of the final version  
806 of the manuscript: all authors.

#### 807 **Competing Interests statement**

808 A.S.B. has received grants from AstraZeneca, Biogen, Bioverativ, Merck, Novartis, and Sanofi. J. D.  
809 sits on the International Cardiovascular and Metabolic Advisory Board for Novartis (since 2010), the  
810 Steering Committee of UK Biobank (since 2011), the MRC International Advisory Group (ING)  
811 member, London (since 2013), the MRC High Throughput Science 'Omics Panel Member, London  
812 (since 2013), the Scientific Advisory Committee for Sanofi (since 2013), the International  
813 Cardiovascular and Metabolism Research and Development Portfolio Committee for Novartis and  
814 the Astra Zeneca Genomics Advisory Board (2018).

815 REFERENCES

- 816 1. Wishart, D. S. Metabolomics for investigating physiological and pathophysiological processes.  
817 *Physiol. Rev.* **99**, 1819–1875 (2019).
- 818 2. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.* **46**,  
819 543–550 (2014).
- 820 3. Draisma, H. H. M. *et al.* Genome-wide association study identifies novel genetic variants  
821 contributing to variation in blood metabolite levels. *Nat. Commun.* **6**, 7208 (2015).
- 822 4. Kettunen, J. *et al.* Genome-wide study for circulating metabolites identifies 62 loci and  
823 reveals novel systemic effects of LPA. *Nat. Commun.* **7**, 11122 (2016).
- 824 5. Illig, T. *et al.* A genome-wide perspective of genetic variation in ... [Nat Genet. 2010] -  
825 PubMed result. *Nat. Genet.* **42**, 137–41 (2010).
- 826 6. Suhre, K. *et al.* Human metabolic individuality in biomedical and pharmaceutical research.  
827 *Nature* **477**, 54–62 (2011).
- 828 7. Rhee, E. P. *et al.* A genome-wide association study of the human metabolome in a  
829 community-based cohort. *Cell Metab.* **18**, 130–43 (2013).
- 830 8. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants associated with  
831 human blood metabolites. *Nat. Genet.* **49**, 568–578 (2017).
- 832 9. Shin, S. Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.* **46**,  
833 543–550 (2014).
- 834 10. Draisma, H. H. M. *et al.* Genome-wide association study identifies novel genetic variants  
835 contributing to variation in blood metabolite levels. *Nat. Commun.* **6**, 7208 (2015).
- 836 11. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to  
837 Common Complex Disease. *Cell* **167**, 1415–1429.e19 (2016).
- 838 12. Bansal, N. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79 (2018).
- 839 13. Learn, D. B., Fried, V. A. & Thomas, E. L. Taurine and hypotaurine content of human  
840 leukocytes. *J. Leukoc. Biol.* **48**, 174–182 (1990).
- 841 14. Yet, I. *et al.* Genetic Influences on Metabolite Levels: A Comparison across Metabolomic  
842 Platforms. *PLoS One* **11**, e0153672 (2016).
- 843 15. Lahiri, S. *et al.* Kinetic characterization of mammalian ceramide synthases: Determination of  
844 Km values towards sphinganine. *FEBS Lett.* **581**, 5289–5294 (2007).
- 845 16. Horowitz, B. *et al.* Asparagine synthetase activity of mouse leukemias. *Science (80-. )*. **160**,  
846 533–535 (1968).
- 847 17. Babu, E. *et al.* Identification of a Novel System L Amino Acid Transporter Structurally Distinct  
848 from Heterodimeric Amino Acid Transporters. *J. Biol. Chem.* **278**, 43838–43845 (2003).
- 849 18. Scott, R. A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in  
850 Europeans. *Diabetes* **66**, 2888–2902 (2017).
- 851 19. Fragkos, K. C. & Forbes, A. Citrulline as a marker of intestinal function and absorption in  
852 clinical settings: A systematic review and meta-analysis. *United Eur. Gastroenterol. J.* **6**, 181–  
853 191 (2018).
- 854 20. Estall, J. L., Koehler, J. A., Yusta, B. & Drucker, D. J. The glucagon-like peptide-2 receptor C  
855 terminus modulates  $\beta$ -arrestin-2 association but is dispensable for ligand-induced  
856 desensitization, endocytosis, and G-protein-dependent effector activation. *J. Biol. Chem.* **280**,  
857 22124–22134 (2005).
- 858 21. Scerri, T. S. *et al.* Genome-wide analyses identify common variants associated with macular  
859 telangiectasia type 2. *Nat. Genet.* **49**, 559–567 (2017).
- 860 22. Gantner, M. L. *et al.* Serine and lipid metabolism in macular disease and peripheral  
861 neuropathy. *N. Engl. J. Med.* **381**, 1422–1433 (2019).
- 862 23. Garrod, A. E. The incidence of alkaptonuria: a study in chemical individuality. *Lancet* **160**,  
863 1616–1620 (1902).
- 864 24. Rath, A. *et al.* Representation of rare diseases in health information systems: The orphanet  
865 approach to serve a wide range of end users. *Hum. Mutat.* **33**, 803–808 (2012).

- 866 25. Stübiger, G. *et al.* Targeted profiling of atherogenic phospholipids in human plasma and  
867 lipoproteins of hyperlipidemic patients using MALDI-QIT-TOF-MS/MS. *Atherosclerosis* **224**,  
868 177–186 (2012).
- 869 26. van der Graaf, A., Kastelein, J. J. P. & Wiegman, A. Heterozygous familial  
870 hypercholesterolaemia in childhood: Cardiovascular risk prevention. *J. Inherit. Metab. Dis.* **32**,  
871 699 (2009).
- 872 27. Lindsay, T. *et al.* Descriptive epidemiology of physical activity energy expenditure in UK adults  
873 (The Fenland study). *Int. J. Behav. Nutr. Phys. Act.* **16**, 126 (2019).
- 874 28. Day, N. *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European  
875 Prospective Investigation of Cancer. *Br. J. Cancer* **80 Suppl 1**, 95–103 (1999).
- 876 29. Moore, C. *et al.* The INTERVAL trial to determine whether intervals between blood donations  
877 can be safely and acceptably decreased to optimise blood supply: Study protocol for a  
878 randomised controlled trial. *Trials* **15**, (2014).
- 879 30. Soininen, P. *et al.* High-throughput serum NMR metabolomics for cost-effective holistic  
880 studies on systemic metabolism. *Analyst* **134**, 1781–5 (2009).
- 881 31. Wittemans, L. B. L. *et al.* Assessing the causal association of glycine with risk of cardio-  
882 metabolic diseases. *Nat. Commun.* **10**, (2019).
- 883 32. Lotta, L. A. *et al.* Genetic Predisposition to an Impaired Metabolism of the Branched-Chain  
884 Amino Acids and Risk of Type 2 Diabetes: A Mendelian Randomisation Analysis. *PLoS Med.*  
885 (2016). doi:10.1371/journal.pmed.1002179
- 886 33. Di Angelantonio, E. *et al.* Efficiency and safety of varying the frequency of whole blood  
887 donation (INTERVAL): a randomised trial of 45 000 donors. *Lancet* **390**, 2360–2371 (2017).
- 888 34. Inouye, M. *et al.* Metabonomic, transcriptomic, and genomic variation of a population cohort.  
889 *Mol. Syst. Biol.* **6**, 441 (2010).
- 890 35. Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-based linkage  
891 analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 892 36. Li, J. & Ji, L. Adjusting multiple testing in multilocus analyses using the eigenvalues of a  
893 correlation matrix. *Heredity (Edinb.)* **95**, 221–227 (2005).
- 894 37. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association  
895 studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- 896 38. Ward, L. D. & Kellis, M. HaploReg: A resource for exploring chromatin states, conservation,  
897 and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.*  
898 **40**, (2012).
- 899 39. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies.  
900 *Nat. Genet.* **48**, 245–252 (2016).
- 901 40. Stacey, D. *et al.* ProGeM: A framework for the prioritization of candidate causal genes at  
902 molecular quantitative trait loci. *Nucleic Acids Res.* **47**, (2019).
- 903 41. Wishart, D. S. *et al.* HMDB 4.0: The human metabolome database for 2018. *Nucleic Acids Res.*  
904 **46**, D608–D617 (2018).
- 905 42. Bateman, A. *et al.* UniProt: The universal protein knowledgebase. *Nucleic Acids Res.* **45**,  
906 D158–D169 (2017).
- 907 43. Rath, A. *et al.* Representation of rare diseases in health information systems: The orphanet  
908 approach to serve a wide range of end users. *Hum. Mutat.* **33**, 803–808 (2012).
- 909 44. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: New perspectives  
910 on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361 (2017).
- 911 45. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-  
912 density imputation and islet-specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
- 913 46. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide  
914 Range of Complex Diseases of Middle and Old Age. *PLoS Med.* **12**, (2015).
- 915 47. Collins, R. What makes UK Biobank special? *The Lancet* (2012). doi:10.1016/S0140-  
916 6736(12)60404-8

- 917 48. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.  
918 *Nucleic Acids Res.* **47**, D766–D773 (2019).
- 919 49. Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis with  
920 multiple genetic variants using summarized data. *Genet. Epidemiol.* **37**, 658–665 (2013).
- 921 50. Burgess, S. & Thompson, S. G. Multivariable Mendelian randomization: The use of pleiotropic  
922 genetic variants to estimate causal effects. *Am. J. Epidemiol.* **181**, 251–260 (2015).
- 923 51. Harrow, J. *et al.* GENCODE: The reference human genome annotation for the ENCODE  
924 project. *Genome Res.* **22**, 1760–1774 (2012).
- 925 52. Lee, J. J. Y., Wasserman, W. W., Hoffmann, G. F., Van Karnebeek, C. D. M. & Blau, N.  
926 Knowledge base and mini-expert platform for the diagnosis of inborn errors of metabolism.  
927 *Genet. Med.* **20**, 151–158 (2018).
- 928 53. Köhler, S. *et al.* The human phenotype ontology in 2017. *Nucleic Acids Res.* **45**, D865–D876  
929 (2017).
- 930 54. Wu, P. *et al.* Mapping ICD-10 and ICD-10-CM codes to phecodes: Workflow development and  
931 initial evaluation. *J. Med. Internet Res.* **21**, (2019).
- 932
- 933

934 **TABLES**

935

936 **Table 1** Results from Mendelian randomisation (MR) analysis between metabolite levels and risk of  
937 macular telangiectasia type 2.

938

<b>Metabolite</b>	<b>Univariable MR</b>	<b>Multivariable MR</b>
<i>Serine</i> (4 SNPs)		
Odds ratio per SD increase	0.06 (0.03; 0.13)	0.10 (0.05; 0.21)
p-value	$9.45 \times 10^{-12}$	$2.95 \times 10^{-9}$
<i>Glycine</i> (15 SNPs)		
Odds ratio per SD increase	0.17 (0.08; 0.37)	0.50 (0.29; 0.87)
p-value	$9.99 \times 10^{-6}$	$1.35 \times 10^{-2}$

939 MR estimates are based on the inverse variance-weighted method using random effects to pool

940 estimates. All single nucleotide polymorphisms (SNPs) significantly associated with either serine or

941 glycine have been included in multivariable MR analysis. SD = standard deviation

942

943



944 **FIGURE LEGENDS**

945

946 **Figure 1A** Overlap among the 174 plasma metabolites investigated in the present study across three  
947 different techniques: Biocrates p180 Kit, Metabolon HD4, and proton nuclear magnetic resonance  
948 spectroscopy (NMR). **B** A three-dimensional Manhattan plot displaying chromosomal position (x-  
949 axis) of significant associations ( $p < 4.9 \times 10^{-10}$ , z-axis) across all metabolites (y-axis). Colours indicate  
950 metabolite groups. **C** A top view of the 3D-Manhattan plot. Dots indicate significantly associated loci.  
951 Colours indicate novelty of metabolite – locus associations. Loci with indication for pleiotropy have  
952 been annotated. Barplot on the right-hand side indicates samples sizes for each metabolite  
953 comparing the present samples sizes with the previous largest genome-wide association study for a  
954 specific metabolite. PCs = phosphatidylcholines, LysoPCs = lysophosphatidylcholines

955

956 **Figure 2A** Distribution of pleiotropy, i.e. number of associated metabolites, among loci identified in  
957 the present study. **B** Distribution of polygenicity of metabolites, i.e. number of identified loci for  
958 each metabolite under investigation. **C** Scatterplot comparing the estimated heritability of each  
959 metabolite against the number of associated loci. Size of the dots indicates samples sizes. **D**  
960 Heritability estimates for single metabolites. Colours indicate the proportion of heritability  
961 attributed to single nucleotide polymorphisms (SNPs) with large effect sizes ( $\beta > 0.25$  per allele). **E** –  
962 **M** SNP – metabolite association with indication of non-additive effects. Beta is an estimate from the  
963 departure of linearity. **N** Barplot showing the increase in heritability and explained phenotypic  
964 variance for each SNP – metabolite pair when including non-additive effects.

965

966 **Figure 3A** Scatterplot comparing the minor allele frequencies (MAF) of associated variants with  
967 effect estimates from linear regression models (N loci=499). Colours indicate possible functional  
968 consequences of each variant: maroon – nonsynonymous variant; blue – in strong LD ( $r^2 > 0.8$ ) with a  
969 nonsynonymous variant and grey otherwise. **B/C** Distribution of effect sizes (B) and allele  
970 frequencies (C) based on the type of single nucleotide polymorphism (SNP) (0 – non-coding or  
971 synonymous, 1 – in strong LD with nonsynonymous, 2 - nonsynonymous). **D** Distribution of  
972 functional annotations of metabolite associated variants (red), trait-associated variants (blue –  
973 continuous, purple – diseases) obtained from the GWAS catalogue, and all SNPs included in the  
974 present genome-wide association studies. The inlet for exonic variants distinguishes between  
975 synonymous (syn) and nonsynonymous variants (nsyn).

976 **Figure 4A** Comparison between the genetically prioritized versus biological knowledge-based  
977 approaches used in the present study to assign candidate genes to metabolite associated single  
978 nucleotide polymorphisms. The Venn-diagram displays the overlap between both approaches. **B**  
979 Enrichment of genetically-prioritized genes among biologically plausible or genes linked to inborn  
980 errors of metabolism (IEM). **C** Proportion of genetically-prioritized genes encoding for either  
981 enzymes or transporters.

982

983 **Figure 5A** Enrichment of associations with type 2 diabetes (T2D: 80,983 cases, 842,909 controls)  
984 among metabolite-associated SNPs. Blue dots indicate metabolite-SNPs and grey dots indicate a  
985 random selection of matched control SNPs. **B** Opposing  $-\log_{10}(p\text{-values})$  from the genome-wide  
986 association study of plasma citrulline with those from the T2D-GWAS for SNPs located around  
987 *GLP2R*. The legend in the upper left gives the posterior probabilities (PP) from statistical  
988 colocalisation analysis.  $H_4 = PP$  for the hypothesis of a shared causal variant. **C** Individual association  
989 summary statistics for all citrulline associated SNPs (coded by the citrulline increasing allele) for T2D.  
990 **D** GLP-2 dose response curves in cAMP assay for GLP2R wild-type and mutant receptors. The dose  
991 response curves of cAMP stimulation by GLP-2 in CHO K1 cells transiently transfected with either  
992 GLP2R wild-type or mutant constructs. Data were normalised to the wild-type maximal and minimal  
993 response, with 100% being GLP-2 maximal stimulation of the wild-type GLP2R, and 0% being wild-  
994 type GLP2R cells with buffer only. Mean  $\pm$  standard errors are presented (n=4). **E-F** Summary of wild-  
995 type and mutant GLP2R beta-arrestin 1 and beta-arrestin 2 responses. Area under the curve (AUC)  
996 summary data (n=3-4) displayed for beta-arrestin 1 recruitment (E) and beta-arrestin 2 recruitment  
997 (F). AUCs were calculated using the 5 minutes prior to ligand addition as the baseline value. Mean  $\pm$   
998 standard errors are presented. Normal distribution of log<sub>10</sub> transformed data was determined by  
999 the D'Agostino & Pearson normality test. Following this statistical significance was assessed by one-  
1000 way ANOVA with post hoc Bonferroni test. \*\*\*p<0.001, \*p<0.05. **G** Schematic sketch for the location  
1001 of the missense variant induces amino acid substitution in the glucagon-like peptide-2 receptor  
1002 (GLP2R).

1003

1004 **Figure 6A** Results from polygenic risk scores (PGS) for each metabolite on risk for macular  
1005 telangiectasia type 2 (MacTel). The dotted line indicates the level of significance after correction for  
1006 multiple testing. **B** Effect estimates of serine-associated genetic variants on the risk for MacTel. **C**  
1007 Comparison of effect sizes for lead variants associated with plasma serine levels and the risk for  
1008 MacTel. **D** Receiver operating characteristic curves (ROC) comparing the discriminative performance

1009 for MacTel using a) sex, the first genetic principal component, and two MacTel variants (*rs73171800*  
1010 and *rs9820286*) not associated with metabolite levels (PMID: 28250457), and b) additionally  
1011 including genetically predicted serine and glycine at individual levels as described in the methods.  
1012 The area under the curve (AUC) is given in the legend.

1013 **Figure 7A** Scheme of the workflow to link common variation in genes causing inborn errors of  
1014 metabolism (IEM) to complex diseases. **7B** Flowchart for the systematic identification of metabolite-  
1015 associated variants to genes and diseases related to inborn errors of metabolism (IEM). **C** P-values  
1016 from phenome-wide association studies among UK Biobank using variants mapping to genes  
1017 knowing to cause IEMs and binary outcomes classified with the ICD-10 code. Colours indicate  
1018 disease classes. The dotted line indicates the significance threshold controlling the false discovery  
1019 rate at 5%. **D** Posterior probabilities (PPs) from statistical colocalisation analysis for each significant  
1020 triplet consisting of a metabolite, a variant, and a ICD-10 code among UK Biobank. The dotted line  
1021 indicates high likelihood (>80%) for one of the four hypothesis tested: H0 – no signal; H1 – signal  
1022 unique to the metabolite; H2 – signal unique to the trait; H3 – two distinct causal variants in the  
1023 same locus and H4 – presence of a shared causal variant between a metabolite and a given trait.

1024