

Links between gut microbiome composition and fatty liver disease in a large population sample

Matti O. Ruuskanen^{1,2*}, Fredrik Åberg^{3,4}, Ville Männistö^{5,6}, Aki S. Havulinna^{2,7}, Guillaume Méric^{8,9}, Yang Liu^{8,10}, Rohit Loomba^{11,12}, Yoshiki Vázquez-Baeza^{13,14}, Anupriya Tripathi^{15,16,17}, Liisa M. Valsta², Michael Inouye^{8,18}, Pekka Jousilahti², Veikko Salomaa², Mohit Jain^{12,19}, Rob Knight^{13,14,20,21}, Leo Lahti²², Teemu J. Niiranen^{1,2,23}

¹Department of Internal Medicine, University of Turku, Turku, Finland

²Department of Public Health Solutions, Finnish Institute for Health and Welfare, Helsinki, Finland

³Transplantation and Liver Surgery Clinic, Helsinki University Hospital, University of Helsinki, Helsinki, Finland

⁴The Transplant Institute, Sahlgrenska University Hospital, Gothenburg, Sweden

⁵Department of Medicine, Kuopio University Hospital, University of Eastern Finland, Kuopio, Finland

⁶Department of Experimental Vascular Medicine, Amsterdam UMC, University of Amsterdam, Amsterdam, The Netherlands

⁷Institute for Molecular Medicine Finland, FIMM - HiLIFE, Helsinki, Finland

⁸Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia

⁹Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria, Australia

¹⁰Department of Clinical Pathology, The University of Melbourne, Melbourne, Victoria, Australia

¹¹Department of Medicine, NAFLD Research Center, La Jolla, CA, USA

¹²Department of Medicine, University of California, San Diego, La Jolla, CA, USA

¹³Jacobs School of Engineering, University of California, San Diego, La Jolla, CA, USA

¹⁴Center for Microbiome Innovation, University of California San Diego, La Jolla, California, USA

¹⁵Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego, La Jolla, California, USA

¹⁶Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego, La Jolla, California, USA

¹⁷Division of Biological Sciences, University of California, San Diego, La Jolla, California, USA

¹⁸Department of Public Health and Primary Care, Cambridge University, Cambridge, United Kingdom

¹⁹Department of Pharmacology, University of California San Diego, La Jolla, California, USA

²⁰Department of Pediatrics, School of Medicine, University of California San Diego, La Jolla, California, USA

²¹Department of Computer Science & Engineering, University of California San Diego, La Jolla, California, USA

²²Department of Future Technologies, University of Turku, Turku, Finland

²³Division of Medicine, Turku University Hospital, Turku, Finland

*Correspondence: Matti Ruuskanen, matti.ruuskanen@utu.fi

Running head: Gut microbiome composition and fatty liver

40 **Abstract**

41 Fatty liver disease is the most common liver disease in the world. It is characterized by a build-
42 up of excess fat in the liver that can lead to cirrhosis and liver failure. The link between fatty
43 liver disease and gut microbiome has been known for at least 80 years. However, this association
44 remains mostly unstudied in the general population because of underdiagnosis and small sample
45 sizes. To address this knowledge gap, we studied the link between the Fatty Liver Index (FLI), a
46 well-established proxy for fatty liver disease, and gut microbiome composition in a
47 representative, ethnically homogeneous population sample in Finland. We based our models on
48 biometric covariates and gut microbiome compositions from shallow metagenome sequencing.
49 Our classification models could discriminate between individuals with a high FLI (≥ 60 ,
50 indicates likely liver steatosis) and low FLI (< 60) in our validation set, consisting of 30% of the
51 data not used in model training, with an average AUC of 0.75. In addition to age and sex, our
52 models included differences in 11 microbial groups from class *Clostridia*, mostly belonging to
53 orders *Lachnospirales* and *Oscillospirales*. Pathway analysis of representative genomes of the
54 FLI-associated taxa in (NCBI) *Clostridium* subclusters IV and XIVa indicated the presence of
55 *e.g.*, ethanol fermentation pathways. Through modeling the fatty liver index, our results provide
56 with high resolution associations between gut microbiota composition and fatty liver in a large
57 representative population cohort and support the role of endogenous ethanol producers in the
58 development of fatty liver.

59

60 **Keywords: Metagenomics, human gut, fatty liver, fatty liver index, population sample**

61 **Introduction**

62 Fatty liver disease affects roughly a quarter of the world's population.¹ It is characterized by
63 accumulation of fat in the liver cells and is intimately linked with pathophysiology of metabolic
64 syndrome.²⁻⁴ Fatty liver disease can be broadly divided into two variants: non-alcoholic fatty
65 liver disease (NAFLD), attributed to high caloric intake, and alcohol associated fatty liver
66 disease, attributed to high alcohol consumption. Even though the rate of progressions and
67 underlying causes of both diseases might be different, they can be broadly sub-divided into those
68 who have fat accumulation in the liver with no or minimal inflammation or those who have
69 additional features of cellular injury and active inflammation with or without fibrosis typically
70 seen in peri-sinusoidal area.⁵ Patients with steatohepatitis may progress to cirrhosis and
71 hepatocellular carcinoma and have increased risk of liver-related morbidity and mortality,
72 globally amounting to hundreds of thousands of deaths.⁶

73

74 The human gut harbors up to 10^{12} microbes per gram of content,⁷ and is intimately connected
75 with the liver. Thus, it is no surprise that gut microbiome composition appears to have a strong
76 connection with liver disease.⁸ Numerous studies over the past 80 years have reported
77 associations between gut microbial composition and liver disease.⁹ For example, gut
78 permeability and overgrowth of bacteria in the small intestine,¹⁰ changes in
79 *Gammaproteobacteria* and *Erysipelotrichi* abundance during choline deficiency,¹¹ elevated
80 abundance of ethanol-producing bacteria,^{12,13} metagenomic signatures of specific bacterial
81 species,^{14,15} have all been linked to NAFLD in small case-control patient samples. However, the
82 microbial signatures often overlap between NAFLD and metabolic diseases, while those of more
83 serious liver disease such as steatohepatitis and cirrhosis are more clear.¹⁶ For example, oral taxa

84 appear to invade the gut in liver cirrhosis,¹⁷ and this phenotype can accurately be detected by
85 analyzing the fecal microbiome composition (AUC = 0.87 in a validation cohort).⁸ Furthermore,
86 we recently demonstrated good prediction accuracy for incident liver disease diagnoses (AUC =
87 0.83 for non-alcoholic liver disease, AUC = 0.96 for alcoholic liver disease, during ~15 years),¹⁸
88 showing that the signatures of serious future liver disease are easy to detect.

89
90 The mechanisms underlying the contribution of gut microbiome content with fatty liver disease
91 are thought to be primarily linked to gut bacterial metabolism. Bacterial metabolites can indeed
92 be translocated from the gut through the intestinal barrier into the portal vein and transported to
93 the liver, where they interact with liver cells, and can lead to inflammation and steatosis.¹⁹ Short-
94 chain fatty acid production, conversion of choline into methylamines, modification of bile acids
95 (BA) into secondary BA, and ethanol production, all of which are mediated by gut bacteria, are
96 also known to be aggravating factors for NAFLD.¹⁹ Recent studies have also suggested that
97 endogenous ethanol production by gut bacteria could lead to an increase in gut membrane
98 permeability.¹³ This can facilitate the translocation of bacterial metabolites and cell components
99 such as lipopolysaccharides from the gut to the liver, leading to further inflammation and
100 possible development of NAFLD.²⁰

101
102 Liver biopsy assessment is the current gold standard for diagnosis of fatty liver disease and its
103 severity,²¹ but it is also impractical and unethical in a population-based setting. Ultrasound and
104 MRI based assessment can help detect presence of fatty liver, however, this data is not available
105 in our cohort. Regardless, recent studies have shown that indices based on anthropometric

106 measurements and standard blood tests can be a reliable tool for non-invasive diagnosis of fatty
107 liver particularly in population-based epidemiologic studies.^{22,23}
108
109 Here, we designed and conducted computational analyses to examine the links between fatty
110 liver and gut microbiome composition in a representative population sample of 7211 extensively
111 phenotyped Finnish individuals.²⁴ Because fatty liver disease is generally underdiagnosed in the
112 general population,²⁵ we used population-wide measurements of BMI, waist circumference,
113 blood triglycerides and gamma-glutamyl-transferase (GGT) to calculate a previously validated
114 Fatty Liver Index (FLI) for each participant as a proxy for fatty liver.²⁶ In parallel, we used
115 shallow shotgun sequencing to analyze gut microbiome composition,²⁷ which also enabled the
116 use of phylogenetic and pathway prediction methods. In this work, we describe high-resolution
117 associations between fatty liver and individual gut microbial taxa and clades, which are
118 generalizable at the population level.

119

120 **Results**

121 *Bacterial community structure is correlated with Fatty Liver Index in a population*

122 *sample*

123 To investigate the link between fatty liver disease (using FLI as a proxy; **Figures 1A, 1B**) and
124 gut microbial composition, we used linear regression (adjusted $R^2 = 0.29$) on the three first
125 principal component (PC) axes of the fecal bacterial beta-diversity (between individuals), sex,
126 age, and alcohol to model FLI. $\log_{10}(\text{FLI})$ significantly correlated with all three bacterial PC
127 axes, sex, age, and alcohol use (all $P < 1 \times 10^{-6}$). Correlations between FLI and archaeal PC axes
128 were not significant ($P > 0.05$). The effect size estimate on $\log_{10}(\text{FLI})$ was a magnitude larger for

129 PC1 (0.11 ± 0.008) than for PC2 (0.04 ± 0.008) and PC3 (-0.06 ± 0.008). The relationships
130 between FLI and the bacterial PC components representing their beta-diversity are visualized for
131 each of the three components in **Figure 1C**. In our analyses, we classified our reads against the
132 Genome Taxonomy Database (GTDB),²⁸ and thus the taxonomy discussed in this study follows
133 the standardized GTDB taxonomy, unless otherwise noted.

134

135 Bacterial clades with the highest positive loadings on PC1 (and therefore associated with higher
136 FLI values) included members of the *Lachnospirales* and *Oscillospirales* taxonomic orders, of
137 the *Bacilli* class, and of the *Ruminococcaceae*, *Bacteroidaceae* and *Lachnospiraceae* families
138 (**Figure S2**). These observations led us to further analyses within a machine learning framework
139 to estimate the relative contributions of individual bacterial taxa to the differences in FLI
140 between study participants.

141

142 ***Bacterial lineages within the NCBI Clostridium subclusters IV and XIVa associate*** 143 ***with FLI***

144 In our machine learning framework, we used the known covariates in addition to individual
145 archaeal and bacterial “balances” as the predicting features. Briefly, each balance represents a
146 single internal node in a phylogenetic tree, and its value is a log-ratio of the abundances of the
147 two clades descending this node (for details, see methods and ref. 29). Continuous FLI and
148 differences between FLI groups (FLI < 60, $N = 4359$ and FLI ≥ 60 , $N = 1910$; see **Figures 1A,**
149 **1B**) were modeled with gradient boosting regression or classification using Leave-One-Group-
150 Out Cross-Validation (LOGOCV) between participants from different regions.

151

152 After feature selection and Bayesian hyperparameter optimization, the correlation between the
153 predictions of the final regression models (age, sex, self-reported alcohol use, and 18 bacterial
154 balances as features; each trained on the data from 5/6 regions) and true values in unseen data
155 from the omitted region averaged $R^2 = 0.30$ (0.26 – 0.33). After feature selection and
156 optimization, the main classification models (age, sex, and 11 bacterial balances as features; each
157 trained on the data from 5/6 regions) averaged AUC = 0.75 (**Table S1**) and AUPRC = 0.56
158 (baseline at 0.30; **Table S2**) on (unseen) test data from the omitted region. Models trained using
159 only the covariates averaged AUC = 0.71 (AUPRC = 0.47) and using only the 11 bacterial
160 balances they averaged AUC = 0.66 (AUPRC = 0.47) on test data. Alternative models were
161 constructed by excluding participants with FLI between 30 and 60 (N = 1583) and discerning
162 between groups of $FLI < 30$ (N = 2776) and $FLI \geq 60$ (N = 1910). These models averaged AUC
163 = 0.80 (AUPRC = 0.75, baseline at 0.41) on their respective test data. They averaged AUC =
164 0.76 (AUPRC = 0.68) when using only the covariates, and AUC = 0.70 (AUPRC = 0.63) when
165 using only the 20 bacterial balances.

166

167 Because training data from all 6 regions was used to prevent overfitting in the selection of core
168 features for all of the models, and similarly in searching for common hyperparameters,
169 participants from the validation region of each model (in the training partition) partly influenced
170 these parameters. Thus, we also constructed classification models discerning between the $FLI <$
171 60 and $FLI \geq 60$ groups, where data of the validation region was completely excluded in the
172 feature selection and hyperparameter optimization of each LOGOCV model. These models,
173 using their individual feature sets and hyperparameters, averaged AUC = 0.75 and AUPRC =
174 0.57 (baseline at 0.30) on test data from their respective validation regions (**Table S3**). Using

175 only covariates, they averaged AUC = 0.71 (AUPRC = 0.47), and AUC = 0.67 (AUPRC = 0.48)
176 with only the bacterial balances.

177

178 To facilitate interpretability of the results, we primarily continued examining the main
179 classification models using a common set of core features. In these models, the median effect
180 sizes of the features on the model predictions at their minimum and maximum values were
181 highest for age, followed by sex, and the 11 balances in the phylogenetic tree (**Figures S2, S3**).
182 All 11 associated balances were in phylum *Firmicutes*, class *Clostridia*, and largely in the NCBI
183 *Clostridium* subclusters IV and XIVa (**Figure 2**). The specific taxa represented standardized
184 GTDB genera (NCBI in brackets) *Negativibacillus* (*Clostridium*), *Clostridium M*
185 (*Lachnoclostridium* / *Clostridium*), *CAG-81* (*Clostridium*), *Dorea* (*Merdimonas* / *Mordavella* /
186 *Dorea* / *Clostridium* / *Eubacterium*), *Faecalicatena* (*Blautia* / *Ruminococcus* / *Clostridium*),
187 *Blautia* (*Blautia*), *Sellimonas* (*Sellimonas* / *Drancourtella*), *Clostridium Q* (*Lachnoclostridium*
188 [*Clostridium*]) and *Tyzzarella* (*Tyzzarella* / *Coprococcus*). Notably, all but one of the features in
189 the main classification models (n226) were identified in the feature selection for the alternative
190 models (constructed otherwise identically, but $FLI < 30$ was compared against $FLI \geq 60$ in
191 different data partitions), together with 10 additional balances (**Figure S4**). Only one of the
192 balances in the alternative models was outside phylum *Firmicutes* (n1712 in *Bacteroidota*), and
193 in addition, 4 balances were outside class *Clostridia* (n481 in *Negativicutes*; n826, n1009 and
194 n918 in *Bacilli*).

195

196 In addition to blood test results, FLI is based on anthropometric markers linked to metabolic
197 syndrome, waist circumference and BMI. This prompted us to attempt to dissect the index and

198 identify which of the covariates and associated microbial balances from the phylogenetic tree can
199 be linked to blood GGT and triglycerides measurements (see **Figure 1B**), and therefore would be
200 more specific to hepatic steatosis and liver damage. To do so, we performed feature selection
201 (similarly to continuous FLI) for GGT and triglycerides measurements in subsets of participants
202 grouped by age, sex, and BMI. The feature selection identified two balances within the NCBI
203 *Clostridia* XIVa subcluster (identified as n336 and n330) which were important for both GGT
204 and triglyceride level prediction, and thus likely specific to liver function (**Figure 2**). Bacterial
205 taxa were positively linked to liver function in these balances, and included (NCBI species)
206 *Clostridium clostridioforme*, *C. bolteae*, *C. citroniae*, *C. saccharolyticum* and *C. symbiosum*.
207

208 ***Ethanol and acetate production pathways are identified in representative bacterial***
209 ***genomes from taxa linked to liver function***

210 The values of predictive balances in the phylogenetic tree cannot be summarized for individual
211 taxa, which means that only a qualitative investigation of the associations between their
212 metabolism and fatty liver was possible in this study. We identified genetic pathways predicted
213 to encode for SCFA (acetate, propanoate, butanoate) and ethanol production, BA metabolism,
214 and choline degradation to trimethylamine in representative genomes from the taxa we identified
215 to be linked to liver function (**Figure S3**). These specific processes were chosen because they
216 have been previously identified to have a mechanistic link to NAFLD (see *e.g.*, ref. 19).

217
218 Acetate and ethanol production pathways appeared to be more abundant in the representative
219 genomes of the taxa which had a positive association with FLI. In the liver function specific
220 clades, n336 and n330, MetaCyc pathways for pyruvate fermentation to ethanol III (PWY-6587)

221 and L-glutamate degradation V (via hydroxyglutarate; P162-PWY; produces acetate and
222 butanoate) were present only in genomes positively associated with FLI. In balance n336, also
223 heterolactic fermentation (P122-PWY; produces ethanol and lactate) was more often encoded in
224 the FLI-associated clade (3/5) than the opposing clade (1/2). In representative genomes from the
225 non-liver-specific balance n355, potential ethanol producers (PWY-6587) were seen in the
226 positively associated clade, but for most balances such trends were not clear in the qualitative
227 analysis. Furthermore, we did not detect any of these pathways in the representative genomes of
228 two individual taxa positively associated with FLI, *Negativibacillus sp000435195* and *Phoceia*
229 *massiliensis* (**Figure S3**).

230

231 **Discussion**

232 The pathophysiology of fatty liver disease in general, and NAFLD in particular, is complex and
233 its clinical diagnosis can be difficult.³⁰ In this study, we leveraged multi-omics data from a
234 large population study (FINRISK02) to identify broad links between the overall gut
235 microbiome composition and fatty liver disease, using FLI as a recognized proxy (**Figure 1C**),
236 and identified specific microbial taxa and lineages predictive of a higher FLI (**Figure 2**).

237 Considering that the predictive ability of FLI for clinically diagnosed NAFLD ranges between
238 AUC = 0.81 – 0.93, in populations of Caucasian ethnicity such as ours,²³ our models were able
239 to reasonably predict the FLI group with AUC = 0.75 (AUPRC = 0.56, baseline at 0.30), while
240 extrapolating to a validation region not used in training of the model.

241

242 Our additional analyses support these results. Excluding participants with intermediate FLI
243 (between 30 – 60) increased the accuracy slightly (to AUC = 0.8 and AUPRC = 0.75, baseline

244 at 0.41). However, discerning between participants with probable fatty liver disease (FLI \geq 60)
245 from others presents a clinically more relevant target for detecting changes in microbiome
246 composition associated with development of the disease. In another set of models, we negated
247 the influence of validation region data in the individual models also for feature selection and
248 hyperparameter optimization during training. This led to individualized sets of features and
249 parameters in the models, but the average performance of the models was almost identical on
250 validation region samples in the test data (AUC = 0.75 and AUPRC 0.57, baseline at 0.30). The
251 aim of our study was to find patterns in microbiome composition which would be generalizable
252 across the 6 sampled geographic regions in Finland and easy to interpret. Thus, we consider the
253 use of all training data to define the common core feature set justified. This goal also guided
254 our overall modeling architecture and likely led to a lower performance than if we instead
255 performed interpolation within a smaller scale (see *e.g.*, ref. 31).

256

257 When interpreting results, several different levels of associations can be considered according
258 to types of fatty liver disease and the gut microbiome composition. Because FLI has been
259 mostly validated with simple steatosis and NAFLD,^{23,26} we can conservatively contextualize
260 our findings with previous associative work that used these diagnoses or clinical
261 manifestations, only.

262

263 *FLI modeling reveals consistent associations between gram-positive Clostridia and* 264 *fatty liver disease*

265 We found significant linear correlations between the first three bacterial PC-axes of our
266 samples (a measure of beta diversity) and FLI (see results and **Figure 1C**). Previous studies

267 have shown differences in beta diversity in relation to NAFLD.³² However, FLI used in our
268 study as a proxy for liver disease also includes features such as BMI and waist circumference,
269 which associate with metabolic syndrome and type 2 diabetes.¹⁶ Links between these diseases
270 and gut microbiome composition are well documented in previous studies.³³ It is thus not
271 surprising that bacterial beta diversity and FLI were correlated, but unfortunately this simple
272 correlation does not enable untangling the relative contributions of fatty liver disease and other
273 metabolic diseases to the differences in bacterial beta diversity.

274

275 Several studies have reported highly specific changes in microbial abundances in relation to
276 NAFLD.^{12,34–36} In summary, while also conflicting results have been reported, generally
277 increases in *Lactobacillus* and *Escherichia* genera, and a decrease in *Coprococcus* genus have
278 been most often associated with a NAFLD diagnosis.³⁷ Furthermore, increased abundance of
279 several gram-positive bacteria belonging to the *Clostridium* genus have often been positively
280 linked with NAFLD.^{14,38} Differences in unconstrained between-samples (beta) diversity have
281 been also documented for persistent NAFLD,³² and along the NAFLD-cirrhosis spectrum.⁸

282

283 In our study, abundances of bacteria from the *Coprococcus* genus were not specifically
284 associated with FLI, although the genus was nested inside our predictive balances. Strikingly,
285 we did not identify any bacterial associations with FLI outside of the *Firmicutes* phylum. A
286 possible reason for this might be the higher relative abundance of phylum *Firmicutes* at high
287 latitudes, where Finland is.³⁹ Among the associations we identified, *Faecalicatena gnavus*
288 (NCBI: *Ruminococcus gnavus*) was positively linked with FLI as part of 3 predictive balances,
289 and associated in previous studies with liver cirrhosis.¹⁷ Interestingly, none of the oral

290 *Firmicutes*, such as *Veillonella*, suggested to invade the gut, were identified in our own
291 analyses. This might be caused by using FLI as a proxy, which is likely not closely associated
292 with advanced liver disease, such as cirrhosis, and thus would target an earlier phase of liver
293 disease progression.

294

295 Two individual taxa, *Negativibacillus sp000435195* and *Phoceia massiliensis*, were highly
296 predictive of FLI group (**Figures 2, S2**), but not of its liver function-specific components. The
297 associations of these taxa with fatty liver disease have not been documented previously.

298 However, a decreasing abundance of both bacteria, *Negativibacillus sp000435195* (NCBI:
299 *Clostridium sp. CAG:169*) and *Phoceia massiliensis* (NCBI: *Phoceia massiliensis*), were seen
300 when the intake of meat and refined cereal was reduced isocalorically in favor of fruit,
301 vegetables, wholegrain cereal, legumes, fish and nuts in overweight and obese subjects in
302 Italy.⁴⁰ While comparisons between these studies are difficult due to annotation, bacteria such
303 as *Faecalicatena gnavus* (NCBI: *Ruminococcus gnavus*) and *Clostridium Q saccharolyticum*
304 (NCBI: *Clostridium saccharolyticum*) were also found to respond negatively to the
305 Mediterranean diet. Together with their positive association with FLI in our study, these
306 observations would warrant further study on these species as plausible biomarkers for healthy
307 diet choices.

308

309 Most taxa in our study with a positive association with FLI belonged to the (broadly defined)
310 *Clostridium* NCBI genus, which supports several previous observations.^{14,38} However,
311 taxonomic standardization according to GTDB has identified the *Clostridium* genus as the most
312 phylogenetically inconsistent of all bacterial genera in the NCBI taxonomy, and divides it into

313 a total of 121 monophyletic genera in 29 distinct families.²⁸ These reassignments, although
314 more accurate and sensible, complicate comparisons to previous research studies. However, our
315 results strongly suggest that this finer taxonomic resolution might robustly reveal novel
316 discoveries. Thus, while (shallow) shotgun metagenomic sequencing is often more costly than
317 amplicon sequencing, this might be justified by the increased resolution which is required to
318 accurately identify specific taxon-based associations (see *e.g.*, refs. 27,41).

319

320 ***Bacterial taxa associated with a high FLI have a genetic potential to exacerbate the***
321 ***development of fatty liver disease***

322 We identified several plausible new associations between individual taxa and clades of bacteria
323 and fatty liver. All taxa were from class *Clostridia*, which are obligate anaerobes. We observed
324 that reference genomes from the bacterial taxa positively associated with FLI in the liver-
325 specific balances harbored several genetic pathways necessary for ethanol production.
326 Specifically, genes predicted to enable the fermentation of pyruvate to ethanol (MetaCyc PWY-
327 6587) appeared to be common. Endogenous production of ethanol has been known to both
328 induce hepatic steatosis and increase intestinal permeability,⁴² and several of the taxa identified
329 in our study have also been experimentally shown to produce ethanol, such as *C. M*
330 *asparagiforme*, *C. M bolteae*, *C. M clostridioforme* / *C. M clostridioforme A*⁴³, and *C. Q*
331 *Saccharolyticum*.⁴⁴ The relative abundances of these putatively ethanol-producing taxa were
332 predictive of FLI groups in previously unseen data. However, the self-reported alcohol
333 consumption from the participants was not among the best predictors for the FLI groups, as it
334 was excluded in the feature selection step.

335

336 All reference genomes from taxa positively associated with FLI in balance n330 harboured
337 genes predicted to encode for the L-glutamate fermentation V (P162-PWY; **Figure S3**)
338 pathway, which results in the production of acetate and butanoate. Glutamate fermentation
339 could lead to increased microbial protein fermentation in the gut, which has been previously
340 been linked with obesity, diabetes and NAFLD.⁴⁵ Recently, the combined intake of fructose
341 and microbial acetate production in the gut was experimentally observed to contribute to
342 lipogenesis in the liver in a mouse model.⁴⁶ Interestingly, *C. Q saccharolyticum* (in our study, a
343 FLI-associated species deriving from balance n330), was experimentally shown to ferment
344 various carbohydrates, including fructose, to acetate, hydrogen, carbon dioxide, and ethanol.⁴⁴
345 Furthermore, while our own pathway analysis did not detect BA modification pathways in the
346 reference genome of *C. Q saccharolyticum*, a strain of this species has been highlighted as a
347 probable contributor to NAFLD development through the synthesis of secondary BA.¹⁵ The
348 links between dietary intake and gene regulation, combined with microbial fermentation in the
349 gut warrant further mechanistic experiments to elucidate their links with fatty liver, and likely
350 other metabolic diseases.

351

352 Intriguingly, NAFLD-associated ethanol producing bacteria in previous cohort studies have all
353 been gram-negatives, such as (NCBI-defined) *Klebsiella pneumoniae*,¹³ and *Escherichia coli*.¹²
354 In our population sample, instead of gram-negatives, bacteria from the *C. M bolteae*, *C. M*
355 *clostridioforme* / *C. M clostridioforme A* and *C. M citroniae* species (linked in our study with
356 FLI as deriving from balance n336) have been described as opportunistic pathogens,⁴⁷ and are
357 hypothesized to exacerbate fatty liver development similarly through endogenous ethanol
358 production. This result suggests that geographical,³¹ and ethnic variability,⁴⁸ might also

359 strongly affect gut microbiome composition and its associations with disease. In addition to
360 putative endogenous ethanol producers, we identified other FLI-associated taxa deriving from
361 balance n330, for which reference genomes harbored a genetic pathway predicted to encode for
362 the ability to ferment L-lysine to acetate and butyrate. While the production of these SCFAs is
363 often considered beneficial for gut health, other metabolism of proteolytic bacteria might
364 negatively contribute to fatty liver disease.⁴⁹

365

366 Through modeling a previously validated risk index for fatty liver, we could associate specific
367 members of the gut microbiome with the disease across geographical regions in this
368 representative sample of the general population in Finland. In addition, sex and age of
369 participants were also strongly predictive of the FLI group in our models (**Figures 2, S2, Table**
370 **S1**). Their similar positive associations with fatty liver disease are known from previous
371 studies.^{50,51} The associated microbial balances could be used to improve the predictions above
372 the baseline of these covariates on 5/6 regions in Finland. For example, in the model cross-
373 validated with Lapland the balances were more predictive of FLI group than the covariates by
374 themselves, and their combination increased the AUC further. Yet, when testing the model
375 where Turku/Loimaa region was used for cross-validation, the microbial balances were slightly
376 predictive of FLI group but failed to improve the AUC over the covariates (**Table S1**). This
377 pattern might stem from the cultural and genetic west-east division in Finland,^{52,53} with a closer
378 proximity of the Helsinki/Vantaa region to eastern regions than Turku/Loimaa, in both terms. It
379 is thus likely that further incorporation and investigation on the use of spatial information in
380 microbiome modeling would elucidate these geographical patterns in taxa-disease associations.
381

382 ***Conclusions***

383 Modeling an established risk index for fatty liver enabled the detection of associations between
384 the disease and gut microbiome composition, even to the level of individual taxa. These taxa
385 and clades were all from the obligately anaerobic gram-positive class *Clostridia*, from several
386 redefined GTDB genera previously included in the polyphyletic NCBI genus *Clostridium*.
387 Many of the representative genomes of taxa positively associated with fatty liver had genomic
388 potential for endogenous ethanol production. This suggests a possible mechanistic link to the
389 pathophysiology of liver disease through increased gut permeability and induction of hepatic
390 steatosis. Further mechanistic links with microbial production of SCFAs, especially acetate,
391 and fatty liver development are also likely. Our models were able to predict the FLI group of
392 participants across geographical regions in Finland, showing that the associations are robust
393 and mostly generalizable in the sampled population.

394

395 **Materials and Methods**

396 ***Survey details and sample collection***

397 Cardiovascular disease risk factors have been monitored in Finland since 1972 by conducting a
398 representative population survey every five years.⁵⁴ In the FINRISK 2002 survey, a stratified
399 random population sample was conducted on six geographical regions in Finland. These are
400 North Karelia and Northern Savo in eastern Finland, Turku and Loimaa regions in southwestern
401 Finland, the cities of Helsinki and Vantaa in the capital region, the provinces of Northern
402 Ostrobothnia and Kainuu in northwestern Finland, and the province of Lapland in northern
403 Finland. The Coordinating Ethics Committee of the Helsinki University Hospital District

404 approved our study protocol (Ref. 558/E3/2001) and all participants have given their written
405 informed consent.

406

407 Briefly, at baseline examination the participants filled out a questionnaire form, and trained
408 nurses carried out a physical examination and blood sampling in local health centers or other
409 survey sites. Data was collected for physiological measures, biomarkers, and dietary,
410 demographic and lifestyle factors. Stool samples were collected by giving willing participants a
411 stool sampling kit with detailed instructions. These samples were mailed overnight between
412 Monday and Thursday under Finnish winter conditions to the laboratory of the Finnish Institute
413 for Health and Welfare, where they were stored at -20°C. In 2017, the samples were shipped still
414 unthawed to University of California San Diego for microbiome sequencing.

415

416 Further details of the FINRISK cohorts and sampling have been extensively covered in previous
417 publications (for details, see refs. 24,55). The Coordinating Ethics Committee of the Helsinki
418 University Hospital District approved our study protocol. All participants have given their
419 written informed consent.

420

421 ***Stool DNA extraction and shallow shotgun metagenome sequencing***

422 A miniaturized version of the Kapa HyperPlus Illumina-compatible library prep kit (Kapa
423 Biosystems) was used for library generation, following the previously published protocol.⁵⁶
424 DNA extracts were normalized to 5 ng total input per sample in an Echo 550 acoustic liquid
425 handling robot (Labcyte Inc). A Mosquito HV liquid-handling robot (TTP Labtech Inc was used
426 for 1/10 scale enzymatic fragmentation, end-repair, and adapter-ligation reactions). Sequencing

427 adapters were based on the iTru protocol,⁵⁷ in which short universal adapter stubs are ligated first
428 and then sample-specific barcoded sequences added in a subsequent PCR step. Amplified and
429 barcoded libraries were then quantified by the PicoGreen assay and pooled in approximately
430 equimolar ratios before being sequenced on an Illumina HiSeq 4000 instrument to an average
431 read count of approximately 900,000 reads per sample.

432

433 *Taxonomic matching and phylogenetic transforms*

434 To improve the taxonomic assignments of our reads, we used a custom index,⁵⁸ based on the
435 Genome Taxonomy Database (GTDB) release 89 taxonomic redefinitions,^{28,59} for read
436 classification with default parameters in Centrifuge 1.0.4.⁶⁰ After read classification, all
437 following steps were performed with R version 3.5.2.⁶¹ To reduce the number of spurious read
438 assignments, and to facilitate more accurate phylogenetic transformations, only reads classified
439 at the species level, matching individual GTDB reference genomes, were retained. Samples with
440 less than 50,000 reads, from pregnant participants or recorded antibiotic use in the past 6 months
441 were removed, resulting in a final number of 6,269 samples. We first filtered taxa not seen with
442 more than 3 counts in at least 1% of samples and those with a coefficient of variation ≤ 3 across
443 all samples, following McMurdie and Holmes⁶², with a slight adaption from 20% of samples to
444 1% of samples, because of our larger sample size. The complete bacterial and archaeal
445 phylogenetic trees of the GTDB release 89 reference genomes, constructed from an alignment of
446 120 bacterial or 122 archaeal marker genes,²⁸ were then combined with our taxa tables. The
447 resulting trees were thus subset only to species which were observed in at least one sample in our
448 data. The read counts were transformed to phylogenetic node balances in both trees with

449 PhILR.²⁹ The default method for PhILR inputs a pseudocount of 1 for taxa absent in an
450 individual sample before the transform.
451
452 In this study, we did not specifically and solely use relative abundances at various taxonomic
453 levels, as is common practice for microbiome studies. Instead, we applied a PhILR
454 transformation to our microbial composition data,²⁹ introducing the concept of microbial
455 “balances”. Indeed, evolutionary relationships of all species harbored in each microbiome
456 sample can be represented on a phylogenetic tree, with species typically shown as external nodes
457 that are related to each other by multiple branches connected by internal nodes. In this context,
458 the value of a given microbial “balance” is defined as the log-ratio of the geometric mean
459 abundance between two groups of microbes descending from the same corresponding internal
460 node on a microbial phylogenetic tree. This phylogenetic transform was used because it i)
461 addresses the compositionality of the metagenomic read data,⁶³ ii) permits simultaneous
462 comparison of all clades without merging the taxa by predefined taxonomic levels, and iii)
463 enables evolutionary insights into the microbial community. The links between microbes and
464 their environment, such as the human gut, is mediated by their functions. Different functions are
465 known to be conserved at different taxonomic resolutions, and most often at multiple different
466 resolutions.⁶⁴ Thus, associations between the microbes and the response variable are likely not
467 best explained by predefined taxonomic levels. In the absence of functional data, concurrently
468 analyzing all clades (partitioned by the nodes in the phylogenetic tree) would likely enable the
469 detection of the associations at the appropriate resolution depending on the function and the local
470 tree topography.
471

472 *Covariates*

473 Because fatty liver disease is underdiagnosed at the population level,²⁵ and our sampling did not
474 have extensive coverage of liver fat measurements, we chose to use the Fatty Liver index as a
475 proxy for fatty liver.²⁶ Furthermore, the index performs well in cohorts of Caucasian ethnicity,
476 such as ours, to diagnose the presence of NAFLD.²³ We calculated FLI after Bedogni et al.²⁶:

$$477 \left(e^{0.953 \cdot \log_e(\text{triglycerides mg/dL}) + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745} \right) /$$

$$478 \left(1 + e^{0.953 \cdot \log_e(\text{triglycerides mg/dL}) + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745} \right) * 100. \text{ We chose}$$

479 the cutoff at $\text{FLI} \geq 60$ to identify participants likely to be diagnosed with hepatic steatosis

480 (positive likelihood ratio = 4.3 and negative likelihood ratio = 0.5, after Bedogni et al.²⁶).

481 Triglycerides, gamma glutamyl transferase (GGT), BMI and waist circumference measurements

482 had near complete coverage for the participants in our data. Self-reported alcohol use was

483 calculated as grams of pure ethanol per week. Cases with missing values were omitted in linear

484 regression models. At least one feature used for FLI calculation was missing for 20 participants

485 (0.3%) and the self-reported alcohol use was missing for 247 participants (3.9%). In the machine

486 learning framework, missing values for FLI and self-reported alcohol use were mean imputed.

487 However, for the feature selection to identify liver function-specific balances, GGT, triglycerides

488 and BMI were not imputed but observations where any of these were missing were simply

489 removed.

490

491 *Beta-diversity and linear modeling of FLI*

492 Beta-diversity was calculated as Euclidian distance of the PhILR balances through Principal

493 Component Analysis (PCA) on bacterial and archaeal balances separately with “rda” in vegan

494 2.5.6.⁶⁵ A linear regression model was constructed for FLI with “lm” in base R,⁶¹ with \log_{10} -

495 transformed FLI as the dependent variable and with first three bacterial PCs, sex, age, and self-
496 reported alcohol use as the independent variables. Archaeal PCs were dropped from the model
497 because none of them were significantly correlated with FLI (all $P > 0.05$). Variation of the
498 samples on the top two bacterial PC axes by their effect sizes in the model were plotted together
499 with a unit vector of $\log_{10}(\text{FLI})$ to show their correlation.

500

501 ***FLI modeling within a machine learning framework***

502 In the machine learning framework, both regression and categorical models were constructed for
503 FLI. The feature selection, hyperparameter optimization and cross-validation methods were
504 identical for both approaches, unless otherwise stated. The continuous or categorical FLI (groups
505 of $\text{FLI} < 60$ and $\text{FLI} \geq 60$) were modeled with xgboost 0.90.0.2,⁶⁶ by using both bacterial and
506 archaeal balances, sex, age, and self-reported alcohol use as preliminary predictor features. We
507 used FLI 60 as the cutoff for ruling in fatty liver (steatosis) for the classification, after Bedogni et
508 al., (2006). The data was first split to 70% train and 30% test sets while preserving sex and
509 region balance. To take into account geographical differences (see *e.g.*, ref. 31) and to find robust
510 patterns across all 6 sampled regions in Finland between the features and FLI group, we used
511 Leave-One-Group-Out Cross-Validation (LOGOCV) inside the 70% train set to construct 6
512 separate models in each optimization step. Because of high dimensionality of the data (3423
513 predictor features) feature selection by filtering was first performed inside the training data,
514 based on random forest permutation as recommended by Bommert et al.⁶⁷ Briefly, permutation
515 importance is based on accuracy, or specifically the difference in accuracy between real and
516 permuted (random) values of the specific variable, averaged in all trees across the whole random
517 forest. The permutation importance in models based on the 6 LOGOCV subsets of the training

518 data were calculated with mlr 2.16.0,⁶⁸ and the simple intersect between the top 50 features in all
519 LOGOCV subsets were retained as the final set of features. Thus, the feature selection was
520 influenced by the training data from all 6 geographical regions, but this only serves to limit the
521 number of chosen features because of the required simple intersect. This approach was used to
522 obtain a set of core predictive features which would have potential for generalizability across the
523 regions. The number of features included in the models by this approach was deemed
524 appropriate, since the relative effect size of the last included predictor was very small (< 0.1
525 change in classification probability across its range).

526

527 Bayesian hyperparameter optimization of the xgboost models was then performed with only the
528 selected features. An optimal set of parameters for the xgboost models were searched over all
529 LOGOCV subsets with “mbo” in mlrMBO 1.1.3,⁶⁹ using 30 preliminary rounds with randomized
530 parameters, followed by 100 optimization rounds. Parameters in the xgboost models and their
531 considered ranges were learning rate (eta) [0.001, 0.3], gamma [0.1, 5], maximum depth of a tree
532 [2, 8], minimum child weight [1, 10], fraction of data subsampled per each iteration [0.2, 0.8],
533 fraction of columns subsampled per tree [0.2, 0.9], and maximum number of iterations (nrounds)
534 [50, 5000]. The parameters recommended by these searchers were as following for regression:
535 eta=0.00889; gamma=2.08; max_depth=2; min_child_weight=8; subsample=0.783;
536 colsample_bytree=0.672; nrounds=1810, and for classification: eta=0.00107; gamma=0.137;
537 max_depth=5; min_child_weight=9; subsample=0.207; colsample_bytree=0.793; nrounds=4328.

538 We used Root-Mean-Square Error (RMSE) for the regression models and Area Under the ROC
539 Curve (AUC) for the classification models to measure model fit on the left-out data (region) in
540 each LOGOCV subset. The final models were trained on the LOGOCV subset training data, the

541 data from one region thus omitted per model, and using the selected features and optimized
542 hyperparameters. Validation of these models was conducted against participants only from the
543 region omitted from each model, in the 30% test data which was not used in model training or
544 optimization. Sensitivity analysis was conducted by using only the predictive covariates (sex and
545 age) or balances separately, with the same hyperparameters, data partitions and final validation
546 as for the full models.

547

548 ***Partial dependence interpretation of the FLI classification models***

549 Because the classification models have a more clinically relevant modeling target for the
550 difference between $FLI < 60$ and $FLI \geq 60$, the latter used to rule in fatty liver,²⁶ we further
551 interpreted the partial dependence of their predictions. Partial dependence of the classification
552 model predictions on individual features was calculated with “partial” in pdp 0.7.0.⁷⁰ The partial
553 dependence of the features on the model predictions was also plotted, overlaying the results from
554 each of the 6 models. For each feature, its relative effect on the model prediction was estimated
555 as medians of the minimum and maximum \hat{y} (output probability of the model for the $FLI \geq 60$
556 class), calculated at the minimum and maximum values of the feature separately in each of the 6
557 models. The relative effects of the balances were then overlaid as a heatmap on a genome
558 cladogram which covers all balances in the model with ggtree 2.1.1.⁷¹

559

560 ***Construction of alternative classification models to discern between***

561 ***FLI < 30 and FLI \geq 60 groups***

562 To assess robustness of the models and how removing the participants with intermediate FLI
563 (between 30 and 60) affects model performance, we removed this group ($N = 1910$) and

564 constructed alternative classification models to discern between the $FLI < 30$ and $FLI \geq 60$
565 groups. Other than removing the intermediate FLI participants and resulting new random split to
566 the train (70%) and test (30%) sets, these models were constructed identically to the main
567 models, including LOGOCV design, feature selection, and hyperparameter optimization. The
568 recommended parameters for this classification task were $\eta=0.00102$; $\gamma=3.7$;
569 $\max_depth=2$; $\min_child_weight=5$; $subsample=0.49$; $colsample_bytree=0.631$; $nrounds=3119$.
570 Interpretation of partial dependence was also performed identically, but only the relative effects
571 of the model features were plotted without a cladogram.

572

573 ***Exclusion of validation region data from feature selection and hyperparameter***
574 ***optimization***

575 Because training data from all 6 regions is used to inform the selection of optimal features and
576 hyperparameters, the validation region data cannot be considered completely independent from
577 the training of the LOGOCV models. Thus, we constructed a set of classification models for the
578 $FLI \geq 60$ and $FLI < 60$ groups, where all validation region samples also in the training data were
579 excluded from the simple intercept of top 50 features in each LOGOCV set and from the
580 subsequent hyperparameter optimization. These models with individualized features and
581 hyperparameters were then tested on the validation region samples in the unseen test data to
582 estimate how model performance was affected. The main test (70%) and train (30%) sets were
583 identical to the main models, but additionally 6 randomized 70/30 splits nested inside the test set
584 (excluding the validation region) were used in hyperparameter optimization to reduce overfitting.
585 Average optimal hyperparameters in the 6 models were $\eta=0.00106$; $\gamma=4.3$;
586 $\max_depth=2$; $\min_child_weight=7$; $subsample=0.36$; $colsample_bytree=0.613$; $nrounds=1772$.

587

588 ***Identification of predictive features specific to liver function***

589 Because the FLI also incorporates BMI and waist circumference, and they strongly contribute to
590 the index,²⁶ we deemed it necessary to further investigate which of the identified balances were
591 specific to liver function. The participants were first grouped by age (< 40, 40 – 60, and 60 <),
592 sex (female or male) and BMI (< 25, 25 – 30, and 30 <) into 18 categories ($N = 105 \sim 711$ per
593 category). We performed feature selection similarly to the FLI models by fitting random forest
594 regressors for GGT and triglycerides with mlr 2.16.0.⁶⁸ This was done separately in each of the
595 18 categories, and in each category, we again used LOGOCV with the regions to obtain 6 runs
596 per category. Finally, the features predictive of GGT or triglycerides in each category were
597 selected as the intersect of top 50 features in the 6 LOGOCV iterations by permutation
598 importance. The intersect of features predictive of GGT or triglycerides in any of the categories
599 and the features predictive of categorical FLI were identified as specific to liver function.

600

601 ***Pathway inference for taxa associated with FLI***

602 Our taxonomic matching of the reads is based on the genomes of GTDB (release 89),²⁸ which are
603 all complete or nearly complete and available in online databases. This enables us to estimate the
604 likely genetic content, and thus, the metabolic potential of the microbes associated with FLI. We
605 use this approach because the sequencing depth of our samples does not allow assembling
606 contigs and (metagenome-assembled) genomes, required for pathway predictions. Because of the
607 compositional phylogenetic transform, among other features of our data, previously developed
608 approaches such as PICRUSt,⁷² could not be used here.

609

610 The genomes of all 336 bacteria under at least one of the predictive balances were downloaded
611 from NCBI. 119 of these genomes were originally not annotated, which is a requirement for
612 pathway prediction. Therefore, Prokka v1.14.6,⁷³ was used to annotate the 119 unannotated
613 genomes as a preliminary step. Pathway predictions were then performed for all 336 genomes
614 with mpwt v0.5.3 multiprocessing tool,⁷⁴ for the PathoLogic pipeline of Pathway Tools 23.0.⁷⁵
615 Pathways for ethanol and short chain fatty acid (acetate, butyrate, propionate) production, bile
616 acid metabolism, and choline degradation to trimethylamine were identified from MetaCyc
617 pathway classifications (see ref. 76, and **Table S4**). The prevalence of these processes was then
618 assessed in the analyzed genomes and summarized per process to consider the possible links of
619 the taxa with fatty liver pathophysiology. Finally, the presence of individual pathways for acetate
620 and ethanol production was also outlined for each genome.

621

622 *Data availability statement*

623 The analysis code written for this study is included with the Supplementary Information. The
624 datasets generated during and analyzed during the current study are not public, but are available
625 based on a written application to the THL Biobank as instructed in: [https://thl.fi/en/web/thl-](https://thl.fi/en/web/thl-biobank/for-researchers)
626 [biobank/for-researchers](https://thl.fi/en/web/thl-biobank/for-researchers)

627

628 *Disclosure of interest*

629 V.S. has consulted for Novo Nordisk and Sanofi and received honoraria from these companies.
630 He also has ongoing research collaboration with Bayer AG, all unrelated to this study. R.L.
631 serves as a consultant or advisory board member for Anylam/Regeneron, Arrowhead
632 Pharmaceuticals, AstraZeneca, Bird Rock Bio, Boehringer Ingelheim, Bristol-Myer Squibb,

633 Celgene, Cirius, CohBar, Conatus, Eli Lilly, Galmed, Gemphire, Gilead, Glympse bio, GNI, GRI
634 Bio, Inipharm, Intercept, Ionis, Janssen Inc., Merck, Metacrine, Inc., NGM Biopharmaceuticals,
635 Novartis, Novo Nordisk, Pfizer, Prometheus, Promethera, Sanofi, Siemens, and Viking
636 Therapeutics. In addition, his institution has received grant support from Allergan, Boehringer-
637 Ingelheim, Bristol-Myers Squibb, Cirius, Eli Lilly and Company, Galectin Therapeutics, Galmed
638 Pharmaceuticals, GE, Genfit, Gilead, Intercept, Grail, Janssen, Madrigal Pharmaceuticals,
639 Merck, NGM Biopharmaceuticals, NuSirt, Pfizer, pH Pharma, Prometheus, and Siemens. He is
640 also co-founder of Liponexus, Inc.

641

642 ***Funding details***

643 This research was supported in part by grants from the Finnish Foundation for Cardiovascular
644 Research, the Emil Aaltonen Foundation, the Paavo Nurmi Foundation, the Urmas Pekkala
645 Foundation, the Finnish Medical Foundation, the Sigrid Juselius Foundation, the Academy of
646 Finland (#321356 to A.H.; #295741, #307127 to L.L.; #321351 to T.N.) and the National
647 Institutes of Health (R01ES027595 to M.J.). R.L. receives funding support from NIEHS
648 (5P42ES010337), NCATS (5UL1TR001442), NIDDK (U01DK061734, R01DK106419,
649 P30DK120515, R01DK121378, R01DK124318), and DOD PRCRP (W81XWH-18-2-0026).
650 Additional support was provided by Illumina, Inc. and Janssen Pharmaceutica through their
651 sponsorship of the Center for Microbiome Innovation at UCSD.

652

653 ***Authors' contributions***

654 M.R., F.Å., V.M., V.S., R.K., L.L and T.N designed the work. A.H., L.V., G.M., P.J., V.S., M.J
655 and R.K. acquired the data. M.R., L.L. and T.N. analyzed the data. M.R. wrote the manuscript in

656 consultation with all authors. M.I., P.J., V.S., R.K., L.L. and T.N. supervised the work. All
657 authors gave final approval of the version to be published.

658

659 *Acknowledgements*

660 We thank all participants of the FINRISK 2002 survey for their contributions to this work, and
661 Tara Schwartz for assistance with laboratory work.

662

663 **References**

- 664 1. Younossi ZM, Koenig AB, Abdelatif D, Fazel Y, Henry L, Wymer M. Global
665 epidemiology of nonalcoholic fatty liver disease—Meta-analytic assessment of prevalence,
666 incidence, and outcomes. *Hepatology* 2016; 64:73–84.
- 667 2. Marchesini G, Bugianesi E, Forlani G, Cerrelli F, Lenzi M, Manini R, Natale S, Vanni E,
668 Villanova N, Melchionda N, et al. Nonalcoholic fatty liver, steatohepatitis, and the
669 metabolic syndrome. *Hepatology* 2003; 37:917–23.
- 670 3. Chalasani N, Younossi Z, Lavine JE, Diehl AM, Brunt EM, Cusi K, Charlton M, Sanyal
671 AJ. The diagnosis and management of non-alcoholic fatty liver disease: Practice Guideline
672 by the American Association for the Study of Liver Diseases, American College of
673 Gastroenterology, and the American Gastroenterological Association. *Hepatology* 2012;
674 55:2005–23.
- 675 4. Yki-Järvinen H. Non-alcoholic fatty liver disease as a cause and a consequence of
676 metabolic syndrome. *Lancet Diabetes Endocrinol* 2014; 2:901–10.
- 677 5. Toshikuni N, Tsutsumi M, Arisawa T. Clinical differences between alcoholic liver disease
678 and nonalcoholic fatty liver disease. *World J Gastroenterol* 2014; 20:8393–406.
- 679 6. Rinella M, Charlton M. The globalization of nonalcoholic fatty liver disease: Prevalence
680 and impact on world health. *Hepatology* 2016; 64:19–22.
- 681 7. Gilbert JA, Blaser MJ, Caporaso JG, Jansson JK, Lynch SV, Knight R. Current
682 understanding of the human microbiome. *Nat Med* 2018; 24:392–400.
- 683 8. Caussy C, Tripathi A, Humphrey G, Bassirian S, Singh S, Faulkner C, Bettencourt R, Rizo
684 E, Richards L, Xu ZZ, et al. A gut microbiome signature for cirrhosis due to nonalcoholic
685 fatty liver disease. *Nature Communications* 2019; 10:1406.

- 686 9. Compare D, Coccoli P, Rocco A, Nardone OM, De Maria S, Cartenì M, Nardone G. Gut–
687 liver axis: The impact of gut microbiota on non alcoholic fatty liver disease. *Nutrition,*
688 *Metabolism and Cardiovascular Diseases* 2012; 22:471–6.
- 689 10. Miele L, Valenza V, Torre GL, Montalto M, Cammarota G, Ricci R, Mascianà R, Forgione
690 A, Gabrieli ML, Perotti G, et al. Increased intestinal permeability and tight junction
691 alterations in nonalcoholic fatty liver disease. *Hepatology* 2009; 49:1877–87.
- 692 11. Spencer MD, Hamp TJ, Reid RW, Fischer LM, Zeisel SH, Fodor AA. Association Between
693 Composition of the Human Gastrointestinal Microbiome and Development of Fatty Liver
694 With Choline Deficiency. *Gastroenterology* 2011; 140:976–86.
- 695 12. Zhu L, Baker SS, Gill C, Liu W, Alkhoury R, Baker RD, Gill SR. Characterization of gut
696 microbiomes in nonalcoholic steatohepatitis (NASH) patients: A connection between
697 endogenous alcohol and NASH. *Hepatology* 2013; 57:601–9.
- 698 13. Yuan J, Chen C, Cui J, Lu J, Yan C, Wei X, Zhao X, Li N, Li S, Xue G, et al. Fatty Liver
699 Disease Caused by High-Alcohol-Producing *Klebsiella pneumoniae*. *Cell Metabolism*
700 2019; 30:675-688.e7.
- 701 14. Loomba R, Seguritan V, Li W, Long T, Klitgord N, Bhatt A, Dulai PS, Caussy C,
702 Bettencourt R, Highlander SK, et al. Gut Microbiome-Based Metagenomic Signature for
703 Non-invasive Detection of Advanced Fibrosis in Human Nonalcoholic Fatty Liver Disease.
704 *Cell Metabolism* 2017; 25:1054-1062.e5.
- 705 15. Jiao N, Wu D, Yang Z, Fang S, Li X, Yuan M, Zhu R, Zhu L. Gut bacteria contributes to
706 NAFLD pathogenesis by promoting secondary bile acids biosynthesis. *The FASEB Journal*
707 2019; 33:126.4-126.4.
- 708 16. Aron-Wisnewsky J, Vigliotti C, Witjes J, Le P, Holleboom AG, Verheij J, Nieuwdorp M,
709 Clément K. Gut microbiota and human NAFLD: disentangling microbial signatures from
710 metabolic disorders. *Nature Reviews Gastroenterology & Hepatology* 2020; 17:279–97.
- 711 17. Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, Guo J, Le Chatelier E, Yao J, Wu L, et al.
712 Alterations of the human gut microbiome in liver cirrhosis. *Nature* 2014; 513:59–64.
- 713 18. Liu Y, Meric G, Havulinna AS, Teo SM, Ruuskanen M, Sanders J, Zhu Q, Tripathi A,
714 Verspoor K, Cheng S, et al. Early prediction of liver disease using conventional risk factors
715 and gut microbiome-augmented gradient boosting [Internet]. *Genetic and Genomic*
716 *Medicine*; 2020 [cited 2020 Jul 28]. Available from:
717 <http://medrxiv.org/lookup/doi/10.1101/2020.06.24.20138933>
- 718 19. Safari Z, Gérard P. The links between the gut microbiome and non-alcoholic fatty liver
719 disease (NAFLD). *Cell Mol Life Sci* 2019; 76:1541–58.
- 720 20. Carpino G, Del Ben M, Pastori D, Carnevale R, Baratta F, Overi D, Francis H, Cardinale V,
721 Onori P, Safarikia S, et al. Increased liver localization of lipopolysaccharides in human and
722 experimental non-alcoholic fatty liver disease. *Hepatology* 2019; :hep.31056.

- 723 21. Li Q, Dhyani M, Grajo JR, Sirlin C, Samir AE. Current status of imaging in nonalcoholic
724 fatty liver disease. *World J Hepatol* 2018; 10:530–42.
- 725 22. Koehler EM, Schouten JNL, Hansen BE, Hofman A, Stricker BH, Janssen HLA. External
726 Validation of the Fatty Liver Index for Identifying Nonalcoholic Fatty Liver Disease in a
727 Population-based Study. *Clinical Gastroenterology and Hepatology* 2013; 11:1201–4.
- 728 23. Vanni E, Bugianesi E. Editorial: utility and pitfalls of Fatty Liver Index in epidemiologic
729 studies for the diagnosis of NAFLD. *Aliment Pharmacol Ther* 2015; 41:406–7.
- 730 24. Salosensaari A, Laitinen V, Havulinna AS, Meric G, Cheng S, Perola M, Valsta L, Alfthan
731 G, Inouye M, Watrous JD, et al. Taxonomic Signatures of Long-Term Mortality Risk in
732 Human Gut Microbiota [Internet]. *Epidemiology*; 2020 [cited 2020 Jan 4]. Available from:
733 <http://medrxiv.org/lookup/doi/10.1101/2019.12.30.19015842>
- 734 25. Alexander M, Loomis AK, Fairburn-Beech J, van der Lei J, Duarte-Salles T, Prieto-
735 Alhambra D, Ansell D, Pasqua A, Lapi F, Rijnbeek P, et al. Real-world data reveal a
736 diagnostic gap in non-alcoholic fatty liver disease. *BMC Medicine* 2018; 16:130.
- 737 26. Bedogni G, Bellentani S, Miglioli L, Masutti F, Passalacqua M, Castiglione A, Tiribelli C.
738 The Fatty Liver Index: a simple and accurate predictor of hepatic steatosis in the general
739 population. *BMC Gastroenterol* 2006; 6:33.
- 740 27. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R,
741 Knights D. Evaluating the Information Content of Shallow Shotgun Metagenomics.
742 *mSystems* [Internet] 2018 [cited 2020 Apr 9]; 3. Available from:
743 <https://msystems.asm.org/content/3/6/e00069-18>
- 744 28. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A,
745 Hugenholtz P. A standardized bacterial taxonomy based on genome phylogeny
746 substantially revises the tree of life. *Nat Biotechnol* 2018; 36:996–1004.
- 747 29. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform
748 enhances analysis of compositional microbiota data. *eLife* 2017; 6:e21887.
- 749 30. Haas JT, Francque S, Staels B. Pathophysiology and Mechanisms of Nonalcoholic Fatty
750 Liver Disease. *Annual Review of Physiology* 2016; 78:181–205.
- 751 31. He Y, Wu W, Zheng H-M, Li P, McDonald D, Sheng H-F, Chen M-X, Chen Z-H, Ji G-Y,
752 Zheng Z-D-X, et al. Regional variation limits applications of healthy gut microbiome
753 reference ranges and disease models. *Nature Medicine* 2018; 24:1532–5.
- 754 32. Kim H-N, Joo E-J, Cheong HS, Kim Y, Kim H-L, Shin H, Chang Y, Ryu S. Gut
755 Microbiota and Risk of Persistent Nonalcoholic Fatty Liver Diseases. *Journal of Clinical
756 Medicine* 2019; 8:1089.
- 757 33. Castaner O, Goday A, Park Y-M, Lee S-H, Magkos F, Shiow S-ATE, Schröder H. The Gut
758 Microbiome Profile in Obesity: A Systematic Review. *Int J Endocrinol* [Internet] 2018

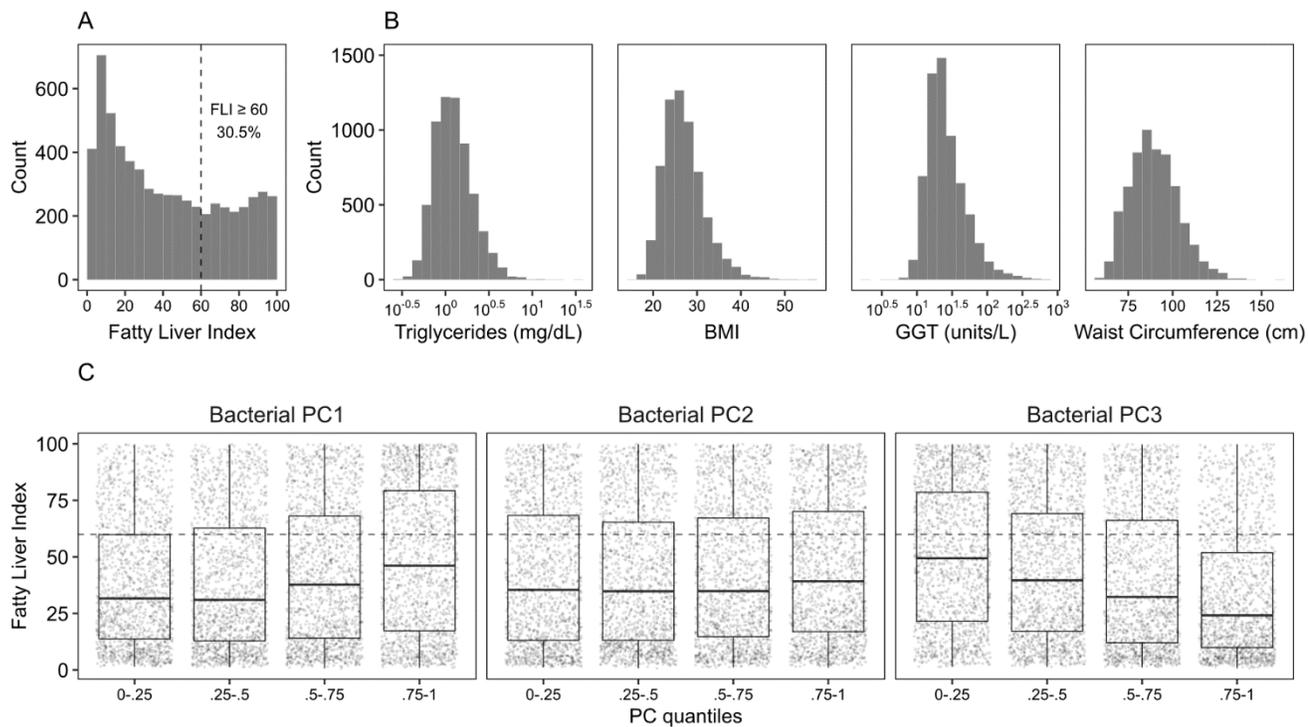
- 759 [cited 2020 Apr 3]; 2018. Available from:
760 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5933040/>
- 761 34. Wigg AJ, Roberts-Thomson IC, Dymock RB, McCarthy PJ, Grose RH, Cummins AG. The
762 role of small intestinal bacterial overgrowth, intestinal permeability, endotoxaemia, and
763 tumour necrosis factor α in the pathogenesis of non-alcoholic steatohepatitis. *Gut* 2001;
764 48:206–11.
- 765 35. Mouzaki M, Comelli EM, Arendt BM, Bonengel J, Fung SK, Fischer SE, McGilvray ID,
766 Allard JP. Intestinal microbiota in patients with nonalcoholic fatty liver disease. *Hepatology*
767 2013; 58:120–7.
- 768 36. Shen F, Zheng R-D, Sun X-Q, Ding W-J, Wang X-Y, Fan J-G. Gut microbiota dysbiosis in
769 patients with non-alcoholic fatty liver disease. *Hepatobiliary & Pancreatic Diseases*
770 *International* 2017; 16:375–81.
- 771 37. Sharpton SR, Ajmera V, Loomba R. Emerging Role of the Gut Microbiome in
772 Nonalcoholic Fatty Liver Disease: From Composition to Function. *Clinical*
773 *Gastroenterology and Hepatology* 2019; 17:296–306.
- 774 38. Jiang W, Wu N, Wang X, Chi Y, Zhang Y, Qiu X, Hu Y, Li J, Liu Y. Dysbiosis gut
775 microbiota associated with inflammation and impaired mucosal immune function in
776 intestine of humans with non-alcoholic fatty liver disease. *Sci Rep* 2015; 5:8096.
- 777 39. Suzuki TA, Worobey M. Geographical variation of human gut microbial composition.
778 *Biology Letters* 2014; 10:20131037.
- 779 40. Meslier V, Laiola M, Roager HM, Filippis FD, Roume H, Quinquis B, Giacco R, Mennella
780 I, Ferracane R, Pons N, et al. Mediterranean diet intervention in overweight and obese
781 subjects lowers plasma cholesterol and causes changes in the gut microbiome and
782 metabolome independently of energy intake. *Gut* [Internet] 2020 [cited 2020 Jun 2];
783 Available from: <https://gut.bmj.com/content/early/2020/02/18/gutjnl-2019-320438>
- 784 41. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Knight R, Knights D. SHOGUN: a
785 modular, accurate, and scalable framework for microbiome quantification. *Bioinformatics*
786 2020; :btaa277.
- 787 42. de Faria Ghetti F, Oliveira DG, de Oliveira JM, de Castro Ferreira LEVV, Cesar DE,
788 Moreira APB. Influence of gut microbiota on the development and progression of
789 nonalcoholic steatohepatitis. *Eur J Nutr* 2018; 57:861–76.
- 790 43. Mohan R, Namsolleck P, Lawson PA, Osterhoff M, Collins MD, Alpert C-A, Blaut M.
791 *Clostridium asparagiforme* sp. nov., isolated from a human faecal sample. *Systematic and*
792 *Applied Microbiology* 2006; 29:292–9.
- 793 44. Murray WD, Khan AW, van den BERG L. *Clostridium saccharolyticum* sp. nov., a
794 Saccharolytic Species from Sewage Sludge. *International Journal of Systematic*
795 *Bacteriology* 1982; 32:132–5.

- 796 45. Diether NE, Willing BP. Microbial Fermentation of Dietary Protein: An Important Factor in
797 Diet–Microbe–Host Interaction. *Microorganisms* [Internet] 2019 [cited 2020 Apr 24]; 7.
798 Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6352118/>
- 799 46. Zhao S, Jang C, Liu J, Uehara K, Gilbert M, Izzo L, Zeng X, Trefely S, Fernandez S, Carrer
800 A, et al. Dietary fructose feeds hepatic lipogenesis via microbiota-derived acetate. *Nature*
801 2020; 579:586–91.
- 802 47. Dehoux P, Marvaud JC, Abouelleil A, Earl AM, Lambert T, Dauga C. Comparative
803 genomics of *Clostridium bolteae* and *Clostridium clostridioforme* reveals species-specific
804 genomic properties and numerous putative antibiotic resistance determinants. *BMC*
805 *Genomics* 2016; 17:819.
- 806 48. Deschasaux M, Bouter KE, Prodan A, Levin E, Groen AK, Herrema H, Tremaroli V,
807 Bakker GJ, Attaye I, Pinto-Sietsma S-J, et al. Depicting the composition of gut microbiota
808 in a population with varied ethnic origins but shared geography. *Nature Medicine* 2018;
809 24:1526–31.
- 810 49. Canfora EE, Meex RCR, Venema K, Blaak EE. Gut microbial metabolites in obesity,
811 NAFLD and T2DM. *Nature Reviews Endocrinology* 2019; 15:261–73.
- 812 50. Cheng H-Y, Wang H-Y, Chang W-H, Lin S-C, Chu C-H, Wang T-E, Liu C-C, Shih S-C.
813 Nonalcoholic Fatty Liver Disease: Prevalence, Influence on Age and Sex, and Relationship
814 with Metabolic Syndrome and Insulin Resistance. *International Journal of Gerontology*
815 2013; 7:194–8.
- 816 51. Lonardo A, Nascimbeni F, Ballestri S, Fairweather D, Win S, Than TA, Abdelmalek MF,
817 Suzuki A. Sex Differences in Nonalcoholic Fatty Liver Disease: State of the Art and
818 Identification of Research Gaps. *Hepatology* 2019; 70:1457–69.
- 819 52. Näyhä S. Geographical variations in cardiovascular mortality in Finland, 1961-1985. *Scand*
820 *J Soc Med Suppl* 1989; 40:1–48.
- 821 53. Kerminen S, Havulinna AS, Hellenthal G, Martin AR, Sarin A-P, Perola M, Palotie A,
822 Salomaa V, Daly MJ, Ripatti S, et al. Fine-Scale Genetic Structure in Finland. *G3: Genes,*
823 *Genomes, Genetics* 2017; 7:3459–68.
- 824 54. Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, Kuulasmaa K,
825 Laatikainen T, Männistö S, Peltonen M, et al. Cohort Profile: The National FINRISK
826 Study. *International Journal of Epidemiology* 2018; 47:696–696i.
- 827 55. Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T, Männistö S,
828 Salomaa V, Sundvall J, Puska P. Forty-year trends in cardiovascular risk factors in Finland.
829 *Eur J Public Health* 2015; 25:539–46.
- 830 56. Sanders JG, Nurk S, Salido RA, Minich J, Xu ZZ, Zhu Q, Martino C, Fedarko M, Arthur
831 TD, Chen F, et al. Optimizing sequencing protocols for leaderboard metagenomics by
832 combining long and short reads. *Genome Biology* 2019; 20:226.

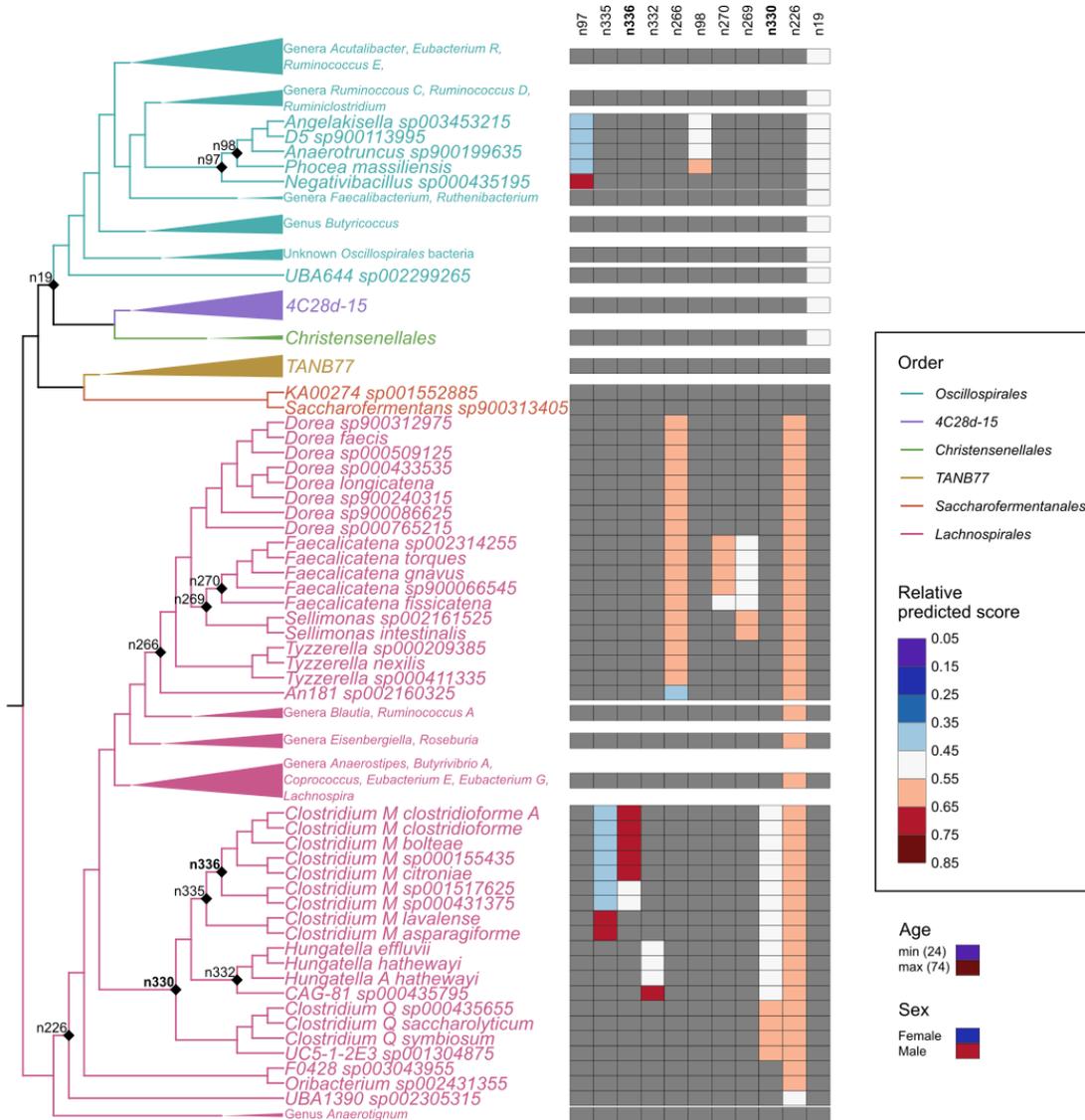
- 833 57. Glenn TC, Nilsen RA, Kieran TJ, Sanders JG, Bayona-Vásquez NJ, Finger JW, Pierson
834 TW, Bentley KE, Hoffberg SL, Louha S, et al. Adapterama I: universal stubs and primers
835 for 384 unique dual-indexed or 147,456 combinatorially-indexed Illumina libraries (iTru &
836 iNext). *PeerJ* 2019; 7:e7755.
- 837 58. Méric G, Wick RR, Watts SC, Holt KE, Inouye M. Correcting index databases improves
838 metagenomic studies. *bioRxiv* 2019; :712166.
- 839 59. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. A complete
840 domain-to-species taxonomy for Bacteria and Archaea. *Nature Biotechnology* 2020; :1–8.
- 841 60. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification
842 of metagenomic sequences. *Genome Res* [Internet] 2016 [cited 2018 May 12]; Available
843 from: <http://genome.cshlp.org/content/early/2016/11/16/gr.210641.116>
- 844 61. R Core Team. R: A language and environment for statistical computing [Internet]. Vienna,
845 Austria: R Foundation for Statistical Computing; 2018 [cited 2019 Mar 4]. Available from:
846 <https://www.R-project.org/>
- 847 62. McMurdie PJ, Holmes S. phyloseq: An R package for reproducible interactive analysis and
848 graphics of microbiome census data. *PLoS ONE* 2013; 8:e61217.
- 849 63. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome Datasets Are
850 Compositional: And This Is Not Optional. *Front Microbiol* [Internet] 2017 [cited 2020 Jul
851 20]; 8. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5695134/>
- 852 64. Louca S, Polz MF, Mazel F, Albright MBN, Huber JA, O'Connor MI, Ackermann M, Hahn
853 AS, Srivastava DS, Crowe SA, et al. Function and functional redundancy in microbial
854 systems. *Nat Ecol Evol* 2018; 2:936–43.
- 855 65. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin PR,
856 O'Hara RB, Simpson GL, Solymos P, et al. vegan: Community Ecology Package [Internet].
857 2018 [cited 2018 Jun 4]. Available from: <https://CRAN.R-project.org/package=vegan>
- 858 66. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd
859 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining -
860 KDD '16* 2016; :785–94.
- 861 67. Bommert A, Sun X, Bischl B, Rahnenführer J, Lang M. Benchmark for filter methods for
862 feature selection in high-dimensional classification data. *Computational Statistics & Data
863 Analysis* 2020; 143:106839.
- 864 68. Bischl B, Lang M, Kotthoff L, Schiffner J, Richter J, Studerus E, Casalicchio G, Jones ZM.
865 mlr: Machine Learning in R. *Journal of Machine Learning Research* 2016; 17:1–5.
- 866 69. Bischl B, Richter J, Bossek J, Horn D, Thomas J, Lang M. mlrMBO: A Modular
867 Framework for Model-Based Optimization of Expensive Black-Box Functions.

- 868 arXiv:170303373 [stat] [Internet] 2018 [cited 2020 Feb 18]; Available from:
869 <http://arxiv.org/abs/1703.03373>
- 870 70. Greenwell BM. pdp: An R Package for Constructing Partial Dependence Plots. *The R*
871 *Journal* 2017; 9:421–36.
- 872 71. Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. ggtree: an r package for visualization and
873 annotation of phylogenetic trees with their covariates and other associated data. *Methods in*
874 *Ecology and Evolution* 2017; 8:28–36.
- 875 72. Douglas GM, Maffei VJ, Zaneveld J, Yurgel SN, Brown JR, Taylor CM, Huttenhower C,
876 Langille MGI. PICRUSt2: An improved and extensible approach for metagenome
877 inference. *bioRxiv* 2019; :672295.
- 878 73. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014; 30:2068–
879 9.
- 880 74. Belcour A, Frioux C, Aite M, Bretaudeau A, Siegel A. Metage2Metabo: metabolic
881 complementarity applied to genomes of large-scale microbiotas for the identification of
882 keystone species. *bioRxiv* 2019; :803056.
- 883 75. Karp PD, Midford PE, Billington R, Kothari A, Krummenacker M, Latendresse M, Ong
884 WK, Subhraveti P, Caspi R, Fulcher C, et al. Pathway Tools version 23.0 update: software
885 for pathway/genome informatics and systems biology. *Briefings in Bioinformatics*
886 2019; :bbz104.
- 887 76. Caspi R, Billington R, Fulcher CA, Keseler IM, Kothari A, Krummenacker M, Latendresse
888 M, Midford PE, Ong Q, Ong WK, et al. The MetaCyc database of metabolic pathways and
889 enzymes. *Nucleic Acids Res* 2018; 46:D633–9.
- 890

891 Figures



892 **Figure 1.** Distribution of FLI (A), its components (B), and FLI in quantiles of the first three PC
893 components of the fecal bacterial composition of the participants (C). The cutoff at FLI = 60
894 used to divide the participants is indicated with a dashed line in panels A and C.



895 **Figure 2.** Relative effects of predictive balances and covariates on the FLI < 60 and FLI ≥ 60
 896 classification model (AUC = 0.75) predictions. Nodes of the balances are indicated in the
 897 cladogram and the relative effect sizes of their clades (opposite sides of each balance) are shown
 898 in the associated heatmap. The relative effect sizes of the covariates (age and sex) are shown
 899 below the legend with a heatmap on the same scale as was used for the balances. The two liver-
 900 specific balances associated with triglyceride and GGT levels are indicated with bold font.
 901 Clades with redundant information have been collapsed but their major genera are indicated. The
 902 complete tree is included in **Figure S3**.