# Genome-wide association study reveals loci with sex-specific effects on 1 plasma bile acids 2

3

Arianna Landini<sup>1,2</sup>, Dariush Ghasemi-Semeskandeh<sup>3,4</sup>, Åsa Johansson<sup>5</sup>, Shahzad Ahmad<sup>6</sup>, 4

Gerhard Liebisch<sup>7</sup>, Carsten Gnewuch<sup>7</sup>, Regeneron Genetics Center<sup>8</sup>, Gannie Tzoneva<sup>8</sup>, Alan 5

R. Shuldiner<sup>8</sup>, Andrew A. Hicks<sup>3</sup>, Peter Pramstaller<sup>3</sup>, Cristian Pattaro<sup>3</sup>, Harry Campbell<sup>2</sup>, Ozren 6

Polašek<sup>9,10</sup>, Nicola Pirastu<sup>11</sup>, Caroline Hayward<sup>1</sup>, Mohsen Ghanbari<sup>6</sup>, Ulf Gyllensten<sup>5</sup>, 7

Christian Fuchsberger<sup>3</sup>, James F. Wilson<sup>\*1,2</sup> & Lucija Klarić<sup>\*1</sup> 8

- 9
- 10 1 MRC Human Genetics Unit, Institute for Genetics and Cancer, University of Edinburgh,
- 11 Edinburgh, United Kingdom
- 2 Centre for Global Health Research, Usher Institute, University of Edinburgh, Edinburgh, 12
- 13 United Kingdom

14 3 Institute for Biomedicine, Eurac Research, Affiliated Institute of the University of Lübeck,

15 Bolzano, Italy

16 4 Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands

17 5 Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala

18 University, Uppsala, Sweden

19 6 Department of Epidemiology, Erasmus MC University Medical Center, Rotterdam, The 20 Netherlands

21 7 Institute of Clinical Chemistry and Laboratory Medicine, University Hospital Regensburg,

- 22 Regensburg, Germany
- 23 8 Regeneron Genetics Center, Tarrytown, NY, USA
- 24 9 Department of Public Health, School of Medicine, University of Split, Split, Croatia
- 25 10 Algebra University College, Zagreb, Croatia
- 26 11 Genomics Research Centre, Human Technopole, Milan, Italy
- 27
- 28 \* Authors contributed equally.
- 29 Correspondence to: J.F.W (jim.wilson@ed.ac.uk) or L.K. (lucija.klaric@ed.ac.uk)

#### 30 Abstract

# 31

Bile acids are essential for food digestion and nutrient absorption, but also act as signalling 32 molecules involved in hepatobiliary diseases, gastrointestinal disorders and carcinogenesis. 33 34 While many studies have focused on the genetic determinants of blood metabolites, research 35 focusing specifically on genetic regulation of bile acids in the general population is currently lacking. Here we investigate the genetic architecture of primary and secondary bile acids in 36 37 blood plasma, reporting associations with both common and rare variants. By performing genome-wide association analysis (GWAS) of plasma blood levels of 18 bile acids (N = 4923) 38 we identify two significantly associated loci, a common variant mapping to SLCO1B1 39 (encoding a liver bilirubin and drug transporter) and a rare variant in *PRKG1* (encoding soluble 40 cyclic GMP-dependent protein kinase). For these loci, in the sex-stratified GWAS (N $^{\uparrow}$  = 820, 41  $N^{\bigcirc}_{=} = 1088$ ), we observe sex-specific effects (*SLCO1B1*  $\beta \stackrel{\frown}{=} = -0.51$ ,  $P = 2.30 \times 10^{-13}$ ,  $\beta \stackrel{\bigcirc}{=} = -$ 42 0.3,  $P = 9.90 \times 10^{-07}$ ; *PRKG1*  $\beta$   $\circlearrowleft$  = -0.18,  $P = 1.80 \times 10^{-01}$ ,  $\beta$   $\heartsuit$  = -0.79,  $P = 8.30 \times 10^{-11}$ ), 43 44 corroborating the contribution of sex to bile acid variability. Using gene-based aggregate tests and whole exome sequencing, we identify rare pLoF and missense variants potentially 45

46 associated with bile acid levels in 3 genes (OR1G1, SART1 and SORCS2), some of which have

47 been linked with liver diseases.

### 48 Introduction

Bile acids (BAs) are synthesised from cholesterol in the liver and subsequently stored in the 49 gallbladder. After ingestion of food, BAs are secreted into the small intestine, where they 50 contribute to the digestion of lipid-soluble nutrients<sup>1</sup>. Approximately 95% of BAs are then re-51 52 absorbed by the intestinal epithelium and transported back to the liver via the portal vein - a process termed "enterohepatic circulation"<sup>2</sup>. Primary bile acids in humans consists of cholic 53 54 acid (CA), chenodeoxycholic acid (CDCA), and their taurine- or glycine-bound derivatives 55 (TCA and TCDCA, GCA and GCDCA). Once secreted in the lower gastrointestinal tract, 56 primary BAs are heavily modified by the gut microbiota to produce a broad range of secondary BAs, with deoxycholic acid (DCA), a CA derivative, and lithocholic acid (LCA), a CDCA 57 derivative, being the most prevalent<sup>2</sup>. Bile acids also act as hormone-like signalling molecules. 58 serving as ligands to nuclear (hormone) receptors. Through activation of these diverse 59 60 signalling pathways, BAs control not only their own transport and metabolism, but also lipid and glucose metabolism, and innate and adaptive immunity<sup>3</sup>. Bile acids are thus involved in 61 62 regulating several physiological systems, such as fat digestion, cholesterol metabolism, vitamin absorption, and liver function<sup>4</sup>. In addition, given their role in coordinating bile homeostasis, 63 biliary physiology and gastrointestinal functions, impaired signalling of BAs is associated with 64 65 development of hepatobiliary diseases, such as cholestatic liver disorders, cholesterol gallstone 66 disease and other gallbladder-related conditions<sup>5</sup>, and of inflammatory bowel disease<sup>6</sup>. Further, bile acids have been implicated in carcinogenesis - specifically oesophageal, gastric, 67 68 hepatocellular, pancreatic, colorectal, breast, prostate and ovarian cancer - both as pro-69 carcinogenic agents and tumour suppressors<sup>7</sup>. Thanks to their role as signalling molecules, BAs have been considered as possible targets for the treatment of metabolic syndrome and various 70 71 metabolic diseases<sup>8</sup>. Further, BAs are able to facilitate and promote drug permeation through biological membranes, making them of general interest for drug formulation and delivery<sup>9</sup>. 72

73 While many studies have focused on the genetic determinants of blood metabolites<sup>10–15</sup>, research focusing specifically on bile acids in a large sample from the general population is 74 currently lacking. Here we investigate the genetic architecture of primary and secondary BAs, 75 76 reporting associations with both common and low-frequency/rare variants. First, we performed 77 a genome-wide association meta-analysis (GWAMA) of plasma blood levels of 18 BA traits (N=4923). For a subset of this sample (female N=1088, male N=820), we perform sex-stratified 78 GWAMA, to describe sex-specific genetic contributions to BA variability. We then explore 79 80 whether complex traits or diseases have a role in influencing BA variability by using Mendelian Randomisation. We finally employ multiple gene-based aggregation tests to investigate rare 81 (MAF < 5%) predicted loss of function (pLoF) and missense variants from whole exome 82 sequencing affecting the 18 BA traits in a subset of our cohorts (N=1006). 83

#### **Results** 84

85

## 86 Loci associated with serum levels of bile acids

87

88 To investigate the genetic control of bile acids, we performed a GWAS meta-analysis on five cohorts of European descent (N = 4923), studying the associations of blood plasma levels of 89 90 18 primary and secondary bile acid traits with HRC-imputed genotypes/whole exome sequence data. Based on the number of below limit-of-detection (LOD) measurements, BAs were 91 92 analysed either as quantitative or binary traits (Supplementary Table 1). In addition, two 93 analysis approaches were carried out in parallel for quantitative traits: in one case, <LOD values were considered as missing, in the other case, they were imputed (Methods). An additive 94 linear model was assumed for each bile acid trait, followed by fixed-effect inverse-variance 95 meta-analysis. Overall, we identified 2 loci that passed the significance threshold (p-96 value  $< 3.57 \times 10^{-9}$ , Bonferroni adjusted for the number of independent bile acid traits) 97 (Figure 1), near the SLCO1B1 and PRKG1 genes. The most strongly associated locus (p=1.14 98 99 x  $10^{-16}$ ), on chromosome 12 near SLCO1B1, showed consistent directionality across 4 of the 5 populations (Table 1), with the effect allele T of the sentinel SNP, rs4149056, associated with 100 decreased serum levels of GDCA (quantitative). In the same locus, we found GLCA and the 101 imputed GDCA trait to be significantly associated with the rs73079476 variant (Supplementary 102 Table 2), in high linkage disequilibrium with the sentinel SNP, rs4149056 ( $r^2 = 0.97$ ). On the 103 other hand, rs146800892, the sentinel SNP on chromosome 10 near PRKG1, has a minor allele 104 105 frequency (MAF) lower than 1% in any cohort but CROATIA-Vis and might thus represent a 106 cohort-specific association with GCA (Supplementary Table 2).



107

Figure 1. Summary Manhattan plot pooling together meta-analysis results obtained 108 across 18 bile acid traits. The pooling was performed by selecting the lowest p value (y-axis) 109 110 from the 18 bile acids for every genomic position (x-axis). The Bonferroni-corrected genome-

- wide significance threshold (horizontal red line) corresponds to  $3.57 \times 10^{-9}$ . For simplicity, 111
- SNPs with p value >  $1 \times 10^{-3}$  are not plotted. P values are derived from the two-sided Wald test 112
- with one degree of freedom. 113

All samples										
Locus	Gene	SNP	EA	OA	EAF	Beta	Р	SE	Ν	Lead BA
12:20994540-21463812	SLCO1B1	rs4149056	Т	C	0.839	-0.25	1.10x10 <sup>-16</sup>	0.03	4547	GDCA
10:53832549-53832549	PRKG1	rs146800892	Т	C	0.988	-0.96	3.30x10 <sup>-09</sup>	0.16	900	GCA
Sex-stratified										
Locus	Gene	SNP	EA	OA	Beta M	Beta F	РМ	P F	N (M/F)	Lead BA
12:20994540-21463812	SLCO1B1	rs73079476	Α	С	-0.51	-0.31	2.30x10 <sup>-13</sup>	9.90x10 <sup>-07</sup>	820/1088	GDCA
10:53832549-53832549	PRKG1	rs117834398	Т	G	-0.18	-0.79	1.80x10 <sup>-01</sup>	8.30x10 <sup>-11</sup>	820/1088	GCA

# 114 Table 1: Loci genome-wide significantly associated with at least one of the 18 bile acid traits in all samples and sex-stratified GWAMA

Each locus is represented by the SNP with the strongest association in the region, according to the p-value rejecting the null hypothesis of no association with at least one of 18 bile acid traits. In all samples analysis, an association was considered significant if the p-value was lower than  $3.57 \times 10^{-9}$ , the genome-wide significance threshold Bonferroni-corrected for the number of independent bile acid traits. In sex-stratified analysis, an association was considered significance threshold Bonferroni-corrected for the number of independent bile acid traits. In sex-stratified analysis, an association was considered significant if the p-value was lower than  $5 \times 10^{-9}$ , the genome-wide significance threshold Bonferroni-corrected for the number of independent bile acid quantitative traits. The two SNPs in the *SLCO1B1* locus are in high LD (LD r<sup>2</sup> = 0.97), while the two SNPs in the *PRKG1* locus represent two distinct signals (LD r<sup>2</sup> < 0.001).

Locus - coded as 'chromosome: locus start-locus end' (GRCh37 human genome build); Gene - suggested candidate gene; SNP - variant with the strongest association in the locus; EA - SNP allele for which the effect estimate is reported; OA - other allele; EAF - frequency of the effect allele; Beta - effect estimate for the SNP and bile acid with the strongest association in the locus; SE - standard error of the effect estimate, P - p-value of the effect estimate (two-sided Wald test with one degree of freedom); N - sample size; Lead BA - bile acid with the strongest association to the reported SNP; M - male specific analysis; F - female specific analysis.

#### 116 Sex-specific associations of bile acid serum levels

117 To investigate whether the genetic component influencing bile acid variation may differ 118 between men and women, we performed sex-specific GWAS meta-analysis of the 14 quantitative (imputed) bile acid traits for ORCADES and CROATIA-Vis cohorts (female 119 120 N=1088, male N=820) and discovered two sex-specific associations. The association of GDCA 121 with rs73079476 from the SLCO1B1 locus was significant in male-only GWAS (beta = -0.51, p-value =  $2.28 \times 10^{-13}$ ) (Table 1, Figure 2A). The signal for the same locus in female-only 122 GWAS, while consistent in terms of directionality, has a smaller effect size than in male-only 123 analysis (beta = -0.31) and does not reach the significance threshold (p-value =  $9.86 \times 10^{-7}$ ). 124 125 despite the slightly higher sample size (Figure 2). This suggests that the genetic effect of SLCO1B1 locus on the serum levels of GDCA is larger in men than women. We also identified 126 127 a sex-specific association of GCA at the *PRKG1* locus. In contrast to *SLCO1B1*, the sentinel 128 SNP in *PRKG1*, rs117834398, has a larger effect in females than in males (female beta = -0.79, male beta = -0.18), and passed the significant threshold only in the female-specific analysis 129 (female p-value =  $8.26 \times 10^{-11}$ , male p-value =  $1.81 \times 10^{-1}$ ) (Table 1, Figure 2B). Interestingly, 130 131 the sentinel SNPs at the *PRKG1* locus for the full meta-analysis and for the female-specific analysis are in linkage equilibrium ( $r^2 < 0.01$ ) and represent two independent associations in that 132 locus. Overall, none of the significant association identified in one sex was replicated in the 133 134 other, suggesting that the genetic contribution to serum BA levels is likely to be different in 135 males and females. We have identified 13 additional associations (p-value  $< 5 \times 10^{-9}$ ). Bonferroni adjusted for the number of independent quantitative bile acid traits) that might have 136 sex-specific effects (Supplementary Table 3, Supplementary Figure 1). However, given the low 137 allele frequencies and allele counts in the two analysed cohorts, further analyses are required 138 139 to replicate these associations.





140 Figure 2. Sex-specific associations. The effect of rs73079476 on chromosome 12 on GDCA 141 bile acid is almost as twice strong in males compared to the effect in females (Panel A). The effect of rs117834398 on GCA bile acid is stronger in females than in males (Panel B). N -142 143 sample size, MAF - minor allele frequency, MAC - minor allele count, CI - confidence 144 interval.

145

## Link with complex traits and diseases 146

Next, we assessed whether variants associated with BA levels have been previously associated 147 with any other biochemical traits and diseases. Using Phenoscanner<sup>16,17</sup> we found that 148 rs4149056, sentinel SNP in *SLCO1B1* locus, and its proxies  $(r^2 > 0.8)$ , were also associated 149 with concentration of bilirubin, non-bile acid metabolites, mean corpuscular haemoglobin, sex 150 hormone binding globulin and estrone conjugates, and various responses to drugs (i.e., statin-151 152 induced myopathy, LDL-cholesterol response to simvastatin and methotrexate clearance in acute lymphoblastic leukaemia) (Supplementary Table 4). To obtain deeper insight into the 153 causal relationship between BAs and diseases, we conducted bi-directional Mendelian 154 Randomisation (MR) analysis. Using the sentinel SNPs associated with GLCA, GDCA and 155 GCA (Table 1) as instrumental variables we tested whether genetically increased levels of BA 156 157 influence levels or risk for 548 biochemical traits and diseases available in the IEU Open GWAS database<sup>18</sup> (Supplementary Table 5). Levels of GLCA and GDCA were significantly 158

 $(p-value < 0.05/(548x3) = 3.04 \times 10^{-5})$  associated with different biochemical measurements, 159 such as levels of sex hormone-binding globulin, testosterone, triglycerides, vitamin D, alanine 160 161 transaminase and galectin-3; with blood traits, such as mean corpuscular haemoglobin and 162 mean corpuscular volume; and with diseases and their risk factors, such as daytime dozing and 163 stroke (Supplementary Table 6). These MR tests were performed using the Wald ratio test 164 utilising only a single instrument, thus the results of causal relationship between BAs and traits/diseases should be interpreted with caution. Yet our results suggest a possible overlap in 165 166 genetic regulation, involving the SLCO1B1 locus. Next, to assess whether complex traits and 167 disease could have an effect on bile acid levels, we performed reverse MR using 548 traits/diseases as exposure and bile acids as outcomes. We found no significant associations, 168 169 suggesting that none of the tested diseases or complex traits have an effect on BA levels 170 (Supplementary Table 7).

171

### 172 Exome-wide rare variant analysis of bile acids

To assess the contribution of low frequency and rare variants to the bile acid genetic 173 architecture, we performed exome-wide gene-based tests across 18 bile acid traits in the 174 175 ORCADES cohort (N = 1006) by testing the aggregated effect of rare (MAF <5%) predicted 176 loss-of-function (pLoF) and non-synonymous missense variants. We identified significant 177 association (p-value  $<1.79 \times 10^{-7}$ ) of rare variants from 3 genes with 2 bile acid traits (quantitative CA and binary THDCA). For these associations, a significant p-value was 178 reported by at least 2 of the 4 aggregation tests used. Rare variants significantly associated with 179 180 quantitative bile acid trait CA are located in the OR1G1 gene, while those associated with binary bile acid trait THDCA are located in SART1 and SORCS2 genes (Table 2, 181 182 Supplementary Table 8). We further identified significant association of rare variants from EPS8L1 gene with quantitative bile acid trait DCA and from EEF2K with binary bile acid trait 183 THDCA (Supplementary Table 8). However, a significant p-value was reported by only one of 184 185 the 4 aggregation tests used. Due to the lack of replication across aggregations tests, we 186 considered these associations as not robust.

187

BA	Trait type	Gene	MAF	Functional consequence	N variants	Aggregation test	Р	AC
CA	Quantitative	OR1G1	<0.01	Missense variants	2	SKAT-O	1.67x10 <sup>-8</sup>	17
THDCA	Binary	SORCS2	< 0.05	pLoF and missense variants	10	SMMAT-E	1.44x10 <sup>-8</sup>	174
THDCA	Binary	SART1	< 0.01	pLoF and missense variants	4	SKAT	1.19x10 <sup>-7</sup>	25

Table 2. Gene-based aggregation analysis results for bile acid traits in ORCADES cohort 188

BA- bile acid trait tested for rare variants association; Trait type – whether BA was analysed as a quantitative or binary trait; Gene gene for which variants were aggregated; MAF - upper bound for minor allele frequency of tested variants; Functional consequence predicted functional consequence for aggregated variants; N variants - number of variants in the mask; Aggregation test - rare-variants aggregation test reporting the lowest p-value out of 4 aggregation tests; P - p-value for the aggregation test; AC - cumulative allele count of all the variants in a mask. Bonferroni-corrected discovery p-value threshold was set to 1.79x10<sup>-7</sup> (0.05/20,000 estimate of number of genes in the human genome/14 number of independent bile acids).

# 189

## 190 Discussion

Bile acids (BAs) are synthesised from cholesterol in the liver and then secreted into the small 191 intestine to emulsify and promote absorption of lipid-soluble nutrients. BAs also act as 192 193 hormone-like signalling molecules and have been linked to regulation of lipid and glucose 194 metabolism, immunity, vitamin absorption, hepatobiliary diseases, inflammatory bowel disease and cancer. Despite the crucial role of BAs on whole-body physiology, their genetic 195 196 architecture has not been extensively investigated in a large sample from the general population. In this study, we performed both pooled and sex-stratified genome-wide 197 198 association meta-analysis of plasma levels of 18 bile acid compounds, including both primary and secondary forms, in 4923 European individuals. 199

200 We identified two secondary bile acids (GDCA and GLCA) significantly associated with a 201 locus encompassing the SLCO1B1 gene. The encoded protein, OATP1B1 (organic anion transporting polypeptide 1B1), is a well-known human hepatocyte transporter mediating the 202 203 uptake of various endogenous compounds such as bile salts, bilirubin glucuronides, thyroid hormones and steroid hormone metabolites, and also clinically frequently used drugs like 204 statins, HIV protease inhibitors, and the anti-cancer agents irinotecan or methotrexate<sup>19–23</sup>. The 205 sentinel SNP of the SLCO1B1 locus, rs4149056, is a missense variant (p.Val174Ala) which 206 has been linked by previous GWA studies to blood concentration of several metabolites, 207 including vitamin  $D^{24}$ , triglycerides<sup>25</sup> and bilirubin<sup>26</sup>, a compound resulting from the 208 breakdown of haem catabolism and excreted as a major component of bile. This same variant 209 210 has also been associated with levels of sex hormone-binding globulin and testosterone<sup>27</sup>. The 211 knock-out of the gene in mice results in abnormal liver physiology and abnormal xenobiotic pharmacokinetic phenotypes (Open Targets<sup>28</sup>). A rare variant from the *PRKG1* locus was 212 significantly associated with levels of glycocholic acid (GCA). PRKG1 encodes a Protein 213 214 Kinase CGMP-Dependent 1, a protein involved in signal transduction and a key mediator of the nitric oxide/cGMP. The sentinel variant in the region, rs146800892, only passes the MAF 215 216 threshold (MAF > 0.01) in the CROATIA-Vis cohort, which is therefore the only cohort contributing to this association. Due to its demographic history and geographic position, 217 CROATIA-Vis is a genetic isolate<sup>29</sup> so it is possible that this variant has increased in frequency 218 compared to a general population<sup>30</sup>. The mechanism of how the variation within this gene could 219 relate to bile acid levels is unclear and would need to be further investigated. 220

221 In the sex-stratified GWAS meta-analysis, we observed sex-specific associations for the two identified loci. Levels of glycodeoxycholic acid (GDCA) are more strongly associated with the 222 variant in SLCO1B1 in men than in women, while female levels of GCA are more strongly 223

224 affected by the variant in PRKG1 than male levels. Later, our Mendelian randomization 225 analysis did not provide evidence that testosterone, oestradiol, sex hormone-binding globulin or other sex-related traits have causal effects on plasma BA levels. While this could be due to 226 a lack of statistical power of our BA meta-analysis, we currently have no evidence to suggest 227 an effect of sex-related hormones on BA levels mediated by genetics. We also detected 228 229 associations with variants from the same gene, *PRKG1*, in the main, non-stratified analysis. 230 However, the two associations (sex-specific and pooled) appear to be independent (LD  $r^2$ 231 <0.001). While the association from the pooled analysis might be either false positive or 232 population-specific, the independent association from the sex-stratified analysis replicates well 233 between two analysed cohorts (CROATIA-Vis and ORCADES).

234 After assaying common variants through GWAS, we performed exome-wide gene-based association tests in a subset of our samples (N = 1006), to investigate the genetic contribution 235 236 of rare and low frequency (MAF <5%) coding variants (pLoF and missense) to bile acid levels. 237 Overall, we identified associations with rare variants from 3 genes, OR1G1, SART1 and 238 SORCS2. OR1G1 is an olfactory receptor gene, whose coded protein receptor interacts with 239 odorant molecules in the nose to initiate a neuronal response triggering the perception of 240 smell<sup>31,32</sup>. In addition to the nasal level, the olfactory receptor coded by *OR1G1* is expressed 241 also by enterochromaffin cells, specialised enteroendocrine cells of the gastrointestinal tract. 242 Braun et al<sup>33</sup>. determined that certain olfactory cues from spices and odorants, such as thymol, present in the luminal environment of the gut may stimulate serotonin release via olfactory 243 receptors present in enterochromaffin cells. Between 90% and 95% of total body serotonin is 244 in fact synthesised by enterochromaffin cells<sup>34</sup>: serotonin controls gut motility and secretion 245 and is implicated in pathologic conditions such as vomiting, diarrhoea, and irritable bowel 246 syndrome<sup>33</sup>. In mice, gut serotonin was shown to stimulate bile acid synthesis and secretion by 247 the liver and gallbladder. Thus, release of serotonin in response to odorant cues increases bile 248 acid turnover<sup>35</sup>. The hypoxia-associated factor (HAF), encoded by SART1 gene and also known 249 250 as SART1(800), is involved in proliferation and hypoxia-related signalling. The protein encoded by SORCS2 is a receptor for the precursor of nerve growth factor, up-regulation of 251 which has been reported for several liver pathologies, such as hepatotoxin- induced fibrosis<sup>36</sup>, 252 ischemia-reperfusion injury<sup>37</sup>, oxidative injury<sup>38</sup>, cholestatic injury<sup>39</sup> and hepatocellular 253 254 carcinoma<sup>36,40,41</sup>. However, due to unavailability of exome sequencing data in other cohorts 255 these associations were not replicated.

Recently, Chen et al.<sup>42</sup> have performed an association analysis on plasma and faecal levels of 256 bile acids in 297 obese individuals. Their study revealed 27 associated loci, including genes 257 involved in transport of GDP-fucose and zinc/manganese and zinc-finger-protein-related 258 259 genes, mostly associated with bile acid levels in stool. In our study we analysed blood plasma in a much larger sample from a general population and discovered only two associated loci. 260 261 Neither of genes identified in our study were reported in Chen et al, suggesting that genetic regulation of bile acids between stool and blood plasma or between obese and general 262 populations might differ significantly. 263

264 We acknowledge several limitations in the present study. We found only a small percentage of BA variability to be affected by genetics, suggesting that a larger sample size is required to 265 further describe BA genetic architecture. BAs are known to be largely influenced by 266 environmental factors, such as sex and gut microbiota. Female sex and oestrogens are 267 considered relevant regulators of BA production and composition<sup>43,44</sup>. In pregnant women, high 268 269 levels of circulating oestrogen are associated with development of cholestasis, characterised by 270 increased serum bile acids, likely via oestrogen reducing the expression of BA receptor and transport proteins<sup>45</sup>. Similarly, age-related differences in hormone levels influence the 271 differential production of BAs in women<sup>46</sup>. The relevant impact of sex on plasma BA levels 272 273 was confirmed by the sex-stratified analysis, where the two significantly associated loci showed to be sex-specific. Similarly, species-composition of gut microbiota has a great impact 274 275 on BAs levels, especially for secondary BAs that are a direct result of microbiome activity. A recent study describing the effect of gut microbiota on the human plasma metabolome reported 276 277 that both primary BA cholic acid (CA) and secondary BA deoxycholic acid (DCA) show a high 278 percentage of variance explained by the microbiota ( $R^2 = 30\%$  and 36\%, respectively), indicating a strong impact on BAs of the variation in microbiota composition<sup>47</sup>. It is important 279 280 to interpret our findings in the context of the tissue in which BA levels were measured, blood 281 plasma. Bile acids are synthesised in the liver and secreted into the intestine, to be then 282 reabsorbed into portal circulation and returned to the liver: plasma BA levels thus reflect the amount of BAs escaping extraction from the portal blood. Therefore, levels of BAs in plasma 283 284 are likely to be influenced by genes other than those encoding the particular anabolic and catabolic enzymes, including those involved in hepatic function and dysfunction. In line with 285 this, the major genetic contributor to blood BA levels in our study are variants from the 286 287 SLCO1B1 gene, encoding the hepatocyte transporter OATP1B1 and important for flux of bile salts, bilirubin glucuronides and various hormone metabolites, rather than genes encoding key 288 289 enzymes of primary BA synthesis, such as CYP7A1 and CYP7B1<sup>48</sup>. Similarly, some of the 290 genes with rare variants associations have been linked to liver diseases, such as liver cancer<sup>49</sup>, 291 and intrahepatic cholestasis of pregnancy $^{50}$ .

292 In conclusion, we explored the genetic architecture of plasma bile acid levels, including both 293 common and rare variants. By performing GWAS meta-analysis (N = 4923), we identified 2 significantly associated loci, mapping to the SLCO1B1 and PRKG1 genes. In the sex-specific 294 295 GWAS meta-analysis we observed that variants in these genes have different impact on bile 296 acid levels in men and women. To assess relationships between genetically increased levels of 297 bile acids and risk for diseases we performed Mendelian randomisation, but did not find any 298 bile acids affecting disease risk, nor the reverse, which however might be affected by the lack 299 of statistical power. Using the gene-based aggregated tests and whole exome sequencing, we further identified rare pLoF and missense variants in 3 genes associated with BAs, OR1G1, 300 301 SART1 and SORCS2, some of which are known to be involved in liver disease. Additional 302 studies with larger sample sizes and of more diverse ancestry will be necessary to validate our findings, further unravel the genetic architecture of bile acid levels, and to understand their 303 relationship with human diseases and complex traits. 304

### 305 **Materials and methods**

306

## Phenotypic data 307

### 308 Bile acids quantification

Bile acid (BA) analysis was performed from plasma or serum (MICROS cohort) samples by 309 310 liquid chromatography-tandem mass spectrometry (LC-MS/MS) as previously described<sup>51</sup>. 311 The HPLC equipment consisted of a 1200 series binary pump (G1312B), a 1200 series isocratic 312 pump (G1310A) and a degasser (G1379B) (Agilent, Waldbronn, Germany) connected to an 313 HTC Pal autosampler (CTC Analytics, Zwingen, CH). A hybrid triple quadrupole linear ion 314 trap mass spectrometer API 4000 Q-Trap equipped with a Turbo V source ion spray operating 315 in negative ESI mode was used for detection (Applied Biosystems, Darmstadt, Germany). High 316 purity nitrogen was produced by a nitrogen generator NGM 22-LC/MS (cmc Instruments, 317 Eschborn, Germany). Gradient chromatographic separation of BAs was performed on a 50 mm 318 × 2.1 mm (i.d.) Macherey-Nagel NUCLEODUR C18 Gravity HPLC column, packed with 1.8 um particles equipped with a 0.5 µm pre-filter (Upchurch Scientific, Oak Harbor, WA, USA). 319 320 The injection volume was 5 µL and the column oven temperature was set to 50 °C. Mobile phase A was methanol/water (1/1, v/v), mobile phase B was 100% methanol, both containing 321 0.1% ammonium hydroxide (25%) and 10 mmol/L ammonium acetate (pH 9). A gradient 322 323 elution was performed with 100% A for 0.5 min, a linear increase to 50% A until 4.5 min, 324 followed by 0% A from 4.6 until 5.5 min and re-equilibration from 5.6 to 6.5 min with 100% 325 A. The flow rate was set to 500 µL/min. To minimize contamination of the mass spectrometer, 326 the column flow was directed only from 1.0 to 5.0 min into the mass spectrometer using a 327 diverter valve. Otherwise, methanol with a flowrate of 250 µL/min was delivered into the mass 328 spectrometer. The turbo ion spray source was operated in the negative ion mode using the following settings: Ion spray voltage = -4500 V, ion source heater temperature = 450 °C, 329 source gas 1 = 40 psi, source gas 2 = 35 psi and curtain gas setting = 20 psi. Analytes were 330 331 monitored in the multiple reaction monitoring (MRM). Quadrupoles Q1 and Q3 were working 332 at unit resolution. Calibration was achieved by the addition of BAs to EDTA-plasma/serum. A 333 combined BA standard solution containing the indicated amounts (0.5 - 70.5 µmol/L) was 334 placed in a 1.5 ml tube and excess solvent was evaporated under reduced pressure before adding 335 EDTA-plasma/serum. Calibration curves were calculated by linear regression without weighting. Data analysis was performed with Analyst Software 1.4.2. (Applied Biosystems, 336 337 Darmstadt, Germany). The data were exported to Excel spreadsheets and further processed by 338 self-programmed Excel macros which sort the results, calculate the analyte/internal standard peak area ratios, generate calibration lines and calculate sample concentrations. For the 339 340 calculation we selected the internal standard with analogous fragmentation and closest 341 retention time to the respective BA species.

342

Pre-processing of bile acid traits 343

344 Prior to genetic analysis, bile acid traits were grouped into three groups based on the percentage 345 of samples with below the limit of detection (<LOD) measurements: high <LOD group (> ~30% of all samples below LOD) and low <LOD group (<  $\sim$ 7% of all samples below LOD) 346 (Supplementary Table 1). Accordingly, different phenotypic pre-processing and different 347 348 analysis strategies were applied to the groups. Bile acids within a high <LOD were considered 349 as binary traits: individuals were categorised based on whether their bile acid levels were 350 effectively measured (category 1) or were below the LOD (category 0). Bile acid traits 351 belonging to this group were THDCA, TUDCA, TCA and GHDCA. All other bile acids were 352 considered as quantitative traits and were  $log_{10}$ -transformed. However, to increase the sample size, in addition to a complete-case analysis (considering as missing all samples with <LOD), 353 we also imputed <LOD measurements. For each bile acid, imputation of <LOD measurements 354 was performed by fitting a truncated normal distribution, with mean and standard deviation of 355 the effectively measured raw data, truncated (as an upper bound) to the lowest measured value 356 357 for the given bile acid. To do so, we used the "rtnorm" function from the MCMCglmm R 358 package<sup>52</sup>. After imputation, measurements were log<sub>10</sub>-transformed.

359

# Genome-wide association analysis 360

Genome-wide association studies (GWAS) were performed in 5 cohorts of European descent, 361 CROATIA-Vis (N=971), ORCADES (Orkney Complex Disease Study) (N=1019), NSPHS 362 (Northern Sweden Population Health Study) (N=718), MICROS (Micro-Isolates in South 363 Tyrol) (N=1336) and ERF (Erasmus Rucphen Family Study) (N=879), for a combined sample 364 size of 4923. Specific sample size for each bile acid molecule, in both meta-analysis and single 365 366 cohort GWAS, can be found in Supplementary Table 10. Bile acid traits were adjusted for age, 367 sex, batch, population structure/cryptic relatedness by including population principal components or applying linear mixed models and using a kinship matrix estimated from 368 369 genotyped data. Within each cohort, residuals of covariate and population structure/relatedness correction were tested for association with Haplotype Reference Consortium (HRC)<sup>53</sup> imputed 370 371 SNP dosages or SNP genotypes from whole genome sequencing, applying an additive genetic 372 model of association. Details of cohorts, individual-level pre-imputation QC, GWAS software and parameters specific for each cohort can be seen in Supplementary Table 11 Single-cohort 373 374 summary statistics were filtered for minor allele frequency (MAF) > 0.01. The genomic control 375 inflation factor ( $\lambda_{GC}$ ) was calculated for each bile acid trait. Cohort-level  $\lambda_{GC}$  overall ranged from 0.9 to 1.1 for quantitative bile acid traits, both imputed and not, suggesting little residual 376 377 influence of population stratification and family structure (Supplementary Table 12). In a few 378 cases, ERF cohort reported somewhat deflated  $\lambda_{GC}$  (GCDCA at 0.884 and GLCA at 0.899). On 379 the other hand, there was considerable inflation for binary bile acid in the case of NSPHS (Supplementary Table 12), with values of  $\lambda_{GC}$  above 1.1, suggesting that population 380 381 structure/cryptic relatedness was not fully controlled for these traits in the NSPHS cohort.

382

### 383 **Meta-analysis**

Prior to meta-analysis, cohort-level GWAS were quality controlled using the EasyQC software 384 package, following the protocol described in Winkler et al.<sup>54</sup> Cohort-level results were 385 corrected for the genomic control inflation factor, then pooled and analysed with METAL 386 v2011-03-25 software<sup>55</sup>, applying the fixed-effect inverse-variance method. The mean genomic 387 control inflation factor after the meta-analysis was 0.991 (range 0.938 - 1.009), suggesting that 388 389 the confounding effects of the family structure were correctly accounted for (Supplementary 390 Table 12). The standard genome-wide significance threshold was Bonferroni corrected for the number of independent bile acid traits, calculated as  $14 (5x10^{-8}/14 = 3.57x10^{-9})$ . The number 391 of independent bile acid traits was estimated as the sum of the number of binary traits (4) and 392 the number of principal components that jointly explained 99% of the total variance of  $log_{10}$ -393 transformed quantitative traits in each cohort (10) (Supplementary Table 13). 394

395

## 396 Sex-stratified GWAS meta-analysis

To identify possible differences in the genetic contribution to bile acid variability between men 397 398 and women, we performed sex-specific GWAS of the 14 quantitative bile acid traits for 399 ORCADES and CROATIA-Vis cohorts. Given that for the sex-stratified GWAS we implicitly 400 halve our sample size, we performed these analyses only on the imputed bile acid traits. The 401 same analysis steps and procedures already described for the full meta-analysis were applied. 402 Bile acid traits were adjusted for age, sex and batch as fixed effects, and relatedness (estimated 403 as the kinship matrix calculated from genotyped data) as a random effect in a linear mixed model, calculated using the 'polygenic' function from the GenABEL R package<sup>56</sup>. Residuals 404 of covariate and relatedness correction were tested for association with HRC-imputed<sup>53</sup> SNP 405 dosages using the RegScan v0.5 software<sup>57</sup>, applying an additive genetic model of association. 406 407 Prior to meta-analysis, SNPs having a difference in allele frequency between the two cohorts higher than ±0.3 or a minor allele count (MAC) lower or equal to 6 were filtered out. Cohort-408 409 level GWAS were corrected for genomic control inflation factor and then meta-analysed (N = 820 for male and N = 1088 for female individuals) using METAL v2011-03-25 software<sup>55</sup>, 410 411 applying the fixed-effect inverse-variance method. The mean  $\lambda_{GC}$  was 0.993 (range 0.978– 412 1.011) for male-specific meta-analysis and 0.996 (range 0.984-1.003) for female-specific

- meta-analysis. The Bonferroni-corrected significance threshold applied is  $5 \times 10^{-9}$ . 413
- 414

# 415 **Phenoscanner and Mendelian Randomisation**

To assess link between bile acids and diseases we explored the overlap of SNPs associated with 416 BAs with complex human traits by using PhenoScanner v1.1 database<sup>16,17</sup>, taking into account 417 significant genetic association ( $p < 5 \times 10^{-9}$ ) at the same or strongly (LD r2 > 0.8) linked SNPs 418 in populations of European ancestry. We then performed bi-directional Mendelian 419 Randomisation (MR) to investigate the effect of 548 complex traits and diseases available in 420 the IEU Open GWAS database<sup>18</sup> (manually curated list of studies from identifiers ebi-a, ieu-a, 421 ieu-b and ukb-a; the complete list reported in the Supplementary Table 5) on BA levels, and 422

vice-versa. The set of genome-wide significant, LD clumped SNPs used as instruments for
complex traits/diseases was extracted from the selected studies by using the
"extract\_instruments" function from the TwoSampleMR 0.5.6 R package<sup>58</sup>. Similarly, sentinel
SNPs from BAs meta-analysis (Supplementary Table 2) were selected as instruments. MR tests
were performed by using fixed effects inverse variance-weighted (IVW) in case of multiple
instruments or Wald Ratio method in case of a single instrument, as implemented in the
TwoSampleMR 0.5.6 R package<sup>58</sup>. Multiple testing correction was controlled for using either

- 430 the Bonferroni correction or false discovery rate (FDR).
- 431

# 432 Whole-exome sequencing data

433 Exome sequencing

434 The "Goldilocks" exome sequence data for ORCADES cohort was prepared at the Regeneron Genetics Center, following the protocol detailed in Van Hout *et al.*<sup>59</sup> for the UK Biobank 435 436 whole-exome sequencing project. In summary, sequencing was performed using S2 flow cells 437 on the Illumina NovaSeq 6000 platform with multiplexed samples. DNAnexus platform<sup>60</sup> was used for processing raw sequencing data. The files were converted to FASTO format and 438 aligned using the BWA-mem<sup>61</sup> to GRCh38 genome reference. The Picard tool<sup>62</sup> was used for 439 identifying and flagging duplicated reads, followed by calling the genotypes for each individual 440 sample using the WeCall variant caller<sup>63</sup>. During quality control, 33 samples genetically 441 identified as duplicates, 3 samples showing disagreement between genetically determined and 442 reported sex, 4 samples with high rates of heterozygosity or contamination, 2 samples having 443 low sequence coverage (less than 80% of targeted bases achieving 20X coverage) and 1 being 444 discordant with genotyping chip were excluded. Finally, the "Goldilocks" dataset was 445 446 generated by (i) filtering out genotypes with read depth lower than 7 reads, (ii) keeping variants 447 having at least one heterozygous variant genotype with allele balance ratio greater than or equal to 15% (AB  $\geq$  0.15) or at least one homozygous variant genotype, and (iii) filtering out variants 448 with more than 10% of missingness and HWE  $p < 10^{-6}$ . Overall, a total of 2,090 ORCADES 449 450 (820 male and 1,270 female) participants passed all exome sequence and genotype quality control thresholds. A pVCF file containing all samples passing quality control was then created 451 using the GLnexus joint genotyping tool<sup>64</sup>. 452

- 453
- 454 <u>Variant annotation</u>

Exome sequencing variants were annotated as described in Van Hout, *et al.*<sup>59</sup> Briefly, they were annotated with the most severe consequence across all protein-coding transcripts using SnpEff<sup>65</sup>. Gene regions were defined based on Ensembl release 85<sup>66</sup>. Predicted loss-of function (pLoF) variants were defined as variants annotated as start lost, stop gained/lost, splice donor/acceptor and frameshift. The deleteriousness of missense variants was based on dbNSFP  $3.2^{67,68}$  and assessed using the following algorithms: (1) SIFT<sup>69</sup>: "D" (Damaging), (2)

461 Polyphen2 HDIV: "D" (Damaging) or "P" (Possibly damaging), (3) Polyphen2 HVAR<sup>70</sup>: "D" (Damaging) or "P" (Possibly damaging), (4) LRT<sup>71</sup>: "D" (Deleterious) and (5) 462 MutationTaster<sup>72</sup>: "A" (Disease causing automatic) or "D" (Disease causing). If not predicted 463 as deleterious by any of the algorithms the missense variants were considered "likely benign". 464 "possibly deleterious" if predicted as deleterious by at least one of the algorithms and "likely 465 466 deleterious" if predicted as deleterious by all five algorithms.

467

# 468 Exome-wide gene-based aggregation analysis of rare variants

469 Generation of gene masks

470 For each gene, the variants were grouped into four categories (masks), based on severity of 471 their functional consequence. The first mask (mask 1) included only pLoF variants. Masks 3 472 and 4 included both pLoF and variants predicted to be deleterious, by 5/5 algorithms (mask 3) or by at least one algorithm (mask 4). The most permissive mask (mask 2) included pLoF and 473 474 all missense variants. These masks were then further split by the frequencies of the minor allele 475  $(MAF \le 5\%, e.g. mask1 maf5; and MAF \le 1\%, e.g. mask1_maf1)$ , resulting in up to 8 burden 476 tests for each gene (Supplementary Table 9).

477

# **ORCADES** gene-based aggregation analysis 478

We performed variant Set Mixed Model Association Tests (SMMAT)<sup>73</sup> on the 18 bile acid 479 480 traits from ORCADES cohort, quantified and pre-processed as previously described, fitting a 481 GLMM adjusting for age, sex, batch, and familial or cryptic relatedness by kinship matrix. The kinship matrix was estimated from the genotyped data using the 'ibs' function from GenABEL 482 R package<sup>56</sup>. The SMMAT framework includes 4 variant aggregate tests: burden test, sequence 483 kernel association test (SKAT), SKAT-O and SMMAT-E, a hybrid test combining the burden 484 485 test and SKAT. The 4 variant aggregate tests were performed on 8 different pools of genetic 486 variants, called "masks", each one including a different set of variants based on both MAF and predicted consequence of variants (e.g., loss of function and missense) (Supplementary Table 487 9), as described above. Discovery significance threshold was Bonferroni corrected for the 488 489 rough estimate of the number of genes in the human genome, 20,000, and the number of independent bile acid traits, 14, calculated as previously described  $(0.05/20000/14 = 1.79 \times 10^{-1})$ 490 491 <sup>7</sup>). A gene association was considered significant if it passed the above reported Bonferroni corrected significance threshold in at least two of the 4 performed variant aggregate tests and 492 493 if the cumulative allele count of the variants included in the gene was equal or higher than 10.

# 494 Code availability

We used publicly available software tools for all analyses. These software tools are listed inthe main text and in the Methods.

497

# 498 Data availability

499 The full summary statistics from GWAS meta-analysis of bile acids will be uploaded to the University of Edinburgh Datashare repository and to GWAS catalog upon manuscript 500 501 acceptance. There is neither Research Ethics Committee approval, nor consent from individual participants, to permit open release of the individual level research data underlying this study. 502 503 The datasets analysed during the current study are therefore not publicly available. Instead, the research data and/or DNA samples for the ORCADES study are available from 504 505 accessQTL@ed.ac.uk on reasonable request, following approval by the QTL Data Access Committee and in line with the consent given by participants. Each approved project is subject 506 507 to a data or materials transfer agreement (D/MTA) or commercial contract. The summary statistics for complex traits and diseases (full list reported in Supplementary Table 5) are 508 509 available in the IEU Open GWAS database https://gwas.mrcieu.ac.uk/.

### 510 References

- 511 1. Lorbek, G., Lewinska, M. & Rozman, D. Cytochrome P450s in the synthesis of 512 cholesterol and bile acids – from mouse models to human diseases. FEBS J. 279, 513 1516-1533 (2012).
- 514 2. Chiang, J. Y. L. Bile Acid Metabolism and Signaling. Compr. Physiol. 3, 1191–1212 515 (2013).
- 516 3. Thomas, C., Pellicciari, R., Pruzanski, M., Auwerx, J. & Schoonjans, K. Targeting 517 bile-acid signalling for metabolic diseases. Nature Reviews Drug Discovery 7, 678-518 693 (2008).
- 519 4. de Aguiar Vallim, T. Q., Tarling, E. J. & Edwards, P. A. Pleiotropic roles of bile acids 520 in metabolism. Cell Metab 17, 657-669 (2013).
- 521 5. Perino, A., Demagny, H., Velazquez-Villegas, L. & Schoonjans, K. Molecular 522 Physiology of Bile Acid Signaling in Health, Disease, and Aging. Physiol Rev 101, 523 683-731 (2021).
- 524 6. Fiorucci, S. et al. Bile Acid Signaling in Inflammatory Bowel Diseases. Dig Dis Sci 525 66, 674-693 (2021).
- 526 7. Rezen, T. et al. The role of bile acids in carcinogenesis. Cell Mol Life Sci 79, 243 527 (2022).
- 528 Danic, M. et al. Pharmacological Applications of Bile Acids and Their Derivatives in 8. 529 the Treatment of Metabolic Syndrome. Front Pharmacol 9, 1382 (2018).
- 530 9. Stojančević, M., Pavlović, N., Goločorbin-Kon, S. & Mikov, M. Application of bile 531 acids in drug formulation and delivery. Front. Life Sci. 7, 112-122
- 532 10. Surendran, P. et al. Rare and common genetic determinants of metabolic individuality 533 and their effects on human health. Nat Med 28, 2321–2332 (2022).
- 534 Bomba, L. et al. Whole-exome sequencing identifies rare genetic variants associated 11. 535 with human plasma metabolites. Am. J. Hum. Genet. 109, 1038–1054 (2022).
- 536 12. Demirkan, A. et al. Insight in genome-wide association of metabolite quantitative traits 537 by exome sequence analyses. PLoS Genet 11, e1004835 (2015).
- 538 13. Kettunen, J. et al. Genome-wide study for circulating metabolites identifies 62 loci and 539 reveals novel systemic effects of LPA. Nat. Commun. 7, 1–9 (2016).
- 540 14. Lotta, L. A. et al. A cross-platform approach identifies genetic regulators of human 541 metabolism and health. Nat. Genet. 53, 54–64 (2021).
- 542 15. Shin, S. Y. et al. An atlas of genetic influences on human blood metabolites. Nat. 543 Genet. 46, 543-550 (2014).
- 544 16. Staley, J. R. et al. PhenoScanner: A database of human genotype-phenotype associations. Bioinformatics 32, 3207–3209 (2016). 545
- 546 17. Kamat, M. A. et al. PhenoScanner V2: an expanded tool for searching human 547 genotype-phenotype associations. *Bioinformatics* 35, 4851-4853 (2019).
- 548 18. Elsworth, B. et al. The MRC IEU OpenGWAS data infrastructure. bioRxiv 549 2020.08.10.244293 (2020). doi:10.1101/2020.08.10.244293
- 550 19. Hagenbuch, B. & Meier, P. J. Organic anion transporting polypeptides of the OATP/ 551 SLC21 family: phylogenetic classification as OATP/ SLCO superfamily, new 552 nomenclature and molecular/functional properties. Pflugers Arch 447, 653-665 553 (2004).
- 554 20. Ho, R. H. & Kim, R. B. Transporters and drug therapy: implications for drug disposition and disease. Clin Pharmacol Ther 78, 260–277 (2005). 555
- 556 21. International Transporter, C. et al. Membrane transporters in drug development. Nat 557 Rev Drug Discov 9, 215–236 (2010).
- Niemi, M., Pasanen, M. K. & Neuvonen, P. J. Organic anion transporting polypeptide 558 22.

It is made available under a CC-BY 4.0 International license .
--

559		1B1: a genetically polymorphic transporter of major importance for hepatic drug
560		uptake. <i>Pharmacol Rev</i> 63, 157–181 (2011).
561	23.	Nies, A. T., Schwab, M. & Keppler, D. Interplay of conjugating enzymes with OATP
562		uptake transporters and ABCC/MRP efflux pumps in the elimination of drugs. Expert
563		<i>Opin Drug Metab Toxicol</i> <b>4</b> , 545–568 (2008).
564	24.	Revez, J. A. et al. Genome-wide association study identifies 143 loci associated with
565		25 hydroxyvitamin D concentration. Nat Commun 11, 1647 (2020).
566	25.	Willer, C. J. et al. Discovery and refinement of loci associated with lipid levels. Nat.
567		Genet. 45, 1274–1285 (2013).
568	26.	Johnson, A. D. <i>et al.</i> Genome-wide association meta-analysis for total serum bilirubin
569		levels. Hum. Mol. Genet. 18, 2700–2710 (2009).
570	27.	Ruth, K. S. <i>et al.</i> Using human genetics to understand the disease impacts of
571	_ / ·	testosterone in men and women. <i>Nat Med</i> <b>26</b> , 252–258 (2020)
572	28	Ochoa, D. <i>et al.</i> Open Targets Platform: supporting systematic drug-target
573	20.	identification and prioritisation <i>Nucleic Acids Res</i> <b>49</b> D1302–D1310 (2021)
574	29	Vitart V <i>et al.</i> 3000 years of solitude: extreme differentiation in the island isolates of
575	27.	Dalmatia Croatia Fur I Hum Ganat 14 478 487 (2006)
576	20	Zuk O at al Sourching for missing horitability: Designing rare variant association
570	30.	studios (2014), doi:10.1072/ppos.1222562111
577	21	Koh M V Lamos Ir D Liu V & Dowis G The hypoxic associated factor
570	51.	Koll, M. T., Lenios JI., K., Liu, A. & Powis, G. The hypoxia-associated factor
5/9		switches cells from HIF-faipha- to HIF-2aipha-dependent signaling promoting stem
580		cell characteristics, aggressive tumor growth and invasion. <i>Cancer Res</i> 71, 4015–4027
581	22	
582	32.	Semenza, G. L. Hypoxia, clonal selection, and the role of HIF-1 in tumor progression.
583	22	Crit Rev Biochem Mol Biol <b>35</b> , /1–103 (2000).
584	33.	Braun, T., Voland, P., Kunz, L., Prinz, C. & Gratzl, M. Enterochromaffin cells of the
585		human gut: sensors for spices and odorants. <i>Gastroenterology</i> <b>132</b> , 1890–1901 (2007).
586	34.	Erspamer, V. Pharmacology of indole-alkylamines. <i>Pharmacol Rev</i> <b>6</b> , 425–487 (1954).
587	35.	Watanabe, H. <i>et al.</i> Peripheral serotonin enhances lipid metabolism by accelerating
588		bile acid turnover. <i>Endocrinology</i> <b>151</b> , 4776–4786 (2010).
589	36.	Oakley, F. et al. Hepatocytes express nerve growth factor during liver injury: evidence
590		for paracrine regulation of hepatic stellate cell apoptosis. <i>Am J Pathol</i> <b>163</b> , 1849–1858
591		(2003).
592	37.	Ohkubo, T. et al. Early induction of nerve growth factor-induced genes after liver
593		resection-reperfusion injury. J Hepatol 36, 210–217 (2002).
594	38.	Valdovinos-Flores, C. & Gonsebatt, M. E. Nerve growth factor exhibits an antioxidant
595		and an autocrine activity in mouse liver that is modulated by buthionine sulfoximine,
596		arsenic, and acetaminophen. Free Radic Res 47, 404-412 (2013).
597	39.	Gigliozzi, A. <i>et al.</i> Nerve growth factor modulates the proliferative capacity of the
598		intrahepatic biliary epithelium in experimental cholestasis. <i>Gastroenterology</i> <b>127</b> ,
599		1198–1209 (2004).
600	40.	Rasi, G. et al. Nerve growth factor involvement in liver cirrhosis and hepatocellular
601		carcinoma. World J Gastroenterol <b>13</b> , 4986–4995 (2007).
602	41.	Tokusashi, Y. <i>et al.</i> Expression of NGF in hepatocellular carcinoma cells with its
603		receptors in non-tumor cell components. Int I Cancer 114 39–45 (2005)
604	42	Chen L et al Genetic and Microbial Associations to Plasma and Fecal Rile Acids in
605	. 2.	Obesity Relate to Plasma Lipids and Liver Fat Content Coll Ron <b>33</b> 108212 (2020)
606	43	Phelos T Snyder F Rodriguez F Child H & Harvey P The influence of
607	-IJ.	hiological sex and sex hormones on hile acid synthesis and cholesterol homeosteric
608		Biol Say Differ 2010 101 10 1 12 (2010)
000		$D(0), D(1, D(1), 201) + 101 + 0, 1^{-1} - 12 (201).$

609 Li-Hawkins, J. et al. Cholic acid mediates negative feedback regulation of bile acid 44. 610 synthesis in mice. J Clin Invest 110, 1191–1200 (2002). 611 45. Abu-Hayyeh, S. et al. Intrahepatic cholestasis of pregnancy levels of sulfated 612 progesterone metabolites inhibit farnesoid X receptor resulting in a cholestatic 613 phenotype. *Hepatology* 57, 716–726 (2013). 614 Frommherz, L. et al. Age-Related Changes of Plasma Bile Acid Concentrations in 46. 615 Healthy Adults--Results from the Cross-Sectional KarMeN Study. PLoS One 11, 616 e0153959 (2016). 617 47. Dekkers, K. F. et al. An online atlas of human plasma metabolite signatures of gut 618 microbiome composition. Nat. Commun. 2022 131 13, 1-12 (2022). 619 48. Russell, D. W. The Enzymes, Regulation, and Genetics of Bile Acid Synthesis. Annu. Rev. Biochem. 72, 137-174 (2003). 620 621 49. Thomas, C. E. et al. Association between Pre-Diagnostic Serum Bile Acids and 622 Hepatocellular Carcinoma: The Singapore Chinese Health Study. Cancers (Basel) 13, 623 (2021). Manzotti, C., Casazza, G., Stimac, T., Nikolova, D. & Gluud, C. Total serum bile acids 624 50. 625 or serum bile acid profile, or both, for the diagnosis of intrahepatic cholestasis of 626 pregnancy. Cochrane Database Syst Rev 7, CD012546 (2019). 627 51. Scherer, M., Gnewuch, C., Schmitz, G. & Liebisch, G. Rapid quantification of bile 628 acids and their conjugates in serum by liquid chromatography-tandem mass spectrometry. J. Chromatogr. B Anal. Technol. Biomed. Life Sci. 877, 3920-3925 629 630 (2009).631 52. Hadfield, J. D. MCMC Methods for Multi-Response Generalized Linear Mixed 632 Models: TheMCMCglmmRPackage. J. Stat. Softw. 33, (2010). 633 McCarthy, S. et al. A reference panel of 64,976 haplotypes for genotype imputation. 53. 634 Nat. Genet. 48, 1279-1283 (2016). Winkler, T. W. et al. Quality control and conduct of genome-wide association meta-635 54. analyses. Nat. Protoc. 9, 1192-1212 (2014). 636 637 Willer, C. J., Li, Y. & Abecasis, G. R. METAL: Fast and efficient meta-analysis of 55. 638 genomewide association scans. Bioinformatics 26, 2190-2191 (2010). 639 Karssen, L. C., van Duijn, C. M. & Aulchenko, Y. S. The GenABEL Project for 56. 640 statistical genomics. F1000Research 5, (2016). 641 Haller, T., Kals, M., Esko, T., Mägi, R. & Fischer, K. RegScan: A GWAS tool for 57. 642 quick estimation of allele effects on continuous traits and their combinations. Brief. 643 Bioinform. 16, 39-44 (2013). 644 58. Hemani, G. et al. The MR-base platform supports systematic causal inference across 645 the human phenome. *Elife* 7, (2018). 646 59. Van Hout, C. V et al. Exome sequencing and characterization of 49,960 individuals in 647 the UK Biobank. Nature 586, 749–756 (2020). Reid, J. G. et al. Launching genomics into the cloud: deployment of Mercury, a next 648 60. 649 generation sequence analysis pipeline. BMC Bioinformatics 15, 30 (2014). 650 61. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler 651 transform. *Bioinformatics* 25, 1754–1760 (2009). 652 "Picard Toolkit." 2019. Broad Institute, GitHub Repository. 62. https://broadinstitute.github.io/picard/; Broad Institute. 653 654 PLC, G. weCall. (2018). 63. 655 Lin, M. F. et al. GLnexus: joint variant calling for large cohort sequencing. bioRxiv 64. 656 (2018). doi:10.1101/343970 Cingolani, P. et al. A program for annotating and predicting the effects of single 657 65. nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster 658

659		strain w1118; iso-2; iso-3. Fly 6, 80–92 (2012).
660	66.	Zerbino, D. R. et al. Ensembl 2018. Nucleic Acids Res. 46, D754–D761 (2018).
661	67.	Liu, X., Wu, C., Li, C. & Boerwinkle, E. dbNSFP v3.0: A One-Stop Database of
662		Functional Predictions and Annotations for Human Nonsynonymous and Splice-Site
663		SNVs. Hum Mutat 37, 235–241 (2016).
664	68.	Liu, X., Jian, X. & Boerwinkle, E. dbNSFP: a lightweight database of human
665		nonsynonymous SNPs and their functional predictions. Hum Mutat 32, 894–899
666		(2011).
667	69.	Vaser, R., Adusumalli, S., Leng, S. N., Sikic, M. & Ng, P. C. SIFT missense
668		predictions for genomes. Nat Protoc 11, 1–9 (2016).
669	70.	Adzhubei, I. A. et al. A method and server for predicting damaging missense
670		mutations. Nat. Methods 7, 248–249 (2010).
671	71.	Chun, S. & Fay, J. C. Identification of deleterious mutations within three human
672		genomes. Genome Res 19, 1553–1561 (2009).
673	72.	Schwarz, J. M., Rodelsperger, C., Schuelke, M. & Seelow, D. MutationTaster
674		evaluates disease-causing potential of sequence alterations. <i>Nat Methods</i> <b>7</b> , 575–576
675		(2010).
676	73.	Chen, H. et al. Efficient Variant Set Mixed Model Association Tests for Continuous
677		and Binary Traits in Large-Scale Whole-Genome Sequencing Studies. Am. J. Hum.
678		Genet. 104, 260–274 (2019).

679

#### 680 Acknowledgments

The Orkney Complex Disease Study (ORCADES) was supported by the Chief Scientist Office 681 of the Scottish Government (CZB/4/276, CZB/4/710), a Royal Society URF to J.F.W., the 682 MRC Human Genetics Unit quinquennial programme "QTL in Health and Disease", Arthritis 683 Research UK and the European Union framework program 6 EUROSPAN project (contract 684 no. LSHG-CT-2006-018947). DNA extractions were performed at the Edinburgh Clinical 685 Research Facility, University of Edinburgh. We would like to acknowledge the invaluable 686 687 contributions of the research nurses in Orkney, the administrative team in Edinburgh and the 688 people of Orkney. For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from this 689 submission. The CROATIA-VIS study in the Croatian island of Vis was supported through the 690 grants from the Medical Research Council UK and Ministry of Science, Education and Sport 691 692 of the Republic of Croatia (number 108-1080315-0302). The authors collectively thank a large 693 number of individuals for their individual help in organising, planning and carrying out the 694 field work related to the project and data management: Professor Pavao Rudan and the staff of 695 the Institute for Anthropological Research in Zagreb, Croatia (organisation of the field work, anthropometric and physiological measurements, and DNA extraction): Professor Ariana 696 697 Vorko-Jovic and the staff and medical students of the Andrija Stampar School of Public Health 698 of the Faculty of Medicine, University of Zagreb, Croatia (questionnaires, genealogical reconstruction and data entry); Dr Branka Salzer from the biochemistry lab "Salzer", Croatia 699 700 (measurements of biochemical traits); local general practitioners and nurses (recruitment and 701 communication with the study population); and the employees of several other Croatian institutions who participated in the field work, including but not limited to the University of 702 703 Rijeka and Split, Croatia; Croatian Institute of Public Health; Institutes of Public Health in Split 704 and Dubrovnik, Croatia. SNP Genotyping of the Vis samples was carried out by the Genetics 705 Core Laboratory at the Wellcome Trust Clinical Research Facility, WGH, Edinburgh. The 706 MICROS (Micro-Isolates in South Tyrol) study is part of the genomic health care program 'GenNova' and was carried out in three villages of the Val Venosta on the populations of 707 708 Stelvio, Vallelunga and Martello. We thank the primary care practitioners Raffaela Stocker, 709 Stefan Waldner, Toni Pizzecco, Josef Plangger, Ugo Marcadent and the personnel of the 710 Hospital of Silandro (Department of Laboratory Medicine) for their participation and 711 collaboration in the research project. In South Tyrol, the study was supported by the Ministry 712 of Health and Department of Educational Assistance, University and Research of the 713 Autonomous Province of Bolzano and the South Tyrolean Sparkasse Foundation. The Northern 714 Swedish Population Health Study (NSPHS) was funded by the Swedish Medical Research Council (project number K2007-66X-20270-01-3), and the Foundation for Strategic Research 715 (SSF). The NSPHS as part of EUROSPAN (European Special Populations Research Network) 716 717 was also supported by European Commission FP6 STRP grant number 01947 (LSHGCT-2006-718 01947). This work was also supported by the Swedish Society for Medical Research (ÅJ). The 719 authors are grateful for the contribution of district nurse Svea Hennix for data collection and 720 Inger Jonasson for logistics and coordination of the health survey. Finally, the authors thank 721 all the community participants for their interest and willingness to contribute to the study. The 722 Erasmus Rucphen Family (ERF) study was supported by grants from The Netherlands

723 Organisation for Scientific Research (NWO), Erasmus MC, the Centre for Medical Systems 724 Biology (CMSB) and the European Community's Seventh Framework Programme (FP7/2007-2013), ENGAGE Consortium, grant agreement HEALTH-F4-2007- 201413. We are grateful 725 to all general practitioners for their contributions. Cornelia van Duijn and Ben Oostra for 726 727 setting-up the ERF study, Petra Veraart for sorting out the genealogy records, Jeannette 728 Vergeer and Peter Snijders for help in retrieving the materials needed to analyse data. We 729 acknowledge support from the European Union's Horizon 2020 research and innovation 730 programme IMforFUTURE (A.L.: H2020-MSCA-ITN/721815); the RCUK Innovation 731 Fellowship from the National Productivity Investment Fund (L.K.: MR/R026408/1) and the MRC Human Genetics Unit programme grant, 'QTL in Health and Disease' (J.F.W. and C.H.: 732 733 MC UU 00007/10).

- 734
- **Ethics** 735

736 All studies were approved by local research ethics committees and all participants have given written informed consent. The ORCADES study was approved by the NHS Orkney Research 737 738 Ethics Committee and the North of Scotland REC. The CROATIA-Vis study was approved by 739 the ethics committee of the medical faculty in Zagreb and the Multi-Centre Research Ethics Committee for Scotland. The Northern Swedish Population Health Study (NSPHS) was 740 741 approved by the local ethics committee at the University of Uppsala (Regionala 742 Etikprövningsnämnden, Uppsala). The MICROS study was approved by the ethical committee 743 of the Autonomous Province of Bolzano, Italy. The ERF study was approved by the Erasmus 744 institutional medical-ethics committee in Rotterdam, The Netherlands.

745

## 746 **Author contributions**

A.L.: Data analysis and interpretation, visualisation, writing—original draft preparation, 747 writing—review and editing. D.G.-S.: Data analysis, writing—review and editing. Å.J.: Data 748 analysis. S.A.: Data analysis. G.L.: Quantification of bile acids, writing-original draft 749 750 preparation. C.G.: Quantification of bile acids, writing-original draft preparation. G.T.: 751 preparation, quality control and annotation of whole-exome sequencing data. A.R.S.: Funding, 752 writing-review and editing. A.A.H.: Genomic and demographic data provider for MICROS 753 cohort. P.P.: Funding, genomic and demographic data provider for MICROS cohort. C.P.: genomic and demographic data provider for MICROS cohort. H.C. Funding. O.P.: Genomic 754 and demographic data provider for CROATIA-Vis cohort. C.H.: Funding, genomic and 755 756 demographic data provider for CROATIA-Vis cohort. N.P.: supervision and data interpretation 757 for bile acid pre-processing and imputation. M.G.: Genomic and demographic data provider 758 for ERF cohort, writing-review and editing. U.G.: Funding, genomic and demographic data provider for NSPHS cohort. C.F.: Genomic and demographic data provider for MICROS 759 760 cohort. J.F.W.: Funding, conceptualisation, genomic and demographic data provider for ORCADES cohort, supervision, data interpretation, writing-original draft preparation, 761

writing—review and editing. L.K.: Conceptualisation, supervision, data interpretation,
writing—original draft preparation, writing—review and editing.

764

# 765 Competing interests

G.T. and A.R.S. are full-time employees of Regeneron Genetics Center and receive salary,stock and stock options as compensation. L.K. is an employee of Humanity Inc., a company

768 developing direct-to-consumer measures of biological ageing. All other authors declare no

competing interests.