

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. Its connection with the gut
42 microbiome has been known for at least 80 years, but this association remains mostly unstudied
43 in the general population because of underdiagnosis and small sample sizes. To address this
44 knowledge gap, we studied the link between the Fatty Liver Index (FLI), a well-established
45 proxy for fatty liver disease, and gut microbiome composition in a representative, ethnically
46 homogeneous population sample of 6,269 Finnish participants. We based our models on
47 biometric covariates and gut microbiome compositions from shallow metagenome sequencing.
48 Our classification models could discriminate between individuals with a high FLI (≥ 60 ,
49 indicates likely liver steatosis) and low FLI (< 60) in internal cross-region validation, consisting
50 of 30% of the data not used in model training, with an average AUC of 0.75 and AUPRC of 0.56
51 (baseline at 0.30). In addition to age and sex, our models included differences in 11 microbial
52 groups from class *Clostridia*, mostly belonging to orders *Lachnospirales* and *Oscillospirales*.
53 Our models were also predictive of the high FLI group in a different Finnish cohort, consisting
54 of 258 participants, with an average AUC of 0.77 and AUPRC of 0.51 (baseline at 0.21).
55 Pathway analysis of representative genomes of the positively FLI-associated taxa in (NCBI)
56 *Clostridium* subclusters IV and XIVa indicated the presence of *e.g.*, ethanol fermentation
57 pathways. These results support several findings from smaller case-control studies, such as the
58 role of endogenous ethanol producers in the development of fatty liver.

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

60 Introduction

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

72

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

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

88

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

100

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

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

118

119 **Results**

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

121 *sample*

122 In our main analyses, we classified our reads against the Genome Taxonomy Database
123 (GTDB)²⁸. This study mainly follows the GTDB taxonomy, unless otherwise noted. The
124 Centrifuge/GTDB microbiome data used in our main analyses was based on archaeal and
125 bacterial phylogenetic “balances”. This method was used to associate larger groups or clades of
126 related organisms with fatty liver disease, and to avoid grouping of taxa on strict hierarchical
127 taxonomic ranks featuring varying ranges of evolutionary divergence.²⁸ Here, we used the PhILR

128 transform, where each balance represents a single internal node in a phylogenetic tree, and its
129 value is a log-ratio of the abundances of the two descending clades (for details, see methods and
130 ref. ²⁹). Positive values of the balance signify that the clade in the numerator is more abundant,
131 and negative values that the clade in the denominator is more abundant. Thus, each association
132 of a balance with the target variable necessarily includes both microbial clades descending from
133 the node, one of them positively and the other negatively associated with the target variable. The
134 clades in the numerator and denominator can be also freely switched by changing the sign of the
135 balance value to retain the equivalence. Notably, we used this feature to show all balance-FLI
136 associations in the positive direction to facilitate the comparison of their effect sizes (in **Figures**
137 **S4, S7, and S9**).

138

139 Because the combined approach of using the GTDB taxonomy and the recently introduced
140 PhILR -phylogenetic transform complicates the comparison of our results to previous studies, we
141 also conducted more traditional statistical analyses with NCBI-annotated data to anchor our
142 results to previous findings on the associations between fatty liver disease and gut microbiome
143 composition. Overall, the Centrifuge/GTDB classification assigned 5.3 billion reads in the 6,269
144 samples (after exclusions in FINRISK 2002) to 23,457 bacterial and 1,248 archaeal taxa, and the
145 SHOGUN/NCBI classification assigned 5.5 billion reads to 5,024 bacterial and 261 archaeal
146 taxa. Starting from high level descriptions of the microbial communities in the high and low FLI
147 groups (< 60 or ≥ 60 FLI; **Figure 1A**), the phylum-level distributions of bacterial and archaeal
148 taxa appeared to be highly similar between the groups (**Figures S1, S2**). However, the proportion
149 of taxa assigned to *Firmicutes* in Centrifuge/GTDB appeared to be slightly higher than in the
150 SHOGUN/NCBI data. Furthermore, only 58% of the number of reads assigned in

151 Centrifuge/GTDB to 6 main archaeal phyla were assigned to a single main archaeal phylum in
152 SHOGUN/NCBI. Alpha diversity (as Shannon diversity) was significantly lower in the high FLI
153 group, in both the SHOGUN/NCBI data (14.7% lower; AIC = 6685; all $P < 1 \times 10^{-6}$) and the
154 Centrifuge/GTDB (13.4% lower; AIC = 6607; all $P < 1 \times 10^{-4}$) data, while adjusting for age, sex,
155 and self-reported alcohol use in both models.

156

157 To further examine the high-level associations between FLI (as a proxy of fatty liver disease)
158 and microbial community composition in FINRISK 2002, we fit a linear regression model on the
159 three first principal component (PC) axes of the fecal bacterial beta diversity (between
160 individuals), sex, age, and alcohol. $\log_{10}(\text{FLI})$ significantly correlated with all three bacterial PC
161 axes, sex, age, and alcohol use in Centrifuge/GTDB data (adjusted $R^2 = 0.29$; all $P < 1 \times 10^{-6}$),
162 and PC1, PC3, sex, age, and alcohol use in SHOGUN/NCBI data (adjusted $R^2 = 0.27$; all $P <$
163 1×10^{-4}). Correlations between FLI and archaeal PC axes were not significant in
164 Centrifuge/GTDB data (at the chosen significance level, $P > 0.001$), and between FLI and
165 bacterial PC2 in SHOGUN/NCBI data ($P > 0.001$). In Centrifuge/GTDB data, the effect size
166 estimate on $\log_{10}(\text{FLI})$ was a magnitude larger for PC1 (0.11 ± 0.008) than for PC2 ($0.04 \pm$
167 0.008) and PC3 (-0.06 ± 0.008). The relationships between FLI and the bacterial PC components
168 representing their beta diversity in Centrifuge/GTDB data are visualized for each of the three
169 components in **Figure 1C**. A comparison of these relationships in Centrifuge/GTDB and
170 SHOGUN/NCBI is included in the SI (**Figure S3**).

171

172 We also further assessed the phylogenetic balances contributing to the PC axes in the
173 Centrifuge/GTDB data. Bacterial clades associated with higher FLI values, on the positive side

174 of the balances contributing to PC1, included members of orders *Lachnospirales* and
175 *Oscillospirales*, class *Bacilli*, and the *Ruminococcaceae*, *Bacteroidaceae* and *Lachnospiraceae*
176 families (**Figure S4**). Several clades had a negative association with FLI, on the negative side of
177 the balances contributing to PC1, such as order *Christensenellales* and genus *Faecalibacterium*.
178 In addition, genus *Bifidobacterium* in PC2, and family *Bifidobacteriaceae* in PC3 had negative
179 associations with continuous FLI.

180

181 ***Several bacterial taxa are differently abundant between the low and high FLI groups***

182 We also assessed significant differences in abundances of individual taxa between the high and
183 low FLI groups in FINRISK 2002. In Centrifuge/GTDB data, we identified 244 taxa (1% of
184 total) with an increased abundance, and 437 taxa (1.9%) with a decreased abundance in the high
185 FLI group (all Q values < 0.001 ; **Table S7**). In SHOGUN/NCBI data, 80 taxa (1.6%) had an
186 increased abundance, and 44 (0.9%) had a decreased abundance in the high FLI group (all Q
187 values < 0.001). While the number of associated taxa was higher in the Centrifuge/GTDB data
188 than SHOGUN/NCBI data, the proportion of significantly associated taxa was similar between
189 the two methods. In both data sets, family *Lachnospiraceae* comprised over 40% of taxa
190 positively associated with the high FLI group and *Bacteroidaceae* were in the top 3 most
191 common families. The negatively associated taxa were much more diverse, but
192 *Ruminococcaceae* and *Oscillospiraceae* were among the top 3 most common families in both
193 data sets (at least $> 6\%$ of all negatively associated taxa).

194

195 *Bacterial lineages within the NCBI Clostridium subclusters IV and XIVa associate*
196 *with FLI*

197 Continuous FLI and differences between FLI groups in the FINRISK 2002 cohort (FLI < 60, $N =$
198 4,359 and FLI \geq 60, $N = 1,910$; see **Figures 1A, 1B, Table S1**) were modeled with gradient
199 boosting regression or classification using Leave-One-Group-Out Cross-Validation (LOGOCV)
200 between participants from different regions. Only the bacterial PhILR transformed
201 Centrifuge/GTDB data were used here, to find robust associations between phylogenetically
202 related bacterial clades and fatty liver disease (instead of single taxa).

203
204 After feature selection and Bayesian hyperparameter optimization, the correlation between the
205 predictions of the final regression models (age, sex, self-reported alcohol use, and 18 bacterial
206 balances as features; each trained on the data from 5/6 regions) and true values in unseen data
207 from the omitted region averaged $R^2 = 0.30$ (0.26 – 0.33). After feature selection and
208 optimization, the main classification models (age, sex, and 11 bacterial balances as features; each
209 trained on the data from 5/6 regions) averaged AUC = 0.75 (**Table S2**) and AUPRC = 0.56
210 (baseline at 0.30; **Table S3**) on (unseen) test data from the omitted region. Models trained using
211 only the covariates averaged AUC = 0.71 (AUPRC = 0.47) and using only the 11 bacterial
212 balances they averaged AUC = 0.66 (AUPRC = 0.47) on test data. Alternative models were
213 constructed by excluding participants with FLI between 30 and 60 ($N = 1,583$) and discerning
214 between groups of FLI < 30 ($N = 2,776$) and FLI \geq 60 ($N = 1,910$). These models averaged AUC
215 = 0.80 (AUPRC = 0.75, baseline at 0.41) on their respective test data (**Tables S2, S3**). They
216 averaged AUC = 0.76 (AUPRC = 0.68) when using only the covariates, and AUC = 0.70
217 (AUPRC = 0.63) when using only the 20 bacterial balances.

218

219 Because training data from all 6 regions was used to prevent overfitting in the selection of core
220 features for all of the models, and similarly in searching for common hyperparameters,
221 participants from the validation region of each model (in the training partition) partly influenced
222 these parameters. Thus, we also constructed classification models discerning between the $FLI <$
223 60 and $FLI \geq 60$ groups, where data of the validation region was completely excluded in the
224 feature selection and hyperparameter optimization of each LOGOCV model. These models,
225 using their individual feature sets and hyperparameters, averaged $AUC = 0.75$ and $AUPRC =$
226 0.57 (baseline at 0.30) on test data from their respective validation regions (**Table S4**). Using
227 only covariates, they averaged $AUC = 0.71$ ($AUPRC = 0.47$), and $AUC = 0.67$ ($AUPRC = 0.48$)
228 with only the bacterial balances.

229

230 Our external validation data consisted of 258 participants after exclusion of pregnant participants
231 or those on antibiotics in the past 6 months, in the FINRISK 2007 population cohort³⁰ (**Table S1**,
232 **Figure S5**). The participants originate from North Karelia and Helsinki/Vantaa regions in
233 Finland, and their samples were processed with the same methodology as was used for FINRISK
234 2002 (with Centrifuge/GTDB approach and PhILR). In this external validation, the 6 full models
235 trained with covariates and the 11 bacterial balances in FINRISK 2002 averaged $AUC = 0.77$
236 ($AUPRC = 0.51$, baseline at 0.21 ; **Table S5**). The covariate-only models averaged $AUC = 0.72$
237 ($AUPRC = 0.40$) and the balance-only models averaged $AUC = 0.69$ ($AUPRC = 0.44$). The
238 receiver operating characteristic and precision-recall curves based on the averaged predictions of
239 the models, tested on this external validation data, also display good predictive ability ($AUC =$
240 0.78 , $AUPRC = 0.51$ with baseline at 0.51 ; **Figure S6**)

241
242 To facilitate interpretability of the results, we continued examining the main classification
243 models using a common set of core features. In these models, the median effect sizes of the
244 features on the model predictions at their minimum and maximum values were highest for age,
245 followed by sex, and the 11 balances in the phylogenetic tree (**Figures S7, S8**). All 11 associated
246 balances were in phylum *Firmicutes*, class *Clostridia*, and largely in the NCBI *Clostridium*
247 subclusters IV and XIVa (**Figure 2**). The specific taxa represented standardized GTDB genera
248 (NCBI in brackets) *Negativibacillus* (*Clostridium*), *Clostridium M* (*Lachnoclostridium* /
249 *Clostridium*), *CAG-81* (*Clostridium*), *Dorea* (*Merdimonas* / *Mordavella* / *Dorea* / *Clostridium* /
250 *Eubacterium*), *Faecalicatena* (*Blautia* / *Ruminococcus* / *Clostridium*), *Blautia* (*Blautia*),
251 *Sellimonas* (*Sellimonas* / *Drancourtella*), *Clostridium Q* (*Lachnoclostridium* [*Clostridium*]) and
252 *Tyzzarella* (*Tyzzarella* / *Coprococcus*). Notably, all but one of the features in the main
253 classification models (n226) were identified in the feature selection for the alternative models
254 (constructed otherwise identically, but $FLI < 30$ was compared against $FLI \geq 60$ in different data
255 partitions), together with 10 additional balances (**Figure S9**). Only one of the balances in the
256 alternative models was outside phylum *Firmicutes* (n1712 in *Bacteroidota*), and in addition, 4
257 balances were outside class *Clostridia* (n481 in *Negativicutes*; n826, n1009 and n918 in *Bacilli*).
258 Also, negative associations with the high FLI group were seen for *An181 sp002160325* in the
259 balance n266, where it is compared against the clade including *Dorea*, *Faecalicatena*,
260 *Sellimonas* and *Tyzzarella* species (**Figures 2, S8**). A higher abundance of the clade including
261 *Angelakisella*, *D5*, *Anaerotruncus* and *Phocea* species (against *Negativibacillus sp00435195* in
262 balance n97) was also negatively associated with high FLI.
263

264 In addition to blood test results, FLI is based on two anthropometric markers linked to metabolic
265 syndrome, waist circumference and BMI. This prompted us to dissect the Fatty Liver Index and
266 identify which of the covariates and associated microbial balances from the phylogenetic tree can
267 be linked to blood GGT and triglycerides measurements (see **Figure 1B**), and therefore would be
268 most specific to hepatic steatosis and liver damage.³¹ To do so, we performed feature selection
269 (similarly to continuous FLI) for GGT and triglycerides measurements in subsets of participants
270 grouped by age, sex, and BMI. The feature selection identified two balances within the NCBI
271 *Clostridia* XIVa subcluster (identified as n336 and n330) which were important for both GGT
272 and triglyceride level prediction, and thus likely specific to liver function (**Figure 2**). Bacterial
273 taxa were positively linked to liver function in these balances, and included (NCBI species)
274 *Clostridium clostridioforme*, *C. bolteae*, *C. citroniae*, *C. saccharolyticum* and *C. symbiosum*. On
275 the opposite, negatively associated side of the balances were, among others, (NCBI species)
276 *Hungatella effluvii*, *H. hathewayi*, and two new GTDB-defined species *Clostridium M*
277 *sp001517625* and *C. M sp000431375*.

278

279 ***Ethanol and acetate production pathways are identified in representative bacterial***
280 ***genomes from taxa linked to high FLI***

281 The values of predictive balances in the phylogenetic tree cannot be summarized for individual
282 taxa, which means that only a qualitative investigation of the associations between their
283 metabolism and fatty liver was possible in this study. We identified genetic pathways predicted
284 to encode for SCFA (acetate, propanoate, butanoate) and ethanol production, BA metabolism,
285 and choline degradation to trimethylamine (TMA) in representative genomes from the taxa we

286 identified to be linked to liver function (**Figure S8**). These processes were chosen because they
287 have been previously identified to have a mechanistic link to NAFLD (see *e.g.*, ref. ¹⁹).
288
289 Acetate and ethanol production pathways appeared to be more common in the representative
290 genomes of the taxa which had a positive association with FLI. In the liver function specific
291 clades, n336 and n330, MetaCyc pathways for pyruvate fermentation to ethanol III (PWY-6587)
292 and L-glutamate degradation V (via hydroxyglutarate; P162-PWY; produces acetate and
293 butanoate) were present only in genomes positively associated with FLI. In balance n336, also
294 heterolactic fermentation (P122-PWY; produces ethanol and lactate) was more often encoded in
295 the clade positively associated with the high FLI group (3/5) than the opposing negatively
296 associated clade (1/2). In representative genomes from the liver-specific balance n336, potential
297 ethanol producers (PWY-6587) were seen in the positively associated clade (*Clostridium M*
298 *clostridoforme A* and *Clostridium M sp000155435*), and not in the negatively associated clade
299 (*Clostridium M sp001517625* and *Clostridium M sp000431375*). However, for most balances
300 such trends were not clear in the qualitative analysis. Furthermore, we did not detect any of these
301 pathways in the representative genomes of two individual taxa positively associated with FLI,
302 *Negativibacillus sp000435195* and *Phoceia massiliensis* (**Figure S8**).

303

304 **Discussion**

305 The pathophysiology of fatty liver disease in general, and NAFLD in particular, is complex and
306 its clinical diagnosis can be difficult.³² In this study, we utilized metagenomic data from a large
307 population cohort (FINRISK 2002³⁰) to identify broad links between the overall gut
308 microbiome composition and fatty liver disease, using FLI as a recognized proxy (**Figure 1C**),

309 and identified specific microbial taxa and lineages positively and negatively associated with the
310 high FLI group (**Figure 2**). It should be noted, that FLI used in our study as a proxy for liver
311 disease also includes features such as BMI and waist circumference, which associate with
312 metabolic syndrome and diabetes.¹⁶ Links between these diseases and gut microbiome
313 composition are well documented in previous studies.³³ However, fatty liver disease is
314 increasingly thought to be a component of the metabolic syndrome,^{4,34} and while diabetes
315 prevalence is higher in the high FLI group in FINRISK 2002, affected participants still consist
316 only 11% of this group (**Table S1**). Furthermore, we would like to emphasize that our results
317 are not suitable for current clinical application, and should be validated by further, preferably
318 mechanistic studies. We also do not know if our results generalize outside the Finnish
319 population, as all participants in this study were exclusively from Finnish cohorts.

320

321 Considering that the predictive ability of FLI for clinically diagnosed NAFLD ranges between
322 AUC = 0.81 – 0.93, in populations of Caucasian ethnicity such as the Finnish population,²³ our
323 models were able to reasonably predict the FLI group with AUC = 0.75 (AUPRC = 0.56,
324 baseline at 0.30), in our internal cross-region validation. Furthermore, the performance of our
325 predictive models was highly similar in an external, Finnish validation cohort (AUC = 0.77,
326 AUPRC = 0.51, baseline at 0.21).

327

328 Our additional analyses support these main results. While a thorough method comparison is
329 beyond the scope of the current study, the results from the two taxa assignments were very
330 similar despite their differences, such as the fourfold higher number of taxa in the
331 Centrifuge/GTDB data. In the machine learning models (performed only with

332 Centrifuge/GTDB data), excluding participants with intermediate FLI (between 30 – 60)
333 increased the accuracy slightly in the internal cross-validation (to AUC = 0.8 and AUPRC =
334 0.75, baseline at 0.41). However, discerning between participants with probable fatty liver
335 disease ($FLI \geq 60$) from others presents a clinically more relevant target for detecting changes
336 in microbiome composition associated with development of the disease. In another set of
337 models, we negated the influence of validation region data in the individual models also for
338 feature selection and hyperparameter optimization during training. This led to individualized
339 sets of features and parameters in the models, but the average performance of the models was
340 almost identical on validation region samples in the internal cross-validation (AUC = 0.75 and
341 AUPRC 0.57, baseline at 0.30). The aim of our study was to find patterns in microbiome
342 composition which would be generalizable across the 6 sampled geographic regions in Finland
343 and easy to interpret. Thus, we consider the use of all training data to define the common core
344 feature set justified. This goal also guided our overall modeling architecture and likely led to a
345 lower performance than if we instead performed interpolation within a smaller scale (see *e.g.*,
346 ref. ³⁵).

347

348 When interpreting our results, several levels of associations can be considered according to
349 types of fatty liver disease and the gut microbiome composition. Because FLI has been mostly
350 validated with simple steatosis and NAFLD,^{23,26} we can conservatively contextualize our
351 findings with previous associative work that used these diagnoses or clinical manifestations,
352 only. The cutoff used in our study at $FLI \geq 60$ has been used to rule in liver steatosis in a
353 Caucasian cohort comparable to ours,²⁶ but also a cutoff at $FLI \geq 48$ has been found appropriate
354 for simple steatosis in a Portuguese cohort.³⁶ Much lower cutoffs ($FLI \geq 20$ to 30) have been

355 used in Asian cohorts.^{37–39} Thus, it is likely that our high FLI groups include most participants
356 with liver steatosis or fibrosis in both FINRISK cohorts, but the low FLI group also likely
357 includes participants with low grade steatosis.

358

359 *Traditional statistical analyses replicate previous findings on gut microbiome*
360 *composition and fatty liver disease when using FLI as a risk index*

361 Among the significant high level FLI-associated differences in the gut microbiomes of the
362 participants in FINRISK 2002, we found a 14.7% lower Shannon alpha diversity in the high
363 FLI group with SHOGUN/NCBI taxa assignments and 13.4% lower diversity with
364 Centrifuge/GTDB assignments. These results are in good accordance with previous results of
365 decreased gut bacterial diversity in patients with biopsy-proven non-alcoholic steatohepatitis
366 (NASH), the most serious form of NAFLD.⁴⁰ In this case-control study, the Shannon diversity
367 of gut microbiomes in NASH patients without liver cirrhosis was on average 7% lower
368 compared to controls, and in patients with cirrhosis, 14% lower. A significantly decreased gut
369 microbiome alpha diversity of similar magnitude was also seen in cohort participants with
370 persistent NAFLD compared to controls.⁴¹

371

372 In both the SHOGUN/NCBI and Centrifuge/GTDB data, we found significant linear
373 correlations between FLI and beta diversity, or two or three main bacterial PC-axes of the
374 samples, respectively (**Figures 1C, S3**). The model fit was slightly better with
375 Centrifuge/GTDB data, which might be due to the higher number of identified taxa, and thus
376 increased taxonomic resolution (although including putative species in GTDB). Our results
377 support previous observations of differences in beta diversity in relation to persistent

378 NAFLD,⁴¹ and along the NAFLD-cirrhosis spectrum.⁸ Through the loadings of the
379 phylogenetic balances on the PC axes in the Centrifuge/GTDB data, we detected several
380 previously known connections between microbial clades and FLI (**Figure S4**). Among others,
381 we observed a positive association between high FLI and family *Lachnospiraceae* and negative
382 associations for order *Christensenellales*, genus *Faecalibacterium*, and genus *Bifidobacterium*.
383 The positive association is supported by previous findings of their connection with obesity,⁴²
384 and the negative associations by connections to lean individuals and healthy gut microbiome
385 composition.^{43–45}
386
387 Our differential abundance analysis also detected a high number of taxa with significantly
388 increased or decreased abundance in the high FLI group. All following results were observed
389 both in the Centrifuge/GTDB and SHOGUN/NCBI data sets, unless otherwise noted. Majority
390 of the taxa with increased abundance in the high FLI group were from family *Lachnospiraceae*,
391 which supports their positive association with NAFLD reported previously in a number of
392 studies,⁴⁶ but also with obesity (**Table S7**).⁴² The increased abundance of genus *Roseburia* has
393 also been highlighted as a characteristic change in gut microbiome related to NAFLD.^{46,47} In
394 the current study, two members of genus *Roseburia* were in the top 10 taxa most strongly
395 associated with high FLI. Furthermore, our results support previous findings on the positive
396 associations of, for example, *Collinsella*,⁴⁰ *Prevotella copri*,⁴⁸ *Dorea*,⁴⁷ with NAFLD. We also
397 detected increases in *Sutterella* and *Streptococcus*, previously associated with cirrhosis.⁴⁹
398 However, we did not find increases in families *Kiloniellaceae* and *Pasteurellaceae*, previously
399 associated with NAFLD.⁴⁶ Among the individual taxa negatively associated with high FLI,
400 families *Ruminococcaceae* and *Oscillospiraceae* (such as genus *Oscillibacter*) were common,

401 which supports previous findings on their connections with NAFLD.^{12,41,46} A high number of
402 putative (GTDB) species were negatively associated with FLI in the Centrifuge/GTDB data,
403 which were understandably not present in the SHOGUN/NCBI data. Many of these were
404 classified in the recently described order *Christensenellales*,²⁸ including families such as CAG-
405 74, associated with healthy participants,⁵⁰ and *Christensenellaceae*, which are widespread,
406 highly heritable, and associated with health.^{44,51}

407

408 While our results from common statistical experiments mainly supported previous findings, we
409 chose to leverage the phylogenetic information included in the GTDB data to find robust
410 associations between larger bacterial clades and fatty liver disease in the Finnish population.
411 This was accomplished by constructing predictive models to classify participants in the FLI
412 groups based on the phylogenetic balances and covariates, subjected to feature selection and
413 geographical cross-validation.

414

415 ***Predictive modeling of FLI reveals consistent associations between gram-positive***
416 ***Clostridia and fatty liver disease***

417 Strikingly, the strongest associations with FLI in our machine learning models were all inside
418 the *Firmicutes* phylum. A possible reason for this might be the higher relative abundance of
419 phylum *Firmicutes* at high latitudes,⁵² where Finland is. Among the associations we identified,
420 *Faecalicatena gnavus* (NCBI: *Ruminococcus gnavus*) was positively linked with FLI as part of
421 3 predictive balances, and associated in previous studies with liver cirrhosis.¹⁷ In their study,
422 oral *Firmicutes*, such as *Veillonella*, were suggested to invade the gut. While our balance-based
423 approach did not detect these taxa, *Megasphaera elsdenii* were positively associated with the

424 high FLI group in our differential abundance analyses (**Table S7**). This might be due to the
425 strict feature selection employed prior to the predictive modeling.
426
427 Two individual taxa, *Negativibacillus sp000435195* and *Phoceia massiliensis*, both had strong
428 positive associations with the high FLI group (**Figure 2**), but the balances including these
429 species were not predictive of the liver function-specific components (triglycerides and GGT).
430 Positive associations of these taxa with fatty liver disease have not been documented
431 previously. However, a decreasing abundance of both bacteria, *Negativibacillus sp000435195*
432 (NCBI: *Clostridium* sp. CAG:169) and *Phoceia massiliensis* (NCBI: *Phoceia massiliensis*), were
433 seen when the intake of meat and refined cereal was reduced isocalorically in favor of fruit,
434 vegetables, wholegrain cereal, legumes, fish and nuts in overweight and obese subjects in
435 Italy.⁵³ While comparisons between these studies are difficult due to differences in taxa
436 annotations, bacteria such as *Faecalicatena gnnavus* (NCBI: *Ruminococcus gnnavus*) and
437 *Clostridium Q saccharolyticum* (NCBI: *Clostridium saccharolyticum*) were also found to
438 respond negatively to the Mediterranean diet. Thus, further study on the connections of these
439 bacteria with gut health and diet is warranted.
440
441 Among the taxa negatively associated with high FLI, *Hungatella* (see balance n332, **Figure 2**)
442 have been previously shown to correlate negatively with the obesity phenotype in mice⁵⁴ and
443 *H. hathewayi* was found to be a common commensal in the gut of healthy volunteers.⁵⁵
444 However, genus *Hungatella* has also been positively associated with concentrations of
445 trimethylamine-N-oxide (TMAO),⁵⁶ a metabolite associated with cardiovascular disease and
446 NAFLD. In our study, on the positively associated side (of balance n332) opposite to genus

447 *Hungatella* was a novel GTDB species, *CAG-81* sp000435795, previously included in NCBI
448 genus *Clostridium*. The *CAG-81* genus was recently positively associated with TMAO levels in
449 urine in a study using the GTDB classification.⁵⁷ While we did not find the pathway for TMA
450 (precursor to TMAO) production in its genome, this would explain the positive association of
451 the *CAG-81* species with high FLI. Furthermore, the previous contradictory results among these
452 taxa could be explained by grouping of putatively TMA producing taxa in *CAG-81* together
453 with the closely related genus *Hungatella*.

454

455 Most taxa in our study with a positive association with FLI belonged to the broadly defined
456 *Clostridium* NCBI genus, which supports several previous observations.^{14,46,58} However,
457 taxonomic standardization according to GTDB has identified the *Clostridium* genus as the most
458 phylogenetically inconsistent of all bacterial genera in the NCBI taxonomy, and divides it into
459 a total of 121 monophyletic genera in 29 distinct families.²⁸ The GTDB reassignment
460 complicates comparisons to previous studies, but it is phylogenetically and biologically
461 sensible, and can thus provide new insights into the microbiomes. Our results also strongly
462 suggest that despite its higher cost compared to metabarcoding, the increased resolution of
463 (shallow) shotgun metagenomic sequencing is highly useful in identifying specific taxon-
464 disease associations (see *e.g.*, refs. ^{27,59}).

465

466 ***Bacterial taxa positively associated with high FLI have a genetic potential to***
467 ***exacerbate the development of fatty liver disease***

468 We identified several plausible new associations between individual taxa and clades of bacteria
469 and fatty liver. All taxa were from class *Clostridia*, which are obligate anaerobes. We observed

470 that reference genomes from the bacterial taxa positively associated with high FLI in the liver-
471 specific balances harbored several genetic pathways necessary for ethanol production.
472 Specifically, genes predicted to enable the fermentation of pyruvate to ethanol (MetaCyc PWY-
473 6587) appeared to be common. Endogenous production of ethanol has been known to both
474 induce hepatic steatosis and increase intestinal permeability,⁶⁰ and several of the taxa
475 associated with the high FLI group have also been experimentally shown to produce ethanol,
476 such as *C. M asparagiforme*, *C. M bolteae*, *C. M clostridioforme* / *C. M clostridioforme A*⁶¹,
477 and *C. Q Saccharolyticum*.⁶² The relative abundances of these putatively ethanol-producing
478 taxa were also predictive of FLI groups in previously unseen data. However, the self-reported
479 alcohol consumption from the participants was not among the best predictors for the FLI
480 groups, as it was excluded in the feature selection step.

481
482 All reference genomes from taxa positively associated with FLI in balance n330 harbored
483 genes predicted to encode for the L-glutamate fermentation V (P162-PWY; **Figure S8**)
484 pathway, which results in the production of acetate and butanoate. Glutamate fermentation
485 could lead to increased microbial protein fermentation in the gut, which has been previously
486 been linked with obesity, diabetes and NAFLD.⁶³ Recently, the combined intake of fructose
487 and microbial acetate production in the gut was experimentally observed to contribute to
488 lipogenesis in the liver in a mouse model.⁶⁴ Interestingly, *C. Q saccharolyticum* (in our study, a
489 taxa positively associated with high FLI deriving from balance n330), was experimentally
490 shown to ferment various carbohydrates, including fructose, to acetate, hydrogen, carbon
491 dioxide, and ethanol.⁶² Furthermore, while our own pathway analysis did not detect BA
492 modification pathways in the reference genome of *C. Q saccharolyticum*, a strain of this

493 species has been highlighted as a probable contributor to NAFLD development through the
494 synthesis of secondary BA.¹⁵ The links between dietary intake and gene regulation, combined
495 with microbial fermentation in the gut warrant further mechanistic experiments to elucidate
496 their links with fatty liver, and likely other metabolic diseases.

497

498 NAFLD-associated ethanol producing bacteria in previous cohort studies have all been gram-
499 negatives, such as (NCBI-defined) *Klebsiella pneumoniae*,¹³ and *Escherichia coli*.¹² In our
500 population sample, instead of gram-negatives, bacteria from the *C. M bolteae*, *C. M*
501 *clostridioforme* / *C. M clostridioforme A* and *C. M citroniae* species (positively associated with
502 high FLI in balance n336) have been described as opportunistic pathogens,⁶⁵ and are
503 hypothesized to exacerbate fatty liver development similarly through endogenous ethanol
504 production. This result suggests that geographical,³⁵ and ethnic variability,⁶⁶ might also
505 strongly affect gut microbiome composition and its associations with disease. In addition to
506 putative endogenous ethanol producers, we identified other taxa positively associated with high
507 FLI in balance n330, for which reference genomes harbored a genetic pathway predicted to
508 encode for the ability to ferment L-lysine to acetate and butyrate. While the production of these
509 SCFAs is often considered beneficial for gut health, other metabolism of proteolytic bacteria
510 might negatively contribute to fatty liver disease.⁶⁷

511

512 Through modeling a previously validated index for fatty liver, FLI, we found replicable
513 associations with specific microbial taxa and likely liver disease of the participants. In addition,
514 sex and age of participants were also strongly predictive of the FLI group in our models
515 (**Figures 2, S7**). Their similar positive associations with fatty liver disease are known from

516 previous studies.^{68,69} The associated microbial balances could be used to improve the
517 predictions above the baseline of these covariates on 5/6 regions in Finland in the main cohort.
518 For example, in the model cross-validated with Lapland the balances were more predictive of
519 FLI group than the covariates by themselves, and their combination increased the AUC further.
520 Yet, when testing the model where Turku/Loimaa region was used for internal cross-validation,
521 the microbial balances were slightly predictive of FLI group but failed to improve the AUC
522 over the covariates (**Table S2**). This pattern might stem from the cultural and genetic west-east
523 division in Finland,^{70,71} with a closer proximity of the Helsinki/Vantaa region to eastern regions
524 than Turku/Loimaa, in both terms. It is thus likely that further incorporation and investigation
525 on the use of spatial information in microbiome modeling would elucidate these geographical
526 patterns in taxa-disease associations.

527

528 Our models were also able to accurately predict the FLI group of participants in the external
529 validation cohort, which were from the North Karelia and Helsinki/Vantaa regions. The
530 observed difficulty to geographically extrapolate taxa-disease associations³⁵ might mean that
531 associations reported in our study are specific to Finland and nearby regions. Notably, many of
532 the positive associations between specific taxa and fatty liver disease have not been reported
533 previously, but the functional potential of these taxa inferred from genomic data is similar to
534 taxa positively associated with NAFLD in previous studies. Thus, the geographical limits of
535 taxa-disease associations reported in studies such as ours warrant further study. Unfortunately,
536 generalization of our own results outside of Finland also remains to be addressed.

537

538 It is likely that not all associations in the current study are related solely to liver steatosis,
539 because FLI is based on measurements related to metabolic syndrome. However, our approach
540 is supported by recent views of NAFLD as the integral liver component of the metabolic
541 syndrome.^{34,72} Indeed, the prevalences of diabetes and cardiovascular disease in both FINRISK
542 2002 and 2007 cohorts are elevated in the high FLI group, although the majority of the high
543 FLI participants did not have either of these diagnoses at the time of sampling (**Table S1**). We
544 also dissected the FLI by dividing participants into age/sex/BMI groups and detected microbial
545 groups specific to the blood work measurements of liver damage, triglycerides and GGT. These
546 associated taxa can thus be thought of as most closely associated with liver function, if such a
547 division is deemed practical.

548

549 ***Conclusions***

550 Modeling an established risk index for fatty liver enabled the detection of associations between
551 the disease and gut microbiome composition, to the level of individual taxa. While utilizing
552 FLI as a proxy, NCBI taxa identified with standard statistical methods were supportive of
553 previously reported differences between NAFLD cases and healthy controls. In our machine
554 learning framework, all clades robustly predictive of the FLI group were from the obligately
555 anaerobic gram-positive class *Clostridia*, representing several redefined GTDB genera
556 previously included in the NCBI genus *Clostridium*. Many of the representative genomes of
557 taxa positively associated with high FLI had genomic potential for endogenous ethanol
558 production. Our results support previous findings on the likely contribution of ethanol and
559 increased gut permeability on the induction of hepatic steatosis. Further support was also found
560 for the involvement TMA and SCFAs, especially acetate, in the likely pathophysiology of fatty

561 liver disease. Our models were able to predict the FLI group of participants in Finland across
562 geographical regions and in an external Finnish cohort, showing that the associations are robust
563 and generalizable in this population. Based on our results, mechanistic connections between
564 specific microbes and fatty liver disease, and the geographical differences in such taxa-disease
565 associations, should be addressed in further studies.

566

567 **Materials and Methods**

568 *Survey details and sample collection*

569 Cardiovascular disease risk factors have been monitored in Finland since 1972 by conducting a
570 representative population survey every five years.³⁰ In the FINRISK 2002 survey, a stratified
571 random population sample was conducted on six geographical regions in Finland. These are
572 North Karelia and Northern Savo in eastern Finland, Turku and Loimaa regions in southwestern
573 Finland, the cities of Helsinki and Vantaa in the capital region, the provinces of Northern
574 Ostrobothnia and Kainuu in northwestern Finland, and the province of Lapland in northern
575 Finland.

576

577 Briefly, at baseline examination the participants filled out a questionnaire form, and trained
578 nurses carried out a physical examination and blood sampling in local health centers or other
579 survey sites. Data was collected for physiological measures, biomarkers, and dietary,
580 demographic and lifestyle factors. Stool samples were collected by giving willing participants a
581 stool sampling kit with detailed instructions. These samples were mailed overnight between
582 Monday and Thursday under Finnish winter conditions to the laboratory of the Finnish Institute

583 for Health and Welfare, where they were stored at -20°C. In 2017, the samples were shipped still
584 unthawed to University of California San Diego for microbiome sequencing.

585

586 Details of the FINRISK cohorts analyzed in this study are included in the supplementary files
587 (**Table S1**). Further details and sampling have also been extensively covered in previous
588 publications (see refs. ^{24,73}). The Coordinating Ethics Committee of the Helsinki University
589 Hospital District approved the study protocol for FINRISK 2002 (Ref. 558/E3/2001), and all
590 participants have given their written informed consent.

591

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

593 DNA extraction was performed according to the Earth Microbiome Project protocols, with the
594 MagAttract PowerSoil DNA kit (Qiagen), as previously described.⁷⁴ A miniaturized version of
595 the Kapa HyperPlus Illumina-compatible library prep kit (Kapa Biosystems) was used for library
596 generation, following the previously published protocol.⁷⁵ DNA extracts were normalized to 5 ng
597 total input per sample in an Echo 550 acoustic liquid handling robot (Labcyte Inc.). A Mosquito
598 HV liquid-handling robot (TTP Labtech Inc.) was used for 1/10 scale enzymatic fragmentation,
599 end-repair, and adapter-ligation reactions. Sequencing adapters were based on the iTru
600 protocol,⁷⁶ in which short universal adapter stubs are ligated first and then sample-specific
601 barcoded sequences added in a subsequent PCR step. Amplified and barcoded libraries were then
602 quantified by the PicoGreen assay and pooled in approximately equimolar ratios before being
603 sequenced on an Illumina HiSeq 4000 instrument to an average read count of approximately
604 900,000 reads per sample.

605

606 *Taxonomic matching and phylogenetic transforms*

607 We quality trimmed the sequences and removed the sequencing adapters with Atropos.⁷⁷ Host
608 reads were removed by mapping the reads against the human genome assembly GRCh38 with
609 Bowtie2.⁷⁸ To improve the taxonomic assignments of our reads, we used a custom index,⁷⁹ based
610 on the Genome Taxonomy Database (GTDB) release 89 taxonomic redefinitions,^{28,80} for read
611 classification with default parameters in Centrifuge 1.0.4.⁸¹ Viral and eukaryotic sequences were
612 removed in this step, as the database contains only bacterial and archaeal reference genomes.
613 After read classification, all following steps were performed with R version 3.5.2,⁸² using
614 phyloseq 1.30.0⁸³ to manage the data. To reduce the number of spurious read assignments, and to
615 facilitate more accurate phylogenetic transformations, only reads classified at the species level,
616 matching individual GTDB reference genomes, were retained. Samples with less than 50,000
617 reads, from pregnant participants or recorded antibiotic use in the past 6 months were removed,
618 resulting in a final number of 6,269 samples. We first filtered taxa not seen with more than 3
619 counts in at least 1% of samples and those with a coefficient of variation ≤ 3 across all samples,
620 following McMurdie and Holmes⁸³, with a slight adaption from 20% of samples to 1% of
621 samples, because of our larger sample size. The complete bacterial and archaeal phylogenetic
622 trees of the GTDB release 89 reference genomes, constructed from an alignment of 120 bacterial
623 or 122 archaeal marker genes,²⁸ were then combined with our taxa tables. The resulting trees
624 were thus subset only to species which were observed in at least one sample in our data. The read
625 counts were transformed to phylogenetic node balances in both trees with PhILR.²⁹ The default
626 method for PhILR inputs a pseudocount of 1 for taxa absent in an individual sample before the
627 transform.
628

629 In this study, we did not specifically and solely use relative abundances at various taxonomic
630 levels, as is common practice for microbiome studies. Instead, we applied a PhILR
631 transformation to our microbial composition data,²⁹ introducing the concept of microbial
632 “balances”. Indeed, evolutionary relationships of all species harbored in each microbiome
633 sample can be represented on a phylogenetic tree, with species typically shown as external nodes
634 that are related to each other by multiple branches connected by internal nodes. In this context,
635 the value of a given microbial “balance” is defined as the log-ratio of the geometric mean
636 abundance between two groups of microbes descending from the same corresponding internal
637 node on a microbial phylogenetic tree. This phylogenetic transform was used because it i)
638 addresses the compositionality of the metagenomic read data,⁸⁴ ii) permits simultaneous
639 comparison of all clades without merging the taxa by predefined taxonomic levels, and iii)
640 enables evolutionary insights into the microbial community. The links between microbes and
641 their environment, such as the human gut, is mediated by their functions. Different functions are
642 known to be conserved at different taxonomic resolutions, and most often at multiple different
643 resolutions.⁸⁵ Thus, associations between the microbes and the response variable are likely not
644 best explained by predefined taxonomic levels. In the absence of functional data, concurrently
645 analyzing all clades (partitioned by the nodes in the phylogenetic tree) would likely enable the
646 detection of the associations at the appropriate resolution depending on the function and the local
647 tree topography.

648

649 To further validate our approach, assess how the use of the GTDB taxonomic redefinitions and
650 custom database affected our results, and to facilitate comparisons with previous results, we
651 annotated our raw reads in FINRISK 2002 samples also with NCBI taxonomy and performed

652 several additional analyses. For this comparison data, after quality trimming the FINRISK 2002
653 reads and removing host sequences as described above, SHOGUN v1.0.5⁵⁹ was used for
654 taxonomy assignments against the NCBI RefSeq version 82 (May 8, 2017) database containing
655 complete bacterial, archaeal, and viral genomes. To facilitate comparisons between different
656 annotations, we subset the samples included in the SHOGUN/NCBI annotated data to those
657 included in the Centrifuge/GTDB data (for exclusion criteria, see above).

658

659 *Covariates*

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

664 $(e^{0.953 \cdot \log_e(\text{triglycerides mg/dL}) + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745}) /$

665 $(1 + e^{0.953 \cdot \log_e(\text{triglycerides mg/dL}) + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745}) * 100$. We chose

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

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

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

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

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

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

672 in FINRISK 2002 (0.3%) and the self-reported alcohol use was missing for 247 participants

673 (3.9%). In the machine learning framework, missing values for FLI and self-reported alcohol use

674 were mean imputed. However, for the feature selection to identify liver function-specific

675 balances, GGT, triglycerides and BMI were not imputed but observations where any of these
676 were missing were simply removed.

677

678 *Taxa composition and alpha diversity*

679 The baseline compositions of the microbial communities in the samples were summarized at
680 phylum level in the different FLI groups (< 60 and ≥ 60 FLI) with the Centrifuge/GTDB data in
681 FINRISK 2002 and 2007 (**Figure S1**), and with SHOGUN/NCBI data in FINRISK 2002 (**Figure**
682 **S2**) by total sum scaling and merging taxa at phylum level, separately for bacteria and archaea.

683

684 Bacterial alpha diversity of each individual sample in FINRISK 2002 was estimated through
685 Shannon diversity as the mean of 10 random rarefactions of raw annotated read counts (see ref.
686 ⁸⁶), separately in both the Centrifuge/GTDB and SHOGUN/NCBI data sets. Associations
687 between the FLI group (< 60 or ≥ 60 FLI) and Shannon diversity in the data sets was modeled
688 using binomial regression and adjusted for age, sex, and self-reported alcohol use, using “glm” in
689 base R.⁸²

690

691 *Beta diversity and linear modeling of FLI*

692 In the Centrifuge/GTDB data, beta diversity was calculated as Euclidian distance of the PhILR
693 balances through Principal Component Analysis (PCA) on bacterial and archaeal balances
694 separately with “rda” in vegan 2.5.6.⁸⁷ To calculate beta diversity with the SHOGUN/GTDB
695 data, raw bacterial taxa counts were centered log-ratio (CLR) transformed with “transform” in
696 microbiome 1.8.0,⁸⁸ and their Euclidian distances were obtained similarly through PCA. Linear
697 regression models were constructed for FLI with “lm” in base R⁸² with Centrifuge/GTDB data,

698 and separately with SHOGUN/NCBI. $\text{Log}_{10}(\text{FLI})$ was used as the dependent variable and the
699 first three bacterial PCs, sex, age, and self-reported alcohol were used as the independent
700 variables. Archaeal PCs were not included in the models because none of them were
701 significantly correlated with FLI in Centrifuge/GTDB data (all $P > 0.001$). To visualize the
702 association between beta diversity and FLI, the FLI of each participant was plotted against its
703 quantiles along the three bacterial PC axes in Centrifuge/GTDB data (**Figure 1C**). A comparison
704 of the associations with the alternative SHOGUN/NCBI annotated data was also included in the
705 SI (**Figure S3**).

706

707 *Differential abundance of individual taxa between the FLI groups*

708 To facilitate comparisons to previous studies, we assessed the associations between the FLI
709 group (< 60 or ≥ 60 FLI) of the participants and Centrifuge/GTDB and SHOGUN/NCBI
710 annotated individual taxa present in the samples. With both data sets, differential abundance of
711 the bacterial taxa between the FLI groups was assessed with the ALDEx2 compositional data
712 analysis tool.⁸⁹ Briefly, significance of the abundance differences between the groups were
713 estimated with a Welch's t-test, and only taxa with (Benjamini Hochberg) false discovery rate -
714 adjusted P values (or Q values) < 0.001 were retained. The associated taxa were then divided in
715 each data set to those positively or negatively associated with the high FLI group and sorted
716 based on effect sizes estimated from the median CLR differences between the groups.

717

718 *FLI modeling within a machine learning framework*

719 In the machine learning framework, both regression and categorical models were constructed for
720 FLI, using only the Centrifuge/GTDB data. The feature selection, hyperparameter optimization

721 and internal cross-validation methods were identical for both approaches, unless otherwise
722 stated. The continuous or categorical FLI (groups of $FLI < 60$ and $FLI \geq 60$) were modeled with
723 xgboost 0.90.0.2,⁹⁰ by using both bacterial and archaeal balances, sex, age, and self-reported
724 alcohol use as preliminary predictor features. We used FLI 60 as the cutoff for ruling in fatty
725 liver (steatosis) for the classification, after Bedogni et al., (2006).²⁶ The data was first split to
726 70% train and 30% test sets while preserving sex and region balance. To take into account
727 geographical differences (see *e.g.*, ref. ³⁵) and to find robust patterns across all 6 sampled regions
728 in Finland between the features and FLI group, we used Leave-One-Group-Out Cross-Validation
729 (LOGOCV) inside the 70% train set to construct 6 separate models in each optimization step.
730 Because of high dimensionality of the data (3,423 predictor features) feature selection by
731 filtering was first performed inside the training data, based on random forest permutation as
732 recommended by Bommert et al.⁹¹ Briefly, permutation importance is based on accuracy, or
733 specifically the difference in accuracy between real and permuted (random) values of the specific
734 variable, averaged in all trees across the whole random forest. The permutation importance in
735 models based on the 6 LOGOCV subsets of the training data were calculated with mlr 2.16.0,⁹²
736 and the simple intersect between the top 50 features in all LOGOCV subsets were retained as the
737 final set of features. Thus, the feature selection was influenced by the training data from all 6
738 geographical regions, but this only serves to limit the number of chosen features because of the
739 required simple intersect. This approach was used to obtain a set of core predictive features
740 which would have potential for generalizability across the regions. The number of features
741 included in the models by this approach was deemed appropriate, since the relative effect size of
742 the last included predictor was very small (< 0.1 change in classification probability across its
743 range).

744

745 Bayesian hyperparameter optimization of the xgboost models was then performed with only the

746 selected features. An optimal set of parameters for the xgboost models were searched over all

747 LOGOCV subsets with “mbo” in mlrMBO 1.1.3,⁹³ using 30 preliminary rounds with randomized

748 parameters, followed by 100 optimization rounds. Parameters in the xgboost models and their

749 considered ranges were learning rate (eta) [0.001, 0.3], gamma [0.1, 5], maximum depth of a tree

750 [2, 8], minimum child weight [1, 10], fraction of data subsampled per each iteration [0.2, 0.8],

751 fraction of columns subsampled per tree [0.2, 0.9], and maximum number of iterations (nrounds)

752 [50, 5000]. The parameters recommended by these searchers were as following for regression:

753 eta=0.00889; gamma=2.08; max_depth=2; min_child_weight=8; subsample=0.783;

754 colsample_bytree=0.672; nrounds=1,810, and for classification: eta=0.00107; gamma=0.137;

755 max_depth=5; min_child_weight=9; subsample=0.207; colsample_bytree=0.793;

756 nrounds=4,328. We used Root-Mean-Square Error (RMSE) for the regression models and Area

757 Under the ROC Curve (AUC) for the classification models to measure model fit on the left-out

758 data (region) in each LOGOCV subset. Receiver operating characteristic and precision-recall

759 curves for these validation metrics were calculated with “evalmod” in precrec v0.11.2.⁹⁴ The

760 final models were trained on the LOGOCV subset training data, the data from one region thus

761 omitted per model, and using the selected features and optimized hyperparameters. Internal

762 validation of these models was conducted against participants only from the region omitted from

763 each model, in the 30% test data which was not used in model training or optimization.

764 Sensitivity analysis was conducted by using only the predictive covariates (sex and age) or

765 balances separately, with the same hyperparameters, data partitions and cross-region internal

766 validation as for the full models.

767

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

769 Because the classification models have a more clinically relevant modeling target for the
770 difference between $FLI < 60$ and $FLI \geq 60$, the latter used to rule in fatty liver,²⁶ we further
771 interpreted the partial dependence of their predictions. Partial dependence of the classification
772 model predictions on individual features was calculated with “partial” in pdp 0.7.0.⁹⁵ The partial
773 dependence of the features on the model predictions was also plotted, overlaying the results from
774 each of the 6 models. For each feature, its relative effect on the model prediction was estimated
775 as medians of the minimum and maximum yhat (output probability of the model for the $FLI \geq 60$
776 class), calculated at the minimum and maximum values of the feature separately in each of the 6
777 models. The relative effects of the balances were then overlaid as a heatmap on a genome
778 cladogram which covers all balances in the model with ggtree 2.1.1.⁹⁶

779

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

781 ***FLI < 30 and FLI ≥ 60 groups***

782 To assess robustness of the models and how removing the participants with intermediate FLI
783 (between 30 and 60) affects model performance, we removed this group ($N = 1910$) and
784 constructed alternative classification models to discern between the $FLI < 30$ and $FLI \geq 60$
785 groups. Other than removing the intermediate FLI participants and resulting new random split to
786 the train (70%) and test (30%) sets, these models were constructed identically to the main
787 models, including LOGOCV design, feature selection, and hyperparameter optimization. The
788 recommended parameters for this classification task were $\eta=0.00102$; $\gamma=3.7$;
789 $\max_depth=2$; $\min_child_weight=5$; $subsample=0.49$; $colsample_bytree=0.631$; $nrounds=3,119$.

790 Interpretation of partial dependence was also performed identically, but only the relative effects
791 of the model features were plotted without a cladogram.

792

793 ***Exclusion of validation region data from feature selection and hyperparameter***
794 ***optimization***

795 Because training data from all 6 regions is used to inform the selection of optimal features and
796 hyperparameters, the validation region data cannot be considered completely independent from
797 the training of the LOGOCV models. Thus, we constructed a set of classification models for the
798 $FLI \geq 60$ and $FLI < 60$ groups, where all validation region samples also in the training data were
799 excluded from the simple intercept of top 50 features in each LOGOCV set and from the
800 subsequent hyperparameter optimization. These models with individualized features and
801 hyperparameters were then tested on the validation region samples in the unseen test data to
802 estimate how model performance was affected. The main test (70%) and train (30%) sets were
803 identical to the main models, but additionally 6 randomized 70/30 splits nested inside the test set
804 (excluding the validation region) were used in hyperparameter optimization to reduce overfitting.

805 Average optimal hyperparameters in the 6 models were $\eta=0.00106$; $\gamma=4.3$;
806 $\max_depth=2$; $\min_child_weight=7$; $subsample=0.36$; $colsample_bytree=0.613$; $nrounds=1,772$.

807

808 ***External validation of the models in a separate population cohort***

809 To further validate our models and results, we leveraged the data from a more recent population
810 cohort in Finland, FINRISK 2007 (see **Table S1**). In this cohort, the choice of participants,
811 sample collection, and related methods for the data used in the current study were similar to
812 FINRISK 2002 to facilitate inter-cohort comparisons, and are reported elsewhere.³⁰ The study

813 protocol of FINRISK 2007 was approved by the Coordinating Ethical Committee of the Hospital
814 District of Helsinki and Uusimaa (Ref. 229/EO/2006). All participants have signed an informed
815 consent.

816

817 Briefly, compared to FINRISK 2002, FINRISK 2007 features a smaller number of participants
818 who donated fecal samples ($N = 258$ after excluding pregnant individuals or antibiotic use in the
819 last 6 months), they were younger on average, and a smaller proportion of them were in the high
820 FLI group. To produce data for the validation, methods and quality control related to DNA
821 extraction, sequencing, taxonomic assignments, and calculation of FLI values were identical to
822 FINRISK 2002 data, as described above. For the phylogenetic transform (performed otherwise
823 identically), only taxa passing the filtering in FINRISK 2002 bacterial data set were retained in
824 FINRISK 2007 and a pseudo-count of 1 was used for taxa unobserved in the new data, to exactly
825 match the node balance names. The FINRISK 2007 data was then subset to the model features of
826 the main classification models (sex, age, and the 11 bacterial balances), and input in each of the 6
827 LOGOCV classification models. The results of these predictions were then compared against the
828 true FLI groups ($FLI \geq 60$ and $FLI < 60$) of the participants (**Table S5**). Receiver operating
829 characteristic and precision-recall curves for the external validation were calculated similarly to
830 the main models for the AUC and AUPRC metrics and plotted after averaging the predictions of
831 the 6 models to obtain single curves (**Figure S6**).

832

833 *Identification of predictive features specific to liver function*

834 Because the FLI also incorporates BMI and waist circumference, and they strongly contribute to
835 the index,²⁶ we deemed it necessary to further investigate which of the identified balances were

836 specific to liver function. The participants were first grouped by age (< 40, 40 – 60, and 60 <),
837 sex (female or male) and BMI (< 25, 25 – 30, and 30 <) into 18 categories ($N = 105 \sim 711$ per
838 category). We performed feature selection similarly to the FLI models by fitting random forest
839 regressors for GGT and triglycerides with mlr 2.16.0.⁹² This was done separately in each of the
840 18 categories, and in each category, we again used LOGOCV with the regions to obtain 6 runs
841 per category. Finally, the features predictive of GGT or triglycerides in each category were
842 selected as the intersect of top 50 features in the 6 LOGOCV iterations by permutation
843 importance. The intersect of features predictive of GGT or triglycerides in any of the categories
844 and the features predictive of categorical FLI were identified as specific to liver function.
845

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

847 Our taxonomic matching of the reads is based on the genomes of GTDB (release 89),²⁸ which are
848 all complete or nearly complete and available in online databases. This enables us to estimate the
849 likely genetic content, and thus, the metabolic potential of the microbes associated with FLI. We
850 use this approach because the sequencing depth of our samples does not allow assembling
851 contigs and (metagenome-assembled) genomes, required for pathway predictions. Because of the
852 compositional phylogenetic transform, among other features of our data, previously developed
853 approaches such as PICRUSt,⁹⁷ could not be used here.

854

855 The genomes of all 336 bacteria under at least one of the predictive balances were downloaded
856 from NCBI. 119 of these genomes were originally not annotated, which is a requirement for
857 pathway prediction. Therefore, Prokka v1.14.6,⁹⁸ was used to annotate the 119 unannotated
858 genomes as a preliminary step. Pathway predictions were then performed for all 336 genomes

859 with mpwt v0.5.3 multiprocessing tool,⁹⁹ for the PathoLogic pipeline of Pathway Tools 23.0.¹⁰⁰
860 Pathways for ethanol and short chain fatty acid (acetate, butyrate, propionate) production, bile
861 acid metabolism, and choline degradation to trimethylamine were identified from MetaCyc
862 pathway classifications (see ref. ¹⁰¹, and **Table S4**). The prevalence of these processes was then
863 assessed in the analyzed genomes and summarized per process to consider the possible links of
864 the taxa with fatty liver pathophysiology. Finally, the presence of individual pathways for acetate
865 and ethanol production was also outlined for each genome.

866

867 ***Data availability statement***

868 The analysis code written for this study is included with the Supplementary Information. The
869 datasets generated during and analyzed during the current study are not public, but are available
870 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)
871 [biobank/for-researchers](https://thl.fi/en/web/thl-biobank/for-researchers)

872

873 ***Disclosure of interest***

874 V.S. has consulted for Novo Nordisk and Sanofi and received honoraria from these companies.
875 He also has ongoing research collaboration with Bayer AG, all unrelated to this study. R.L.
876 serves as a consultant or advisory board member for Anylam/Regeneron, Arrowhead
877 Pharmaceuticals, AstraZeneca, Bird Rock Bio, Boehringer Ingelheim, Bristol-Myer Squibb,
878 Celgene, Cirius, CohBar, Conatus, Eli Lilly, Galmed, Gemphire, Gilead, Glympse bio, GNI, GRI
879 Bio, Inipharm, Intercept, Ionis, Janssen Inc., Merck, Metacrine, Inc., NGM Biopharmaceuticals,
880 Novartis, Novo Nordisk, Pfizer, Prometheus, Promethera, Sanofi, Siemens, and Viking
881 Therapeutics. In addition, his institution has received grant support from Allergan, Boehringer-

882 Ingelheim, Bristol-Myers Squibb, Cirius, Eli Lilly and Company, Galectin Therapeutics, Galmed
883 Pharmaceuticals, GE, Genfit, Gilead, Intercept, Grail, Janssen, Madrigal Pharmaceuticals,
884 Merck, NGM Biopharmaceuticals, NuSirt, Pfizer, pH Pharma, Prometheus, and Siemens. He is
885 also co-founder of Liponexus, Inc.

886

887 ***Funding details***

888 This research was supported in part by grants from the Finnish Foundation for Cardiovascular
889 Research, the Emil Aaltonen Foundation, the Paavo Nurmi Foundation, the Urmas Pekkala
890 Foundation, the Finnish Medical Foundation, the Sigrid Juselius Foundation, the Academy of
891 Finland (#321356 to A.H.; #295741, #307127 to L.L.; #321351 to T.N.) and the National
892 Institutes of Health (R01ES027595 to M.J.). R.L. receives funding support from NIEHS
893 (5P42ES010337), NCATS (5UL1TR001442), NIDDK (U01DK061734, R01DK106419,
894 P30DK120515, R01DK121378, R01DK124318), and DOD PRCRP (W81XWH-18-2-0026).
895 Additional support was provided by Illumina, Inc. and Janssen Pharmaceutica through their
896 sponsorship of the Center for Microbiome Innovation at UCSD.

897

898 ***Authors' contributions***

899 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
900 and R.K. acquired the data. M.R., L.L. and T.N. analyzed the data. M.R. wrote the manuscript in
901 consultation with all authors. M.I., P.J., V.S., R.K., L.L. and T.N. supervised the work. All
902 authors gave final approval of the version to be published.

903

904 *Acknowledgements*

905 We thank all participants of the FINRISK 2002 and FINRISK 2007 surveys for their
906 contributions to this work, and Tara Schwartz for assistance with laboratory work. We also thank
907 the editor and both anonymous reviewers for their constructive criticism.

908

909 **References**

- 910 1. Younossi ZM, Koenig AB, Abdelatif D, Fazel Y, Henry L, Wymer M. Global
911 epidemiology of nonalcoholic fatty liver disease—Meta-analytic assessment of prevalence,
912 incidence, and outcomes. *Hepatology* 2016; 64:73–84.
- 913 2. Marchesini G, Bugianesi E, Forlani G, Cerrelli F, Lenzi M, Manini R, Natale S, Vanni E,
914 Villanova N, Melchionda N, et al. Nonalcoholic fatty liver, steatohepatitis, and the
915 metabolic syndrome. *Hepatology* 2003; 37:917–23.
- 916 3. Chalasani N, Younossi Z, Lavine JE, Diehl AM, Brunt EM, Cusi K, Charlton M, Sanyal
917 AJ. The diagnosis and management of non-alcoholic fatty liver disease: Practice Guideline
918 by the American Association for the Study of Liver Diseases, American College of
919 Gastroenterology, and the American Gastroenterological Association. *Hepatology* 2012;
920 55:2005–23.
- 921 4. Yki-Järvinen H. Non-alcoholic fatty liver disease as a cause and a consequence of
922 metabolic syndrome. *Lancet Diabetes Endocrinol* 2014; 2:901–10.
- 923 5. Toshikuni N, Tsutsumi M, Arisawa T. Clinical differences between alcoholic liver disease
924 and nonalcoholic fatty liver disease. *World J Gastroenterol* 2014; 20:8393–406.
- 925 6. Rinella M, Charlton M. The globalization of nonalcoholic fatty liver disease: Prevalence
926 and impact on world health. *Hepatology* 2016; 64:19–22.
- 927 7. Gilbert JA, Blaser MJ, Caporaso JG, Jansson JK, Lynch SV, Knight R. Current
928 understanding of the human microbiome. *Nat Med* 2018; 24:392–400.
- 929 8. Caussy C, Tripathi A, Humphrey G, Bassirian S, Singh S, Faulkner C, Bettencourt R, Rizo
930 E, Richards L, Xu ZZ, et al. A gut microbiome signature for cirrhosis due to nonalcoholic
931 fatty liver disease. *Nature Communications* 2019; 10:1406.
- 932 9. Compare D, Coccoli P, Rocco A, Nardone OM, De Maria S, Carteni M, Nardone G. Gut–
933 liver axis: The impact of gut microbiota on non alcoholic fatty liver disease. *Nutrition,
934 Metabolism and Cardiovascular Diseases* 2012; 22:471–6.

- 935 10. Miele L, Valenza V, Torre GL, Montalto M, Cammarota G, Ricci R, Mascianà R, Forgione
936 A, Gabrieli ML, Perotti G, et al. Increased intestinal permeability and tight junction
937 alterations in nonalcoholic fatty liver disease. *Hepatology* 2009; 49:1877–87.
- 938 11. Spencer MD, Hamp TJ, Reid RW, Fischer LM, Zeisel SH, Fodor AA. Association Between
939 Composition of the Human Gastrointestinal Microbiome and Development of Fatty Liver
940 With Choline Deficiency. *Gastroenterology* 2011; 140:976–86.
- 941 12. Zhu L, Baker SS, Gill C, Liu W, Alkhoury R, Baker RD, Gill SR. Characterization of gut
942 microbiomes in nonalcoholic steatohepatitis (NASH) patients: A connection between
943 endogenous alcohol and NASH. *Hepatology* 2013; 57:601–9.
- 944 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
945 Disease Caused by High-Alcohol-Producing *Klebsiella pneumoniae*. *Cell Metabolism*
946 2019; 30:675-688.e7.
- 947 14. Loomba R, Seguritan V, Li W, Long T, Klitgord N, Bhatt A, Dulai PS, Caussy C,
948 Bettencourt R, Highlander SK, et al. Gut Microbiome-Based Metagenomic Signature for
949 Non-invasive Detection of Advanced Fibrosis in Human Nonalcoholic Fatty Liver Disease.
950 *Cell Metabolism* 2017; 25:1054-1062.e5.
- 951 15. Jiao N, Wu D, Yang Z, Fang S, Li X, Yuan M, Zhu R, Zhu L. Gut bacteria contributes to
952 NAFLD pathogenesis by promoting secondary bile acids biosynthesis. *The FASEB Journal*
953 2019; 33:126.4-126.4.
- 954 16. Aron-Wisnewsky J, Vigliotti C, Witjes J, Le P, Holleboom AG, Verheij J, Nieuwdorp M,
955 Clément K. Gut microbiota and human NAFLD: disentangling microbial signatures from
956 metabolic disorders. *Nature Reviews Gastroenterology & Hepatology* 2020; 17:279–97.
- 957 17. Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, Guo J, Le Chatelier E, Yao J, Wu L, et al.
958 Alterations of the human gut microbiome in liver cirrhosis. *Nature* 2014; 513:59–64.
- 959 18. Liu Y, Meric G, Havulinna AS, Teo SM, Ruuskanen M, Sanders J, Zhu Q, Tripathi A,
960 Verspoor K, Cheng S, et al. Early prediction of liver disease using conventional risk factors
961 and gut microbiome-augmented gradient boosting [Internet]. *Genetic and Genomic*
962 *Medicine*; 2020 [cited 2020 Jul 28]. Available from:
963 <http://medrxiv.org/lookup/doi/10.1101/2020.06.24.20138933>
- 964 19. Safari Z, Gérard P. The links between the gut microbiome and non-alcoholic fatty liver
965 disease (NAFLD). *Cell Mol Life Sci* 2019; 76:1541–58.
- 966 20. Carpino G, Del Ben M, Pastori D, Carnevale R, Baratta F, Overi D, Francis H, Cardinale V,
967 Onori P, Safarikia S, et al. Increased liver localization of lipopolysaccharides in human and
968 experimental non-alcoholic fatty liver disease. *Hepatology* 2019; :hep.31056.
- 969 21. Li Q, Dhyani M, Grajo JR, Sirlin C, Samir AE. Current status of imaging in nonalcoholic
970 fatty liver disease. *World J Hepatol* 2018; 10:530–42.

- 971 22. Koehler EM, Schouten JNL, Hansen BE, Hofman A, Stricker BH, Janssen HLA. External
972 Validation of the Fatty Liver Index for Identifying Nonalcoholic Fatty Liver Disease in a
973 Population-based Study. *Clinical Gastroenterology and Hepatology* 2013; 11:1201–4.
- 974 23. Vanni E, Bugianesi E. Editorial: utility and pitfalls of Fatty Liver Index in epidemiologic
975 studies for the diagnosis of NAFLD. *Aliment Pharmacol Ther* 2015; 41:406–7.
- 976 24. Salosensaari A, Laitinen V, Havulinna AS, Meric G, Cheng S, Perola M, Valsta L, Alfthan
977 G, Inouye M, Watrous JD, et al. Taxonomic Signatures of Long-Term Mortality Risk in
978 Human Gut Microbiota [Internet]. *Epidemiology*; 2020 [cited 2020 Jan 4]. Available from:
979 <http://medrxiv.org/lookup/doi/10.1101/2019.12.30.19015842>
- 980 25. Alexander M, Loomis AK, Fairburn-Beech J, van der Lei J, Duarte-Salles T, Prieto-
981 Alhambra D, Ansell D, Pasqua A, Lapi F, Rijnbeek P, et al. Real-world data reveal a
982 diagnostic gap in non-alcoholic fatty liver disease. *BMC Medicine* 2018; 16:130.
- 983 26. Bedogni G, Bellentani S, Miglioli L, Masutti F, Passalacqua M, Castiglione A, Tiribelli C.
984 The Fatty Liver Index: a simple and accurate predictor of hepatic steatosis in the general
985 population. *BMC Gastroenterol* 2006; 6:33.
- 986 27. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R,
987 Knights D. Evaluating the Information Content of Shallow Shotgun Metagenomics.
988 *mSystems* [Internet] 2018 [cited 2020 Apr 9]; 3. Available from:
989 <https://msystems.asm.org/content/3/6/e00069-18>
- 990 28. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil P-A,
991 Hugenholtz P. A standardized bacterial taxonomy based on genome phylogeny
992 substantially revises the tree of life. *Nat Biotechnol* 2018; 36:996–1004.
- 993 29. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform
994 enhances analysis of compositional microbiota data. *eLife* 2017; 6:e21887.
- 995 30. Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, Kuulasmaa K,
996 Laatikainen T, Männistö S, Peltonen M, et al. Cohort Profile: The National FINRISK
997 Study. *International Journal of Epidemiology* 2018; 47:696–696i.
- 998 31. Banderas DZ, Escobedo J, Gonzalez E, Liceaga MG, Ramírez JC, Castro MG. γ -Glutamyl
999 transferase: a marker of nonalcoholic fatty liver disease in patients with the metabolic
1000 syndrome. *European Journal of Gastroenterology & Hepatology* 2012; 24:805–810.
- 1001 32. Haas JT, Francque S, Staels B. Pathophysiology and Mechanisms of Nonalcoholic Fatty
1002 Liver Disease. *Annual Review of Physiology* 2016; 78:181–205.
- 1003 33. Castaner O, Goday A, Park Y-M, Lee S-H, Magkos F, Shiow S-ATE, Schröder H. The Gut
1004 Microbiome Profile in Obesity: A Systematic Review. *Int J Endocrinol* [Internet] 2018
1005 [cited 2020 Apr 3]; 2018. Available from:
1006 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5933040/>

- 1007 34. Eslam M, Sanyal AJ, George J, Sanyal A, Neuschwander-Tetri B, Tiribelli C, Kleiner DE,
1008 Brunt E, Bugianesi E, Yki-Järvinen H, et al. MAFLD: A Consensus-Driven Proposed
1009 Nomenclature for Metabolic Associated Fatty Liver Disease. *Gastroenterology* 2020;
1010 158:1999-2014.e1.
- 1011 35. He Y, Wu W, Zheng H-M, Li P, McDonald D, Sheng H-F, Chen M-X, Chen Z-H, Ji G-Y,
1012 Zheng Z-D-X, et al. Regional variation limits applications of healthy gut microbiome
1013 reference ranges and disease models. *Nature Medicine* 2018; 24:1532–5.
- 1014 36. Carvalhana S, Leitão J, Alves AC, Bourbon M, Cortez-Pinto H. How good is controlled
1015 attenuation parameter and fatty liver index for assessing liver steatosis in general
1016 population: correlation with ultrasound. *Liver International* 2014; 34:e111–7.
- 1017 37. Dehnavi Z, Razmpour F, Naseri MB, Nematy M, Alamdaran SA, Vatanparast HA, Nezhad
1018 MA, Abbasi B, Ganji A. Fatty liver index (FLI) in predicting non-alcoholic fatty liver
1019 disease (NAFLD). *Hepatitis Monthly* 2018; 18.
- 1020 38. Yang B-L, Wu W-C, Fang K-C, Wang Y-C, Huo T-I, Huang Y-H, Yang H-I, Su C-W, Lin
1021 H-C, Lee F-Y, et al. External Validation of Fatty Liver Index for Identifying
1022 Ultrasonographic Fatty Liver in a Large-Scale Cross-Sectional Study in Taiwan. *PLOS*
1023 *ONE* 2015; 10:e0120443.
- 1024 39. Huang X, Xu M, Chen Y, Peng K, Huang Y, Wang P, Ding L, Lin L, Xu Y, Chen Y, et al.
1025 Validation of the Fatty Liver Index for Nonalcoholic Fatty Liver Disease in Middle-Aged
1026 and Elderly Chinese. *Medicine (Baltimore)* 2015; 94:e1682.
- 1027 40. Astbury S, Atallah E, Vijay A, Aithal GP, Grove JI, Valdes AM. Lower gut microbiome
1028 diversity and higher abundance of proinflammatory genus *Collinsella* are associated with
1029 biopsy-proven nonalcoholic steatohepatitis. *Gut Microbes* 2020; 11:569–80.
- 1030 41. Kim H-N, Joo E-J, Cheong HS, Kim Y, Kim H-L, Shin H, Chang Y, Ryu S. Gut
1031 Microbiota and Risk of Persistent Nonalcoholic Fatty Liver Diseases. *Journal of Clinical*
1032 *Medicine* 2019; 8:1089.
- 1033 42. de la Cuesta-Zuluaga J, Corrales-Agudelo V, Velásquez-Mejía EP, Carmona JA, Abad JM,
1034 Escobar JS. Gut microbiota is associated with obesity and cardiometabolic disease in a
1035 population in the midst of Westernization. *Scientific Reports* 2018; 8:11356.
- 1036 43. O’Callaghan A, van Sinderen D. Bifidobacteria and Their Role as Members of the Human
1037 Gut Microbiota. *Front Microbiol* [Internet] 2016 [cited 2020 Nov 9]; 7. Available from:
1038 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4908950/>
- 1039 44. Waters JL, Ley RE. The human gut bacteria Christensenellaceae are widespread, heritable,
1040 and associated with health. *BMC Biology* 2019; 17:83.
- 1041 45. Ferreira-Halder CV, Faria AV de S, Andrade SS. Action and function of *Faecalibacterium*
1042 *prausnitzii* in health and disease. *Best Pract Res Clin Gastroenterol* 2017; 31:643–8.

- 1043 46. Guohong-Liu, Qingxi-Zhao, Hongyun-Wei. Characteristics of intestinal bacteria with fatty
1044 liver diseases and cirrhosis. *Annals of Hepatology* 2019; 18:796–803.
- 1045 47. Raman M, Ahmed I, Gillevet PM, Probert CS, Ratcliffe NM, Smith S, Greenwood R,
1046 Sikaroodi M, Lam V, Crotty P, et al. Fecal Microbiome and Volatile Organic Compound
1047 Metabolome in Obese Humans With Nonalcoholic Fatty Liver Disease. *Clinical*
1048 *Gastroenterology and Hepatology* 2013; 11:868-875.e3.
- 1049 48. Dong TS, Katzka W, Lagishetty V, Luu K, Hauer M, Pisegna J, Jacobs JP. A Microbial
1050 Signature Identifies Advanced Fibrosis in Patients with Chronic Liver Disease Mainly Due
1051 to NAFLD. *Scientific Reports* 2020; 10:2771.
- 1052 49. Bajaj JS, Idilman R, Mabudian L, Hood M, Fagan A, Turan D, White MB, Karakaya F,
1053 Wang J, Atalay R, et al. Diet affects gut microbiota and modulates hospitalization risk
1054 differentially in an international cirrhosis cohort. *Hepatology* 2018; 68:234–47.
- 1055 50. Bowerman KL, Rehman SF, Vaughan A, Lachner N, Budden KF, Kim RY, Wood DLA,
1056 Gellatly SL, Shukla SD, Wood LG, et al. Disease-associated gut microbiome and
1057 metabolome changes in patients with chronic obstructive pulmonary disease. *Nature*
1058 *Communications* 2020; 11:5886.
- 1059 51. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, Beaumont M,
1060 Van Treuren W, Knight R, Bell JT, et al. Human Genetics Shape the Gut Microbiome. *Cell*
1061 2014; 159:789–99.
- 1062 52. Suzuki TA, Worobey M. Geographical variation of human gut microbial composition.
1063 *Biology Letters* 2014; 10:20131037.
- 1064 53. Meslier V, Laiola M, Roager HM, Filippis FD, Roume H, Quinquis B, Giacco R, Mennella
1065 I, Ferracane R, Pons N, et al. Mediterranean diet intervention in overweight and obese
1066 subjects lowers plasma cholesterol and causes changes in the gut microbiome and
1067 metabolome independently of energy intake. *Gut* [Internet] 2020 [cited 2020 Jun 2];
1068 Available from: <https://gut.bmj.com/content/early/2020/02/18/gutjnl-2019-320438>
- 1069 54. Cui C, Li Y, Gao H, Zhang H, Han J, Zhang D, Li Y, Zhou J, Lu C, Su X. Modulation of
1070 the gut microbiota by the mixture of fish oil and krill oil in high-fat diet-induced obesity
1071 mice. *PLOS ONE* 2017; 12:e0186216.
- 1072 55. Manzoor SE, McNulty CAM, Nakiboneka-Ssenabulya D, Lecky DM, Hardy KJ, Hawkey
1073 PM. Investigation of community carriage rates of *Clostridium difficile* and *Hungatella*
1074 *hathewayi* in healthy volunteers from four regions of England. *Journal of Hospital Infection*
1075 2017; 97:153–5.
- 1076 56. Genoni A, Christophersen CT, Lo J, Coghlan M, Boyce MC, Bird AR, Lyons-Wall P,
1077 Devine A. Long-term Paleolithic diet is associated with lower resistant starch intake,
1078 different gut microbiota composition and increased serum TMAO concentrations. *Eur J*
1079 *Nutr* 2020; 59:1845–58.

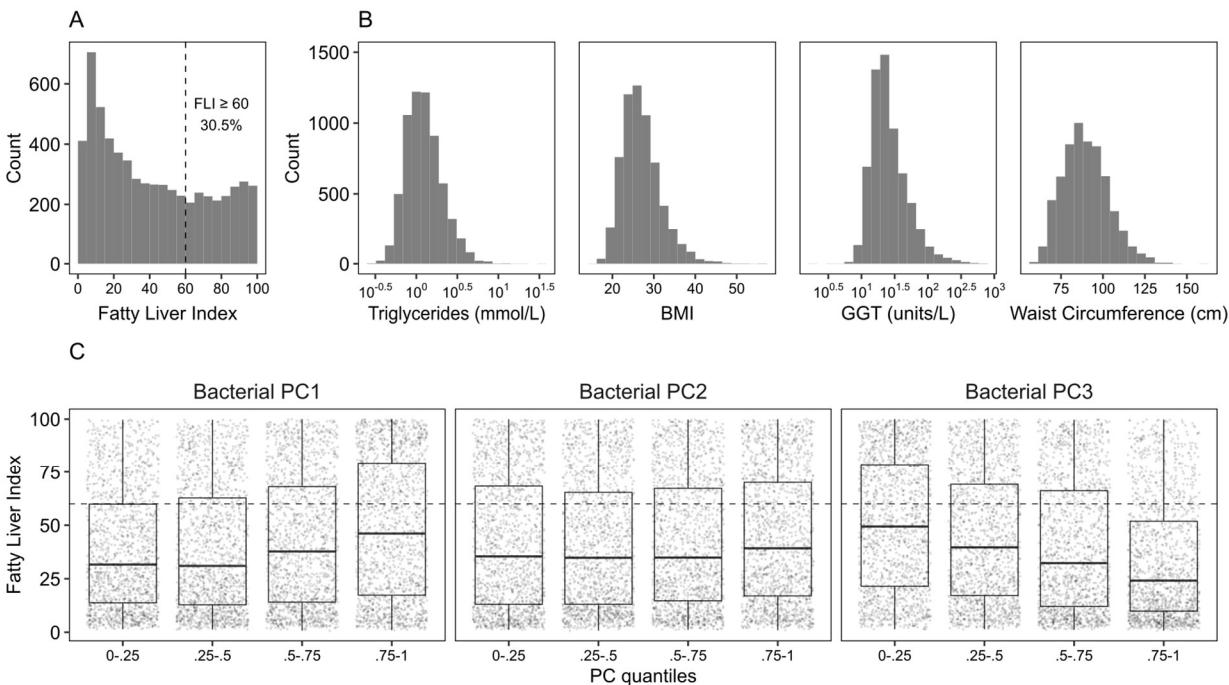
- 1080 57. Burton KJ, Krüger R, Scherz V, Mürger LH, Picone G, Vionnet N, Bertelli C, Greub G,
1081 Capozzi F, Vergères G. Trimethylamine-N-Oxide Postprandial Response in Plasma and
1082 Urine Is Lower After Fermented Compared to Non-Fermented Dairy Consumption in
1083 Healthy Adults. *Nutrients* 2020; 12:234.
- 1084 58. Jiang W, Wu N, Wang X, Chi Y, Zhang Y, Qiu X, Hu Y, Li J, Liu Y. Dysbiosis gut
1085 microbiota associated with inflammation and impaired mucosal immune function in
1086 intestine of humans with non-alcoholic fatty liver disease. *Sci Rep* 2015; 5:8096.
- 1087 59. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Knight R, Knights D. SHOGUN: a
1088 modular, accurate, and scalable framework for microbiome quantification. *Bioinformatics*
1089 2020; :btaa277.
- 1090 60. de Faria Ghetti F, Oliveira DG, de Oliveira JM, de Castro Ferreira LEVV, Cesar DE,
1091 Moreira APB. Influence of gut microbiota on the development and progression of
1092 nonalcoholic steatohepatitis. *Eur J Nutr* 2018; 57:861–76.
- 1093 61. Mohan R, Namsolleck P, Lawson PA, Osterhoff M, Collins MD, Alpert C-A, Blaut M.
1094 *Clostridium asparagiforme* sp. nov., isolated from a human faecal sample. *Systematic and*
1095 *Applied Microbiology* 2006; 29:292–9.
- 1096 62. Murray WD, Khan AW, van den BERG L. *Clostridium saccharolyticum* sp. nov., a
1097 Saccharolytic Species from Sewage Sludge. *International Journal of Systematic*
1098 *Bacteriology* 1982; 32:132–5.
- 1099 63. Diether NE, Willing BP. Microbial Fermentation of Dietary Protein: An Important Factor in
1100 Diet–Microbe–Host Interaction. *Microorganisms* [Internet] 2019 [cited 2020 Apr 24]; 7.
1101 Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6352118/>
- 1102 64. Zhao S, Jang C, Liu J, Uehara K, Gilbert M, Izzo L, Zeng X, Trefely S, Fernandez S, Carrer
1103 A, et al. Dietary fructose feeds hepatic lipogenesis via microbiota-derived acetate. *Nature*
1104 2020; 579:586–91.
- 1105 65. Dehoux P, Marvaud JC, Abouelleil A, Earl AM, Lambert T, Dauga C. Comparative
1106 genomics of *Clostridium bolteae* and *Clostridium clostridioforme* reveals species-specific
1107 genomic properties and numerous putative antibiotic resistance determinants. *BMC*
1108 *Genomics* 2016; 17:819.
- 1109 66. Deschasaux M, Bouter KE, Prodan A, Levin E, Groen AK, Herrema H, Tremaroli V,
1110 Bakker GJ, Attaye I, Pinto-Sietsma S-J, et al. Depicting the composition of gut microbiota
1111 in a population with varied ethnic origins but shared geography. *Nature Medicine* 2018;
1112 24:1526–31.
- 1113 67. Canfora EE, Meex RCR, Venema K, Blaak EE. Gut microbial metabolites in obesity,
1114 NAFLD and T2DM. *Nature Reviews Endocrinology* 2019; 15:261–73.
- 1115 68. Cheng H-Y, Wang H-Y, Chang W-H, Lin S-C, Chu C-H, Wang T-E, Liu C-C, Shih S-C.
1116 Nonalcoholic Fatty Liver Disease: Prevalence, Influence on Age and Sex, and Relationship

- 1117 with Metabolic Syndrome and Insulin Resistance. *International Journal of Gerontology*
1118 2013; 7:194–8.
- 1119 69. Lonardo A, Nascimbeni F, Ballestri S, Fairweather D, Win S, Than TA, Abdelmalek MF,
1120 Suzuki A. Sex Differences in Nonalcoholic Fatty Liver Disease: State of the Art and
1121 Identification of Research Gaps. *Hepatology* 2019; 70:1457–69.
- 1122 70. Näyhä S. Geographical variations in cardiovascular mortality in Finland, 1961-1985. *Scand*
1123 *J Soc Med Suppl* 1989; 40:1–48.
- 1124 71. Kerminen S, Havulinna AS, Helleenthal G, Martin AR, Sarin A-P, Perola M, Palotie A,
1125 Salomaa V, Daly MJ, Ripatti S, et al. Fine-Scale Genetic Structure in Finland. *G3: Genes,*
1126 *Genomes, Genetics* 2017; 7:3459–68.
- 1127 72. Reccia I, Kumar J, Akladios C, Viridis F, Pai M, Habib N, Spalding D. Non-alcoholic fatty
1128 liver disease: A sign of systemic disease. *Metabolism* 2017; 72:94–108.
- 1129 73. Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T, Männistö S,
1130 Salomaa V, Sundvall J, Puska P. Forty-year trends in cardiovascular risk factors in Finland.
1131 *Eur J Public Health* 2015; 25:539–46.
- 1132 74. Marotz L, Schwartz T, Thompson L, Humphrey G, Gogul G, Gaffney J, Amir A, Knight R.
1133 Earth Microbiome Project (EMP) high throughput (HTP) DNA extraction protocol v1
1134 (protocols.io.pdmdi46) [Internet]. 2018 [cited 2020 Nov 10]; Available from:
1135 [https://www.protocols.io/view/earth-microbiome-project-emp-high-throughput-htp-d-](https://www.protocols.io/view/earth-microbiome-project-emp-high-throughput-htp-d-pdmdi46)
1136 [pdmdi46](https://www.protocols.io/view/earth-microbiome-project-emp-high-throughput-htp-d-pdmdi46)
- 1137 75. Sanders JG, Nurk S, Salido RA, Minich J, Xu ZZ, Zhu Q, Martino C, Fedarko M, Arthur
1138 TD, Chen F, et al. Optimizing sequencing protocols for leaderboard metagenomics by
1139 combining long and short reads. *Genome Biology* 2019; 20:226.
- 1140 76. Glenn TC, Nilsen RA, Kieran TJ, Sanders JG, Bayona-Vásquez NJ, Finger JW, Pierson
1141 TW, Bentley KE, Hoffberg SL, Louha S, et al. Adapterama I: universal stubs and primers
1142 for 384 unique dual-indexed or 147,456 combinatorially-indexed Illumina libraries (iTru &
1143 iNext). *PeerJ* 2019; 7:e7755.
- 1144 77. Didion JP, Martin M, Collins FS. Atropos: specific, sensitive, and speedy trimming of
1145 sequencing reads. *PeerJ* [Internet] 2017 [cited 2020 Nov 10]; 5. Available from:
1146 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5581536/>
- 1147 78. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 2012;
1148 9:357–9.
- 1149 79. Méric G, Wick RR, Watts SC, Holt KE, Inouye M. Correcting index databases improves
1150 metagenomic studies. *bioRxiv* 2019; :712166.
- 1151 80. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. A complete
1152 domain-to-species taxonomy for Bacteria and Archaea. *Nature Biotechnology* 2020; :1–8.

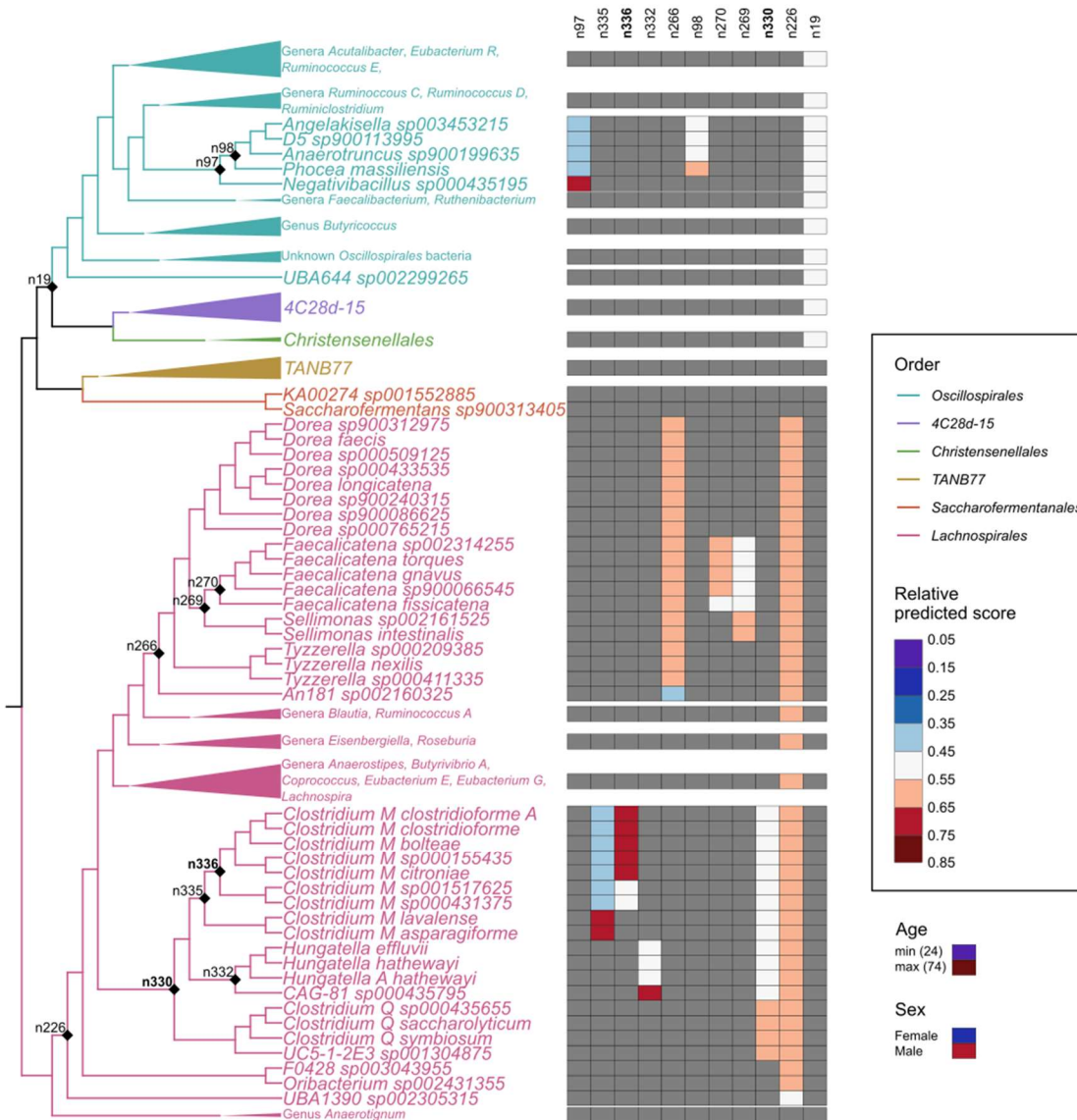
- 1153 81. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification
1154 of metagenomic sequences. *Genome Res* [Internet] 2016 [cited 2018 May 12]; Available
1155 from: <http://genome.cshlp.org/content/early/2016/11/16/gr.210641.116>
- 1156 82. R Core Team. R: A language and environment for statistical computing [Internet]. Vienna,
1157 Austria: R Foundation for Statistical Computing; 2018 [cited 2019 Mar 4]. Available from:
1158 <https://www.R-project.org/>
- 1159 83. McMurdie PJ, Holmes S. phyloseq: An R package for reproducible interactive analysis and
1160 graphics of microbiome census data. *PLoS ONE* 2013; 8:e61217.
- 1161 84. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome Datasets Are
1162 Compositional: And This Is Not Optional. *Front Microbiol* [Internet] 2017 [cited 2020 Jul
1163 20]; 8. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5695134/>
- 1164 85. Louca S, Polz MF, Mazel F, Albright MBN, Huber JA, O'Connor MI, Ackermann M, Hahn
1165 AS, Srivastava DS, Crowe SA, et al. Function and functional redundancy in microbial
1166 systems. *Nat Ecol Evol* 2018; 2:936–43.
- 1167 86. Ruuskanen MO, St. Pierre KA, St. Louis VL, Aris-Brosou S, Poulain AJ. Physicochemical
1168 Drivers of Microbial Community Structure in Sediments of Lake Hazen, Nunavut, Canada.
1169 *Front Microbiol* [Internet] 2018 [cited 2018 Jun 6]; 9. Available from:
1170 <https://www.frontiersin.org/articles/10.3389/fmicb.2018.01138/full>
- 1171 87. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin PR,
1172 O'Hara RB, Simpson GL, Solymos P, et al. vegan: Community Ecology Package [Internet].
1173 2018 [cited 2018 Jun 4]. Available from: <https://CRAN.R-project.org/package=vegan>
- 1174 88. Lahti L, Shetty S. microbiome R package [Internet]. 2019 [cited 2020 Dec 14]. Available
1175 from: <http://microbiome.github.io>
- 1176 89. Fernandes AD, Reid JN, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB. Unifying
1177 the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA
1178 gene sequencing and selective growth experiments by compositional data analysis.
1179 *Microbiome* 2014; 2:15.
- 1180 90. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd
1181 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining -
1182 KDD '16* 2016; :785–94.
- 1183 91. Bommert A, Sun X, Bischl B, Rahnenführer J, Lang M. Benchmark for filter methods for
1184 feature selection in high-dimensional classification data. *Computational Statistics & Data
1185 Analysis* 2020; 143:106839.
- 1186 92. Bischl B, Lang M, Kotthoff L, Schiffner J, Richter J, Studerus E, Casalicchio G, Jones ZM.
1187 mlr: Machine Learning in R. *Journal of Machine Learning Research* 2016; 17:1–5.

- 1188 93. Bischl B, Richter J, Bossek J, Horn D, Thomas J, Lang M. mlrMBO: A Modular
1189 Framework for Model-Based Optimization of Expensive Black-Box Functions.
1190 arXiv:170303373 [stat] [Internet] 2018 [cited 2020 Feb 18]; Available from:
1191 <http://arxiv.org/abs/1703.03373>
- 1192 94. Saito T, Rehmsmeier M. Precrec: fast and accurate precision–recall and ROC curve
1193 calculations in R. *Bioinformatics* 2017; 33:145–7.
- 1194 95. Greenwell BM. pdp: An R Package for Constructing Partial Dependence Plots. *The R*
1195 *Journal* 2017; 9:421–36.
- 1196 96. Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. ggtree: an r package for visualization and
1197 annotation of phylogenetic trees with their covariates and other associated data. *Methods in*
1198 *Ecology and Evolution* 2017; 8:28–36.
- 1199 97. Douglas GM, Maffei VJ, Zaneveld J, Yurgel SN, Brown JR, Taylor CM, Huttenhower C,
1200 Langille MGI. PICRUSt2: An improved and extensible approach for metagenome
1201 inference. *bioRxiv* 2019; :672295.
- 1202 98. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014; 30:2068–
1203 9.
- 1204 99. Belcour A, Frioux C, Aite M, Bretaudeau A, Siegel A. Metage2Metabo: metabolic
1205 complementarity applied to genomes of large-scale microbiotas for the identification of
1206 keystone species. *bioRxiv* 2019; :803056.
- 1207 100. Karp PD, Midford PE, Billington R, Kothari A, Krummenacker M, Latendresse M, Ong
1208 WK, Subhraveti P, Caspi R, Fulcher C, et al. Pathway Tools version 23.0 update: software
1209 for pathway/genome informatics and systems biology. *Briefings in Bioinformatics*
1210 2019; :bbz104.
- 1211 101. Caspi R, Billington R, Fulcher CA, Keseler IM, Kothari A, Krummenacker M, Latendresse
1212 M, Midford PE, Ong Q, Ong WK, et al. The MetaCyc database of metabolic pathways and
1213 enzymes. *Nucleic Acids Res* 2018; 46:D633–9.
- 1214

1215 **Figures**



1216 **Figure 1.** Relative of FLI (A), its components (B), and FLI in quantiles of the first three PC
1217 components of the fecal bacterial composition of the participants (C). The cutoff at FLI = 60
1218 used to divide the participants is indicated with a dashed line in panels A and C.



1219 **Figure 2.** Relative effects of predictive balances and covariates on the FLI < 60 and FLI ≥ 60
 1220 classification model (AUC = 0.75) predictions. Nodes of the balances are indicated in the
 1221 cladogram and the relative effect sizes of their clades (opposite sides of each balance) are shown
 1222 in the associated heatmap. The relative effect sizes of the covariates (age and sex) are shown
 1223 below the legend with a heatmap on the same scale as was used for the balances. The two liver-
 1224 specific balances associated with triglyceride and GGT levels are indicated with bold font.
 1225 Clades with redundant information have been collapsed but their major genera are indicated. The
 1226 complete tree is included in **Figure S8**.