


















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

Matti O. Ruuskanen ^{a,b}, Fredrik Åberg ^{c,d}, Ville Männistö ^{e,f}, Aki S. Havulinna ^{b,g}, Guillaume Méric ^{h,i}, Yang Liu ^{h,j}, Rohit Loomba ^{k,l}, Yoshiki Vázquez-Baeza ^m, Anupriya Tripathi ^{n,o,p}, Liisa M. Valsta ^b, Michael Inouye ^{h,q}, Pekka Jousilahti ^b, Veikko Salomaa ^b, Mohit Jain ^{l,r}, Rob Knight ^{m,s,t}, Leo Lahti ^u, and Teemu J. Niiranen ^{a,b,v}

^aDepartment of Internal Medicine, University of Turku, Turku, Finland; ^bDepartment of Public Health Solutions, Finnish Institute for Health and Welfare, Helsinki, Finland; ^cTransplantation and Liver Surgery Clinic, Helsinki University Hospital, University of Helsinki, Helsinki, Finland; ^dTransplant Institute, Sahlgrenska University Hospital, Gothenburg, Sweden; ^eDepartment of Medicine, Kuopio University Hospital, University of Eastern Finland, Kuopio, Finland; ^fDepartment of Experimental Vascular Medicine, Amsterdam UMC, University of Amsterdam, Amsterdam, The Netherlands; ^gInstitute for Molecular Medicine Finland, FIMM - HiLIFE, Helsinki, Finland; ^hCambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia; ⁱDepartment of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria, Australia; ^jDepartment of Clinical Pathology, The University of Melbourne, Melbourne, Victoria, Australia; ^kDepartment of Medicine, NAFLD Research Center, La Jolla, CA, USA; ^lDepartment of Medicine, University of California San Diego, La Jolla, CA, USA; ^mCenter for Microbiome Innovation, Jacobs School of Engineering, University of California San Diego, La Jolla, CA, USA; ⁿCollaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA; ^oSkaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA; ^pDivision of Biological Sciences, University of California San Diego, La Jolla, CA, USA; ^qDepartment of Public Health and Primary Care, Cambridge University, Cambridge, UK; ^rDepartment of Pharmacology, University of California San Diego, La Jolla, California, USA; ^sDepartment of Pediatrics, School of Medicine, University of California San Diego, La Jolla, California, USA; ^tDepartment of Computer Science & Engineering, University of California San Diego, La Jolla, California, USA; ^uDepartment of Computing, University of Turku, Turku, Finland; ^vDivision of Medicine, Turku University Hospital, Turku, Finland

ABSTRACT

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

ARTICLE HISTORY

Received 17 August 2020
Revised 14 January 2021
Accepted 28 January 2021

KEYWORDS


Metagenomics; human gut; fatty liver; fatty liver index; population sample

Introduction

Fatty liver disease affects roughly a quarter of the world's population.¹ It is characterized by the accumulation of fat in the liver cells and is intimately linked with the pathophysiology of metabolic syndrome.²⁻⁴ Fatty liver disease can be broadly divided into two variants: nonalcoholic fatty liver

disease (NAFLD), attributed to high caloric intake, and alcohol-associated fatty liver disease, attributed to high alcohol consumption. Even though the rate of progressions and underlying causes of both diseases might be different, they can be broadly subdivided into those who have fat accumulation in the liver with no or minimal inflammation or those who have additional features of cellular injury and

CONTACT Matti O. Ruuskanen  matti.ruuskanen@utu.fi  FI-20014, Turun yliopisto, Finland

 Supplemental data for this article can be accessed on the [publisher's website](#).

© 2021 The Author(s). Published with license by Taylor & Francis Group, LLC.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

active inflammation with or without fibrosis typically seen in the peri-sinusoidal area.⁵ Patients with steatohepatitis may progress to cirrhosis and hepatocellular carcinoma and have an increased risk of liver-related morbidity and mortality, globally amounting to hundreds of thousands of deaths.⁶

The human gut harbors up to 10^{12} microbes per gram of content,⁷ and is intimately connected with the liver. Thus, it is no surprise that gut microbiome composition appears to have a strong connection with liver disease.⁸ Numerous studies over the past 80 y have reported associations between gut microbial composition and liver disease.⁹ For example, gut permeability and overgrowth of bacteria in the small intestine,¹⁰ changes in *Gammaproteobacteria* and *Erysipelotrichi* abundance during choline deficiency,¹¹ an elevated abundance of ethanol-producing bacteria,^{12,13} metagenomic signatures of specific bacterial species,^{14,15} have all been linked to NAFLD in small case-control patient samples. However, the microbial signatures often overlap between NAFLD and metabolic diseases, while those of more serious liver disease, such as steatohepatitis and cirrhosis are more clear.¹⁶ For example, oral taxa appear to invade the gut in liver cirrhosis,¹⁷ and this phenotype can accurately be detected by analyzing the fecal microbiome composition (AUC = 0.87 in a validation cohort).⁸ Furthermore, we recently demonstrated good prediction accuracy for incident liver disease diagnoses (AUC = 0.83 for nonalcoholic liver disease, AUC = 0.96 for alcoholic liver disease, during ~15 y),¹⁸ showing that the signatures of serious future liver disease are easy to detect.

The mechanisms underlying the contribution of gut microbiome content with fatty liver disease are thought to be primarily linked to gut bacterial metabolism. Bacterial metabolites can indeed be translocated from the gut through the intestinal barrier into the portal vein and transported to the liver, where they interact with liver cells, and can lead to inflammation and steatosis.¹⁹ Short-chain fatty acid production, conversion of choline into methylamines, modification of bile acids (BA) into secondary BA, and ethanol production, all of which are mediated by gut bacteria, are also known to be aggravating factors for NAFLD.¹⁹ Recent studies have also suggested that endogenous ethanol production by gut bacteria could lead to an increase in gut membrane permeability.¹³ This can facilitate

the translocation of bacterial metabolites and cell components, such as lipopolysaccharides from the gut to the liver, leading to further inflammation and possible development of NAFLD.²⁰

Liver biopsy assessment is the current gold standard for diagnosis of fatty liver disease and its severity,²¹ but it is also impractical and unethical in a population-based setting. Ultrasound and MRI based assessment can help detect the presence of fatty liver, however, this data is not available in our cohort. Regardless, recent studies have shown that indices based on anthropometric measurements and standard blood tests can be a reliable tool for noninvasive diagnosis of fatty liver, particularly in population-based epidemiologic studies.^{22,23}

Here, we designed and conducted computational analyses to examine the links between fatty liver and gut microbiome composition in a representative population sample of 7,211 extensively phenotyped Finnish individuals.²⁴ Because the fatty liver disease is generally underdiagnosed in the general population,²⁵ we used population-wide measurements of BMI, waist circumference, blood triglycerides, and gamma-glutamyl transferase (GGT) to calculate a previously validated Fatty Liver Index (FLI) for each participant as a proxy for fatty liver.²⁶ In parallel, we used shallow shotgun sequencing to analyze gut microbiome composition,²⁷ which also enabled the use of phylogenetic and pathway prediction methods. In this work, we describe high-resolution associations between fatty liver and individual gut microbial taxa and clades, which are replicable in an external Finnish cohort, and thus generalizable in the Finnish population.

Results

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

In our main analyses, we classified our reads against the Genome Taxonomy Database (GTDB).²⁸ This study mainly follows the GTDB taxonomy, unless otherwise noted. The Centrifuge/GTDB microbiome data used in our main analyses were based on archaeal and bacterial phylogenetic “balances”. This method was used to associate larger groups or clades of related organisms with fatty liver disease, and to avoid grouping of taxa on strict hierarchical

taxonomic ranks featuring varying ranges of evolutionary divergence.²⁸ Here, we used the PhILR transform, where each balance represents a single internal node in a phylogenetic tree, and its value is a log-ratio of the abundances of the two descending clades (for details, see methods and ref.²⁹). Positive values of the balance signify that the clade in the numerator is more abundant, and negative values that the clade in the denominator is more abundant. Thus, each association of a balance with the target variable necessarily includes both microbial clades descending from the node, one of them positively and the other negatively associated with the target variable. The clades in the numerator and denominator can be also freely switched by changing the sign of the balance value to retain the equivalence. Notably, we used this feature to show all balance-FLI associations in the positive direction to facilitate the comparison of their effect sizes (in **Figures S4, S7, and S9**).

Because the combined approach of using the GTDB taxonomy and the recently introduced PhILR phylogenetic transform complicates the comparison of our results in previous studies, we also conducted more traditional statistical analyses

with NCBI-annotated data to anchor our results in previous findings on the associations between fatty liver disease and gut microbiome composition. Overall, the Centrifuge/GTDB classification assigned 5.3 billion reads in the 6,269 samples (after exclusions in FINRISK 2002) to 23,457 bacterial and 1,248 archaeal taxa, and the SHOGUN/NCBI classification assigned 5.5 billion reads to 5,024 bacterial and 261 archaeal taxa. Starting from high level descriptions of the microbial communities in the high and low FLI groups (<60 or ≥ 60 FLI; **Figure 1a**), the phylum-level distributions of bacterial and archaeal taxa appeared to be highly similar between the groups (**Figures S1, S2**). However, the proportion of taxa assigned to *Firmicutes* in Centrifuge/GTDB appeared to be slightly higher than in the SHOGUN/NCBI data. Furthermore, only 58% of the number of reads assigned in Centrifuge/GTDB to 6 main archaeal phyla were assigned to a single main archaeal phylum in SHOGUN/NCBI. Alpha diversity (as Shannon diversity) was significantly lower in the high FLI group, in both the SHOGUN/NCBI data (14.7% lower; AIC = 6,685; all $P < 1 \times 10^{-6}$) and the Centrifuge/GTDB (13.4% lower; AIC = 6,607; all

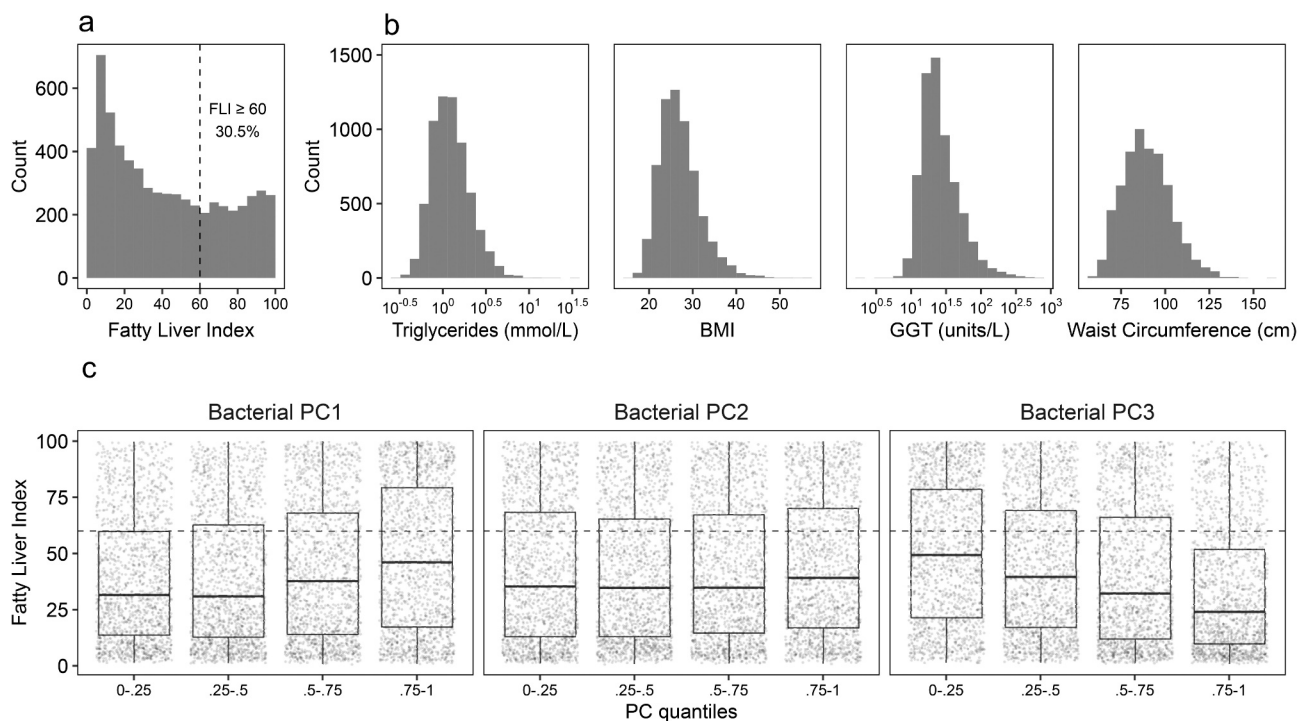


Figure 1. Distribution of FLI (a), its components (b), and FLI in quantiles of the first three PC components of the fecal bacterial composition of the participants (c). The cutoff at FLI = 60 used to divide the participants is indicated with a dashed line in panels a and c.

$P < 1 \times 10^{-4}$) data while adjusting for age, sex, and self-reported alcohol use in both models.

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

We also further assessed the phylogenetic balances contributing to the PC axes in the Centrifuge/GTDB data. Bacterial clades associated with higher FLI values, on the positive side of the balances contributing to PC1, included members of orders *Lachnospirales* and *Oscillospirales*, class *Bacilli*, and the *Ruminococcaceae*, *Bacteroidaceae*, and *Lachnospiraceae* families ([Figure S4](#)). Several clades had a negative association with FLI, on the negative side of the balances contributing to PC1, such as order *Christensenellales* and genus *Faecalibacterium*. In addition, genus *Bifidobacterium* in PC2, and family *Bifidobacteriaceae* in PC3 had negative associations with continuous FLI.

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

We also assessed significant differences in abundances of individual taxa between the high and low FLI groups in FINRISK 2002. In Centrifuge/

GTDB data, we identified 244 taxa (1% of total) with an increased abundance and 437 taxa (1.9%) with a decreased abundance in the high FLI group (all Q values < 0.001 ; [Table S7](#)). In SHOGUN/NCBI data, 80 taxa (1.6%) had an increased abundance, and 44 (0.9%) had a decreased abundance in the high FLI group (all Q values < 0.001). While the number of associated taxa was higher in the Centrifuge/GTDB data than SHOGUN/NCBI data, the proportion of significantly associated taxa was similar between the two methods. In both data sets, family *Lachnospiraceae* comprised over 40% of taxa positively associated with the high FLI group and *Bacteroidaceae* were in the top 3 most common families. The negatively associated taxa were much more diverse, but *Ruminococcaceae* and *Oscillospiraceae* were among the top three most common families in both data sets (at least $> 6\%$ of all negatively associated taxa).

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

Continuous FLI and differences between FLI groups in the FINRISK 2002 cohort (FLI < 60 , $N = 4,359$ and FLI ≥ 60 , $N = 1,910$; see [Figure 1a](#), [Figure 1b](#), [Table S1](#)) were modeled with gradient boosting regression or classification using Leave-One-Group-Out Cross-Validation (LOGOCV) between participants from different regions. Only the bacterial PhILR transformed Centrifuge/GTDB data were used here, to find robust associations between phylogenetically related bacterial clades and fatty liver disease (instead of single taxa).

After feature selection and Bayesian hyperparameter optimization, the correlation between the predictions of the final regression models (age, sex, self-reported alcohol use, and 18 bacterial balances as features; each trained on the data from 5/6 regions) and true values in unseen data from the omitted region averaged $R^2 = 0.30$ (0.26–0.33). After feature selection and optimization, the main classification models (age, sex, and 11 bacterial balances as features; each trained on the data from 5/6 regions) averaged AUC = 0.75 ([Table S2](#)) and AUPRC = 0.56 (baseline at 0.30; [Table S3](#)) on (unseen) test data from the omitted region. Models trained using only the covariates averaged AUC = 0.71 (AUPRC = 0.47) and using only the 11

bacterial balances they averaged AUC = 0.66 (AUPRC = 0.47) on test data. Alternative models were constructed by excluding participants with FLI between 30 and 60 (N = 1,583) and discerning between groups of FLI < 30 (N = 2,776) and FLI ≥ 60 (N = 1,910). These models averaged AUC = 0.80 (AUPRC = 0.75, baseline at 0.41) on their respective test data (**Tables S2, S3**). They averaged AUC = 0.76 (AUPRC = 0.68) when using only the covariates, and AUC = 0.70 (AUPRC = 0.63) when using only the 20 bacterial balances.

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

Our external validation data consisted of 258 participants after exclusion of pregnant participants or those on antibiotics in the past 6 months, in the FINRISK 2007 population cohort (**Table S1, Figure S5**).³⁰ The participants originate from North Karelia and Helsinki/Vantaa regions in Finland, and their samples were processed with the same methodology as was used for FINRISK 2002 (with Centrifuge/GTDB approach and PhILR). In this external validation, the six full models trained with covariates and the 11 bacterial balances in FINRISK 2002 averaged AUC = 0.77 (AUPRC = 0.51, baseline at 0.21; **Table S5**). The covariate-only models averaged AUC = 0.72 (AUPRC = 0.40) and the balance-only models averaged AUC = 0.69 (AUPRC = 0.44). The receiver operating characteristic and precision–recall curves based on the averaged predictions of the models, tested on these external validation data, also display

good predictive ability (AUC = 0.78, AUPRC = 0.51 with baseline at 0.51; **Figure S6**)

To facilitate the interpretability of the results, we continued examining the main classification models using a common set of core features. In these models, the median effect sizes of the features on the model predictions at their minimum and maximum values were highest for age, followed by sex, and the 11 balances in the phylogenetic tree (**Figures S7, S8**). All 11 associated balances were in phylum *Firmicutes*, class *Clostridia*, and largely in the NCBI *Clostridium* subclusters IV and XIVa (**Figure 2**). The specific taxa represented standardized GTDB genera

(NCBI in brackets) *Negativibacillus* (*Clostridium*), *Clostridium M* (*Lachnoclostridium/Clostridium*), CAG-81 (*Clostridium*), *Dorea* (*Merdimonas/Mordavella/Dorea/Clostridium/Eubacterium*), *Faecalicatena* (*Blautia/Ruminococcus/Clostridium*), *Blautia* (*Blautia*), *Sellimonas* (*Sellimonas/Drancourtella*), *Clostridium Q* (*Lachnoclostridium [Clostridium]*), and *Tyzzarella* (*Tyzzarella/Coproccoccus*).

Notably, all but one of the features in the main classification models (n226) were identified in the feature selection for the alternative models (constructed otherwise identically, but FLI < 30 was compared against FLI ≥ 60 in different data partitions), together with 10 additional balances (**Figure S9**). Only one of the balances in the alternative models was outside phylum *Firmicutes* (n1712 in *Bacteroidota*), and in addition, four balances were outside class *Clostridia* (n481 in *Negativicutes*; n826, n1009, and n918 in *Bacilli*).

Also, negative associations with the high FLI group were seen for *An181 sp002160325* in the balance n266, where it is compared against the clade including *Dorea*, *Faecalicatena*, *Sellimonas*, and *Tyzzarella* species (**Figure 2, S8**). A higher abundance of the clade including *Angelakisella*, *D5*, *Anaerotruncus*, and *Phoceia* species (against *Negativibacillus sp00435195* in balance n97) was also negatively associated with high FLI.

In addition to blood test results, FLI is based on two anthropometric markers linked to metabolic syndrome, waist circumference, and BMI. This prompted us to dissect the Fatty Liver Index and identify which of the covariates and associated microbial balances from the phylogenetic tree can be linked to blood GGT and triglyceride

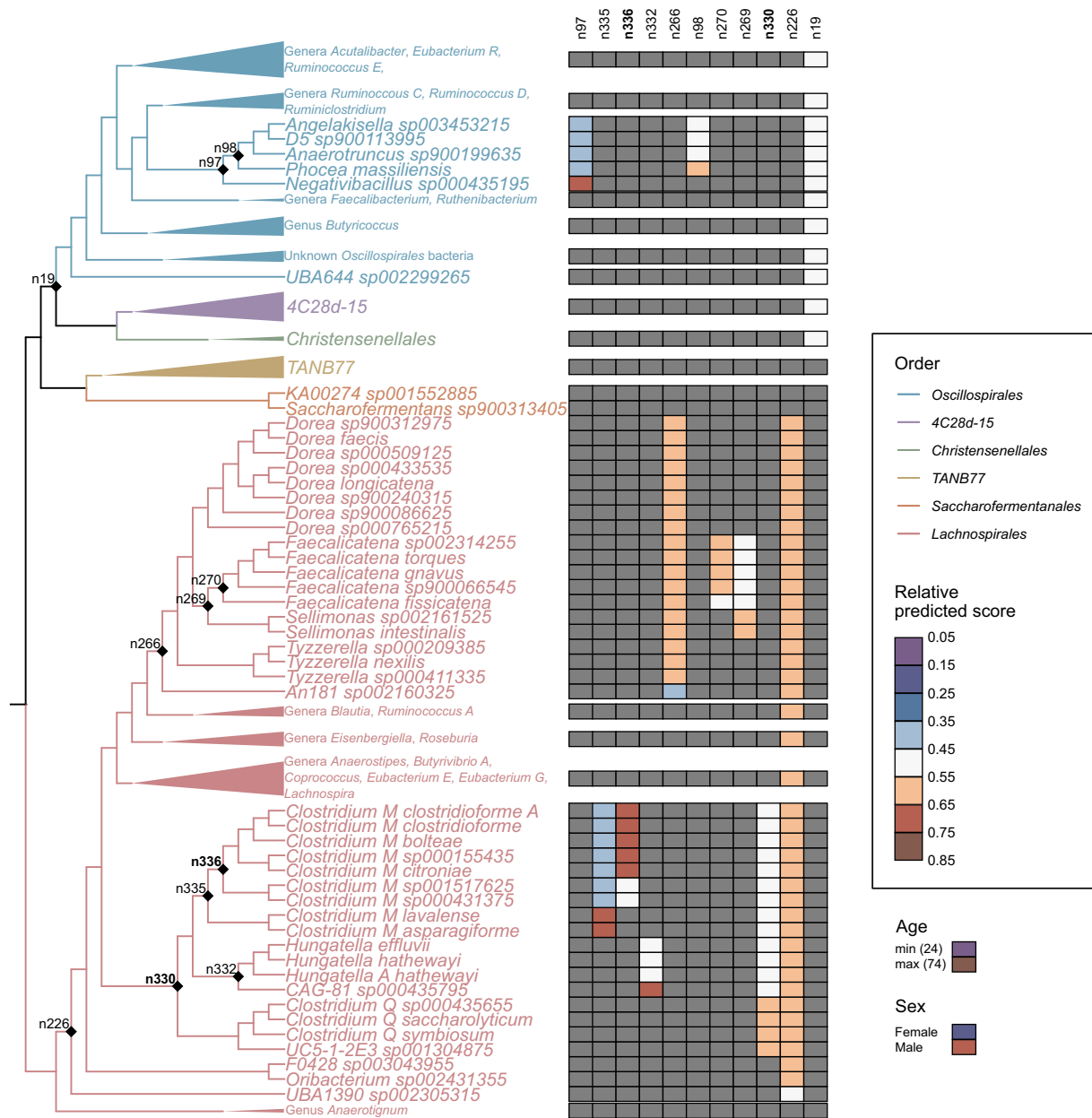


Figure 2. Relative effects of predictive balances and covariates on the FLI < 60 and FLI ≥ 60 classification model (AUC = 0.75) predictions. Nodes of the balances are indicated in the cladogram and the relative effect sizes of their clades (opposite sides of each balance) are shown in the associated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

measurements (see **Figure 1b**), and therefore would be most specific to hepatic steatosis and liver damage.³¹ To do so, we performed feature selection (similarly to continuous FLI) for GGT and triglyceride measurements in subsets of participants grouped by age, sex, and BMI. The feature selection identified two balances within the NCBI *Clostridia* XIVa sub-cluster (identified as n336 and n330) which were

important for both GGT and triglyceride level prediction, and thus likely specific to liver function (**Figure 2**). Bacterial taxa were positively linked to liver function in these balances and included (NCBI species) *Clostridium clostridioforme*, *C. bolteae*, *C. citroniae*, *C. saccharolyticum*, and *C. symbiosum*. On the opposite, negatively associated side of the balances was, among others (NCBI species)

Hungatella effluvii, *H. hathewayi*, and two new GTDB-defined species *Clostridium M sp001517625* and *C. M sp000431375*.

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

The values of predictive balances in the phylogenetic tree cannot be summarized for individual taxa, which means that only a qualitative investigation of the associations between their metabolism and fatty liver was possible in this study. We identified genetic pathways predicted to encode for SCFA (acetate, propanoate, butanoate) and ethanol production, BA metabolism, and choline degradation to trimethylamine (TMA) in representative genomes from the taxa we identified to be linked to liver function (**Figure S8**). These processes were chosen because they have been previously identified to have a mechanistic link to NAFLD (see, e.g., ref.¹⁹).

Acetate and ethanol production pathways appeared to be more common in the representative genomes of the taxa which had a positive association with FLI. In the liver function-specific clades, n336 and n330, MetaCyc pathways for pyruvate fermentation to ethanol III (PWY-6587) and L-glutamate degradation V (via hydroxyglutarate; P162-PWY; produces acetate and butanoate) were present only in genomes positively associated with FLI. In balance n336, also heterolactic fermentation (P122-PWY; produces ethanol and lactate) was more often encoded in the clade positively associated with the high FLI group (3/5) than the opposing negatively associated clade (1/2). In representative genomes from the liver-specific balance n336, potential ethanol producers (PWY-6587) were seen in the positively associated clade (*Clostridium M clostridoforme A* and *Clostridium M sp000155435*), and not in the negatively associated clade (*Clostridium M sp001517625* and *Clostridium M sp000431375*). However, for most balances, such trends were not clear in the qualitative analysis. Furthermore, we did not detect any of these pathways in the representative genomes of two individual taxa positively associated with FLI, *Negativibacillus sp000435195* and *Phoceia massiliensis* (**Figure S8**).

Discussion

The pathophysiology of fatty liver disease in general, and NAFLD in particular, is complex and its clinical diagnosis can be difficult.³² In this study, we utilized metagenomic data from a large population cohort (FINRISK 2002),³⁰ to identify broad links between the overall gut microbiome composition and fatty liver disease, using FLI as a recognized proxy (**Figure 1c**), and identified specific microbial taxa and lineages positively and negatively associated with the high FLI group (**Figure 2**). It should be noted that FLI used in our study as a proxy for liver disease also includes features such as BMI and waist circumference, which associate with metabolic syndrome and diabetes.¹⁶ Links between these diseases and gut microbiome composition are well documented in previous studies.³³ However, fatty liver disease is increasingly thought to be a component of the metabolic syndrome,^{4,34} and while diabetes prevalence is higher in the high FLI group in FINRISK 2002, affected participants still consist only 11% of this group (**Table S1**). Furthermore, we would like to emphasize that our results are not suitable for current clinical application, and should be validated by further, preferably mechanistic studies. We also do not know if our results generalize outside the Finnish population, as all participants in this study were exclusively from Finnish cohorts.

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

Our additional analyses support these main results. While a thorough method comparison is beyond the scope of the current study, the results from the two taxa assignments were very similar despite their differences, such as the fourfold higher number of taxa in the Centrifuge/GTDB data. In the machine learning models (performed only with Centrifuge/GTDB data), excluding participants with intermediate FLI (between 30 and 60) increased

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

When interpreting our results, several levels of associations can be considered according to types of fatty liver disease and the gut microbiome composition. Because FLI has been mostly validated with simple steatosis and NAFLD,^{23,26} we can conservatively contextualize our findings with previous associative work that used these diagnoses or clinical manifestations, only. The cutoff used in our study at FLI \geq 60 has been used to rule in liver steatosis in a Caucasian cohort comparable to ours,²⁶ but also a cutoff at FLI \geq 48 has been found appropriate for simple steatosis in a Portuguese cohort.³⁶ Much lower cutoffs (FLI \geq 20 to 30) have been used in Asian cohorts.^{37–39} Thus, it is likely that our high FLI groups include most participants with liver steatosis or fibrosis in both FINRISK cohorts, but the low FLI group also likely includes participants with low-grade steatosis.

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

Among the significantly high level FLI-associated differences in the gut microbiomes of the

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

In both the SHOGUN/NCBI and Centrifuge/GTDB data, we found significant linear correlations between FLI and beta diversity or two or three main bacterial PC-axes of the samples, respectively (Figure 1c, S3). The model fit was slightly better with Centrifuge/GTDB data, which might be due to the higher number of identified taxa, and thus increased taxonomic resolution (although including putative species in GTDB). Our results support previous observations of differences in beta diversity in relation to persistent NAFLD,⁴¹ and along the NAFLD-cirrhosis spectrum.⁸ Through the loadings of the phylogenetic balances on the PC axes in the Centrifuge/GTDB data, we detected several previously known connections between microbial clades and FLI (Figure S4). Among others, we observed a positive association between high FLI and family *Lachnospiraceae* and negative associations for order *Christensenellales*, genus *Faecalibacterium*, and genus *Bifidobacterium*. The positive association is supported by previous findings of their connection with obesity,⁴² and the negative associations by connections to lean individuals and healthy gut microbiome composition.^{43–45}

Our differential abundance analysis also detected a high number of taxa with significantly increased or decreased abundance in the high FLI group. All following results were observed both in the Centrifuge/GTDB and SHOGUN/NCBI data sets, unless otherwise noted. The majority of the taxa with increased abundance in the high FLI group was from family *Lachnospiraceae*, which supports their positive association with NAFLD reported previously in a number of studies,⁴⁶ but also with obesity (Table S7).⁴² The

increased abundance of genus *Roseburia* has also been highlighted as a characteristic change in the gut microbiome related to NAFLD.^{46,47} In the current study, two members of the genus *Roseburia* were in the top 10 taxa most strongly associated with high FLI. Furthermore, our results support previous findings on the positive associations of, for example, *Collinsella*,⁴⁰ *Prevotella copri*,⁴⁸ *Dorea*,⁴⁷ with NAFLD. We also detected increases in *Sutterella* and *Streptococcus*, previously associated with cirrhosis.⁴⁹ However, we did not find increases in families *Kiloniellaceae* and *Pasteurellaceae*, previously associated with NAFLD.⁴⁶ Among the individual taxa negatively associated with high FLI, families *Ruminococcaceae* and *Oscillospiraceae* (such as genus *Oscillibacter*) were common, which supports previous findings on their connections with NAFLD.^{12,41,46} A high number of putative (GTDB) species were negatively associated with FLI in the Centrifuge/GTDB data, which were understandably not present in the SHOGUN/NCBI data. Many of these were classified in the recently described order *Christensenellales*,²⁸ including families such as *CAG-74*, associated with healthy participants,⁵⁰ and *Christensenellaceae*, which are widespread, highly heritable, and associated with health.^{44,51}

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

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

Strikingly, the strongest associations with FLI in our machine learning models were all inside the *Firmicutes* phylum. A possible reason for this might be the higher relative abundance of phylum *Firmicutes* at high latitudes,⁵² where Finland is. Among the associations we identified, *Faecalicatena gnavus* (NCBI: *Ruminococcus gnavus*) was positively linked with FLI as part of three predictive balances

and associated in previous studies with liver cirrhosis.¹⁷ In their study, oral *Firmicutes*, such as *Veillonella*, were suggested to invade the gut. While our balance-based approach did not detect these taxa, *Megasphaera elsdenii* was positively associated with the high FLI group in our differential abundance analyses (Table S7). This might be due to the strict feature selection employed before the predictive modeling.

Two individual taxa, *Negativibacillus* sp000435195 and *Phoceae massiliensis*, both had strong positive associations with the high FLI group (Figure 2), but the balances including these species were not predictive of the liver function-specific components (triglycerides and GGT). Positive associations of these taxa with fatty liver disease have not been documented previously. However, a decreasing abundance of both bacteria, *Negativibacillus* sp000435195 (NCBI: *Clostridium* sp. CAG:169) and *Phoceae massiliensis* (NCBI: *Phoceae massiliensis*) were seen when the intake of meat and refined cereal was reduced isocalorically in favor of fruit, vegetables, wholegrain cereal, legumes, fish and nuts in overweight and obese subjects in Italy.⁵³ While comparisons between these studies are difficult due to differences in taxa annotations, bacteria, such as *Faecalicatena gnavus* (NCBI: *Ruminococcus gnavus*) and *Clostridium Q saccharolyticum* (NCBI: *Clostridium saccharolyticum*), were also found to respond negatively to the Mediterranean diet. Thus, further study on the connections of these bacteria with gut health and diet is warranted.

Among the taxa negatively associated with high FLI, *Hungatella* (see balance n332, Figure 2) has been previously shown to correlate negatively with the obesity phenotype in mice,⁵⁴ and *H. hathewayi* was found to be a common commensal in the gut of healthy volunteers.⁵⁵ However, genus *Hungatella* has also been positively associated with concentrations of trimethylamine-N-oxide (TMAO),⁵⁶ a metabolite associated with cardiovascular disease and NAFLD. In our study, on the positively associated side (of balance n332) opposite to genus *Hungatella* was a novel GTDB species, *CAG-81* sp000435795, previously included in NCBI genus *Clostridium*. The *CAG-81* genus was recently positively associated with TMAO levels in urine in a study using the GTDB classification.⁵⁷ While we did not find the pathway for TMA (precursor to TMAO) production

in its genome, this would explain the positive association of the *CAG-81* species with high FLI. Furthermore, the previous contradictory results among these taxa could be explained by grouping of putatively TMA producing taxa in *CAG-81* together with the closely related genus *Hungatella*.

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

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

We identified several plausible new associations between individual taxa and clades of bacteria and fatty liver. All taxa were from class *Clostridia*, which are obligate anaerobes. We observed that reference genomes from the bacterial taxa positively associated with high FLI in the liver-specific balances harbored several genetic pathways necessary for ethanol production. Specifically, genes predicted to enable the fermentation of pyruvate to ethanol (MetaCyc PWY-6587) appeared to be common. Endogenous production of ethanol has been known to both induce hepatic steatosis and increase intestinal permeability,⁶⁰ and several of the taxa associated with the high FLI group have also been experimentally shown to produce ethanol, such as *C. M asparagiforme*, *C. M bolteae*, *C. M clostridioforme*/*C. M clostridioforme A*,⁶¹ and *C. Q Saccharolyticum*.⁶² The relative abundances of these putatively ethanol-producing taxa were also predictive of FLI groups in previously unseen data. However, the self-reported alcohol consumption

from the participants was not among the best predictors for the FLI groups, as it was excluded in the feature selection step.

All reference genomes from taxa positively associated with FLI in balance n330 harbored genes predicted to encode for the L-glutamate fermentation V (P162-PWY; **Figure S8**) pathway, which results in the production of acetate and butanoate. Glutamate fermentation could lead to increased microbial protein fermentation in the gut, which has been previously linked with obesity, diabetes, and NAFLD.⁶³ Recently, the combined intake of fructose and microbial acetate production in the gut was experimentally observed to contribute to lipogenesis in the liver in a mouse model.⁶⁴ Interestingly, *C. Q saccharolyticum* (in our study, a taxon positively associated with high FLI deriving from balance n330) was experimentally shown to ferment various carbohydrates, including fructose, to acetate, hydrogen, carbon dioxide, and ethanol.⁶² Furthermore, while our own pathway analysis did not detect BA modification pathways in the reference genome of *C. Q saccharolyticum*, a strain of this species has been highlighted as a probable contributor to NAFLD development through the synthesis of secondary BA.¹⁵ The links between dietary intake and gene regulation, combined with microbial fermentation in the gut warrant further mechanistic experiments to elucidate their links with fatty liver and likely other metabolic diseases.

NAFLD-associated ethanol-producing bacteria in previous cohort studies have all been gram-negatives, such as (NCBI-defined) *Klebsiella pneumoniae*,¹³ and *Escherichia coli*.¹² In our population sample, instead of gram-negatives, bacteria from the *C. M bolteae*, *C. M clostridioforme*/*C. M clostridioforme A* and *C. M citroniae* species (positively associated with high FLI in balance n336) have been described as opportunistic pathogens,⁶⁵ and are hypothesized to exacerbate fatty liver development similarly through endogenous ethanol production. This result suggests that geographical,³⁵ and ethnic variability,⁶⁶ might also strongly affect gut microbiome composition and its associations with disease. In addition to putative endogenous ethanol producers, we identified other taxa positively associated with high FLI in balance n330, for which reference genomes harbored a genetic pathway predicted

to encode for the ability to ferment L-lysine to acetate and butyrate. While the production of these SCFAs is often considered beneficial for gut health, other metabolism of proteolytic bacteria might negatively contribute to fatty liver disease.⁶⁷

Through modeling a previously validated index for fatty liver, FLI, we found replicable associations with specific microbial taxa and likely liver disease of the participants. In addition, the sex and age of participants were also strongly predictive of the FLI group in our models (Figure 2, S7). Their similar positive associations with fatty liver disease are known from previous studies.^{68,69} The associated microbial balances could be used to improve the predictions above the baseline of these covariates on 5/6 regions in Finland in the main cohort. For example, in the model cross-validated with Lapland, the balances were more predictive of the FLI group than the covariates by themselves, and their combination increased the AUC further. Yet, when testing the model where Turku/Loimaa region was used for internal cross-validation, the microbial balances were slightly predictive of the FLI group but failed to improve the AUC over the covariates (Table S2). This pattern might stem from the cultural and genetic west-east division in Finland,^{70,71} with a closer proximity of the Helsinki/Vantaa region to eastern regions than Turku/Loimaa, in both terms. It is thus likely that further incorporation and investigation on the use of spatial information in microbiome modeling would elucidate these geographical patterns in taxa-disease associations.

Our models were also able to accurately predict the FLI group of participants in the external validation cohort, which were from the North Karelia and Helsinki/Vantaa regions. The observed difficulty to geographically extrapolate taxa-disease associations,³⁵ might mean that associations reported in our study are specific to Finland and nearby regions. Notably, many of the positive associations between specific taxa and fatty liver disease have not been reported previously, but the functional potential of these taxa inferred from genomic data is similar to taxa positively associated with NAFLD in previous studies. Thus, the geographical limits of taxa-disease associations reported in studies, such as ours warrant further study. Unfortunately, the generalization of our own results outside of Finland also remains to be addressed.

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

Conclusions

Modeling an established risk index for fatty liver enabled the detection of associations between the disease and gut microbiome composition, to the level of individual taxa. While utilizing FLI as a proxy, NCBI taxa identified with standard statistical methods were supportive of previously reported differences between NAFLD cases and healthy controls. In our machine learning framework, all clades robustly predictive of the FLI group were from the obligately anaerobic gram-positive class *Clostridia*, representing several redefined GTDB genera previously included in the NCBI genus *Clostridium*. Many of the representative genomes of taxa positively associated with high FLI had the genomic potential for endogenous ethanol production. Our results support previous findings on the likely contribution of ethanol and increased gut permeability on the induction of hepatic steatosis. Further support was also found for the involvement of TMA and SCFAs, especially acetate, in the likely pathophysiology of fatty liver disease. Our models were able to predict the FLI group of participants in Finland across geographical regions and in an external Finnish cohort, showing that the associations are robust and generalizable in this population. Based on our results, mechanistic connections between specific microbes and fatty liver disease, and the geographical differences in such taxa-

disease associations should be addressed in further studies.

Materials and Methods

Survey details and sample collection

Cardiovascular disease risk factors have been monitored in Finland since 1972 by conducting a representative population survey every 5 y.³⁰ In the FINRISK 2002 survey, a stratified random population sample was conducted on six geographical regions in Finland. These are North Karelia and Northern Savo in eastern Finland, Turku and Loimaa regions in southwestern Finland, the cities of Helsinki and Vantaa in the capital region, the provinces of Northern Ostrobothnia and Kainuu in northwestern Finland, and the province of Lapland in northern Finland.

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

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

Stool DNA extraction and shallow shotgun metagenome sequencing

DNA extraction was performed according to the Earth Microbiome Project protocols, with the

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

Taxonomic matching and phylogenetic transforms

We quality trimmed the sequences and removed the sequencing adapters with Atropos.⁷⁷ Host reads were removed by mapping the reads against the human genome assembly GRCh38 with Bowtie2.⁷⁸ To improve the taxonomic assignments of our reads, we used a custom index,⁷⁹ based on the Genome Taxonomy Database (GTDB) release 89 taxonomic redefinitions,^{28,80} for reading classification with default parameters in Centrifuge 1.0.4.⁸¹ Viral and eukaryotic sequences were removed in this step, as the database contains only bacterial and archaeal reference genomes. After read the classification, all following steps were performed with R version 3.5.2,⁸² using phyloseq 1.30.0,⁸³ to manage the data. To reduce the number of spurious read assignments, and to facilitate more accurate phylogenetic transformations, only reads classified at the species level, matching individual GTDB reference genomes, were retained. Samples with less than 50,000 reads, from pregnant participants or recorded antibiotic use in the past 6 months were removed, resulting in a final number of 6,269 samples. We first filtered taxa not seen with more than 3 counts in at least 1% of the samples and those with a coefficient of variation ≤ 3 across all samples, following McMurdie and Holmes,⁸³ with a slight

adaption from 20% of the samples to 1% of the samples, because of our larger sample size. The complete bacterial and archaeal phylogenetic trees of the GTDB release 89 reference genomes, constructed from an alignment of 120 bacterial or 122 archaeal marker genes,²⁸ were then combined with our taxa tables. The resulting trees were thus subset only to species that were observed in at least one sample in our data. The read counts were transformed to phylogenetic node balances in both trees with PhILR.²⁹ The default method for PhILR inputs a pseudocount of 1 for taxa absent in an individual sample before the transform.

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

To further validate our approach, assess how the use of the GTDB taxonomic redefinitions and custom database affected our results, and to facilitate comparisons with previous results, we annotated our raw reads in FINRISK 2002 samples also with NCBI taxonomy and performed several additional analyses. For these comparison data, after quality trimming the FINRISK 2002 reads and removing host sequences as described above, SHOGUN v1.0.5,⁵⁹ was used for taxonomy assignments against the NCBI RefSeq version 82 (May 8, 2017) database containing complete bacterial, archaeal, and viral genomes. To facilitate comparisons between different annotations, we subset the samples included in the SHOGUN/NCBI annotated data to those included in the Centrifuge/GTDB data (for exclusion criteria, see above).

Covariates

Because the fatty liver disease is underdiagnosed at the population level,²⁵ and our sampling did not have extensive coverage of liver fat measurements, we chose to use the Fatty Liver index as a proxy for fatty liver.²⁶ Furthermore, the index performs well in cohorts of Caucasian ethnicity, such as ours, to diagnose the presence of NAFLD.²³ We calculated FLI after Bedogni et al.²⁶ ($e^{0.953 \cdot \log_e(\text{triglyceridesmg/dL})} + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745$) / ($1 + e^{0.953 \cdot \log_e(\text{triglyceridesmg/dL})} + 0.139 \cdot \text{BMI} + 0.718 \cdot \log_e(\text{GGT}) + 0.053 \cdot \text{waist circumference} - 15.745$) * 100. We chose the cutoff at FLI ≥ 60 to identify participants likely to be diagnosed with hepatic steatosis (positive likelihood ratio = 4.3 and negative likelihood ratio = 0.5, after Bedogni et al.²⁶). Triglycerides, gamma-glutamyl transferase (GGT), BMI, and waist circumference measurements had near complete coverage for the participants in our data. Self-reported alcohol use was calculated as grams of pure ethanol per week. Cases with missing values were omitted in linear regression models. At least one feature used for FLI calculation was missing for 20 participants in FINRISK 2002 (0.3%) and the self-reported alcohol use was missing for 247 participants (3.9%). In the machine learning framework, missing values for FLI and self-reported alcohol use were mean imputed. However, for the feature selection to

identify liver function-specific balances, GGT, triglycerides, and BMI were not imputed but observations where any of these were missing were simply removed.

Taxa composition and alpha diversity

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

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

Beta diversity and linear modeling of FLI

In the Centrifuge/GTDB data, beta diversity was calculated as Euclidian distance of the PhILR balances through Principal Component Analysis (PCA) on bacterial and archaeal balances separately with “rda” in vegan 2.5.6.⁸⁷ To calculate beta diversity with the SHOGUN/GTDB data, raw bacterial taxa counts were centered log-ratio (CLR) transformed with “transform” in microbiome 1.8.0,⁸⁸ and their Euclidian distances were obtained similarly through PCA. Linear regression models were constructed for FLI with “lm” in base R,⁸² with Centrifuge/GTDB data, and separately with SHOGUN/NCBI. Log₁₀ (FLI) was used as the dependent variable and the first three bacterial PCs, sex, age, and self-reported alcohol were used as the independent variables. Archaeal PCs were not included in the models because none of them was significantly correlated with FLI in Centrifuge/GTDB data (all $P > 0.001$). To visualize the association between beta diversity and FLI, the FLI of each participant was plotted against its quantiles along the three bacterial PC

axes in Centrifuge/GTDB data (**Figure 1c**). A comparison of the associations with the alternative SHOGUN/NCBI annotated data was also included in the SI (**Figure S3**).

Differential abundance of individual taxa between the FLI groups

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

FLI modeling within a machine learning framework

In the machine learning framework, both regression and categorical models were constructed for FLI, using only the Centrifuge/GTDB data. The feature selection, hyperparameter optimization, and internal cross-validation methods were identical for both approaches, unless otherwise stated. The continuous or categorical FLI (groups of FLI < 60 and FLI \geq 60) were modeled with xgboost 0.90.0.2,⁹⁰ by using both bacterial and archaeal balances, sex, age, and self-reported alcohol use as preliminary predictor features. We used FLI 60 as the cutoff for ruling in fatty liver (steatosis) for the classification, after Bedogni et al.²⁶ The data were first split into 70% train and 30% test sets while preserving sex and region balance. To take into account geographical differences (see, e.g., ref.³⁵) and to find robust patterns across all six sampled regions in Finland between the features and FLI group, we used Leave-One-Group-Out Cross-Validation (LOGOCV) inside the 70% train set to construct 6 separate

models in each optimization step. Because of the high dimensionality of the data (3,423 predictor features), feature selection by filtering was first performed inside the training data, based on random forest permutation as recommended by Bommert et al.⁹¹ Briefly, permutation importance is based on accuracy, or specifically, the difference in accuracy between real and permuted (random) values of the specific variable, averaged in all trees across the whole random forest. The permutation importance in models based on the six LOGOCV subsets of the training data was calculated with mlr 2.16.0,⁹² and the simple intersect between the top 50 features in all LOGOCV subsets was retained as the final set of features. Thus, the feature selection was influenced by the training data from all six geographical regions, but this only serves to limit the number of chosen features because of the required simple intersect. This approach was used to obtain a set of core predictive features that would have the potential for generalizability across the regions. The number of features included in the models by this approach was deemed appropriate, since the relative effect size of the last included predictor was very small (<0.1 change in classification probability across its range).

Bayesian hyperparameter optimization of the xgboost models was then performed with only the selected features. An optimal set of parameters for the xgboost models was searched over all LOGOCV subsets with “mbo” in mlrMBO 1.1.3,⁹³ using 30 preliminary rounds with randomized parameters, followed by 100 optimization rounds. Parameters in the xgboost models and their considered ranges were learning rate (eta) [0.001, 0.3], gamma [0.1, 5], maximum depth of a tree [2, 8], minimum child weight [1, 10], fraction of data subsampled per each iteration [0.2, 0.8], fraction of columns subsampled per tree [0.2, 0.9], and maximum number of iterations (nrounds) [50, 5000]. The parameters recommended by these searchers were as following for regression: eta = 0.00889; gamma = 2.08; max_depth = 2; min_child_weight = 8; subsample = 0.783; colsample_bytree = 0.672; nrounds = 1,810, and for classification: eta = 0.00107; gamma = 0.137; max_depth = 5; min_child_weight = 9; subsample = 0.207; colsample_bytree = 0.793; nrounds = 4,328. We used Root-Mean-Square Error (RMSE) for the regression models and Area Under the ROC Curve (AUC) for the classification models to

measure model fit on the left-out data (region) in each LOGOCV subset. Receiver operating characteristic and precision–recall curves for these validation metrics were calculated with “evalmod” in precrec v0.11.2.⁹⁴ The final models were trained on the LOGOCV subset training data, the data from one region thus omitted per model, and using the selected features and optimized hyperparameters. Internal validation of these models was conducted against participants only from the region omitted from each model, in the 30% test data which was not used in model training or optimization. Sensitivity analysis was conducted by using only the predictive covariates (sex and age) or balances separately, with the same hyperparameters, data partitions, and cross-region internal validation as for the full models.

Partial dependence interpretation of the FLI classification models

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

Construction of alternative classification models to discern between FLI < 30 and FLI ≥ 60 groups

To assess the robustness of the models and how removing the participants with intermediate FLI (between 30 and 60) affects model performance, we removed this group (N = 1,910) and constructed alternative classification models to discern between

the $FLI < 30$ and $FLI \geq 60$ groups. Other than removing the intermediate FLI participants and resulting new random split to the train (70%) and test (30%) sets, these models were constructed identically to the main models, including LOGOCV design, feature selection, and hyperparameter optimization. The recommended parameters for this classification task were $\eta = 0.00102$; $\gamma = 3.7$; $\max_depth = 2$; $\min_child_weight = 5$; $subsample = 0.49$; $colsample_bytree = 0.631$; $nrounds = 3,119$. Interpretation of partial dependence was also performed identically, but only the relative effects of the model features were plotted without a cladogram.

Exclusion of validation region data from feature selection and hyperparameter optimization

Because training data from all six regions are used to inform the selection of optimal features and hyperparameters, the validation region data cannot be considered completely independent of the training of the LOGOCV models. Thus, we constructed a set of classification models for the $FLI \geq 60$ and $FLI < 60$ groups, where all validation region samples also in the training data were excluded from the simple intercept of top 50 features in each LOGOCV set and from the subsequent hyperparameter optimization. These models with individualized features and hyperparameters were then tested on the validation region samples in the unseen test data to estimate how model performance was affected. The main test (70%) and train (30%) sets were identical to the main models, but additionally, 6 randomized 70/30 splits nested inside the test set (excluding the validation region) were used in hyperparameter optimization to reduce overfitting. Average optimal hyperparameters in the six models were $\eta = 0.00106$; $\gamma = 4.3$; $\max_depth = 2$; $\min_child_weight = 7$; $subsample = 0.36$; $colsample_bytree = 0.613$; $nrounds = 1,772$.

External validation of the models in a separate population cohort

To further validate our models and results, we leveraged the data from a more recent population cohort in Finland, FINRISK 2007 (see **Table S1**). In this cohort, the choice of participants, sample collection,

and related methods for the data used in the current study were similar to FINRISK 2002 to facilitate inter-cohort comparisons, and are reported elsewhere.³⁰ The study protocol of FINRISK 2007 was approved by the Coordinating Ethical Committee of the Hospital District of Helsinki and Uusimaa (Ref. 229/EO/2006). All participants have signed an informed consent.

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

Identification of predictive features specific to liver function

Because the FLI also incorporates BMI and waist circumference, and they strongly contribute to the index,²⁶ we deemed it necessary to further investigate which of the identified balances were specific to liver function. The participants were first grouped by age (<40 , $40-60$, and $60<$), sex (female or male) and BMI (<25 , $25-30$, and $30<$) into 18 categories ($N = 105-711$ per category). We performed feature selection similarly to the FLI models by fitting

random forest regressors for GGT and triglycerides with mlr 2.16.0.⁹² This was done separately in each of the 18 categories, and in each category, we again used LOGOCV with the regions to obtain 6 runs per category. Finally, the features predictive of GGT or triglycerides in each category were selected as the intersect of top 50 features in the 6 LOGOCV iterations by permutation importance. The intersect of features predictive of GGT or triglycerides in any of the categories and the features predictive of categorical FLI were identified as specific to liver function.

Pathway inference for taxa associated with FLI

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

The genomes of all 336 bacteria under at least one of the predictive balances were downloaded from NCBI. One hundred and nineteen of these genomes were originally not annotated, which is a requirement for pathway prediction. Therefore, Prokka v1.14.6,⁹⁸ was used to annotate the 119 unannotated genomes as a preliminary step. Pathway predictions were then performed for all 336 genomes with mpwt v0.5.3 multiprocessing tool,⁹⁹ for the PathoLogic pipeline of Pathway Tools 23.0.¹⁰⁰ Pathways for ethanol and short chain fatty acid (acetate, butyrate, propionate) production, bile acid metabolism, and choline degradation to trimethylamine were identified from MetaCyc pathway classifications (see ref.¹⁰¹ and **Table S4**). The prevalence of these processes was then assessed in the analyzed genomes and summarized per process to consider the possible links of the taxa with fatty liver pathophysiology. Finally, the presence of individual pathways for acetate and ethanol production was also outlined for each genome.

Acknowledgments

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

Data availability statement

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

Disclosure of interest

V.S. has consulted for Novo Nordisk and Sanofi and received honoraria from these companies. He also has an ongoing research collaboration with Bayer AG, all unrelated to this study. R.L. serves as a consultant or advisory board member for Anylam/Regeneron, Arrowhead Pharmaceuticals, AstraZeneca, Bird Rock Bio, Boehringer Ingelheim, Bristol-Myer Squibb, Celgene, Cirus, CohBar, Conatus, Eli Lilly, Galmed, Gemphire, Gilead, Glympse bio, GNI, GRI Bio, Inipharm, Intercept, Ionis, Janssen Inc., Merck, Metacrine, Inc., NGM Biopharmaceuticals, Novartis, Novo Nordisk, Pfizer, Prometheus, Promethera, Sanofi, Siemens, and Viking Therapeutics. In addition, his institution has received grant support from Allergan, Boehringer-Ingelheim, Bristol-Myers Squibb, Cirus, Eli Lilly and Company, Galectin Therapeutics, Galmed Pharmaceuticals, GE, Genfit, Gilead, Intercept, Grail, Janssen, Madrigal Pharmaceuticals, Merck, NGM Biopharmaceuticals, NuSirt, Pfizer, pH Pharma, Prometheus, and Siemens. He is also co-founder of Liponex, Inc.

Funding

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

ORCID

Matti O. Ruuskanen  <http://orcid.org/0000-0003-4221-2880>
 Fredrik Åberg  <http://orcid.org/0000-0002-3833-0705>
 Ville Männistö  <http://orcid.org/0000-0002-0735-400X>
 Aki S. Havulinna  <http://orcid.org/0000-0002-4787-8959>
 Guillaume Méric  <http://orcid.org/0000-0001-6288-9958>
 Rohit Loomba  <http://orcid.org/0000-0002-4845-9991>
 Yoshiki Vázquez-Baeza  <http://orcid.org/0000-0001-6014-2009>
 Anupriya Tripathi  <http://orcid.org/0000-0001-8912-9684>
 Michael Inouye  <http://orcid.org/0000-0001-9413-6520>
 Veikko Salomaa  <http://orcid.org/0000-0001-7563-5324>
 Mohit Jain  <http://orcid.org/0000-0001-8628-2069>
 Rob Knight  <http://orcid.org/0000-0002-0975-9019>
 Leo Lahti  <http://orcid.org/0000-0001-5537-637X>
 Teemu J. Niiranen  <http://orcid.org/0000-0002-7394-7487>

Authors' contributions

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

References

1. Younossi ZM, Koenig AB, Abdelatif D, Fazel Y, Henry L, Wymer M. Global epidemiology of nonalcoholic fatty liver disease-Meta-analytic assessment of prevalence, incidence, and outcomes. *Hepatology*. 2016;64(1):73–84. doi:10.1002/hep.28431.
2. Marchesini G, Bugianesi E, Forlani G, Cerrelli F, Lenzi M, Manini R, Natale S, Vanni E, Villanova N, Melchionda N, et al. Nonalcoholic fatty liver, steatohepatitis, and the metabolic syndrome. *Hepatology*. 2003;37(4):917–923. doi:10.1053/jhep.2003.50161.
3. Chalasani N, Younossi Z, Lavine JE, Diehl AM, Brunt EM, Cusi K, Charlton M, Sanyal AJ. The diagnosis and management of non-alcoholic fatty liver disease: practice guideline by the American association for the study of liver diseases, American College of Gastroenterology, and the American Gastroenterological Association. *Hepatology*. 2012;55(6):2005–2023. doi:10.1002/hep.25762.
4. Yki-Järvinen H. Non-alcoholic fatty liver disease as a cause and a consequence of metabolic syndrome. *Lancet Diabetes Endocrinol*. 2014;2(11):901–910. doi:10.1016/S2213-8587(14)70032-4.
5. Toshikuni N. Clinical differences between alcoholic liver disease and nonalcoholic fatty liver disease. *World J Gastroenterol*. 2014;20(26):8393–8406. doi:10.3748/wjg.v20.i26.8393.
6. Rinella M, Charlton M. The globalization of nonalcoholic fatty liver disease: prevalence and impact on world health. *Hepatology*. 2016;64(1):19–22. doi:10.1002/hep.28524.
7. Gilbert JA, Blaser MJ, Caporaso JG, Jansson JK, Lynch SV, Knight R. Current understanding of the human microbiome. *Nat Med*. 2018;24(4):392–400. doi:10.1038/nm.4517.
8. Caussy C, Tripathi A, Humphrey G, Bassirian S, Singh S, Faulkner C, Bettencourt R, Rizo E, Richards L, Xu ZZ, et al. A gut microbiome signature for cirrhosis due to nonalcoholic fatty liver disease. *Nat Commun*. 2019;10(1):1406. doi:10.1038/s41467-019-09455-9.
9. Compare D, Coccoli P, Rocco A, Nardone OM, De Maria S, Carteni M, Nardone G. Gut–liver axis: the impact of gut microbiota on non-alcoholic fatty liver disease. *Nutr Metab Cardiovasc Dis*. 2012;22(6):471–476. doi:10.1016/j.numecd.2012.02.007.
10. Miele L, Valenza V, Torre GL, Montalto M, Cammarota G, Ricci R, Mascianà R, Forgione A, Gabrieli ML, Perotti G, et al. Increased intestinal permeability and tight junction alterations in nonalcoholic fatty liver disease. *Hepatology*. 2009;49(6):1877–1887. doi:10.1002/hep.22848.
11. Spencer MD, Hamp TJ, Reid RW, Fischer LM, Zeisel SH, Fodor AA. Association between composition of the human gastrointestinal microbiome and development of fatty liver with choline deficiency. *Gastroenterology*. 2011;140(3):976–986. doi:10.1053/j.gastro.2010.11.049.
12. Zhu L, Baker SS, Gill C, Liu W, Alkhoury R, Baker RD, Gill SR. Characterization of gut microbiomes in non-alcoholic steatohepatitis (NASH) patients: a connection between endogenous alcohol and NASH. *Hepatology*. 2013;57(2):601–609. doi:10.1002/hep.26093.
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 disease caused by high-alcohol-producing *Klebsiella pneumoniae*. *Cell Metab*. 2019;30(4):675–688.e7. doi:10.1016/j.cmet.2019.08.018.
14. Loomba R, Seguritan V, Li W, Long T, Klitgord N, Bhatt A, Dulai PS, Caussy C, Bettencourt R, Highlander SK, et al. Gut microbiome-based metagenomic signature for non-invasive detection of advanced fibrosis in human non-alcoholic fatty liver disease. *Cell Metab*. 2017;25(5):1054–1062.e5. doi:10.1016/j.cmet.2017.04.001.
15. Jiao N, Wu D, Yang Z, Fang S, Li X, Yuan M, Zhu R, Zhu L. Gut bacteria contributes to NAFLD pathogenesis by promoting secondary bile acids biosynthesis. *The FASEB Journal*. 2019;33:4–126
16. Aron-Wisniewsky J, Vigliotti C, Witjes J, Le P, Holleboom AG, Verheij J, Nieuwdorp M, Clément K. Gut microbiota and human NAFLD: disentangling microbial signatures from metabolic disorders. *Nat Rev Gastroenterol Hepatol*. 2020;17(5):279–297. doi:10.1038/s41575-020-0269-9.
17. Qin N, Yang F, Li A, Prifti E, Chen Y, Shao L, Guo J, Le Chatelier E, Yao J, Wu L, et al. Alterations of the human gut microbiome in liver cirrhosis. *Nature*. 2014;513(7516):59–64. doi:10.1038/nature13568.

18. Liu Y, Meric G, Havulinna AS, Teo SM, Ruuskanen M, Sanders J, Zhu Q, Tripathi A, Verspoor K, Cheng S, et al. Early prediction of liver disease using conventional risk factors and gut microbiome-augmented gradient boosting [Internet]. *Genet Genomic Med.* 2020;06(24). [cited 2020 Jul 28]. Available from <https://www.medrxiv.org/content/10.1101/2020.06.24.20138933v1>. doi:10.1101/2020.06.24.20138933.
19. Safari Z, Gérard P. The links between the gut microbiome and non-alcoholic fatty liver disease (NAFLD). *Cell Mol Life Sci.* 2019;76(8):1541–1558. doi:10.1007/s00018-019-03011-w.
20. Carpino G, Del Ben M, Pastori D, Carnevale R, Baratta F, Overi D, Francis H, Cardinale V, Onori P, Safarikia S, et al. Increased liver localization of lipopolysaccharides in human and experimental NAFLD. *Hepatology.* 2019;72(2):470–485. doi:10.1002/hep.31056.
21. Li Q, Dhyani M, Grajo JR, Sirlin C, Samir AE. Current status of imaging in nonalcoholic fatty liver disease. *World J Hepatol.* 2018;10(8):530–542. doi:10.4254/wjh.v10.i8.530.
22. Koehler EM, Schouten JNL, Hansen BE, Hofman A, Stricker BH, Janssen HLA. External validation of the fatty liver index for identifying nonalcoholic fatty liver disease in a population-based study. *Clin Gastroenterol Hepatol.* 2013;11(9):1201–1204. doi:10.1016/j.cgh.2012.12.031.
23. Vanni E, Bugianesi E. Editorial: utility and pitfalls of fatty liver index in epidemiologic studies for the diagnosis of NAFLD. *Aliment Pharmacol Ther.* 2015;41(4):406–407. doi:10.1111/apt.13063.
24. Salosensaari A, Laitinen V, Havulinna AS, Meric G, Cheng S, Perola M, Valsta L, Alftan G, Inouye M, Watrous JD, et al. Taxonomic signatures of long-term mortality risk in human gut microbiota [Internet]. *Epidemiology.* 2020;12(30): [cited 2020 Jan 4]. Available from <https://www.medrxiv.org/content/10.1101/2019.12.30.19015842v2>. doi:10.1101/2019.12.30.19015842
25. Alexander M, Loomis AK, Fairburn-Beech J, van der Lei J, Duarte-Salles T, Prieto-Alhambra D, Ansell D, Pasqua A, Lapi F, Rijnbeek P, et al. Real-world data reveal a diagnostic gap in non-alcoholic fatty liver disease. *BMC Med.* 2018;16(1):130. doi:10.1186/s12916-018-1103-x.
26. Bedogni G, Bellentani S, Miglioli L, Masutti F, Passalacqua M, Castiglione A, Tiribelli C. The fatty liver index: a simple and accurate predictor of hepatic steatosis in the general population. *BMC Gastroenterol.* 2006;6(1):33. doi:10.1186/1471-230X-6-33.
27. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R, Knights D. Evaluating the information content of shallow shotgun metagenomics. *mSystems* [Internet]. 2018; 3. cited 2020 Apr 9 Available from <https://msystems.asm.org/content/3/6/e00069-18>
28. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A, Hugenholtz P. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol.* 2018;36(10):996–1004. doi:10.1038/nbt.4229.
29. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform enhances analysis of compositional microbiota data. *eLife.* 2017;6:e21887. doi:10.7554/eLife.21887.
30. Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, Kuulasmaa K, Laatikainen T, Männistö S, Peltonen M, et al. Cohort profile: the National FINRISK study. *Int J Epidemiol.* 2018;47(3):696. i. doi:10.1093/ije/dyx239.
31. Banderas DZ, Escobedo J, Gonzalez E, Liceaga MG, Ramirez JC, Castro MG. γ -Glutamyl transferase: a marker of nonalcoholic fatty liver disease in patients with the metabolic syndrome. *Eur J Gastroenterol Hepatol.* 2012;24(7):805–810. doi:10.1097/MEG.0b013e328354044a.
32. Haas JT, Francque S, Staels B. Pathophysiology and mechanisms of nonalcoholic fatty liver disease. *Annu Rev Physiol.* 2016;78(1):181–205. doi:10.1146/annurev-physiol-021115-105331.
33. Castaner O, Goday A, Park Y-M, Lee S-H, Magkos F, Shio S-ATE, Schröder H. The gut microbiome profile in obesity: a systematic review. *Int J Endocrinol.* 2018; 2018:4095789. doi:10.1155/2018/4095789.
34. Eslam M, Sanyal AJ, George J, Sanyal A, Neuschwander-Tetri B, Tiribelli C, Kleiner DE, Brunt E, Bugianesi E, Yki-Järvinen H, et al. MAFLD: a consensus-driven proposed nomenclature for metabolic associated fatty liver disease. *Gastroenterology.* 2020;158(7):1999–2014.e1. doi:10.1053/j.gastro.2019.11.312.
35. He Y, Wu W, Zheng H-M, Li P, McDonald D, Sheng H-F, Chen M-X, Chen Z-H, Ji G-Y, Zheng ZDX, et al. Regional variation limits applications of healthy gut microbiome reference ranges and disease models. *Nat Med.* 2018;24(10):1532–1535. doi:10.1038/s41591-018-0164-x.
36. Carvalhana S, Leitão J, Alves AC, Bourbon M, Cortez-Pinto H. How good is controlled attenuation parameter and fatty liver index for assessing liver steatosis in general population: correlation with ultrasound. *Liver Int.* 2014;34(6):e111–7. doi:10.1111/liv.12305.
37. Leathers JS, Balderramo D, Prieto J, Diehl F, Gonzalez-Ballerga E, Ferreiro MR, Carrera E, Barreyro F, Diaz-Ferrer J, Singh D. PIB: a score to select sorafenib treatment candidates for hepatocellular carcinoma in resource-limited settings. *Hepat Mon.* 2018;18(10):18. doi:10.5812/hepatmon.82345.
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 H-C, Lee F-Y, et al. External validation of fatty liver index for identifying ultrasonographic fatty liver in a large-scale cross-sectional study in Taiwan. *PLoS One.* 2015;10(3): e0120443. doi:10.1371/journal.pone.0120443.
39. Huang X, Xu M, Chen Y, Peng K, Huang Y, Wang P, Ding L, Lin L, Xu Y, Chen Y, et al. Validation of the fatty

- liver index for nonalcoholic fatty liver disease in middle-aged and elderly Chinese. *Medicine (Baltimore)*. 2015;94(40):e1682. doi:10.1097/MD.0000000000001682.
40. Astbury S, Atallah E, Vijay A, Aithal GP, Grove JL, Valdes AM. Lower gut microbiome diversity and higher abundance of proinflammatory genus *Collinsella* are associated with biopsy-proven nonalcoholic steatohepatitis. *Gut Microbes*. 2020;11(3):569–580. doi:10.1080/19490976.2019.1681861.
 41. Kim H-N, Joo E-J, Cheong HS, Kim Y, Kim H-L, Shin H, Chang Y, Ryu RS. Gut microbiota and risk of persistent nonalcoholic fatty liver diseases. *J Clin Med*. 2019;8(8):1089. doi:10.3390/jcm8081089.
 42. de la Cuesta-zuluaga J, Corrales-Agudelo V, Velásquez-Mejía EP, Carmona JA, Abad JM, Escobar JS. Gut microbiota is associated with obesity and cardiometabolic disease in a population in the midst of Westernization. *Sci Rep*. 2018;8(1):11356. doi:10.1038/s41598-018-29687-x.
 43. O’Callaghan A, van Sinderen D. Bifidobacteria and their role as members of the human gut microbiota. *Front Microbiol* 2016;7:925. doi:10.3389/fmicb.2016.00925.
 44. Waters JL, Ley RE. The human gut bacteria Christensenellaceae are widespread, heritable, and associated with health. *BMC Biol*. 2019;17(1):83. doi:10.1186/s12915-019-0699-4.
 45. Ferreira-Halder CV, Faria de S AV, Andrade SS. Action and function of *Faecalibacterium prausnitzii* in health and disease. *Best Pract Res Clin Gastroenterol*. 2017;31(6):643–648. doi:10.1016/j.bpg.2017.09.011.
 46. Liu G-L, Qingxi Q-Z, Hongyun H-W. Characteristics of intestinal bacteria with fatty liver diseases and cirrhosis. *Ann Hepatol*. 2019;18(6):796–803. doi:10.1016/j.aohep.2019.06.020.
 47. Raman M, Ahmed I, Gillevet PM, Probert CS, Ratcliffe NM, Smith S, Greenwood R, Sikaroodi M, Lam V, Crotty P, et al. Fecal microbiome and volatile organic compound metabolome in obese humans with nonalcoholic fatty liver disease. *Clin Gastroenterol Hepatol*. 2013;11(7):e3. doi:10.1016/j.cgh.2013.02.015.
 48. Dong TS, Katzka W, Lagishetty V, Luu K, Hauer M, Piseigna J, Jacobs JP. A microbial signature identifies advanced fibrosis in patients with chronic liver disease mainly due to NAFLD. *Sci Rep*. 2020;10(1):2771. doi:10.1038/s41598-020-59535-w.
 49. Bajaj JS, Idilman R, Mabudian L, Hood M, Fagan A, Turan D, White MB, Karakaya F, Wang J, Atalay R, et al. Diet affects gut microbiota and modulates hospitalization risk differentially in an international cirrhosis cohort. *Hepatology*. 2018;68(1):234–247. doi:10.1002/hep.29791.
 50. Bowerman KL, Rehman SF, Vaughan A, Lachner N, Budden KF, Kim RY, Wood DLA, Gellatly SL, Shukla SD, Wood LG, et al. Disease-associated gut microbiome and metabolome changes in patients with chronic obstructive pulmonary disease. *Nat Commun*. 2020;11(1):5886. doi:10.1038/s41467-020-19701-0.
 51. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, Beaumont M, Van Treuren W, Knight R, Bell JT, et al. Human genetics shape the gut microbiome. *Cell*. 2014;159(4):789–799. doi:10.1016/j.cell.2014.09.053.
 52. Suzuki TA, Worobey M. Geographical variation of human gut microbial composition. *Biol Lett*. 2014;10(2):20131037. doi:10.1098/rsbl.2013.1037.
 53. Meslier V, Laiola M, Roager HM, Filippis FD, Roume H, Quinquis B, Giacco R, Mennella I, Ferracane R, Pons N, et al. Mediterranean diet intervention in overweight and obese subjects lowers plasma cholesterol and causes changes in the gut microbiome and metabolome independently of energy intake. *Gut* [Internet] cited 2020 Jun 2]; Available from. 2020;69(7):1258–1268. doi:10.1136/gutjnl-2019-320438
 54. Cui C, Li Y, Gao H, Zhang H, Han J, Zhang D, Li Y, Zhou J, Lu C, Su X. Modulation of the gut microbiota by the mixture of fish oil and krill oil in high-fat diet-induced obesity mice. *PLoS One*. 2017;12(10):e0186216. doi:10.1371/journal.pone.0186216.
 55. Manzoor SE, McNulty CAM, Nakiboneka-Ssenabulya D, Lecky DM, Hardy KJ, Hawkey PM. Investigation of community carriage rates of *Clostridium difficile* and *Hungatella hathewayi* in healthy volunteers from four regions of England. *J Hosp Infect*. 2017;97(2):153–155. doi:10.1016/j.jhin.2017.05.014.
 56. Genoni A, Christophersen CT, Lo J, Coghlan M, Boyce MC, Bird AR, Lyons-Wall P, Devine DA. Long-term Paleolithic diet is associated with lower resistant starch intake, different gut microbiota composition and increased serum TMAO concentrations. *Eur J Nutr*. 2020;59(5):1845–1858. doi:10.1007/s00394-019-02036-y.
 57. Burton KJ, Krüger R, Scherz V, Münger LH, Picone G, Vionnet N, Bertelli C, Greub G, Capozzi F, Vergères G. Trimethylamine-N-Oxide postprandial response in plasma and urine is lower after fermented compared to non-fermented dairy consumption in healthy adults. *Nutrients*. 2020;12(1):234. doi:10.3390/nu12010234.
 58. Jiang W, Wu N, Wang X, Chi Y, Zhang Y, Qiu X, Hu Y, Li J, Liu Y. Dysbiosis gut microbiota associated with inflammation and impaired mucosal immune function in intestine of humans with non-alcoholic fatty liver disease. *Sci Rep*. 2015;5(1):8096. doi:10.1038/srep08096.
 59. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Knight R, Knights D. SHOGUN: a modular, accurate and scalable framework for microbiome quantification. *Bioinformatics*. 2020;36(13):btaa277. doi:10.1093/bioinformatics/btaa277.
 60. de Faria Ghetti F, Oliveira DG, de Oliveira JM, de Castro Ferreira LEVV, Cesar DE, Moreira APB. Influence of gut microbiota on the development and progression of nonalcoholic steatohepatitis. *Eur J Nutr*. 2018;57(3):861–876. doi:10.1007/s00394-017-1524-x.

61. Mohan R, Namsolleck P, Lawson PA, Osterhoff M, Collins MD, Alpert C-A, Blaut M. *Clostridium asparagiforme* sp. nov., isolated from a human faecal sample. *Syst Appl Microbiol.* 2006;29(4):292–299. doi:10.1016/j.syapm.2005.11.001.
62. Murray WD, Khan AW, van den BERG L. *Clostridium saccharolyticum* sp. nov., a saccharolytic species from sewage sludge. *Int J Syst Bacteriol.* 1982;32(1):132–135. doi:10.1099/00207713-32-1-132.
63. Diether NE, Willing BP. Microbial fermentation of dietary protein: an important factor in diet–microbe–host interaction. *Microorganisms.* 2019; 7(1)19: doi:10.3390/microorganisms7010019.
64. Zhao S, Jang C, Liu J, Uehara K, Gilbert M, Izzo L, Zeng X, Trefely S, Fernandez S, Carrer A, et al. Dietary fructose feeds hepatic lipogenesis via microbiota-derived acetate. *Nature.* 2020;579(7800):586–591. doi:10.1038/s41586-020-2101-7.
65. Dehoux P, Marvaud JC, Abouelleil A, Earl AM, Lambert T, Dauga C. Comparative genomics of *Clostridium bolteae* and *Clostridium clostridioforme* reveals species-specific genomic properties and numerous putative antibiotic resistance determinants. *BMC Genomics.* 2016;17(1):819. doi:10.1186/s12864-016-3152-x.
66. Deschasaux M, Bouter KE, Prodan A, Levin E, Groen AK, Herrema H, Tremaroli V, Bakker GJ, Attaye I, Pinto-Sietsma S-J, et al. Depicting the composition of gut microbiota in a population with varied ethnic origins but shared geography. *Nat Med.* 2018;24(10):1526–1531. doi:10.1038/s41591-018-0160-1.
67. Canfora EE, Meex RCR, Venema K, Blaak EE. Gut microbial metabolites in obesity, NAFLD and T2DM. *Nat Rev Endocrinol.* 2019;15(5):261–273. doi:10.1038/s41574-019-0156-z.
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. Nonalcoholic fatty liver disease: prevalence, influence on age and sex, and relationship with metabolic syndrome and insulin resistance. *Int J Gerontol.* 2013;7(4):194–198. doi:10.1016/j.ijge.2013.03.008.
69. Lonardo A, Nascimbeni F, Ballestri S, Fairweather D, Win S, Than TA, Abdelmalek MF, Suzuki A. Sex differences in nonalcoholic fatty liver disease: state of the art and identification of research gaps. *Hepatology.* 2019;70(4):1457–1469. doi:10.1002/hep.30626.
70. Näyhä S. Geographical variations in cardiovascular mortality in Finland, 1961–1985. *Scandinavian Journal of Social Medicine. Supplementum.* 1989;40:1–48.
71. Kerminen S, Havulinna AS, Hellenthal G, Martin AR, Sarin A-P, Perola M, Palotie A, Salomaa V, Daly MJ, Ripatti S, et al. Fine-scale genetic structure in Finland. *G3: Genes|Genomes|Genetics.* 2017;7(10):3459–3468. doi:10.1534/g3.117.300217.
72. Reccia I, Kumar J, Akladios C, Viridis F, Pai M, Habib N, Spalding D. Non-alcoholic fatty liver disease: a sign of systemic disease. *Metabolism.* 2017;72:94–108.
73. Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T, Männistö S, Salomaa V, Sundvall J, Puska P. Forty-year trends in cardiovascular risk factors in Finland. *Eur J Public Health.* 2015;25(3):539–546. doi:10.1093/eurpub/cku174.
74. Marotz L, Schwartz T, Thompson L, Humphrey G, Gogul G, Gaffney J, Amir A, Knight R Earth microbiome project (EMP) high throughput (HTP) DNA extraction protocol v1 (protocols.io.pdmd46) [Internet]. 2018 [cited 2020 Nov 10]; Available from: <https://www.protocols.io/view/earth-microbiome-project-emp-high-throughput-htp-d-pdmd46>
75. Sanders JG, Nurk S, Salido RA, Minich J, Xu ZZ, Zhu Q, Martino C, Fedarko M, Arthur TD, Chen F, et al. Optimizing sequencing protocols for leaderboard metagenomics by combining long and short reads. *Genome Biol.* 2019;20(1):226. doi:10.1186/s13059-019-1834-9.
76. Glenn TC, Nilsen RA, Kieran TJ, Sanders JG, Bayona-Vásquez NJ, Finger JW, Pierson TW, Bentley KE, Hoffberg SL, Louha S, et al. Adapterama I: universal stubs and primers for 384 unique dual-indexed or 147,456 combinatorially-indexed Illumina libraries (iTru & iNext). *PeerJ.* 2019;7:e7755. doi:10.7717/peerj.7755.
77. Didion JP, Martin M, Collins FS. Atropos: specific, sensitive, and speedy trimming of sequencing reads. *PeerJ.* 2017;5:e3720. doi:10.7717/peerj.3720.
78. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 2012;9(4):357–359. doi:10.1038/nmeth.1923.
79. Méric G, Wick RR, Watts SC, Holt KE, Inouye M. Correcting index databases improves metagenomic studies. *bioRxiv.* 2019;712166. Available from <https://www.biorxiv.org/content/10.1101/712166v1>. doi:10.1101/712166.
80. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. A complete domain-to-species taxonomy for bacteria and archaea. *Nat Biotechnol.* 2020;38(9):1079–86. doi:10.1038/s41587-020-0501-8.
81. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Research.* 2016 [cited 2018 May 12];26(12):1721–1729. doi:10.1101/gr.210641.116.
82. R Core Team. R: a language and environment for statistical computing [Internet]. Vienna, Austria: R Foundation for Statistical Computing; 2018 [cited 2019 Mar 4]. Available from: <https://www.R-project.org/>
83. McMurdie PJ, Holmes S, Watson M. phyloseq: an R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS ONE.* 2013;8(4):e61217. doi:10.1371/journal.pone.0061217.
84. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are compositional: and this is not optional. *Front Microbiol.* 2017;8:2224. doi:10.3389/fmicb.2017.02224.

85. Louca S, Polz MF, Mazel F, Albright MBN, Huber JA, O'Connor MI, Ackermann M, Hahn AS, Srivastava DS, Crowe SA, et al. Function and functional redundancy in microbial systems. *Nat Ecol Evol.* 2018;2(6):936–943. doi:10.1038/s41559-018-0519-1.
86. Ruuskanen MO, St. Pierre KA, St. Louis VL, Aris-Brosou S, Poulain AJ. Physicochemical drivers of microbial community structure in sediments of lake hazen, Nunavut, Canada. *Front Microbiol* cited 2018 Jun 6]; 9. Available from. 2018;9:1138. doi:10.3389/fmicb.2018.01138
87. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin PR, O'Hara RB, Simpson GL, Solymos P, et al. *vegan*: community ecology package [Internet]. 2018 [cited 2018 Jun 4]. Available from: <https://CRAN.R-project.org/package=vegan>
88. Lahti L, Shetty S. *microbiome* R package [Internet]. 2019 [cited 2020 Dec 14]. Available from: <http://microbiome.github.io>
89. Fernandes AD, Reid JN, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB. Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome.* 2014;2(1):15. doi:10.1186/2049-2618-2-15.
90. Chen T, Guestrin C. XGBoost: a scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16; San Francisco (CA); 2016;785–794.* <https://dl.acm.org/doi/proceedings/10.1145/2939672>.
91. Bommert A, Sun X, Bischl B, Rahnenführer J, Lang M. Benchmark for filter methods for feature selection in high-dimensional classification data. *Comput Stat Data Anal.* 2020;143:106839. doi:10.1016/j.csda.2019.106839.
92. Bischl B, Lang M, Kotthoff L, Schiffner J, Richter J, Studerus E, Casalicchio G, Jones ZM. *mlr*: machine learning in R. *J Mach Learn Res.* 2016;17:1–5.
93. Bischl B, Richter J, Bossek J, Horn D, Thomas J, Lang M. *mlrMBO*: a modular framework for model-based optimization of expensive black-box functions. arXiv:170303373 [stat] [Internet] 2018 [cited 2020 Feb 18]; Available from: <http://arxiv.org/abs/1703.03373>
94. Saito T, Rehmsmeier M. *Precrec*: fast and accurate precision–recall and ROC curve calculations in R. *Bioinformatics.* 2017;33(1):145–147. doi:10.1093/bioinformatics/btw570.
95. Greenwell BM. *pdp*: an R package for constructing partial dependence plots. *R J.* 2017;9(1):421–436. doi:10.32614/RJ-2017-016.
96. Yu G, Smith DK, Zhu H, Guan Y, Lam TTY, McInerny G. *ggtree* : an r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol Evol.* 2017;8(1):28–36. doi:10.1111/2041-210X.12628.
97. Douglas GM, Maffei VJ, Zaneveld J, Yurgel SN, Brown JR, Taylor CM, Huttenhower C, Langille MGL. *PICRUSt2* for prediction of metagenome functions. *Nat Biotechnol.* 2020;38(6):685–688. doi:10.1038/s41587-020-0548-6
98. Seemann T. *Prokka*: rapid prokaryotic genome annotation. *Bioinformatics.* 2014;30(14):2068–2069. doi:10.1093/bioinformatics/btu153.
99. Belcour A, Frioux C, Aite M, Bretaudeau A, Hildebrand F, Siegel A. *Metage2Metabo*, microbiota-scale metabolic complementarity for the identification of key species. *eLife.* 2020;9:e61968. doi:10.7554/eLife.61968
100. Karp PD, Billington R, Caspi R, Fulcher A, Latendresse M, Kothari A, Keseler IM, Krummenacker M, Midford PE, Ong Q, et al. The *BioCyc* collection of microbial genomes and metabolic pathways. *Brief Bioinform.* 2019;20(4):bbz104. doi:10.1093/bib/bbx085.
101. Caspi R, Billington R, Fulcher CA, Keseler IM, Kothari A, Krummenacker M, Latendresse M, Midford PE, Ong Q, Ong WK, et al. The *MetaCyc* database of metabolic pathways and enzymes. *Nucleic Acids Res.* 2018;46(D1):D633–9. doi:10.1093/nar/gkx935.