

1 **Combined effects of host genetics and diet on human gut microbiota and**
2 **incident disease in a single population cohort**

3
4 Youwen Qin^{1,2}, Aki S. Havulinna³, Yang Liu^{1,4}, Pekka Jousilahti³, Scott C. Ritchie^{1,5-7}, Alex
5 Tokolyi⁸, Jon G. Sanders^{9,10}, Liisa Valsta³, Marta Brożyna¹, Qiyun Zhu¹¹, Anupriya Tripathi^{11,12},
6 Yoshiki Vazquez-Baeza^{13,14}, Rohit Loomba¹⁵, Susan Cheng¹⁶, Mohit Jain^{11,13}, Teemu Niiranen^{3,17},
7 Leo Lahti¹⁸, Rob Knight^{11,13,14}, Veikko Salomaa³, Michael Inouye^{1,2,5-7,19-21*}§, Guillaume Méric^{1,22*}§
8

9 ¹Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia;
10 ²School of BioSciences, The University of Melbourne, Melbourne, Victoria, Australia; ³Department of Public
11 Health Solutions, Finnish Institute for Health and Welfare, Helsinki, Finland; ⁴Department of Clinical Pathology,
12 The University of Melbourne, Melbourne, Victoria, Australia; ⁵Cambridge Baker Systems Genomics Initiative,
13 Department of Public Health and Primary Care, University of Cambridge, UK; ⁶British Heart Foundation Centre of
14 Research Excellence, University of Cambridge, UK; ⁷National Institute for Health Research Cambridge Biomedical
15 Research Centre, University of Cambridge and Cambridge University Hospitals, Cambridge, UK; ⁸Wellcome
16 Sanger Institute, Wellcome Genome Campus, Hinxton, UK; ⁹Department of Ecology and Evolutionary Biology,
17 Cornell University, Ithaca, NY, USA; ¹⁰Cornell Institute for Host-Microbe Interaction and Disease, Cornell
18 University, Ithaca, NY, USA; ¹¹Department of Pediatrics, School of Medicine, University of California San Diego,
19 La Jolla, CA, USA; ¹²Division of Biological Sciences, University of California San Diego, La Jolla, California,
20 USA; ¹³Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, USA; ¹⁴Department
21 of Computer Science & Engineering, Jacobs School of Engineering, University of California San Diego, La Jolla,
22 CA, USA; ¹⁵NAFLD Research Center, Department of Medicine, University of California San Diego, La Jolla, CA,
23 USA; ¹⁶Smidt Heart Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA; ¹⁷Department of Medicine,
24 Turku University Hospital and University of Turku, Turku, Finland; ¹⁸Department of Future Technologies,
25 University of Turku, Turku, Finland; ¹⁹British Heart Foundation Cardiovascular Epidemiology Unit, Department of
26 Public Health and Primary Care, University of Cambridge, UK; ²⁰Health Data Research UK Cambridge, Wellcome
27 Genome Campus & University of Cambridge, UK; ²¹The Alan Turing Institute, London, UK; ²²Department of
28 Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria, Australia.
29

30 § These authors contributed equally

31 *Corresponding authors: Michael Inouye: mi336@medschl.cam.ac.uk; Guillaume Méric:
32 guillaume.meric@baker.edu.au.
33

34 **Human genetic variation affects the gut microbiota through a complex combination of**
35 **environmental and host factors. Here, we characterize genetic variations associated with**
36 **microbial abundances in a single large-scale population-based cohort of 5,959 genotyped**
37 **individuals with matched gut microbial metagenomes, dietary and health records**
38 **(prevalent and follow-up). We identified 567 independent SNP-taxon associations.**
39 **Variants at the *LCT* locus associated with *Bifidobacterium* and other taxa, but they differed**
40 **according to dairy intake. Furthermore, levels of *Faecalicatena lactaris* associated with**
41 ***ABO*, and suggested preferential utilization of secreted blood antigens as energy source in**
42 **the gut. *Enterococcus faecalis* levels associated with variants in the *MED13L* locus, which**
43 **has been linked to colorectal cancer. Mendelian randomization analysis indicated a**
44 **potential causal effect of *Morganella* on major depressive disorder, consistent with**
45 **observational incident disease analysis. Overall, we identify and characterize the intricate**
46 **nature of host-microbiota interactions and their association with disease.**
47

48 Humans have co-evolved with the microbial communities that colonize them, resulting in a
49 complex assembly of thousands of microbial species mutualistically living in their
50 gastrointestinal tract. A fine-tuned interplay between microbial and human physiologies can
51 impact multiple aspects of development and health to the point that dysbiosis is often associated
52 with disease^{1,2}. As such, increasing evidence points to the influence of human genetic variation
53 on the composition and modulation of their gut microbiota.

54
55 Past genetic studies have collectively revealed important host-microbe interactions^{3–13}. Previous
56 twin studies detected significant heritability signal from the presence and abundance of only a
57 few microbial taxa, such as some *Firmicutes*¹⁴, suggesting a strong transientness and variability
58 in gut microbial composition, as well as an important influence from external factors^{5,14–17}.
59 Nonetheless, a well-described association between *Bifidobacterium* levels and *LCT-MCM6*,
60 governing the phenotype of lactase persistence throughout adulthood in Europeans, was
61 uncovered in 2015³ and subsequently replicated by later studies^{5,6,8–11}, suggesting a very strong
62 influence of the evolution of dairy diet in modern humans on their gut bacteria. Additionally,
63 genes involved in immune and metabolic processes⁸ but also disease¹⁸ were also associated with
64 gut microbial variation. Despite several promising findings, reproducibility across studies
65 varying in sampling and methods is generally poor, and most previously reported associations
66 lose significance after multiple testing correction¹⁹. The individual gut microbiota is largely
67 influenced by environmental variables, mostly diet and medication^{20–22}, which could explain a
68 larger proportion of microbiome variance than identifiable host genetic factors^{8,9}. Biological
69 factors could also influence the cross-study reproducibility of results. GWAS would typically
70 not reproducibly identify genetic associations with taxa harboring microbial functions
71 potentially shared by multiple unrelated species^{23,24}. Indeed, a certain degree of functional
72 redundancy has been observed in human gut microbial communities²⁴, which is believed to play
73 a role in the resistance and resilience to perturbations^{25–27}. However, both assembly and
74 functioning in human gut microbial communities seem to be driven by the presence of a few
75 particular and identifiable keystone taxa²⁸, which exert key ecological and modulatory roles on
76 gut microbial composition independently of their abundance^{29,30}. Such taxa are relatively
77 prevalent across individuals and thought to be part of the human “core” microbiota^{29,30}, which
78 makes them potentially identifiable through GWAS.

79
80 Increasing sample size in studied populations could yield novel and robustly associated results,
81 and alleviate the effect of confounding technical or biological factors. This could be achieved
82 either by performing meta-analyses of GWAS conducted in various populations¹¹, or by using
83 larger cohort datasets. In this study, we used a large single homogenous population cohort with
84 matching human genotypes and shotgun fecal metagenomes (N=5959; FINRISK 2002 (FR02))
85 to identify novel genome-wide associations between human genotypes and gut microbial
86 abundances (**Extended Data Figure 1**). We further leveraged additional and extensive health
87 registry and dietary individual data to investigate the effects of diet and genotype on particular
88 host-microbial associations, and to predict incident disease linked to gut microbial variation.

89 **Results**

90 **Genome-wide association analysis of gut microbial taxa**

91
92
93
94 Genome-wide association tests were applied to 2,801 microbial taxa and 7,967,866 human
95 genetic variants from 5,959 individuals enrolled in the FR02 cohort (**Supplementary Table 9**),
96 which includes all taxa discovered to be prevalent in >25% of the cohort (**Methods**). Using a
97 genome-wide significance threshold ($p < 5.0 \times 10^{-8}$), a total of 471 distinct Genome Taxonomy
98 Database (GTDB) taxa, which represented 17% of all tested taxa and included 11 phyla, 19

99 classes, 24 orders, 62 families, 146 genera and 209 species, were found to be associated with at
100 least one genetic variant (**Figure 1, Supplementary Table 1**). Conditional analysis found 567
101 independent SNP-taxon associations at genome-wide significance in 411 loci (**Supplementary**
102 **Table 1**). Heritability across the 2,801 taxa ranged between $h^2=0.001$ to 0.214, with the highest
103 values observed for taxa belonging to the *Firmicutes* and *Firmicutes_A* GTDB phyla, both of
104 which encompassed half (237/471, 50.3%) of all associated taxa with genetic variation
105 (**Extended Data Figure 2**). There was no difference in SNP heritability between groups of
106 associated or non-associated taxa at genome-wide significance ($p=0.23$). Adjusting for antibiotic
107 prescription did not change any study-wide significant associations, and only 32 out of 567
108 genome-wide associations moved slightly above $p=5\times 10^{-8}$, which is likely by chance given
109 inclusion of any additional covariate (**Supplementary Table 10**). After adjustment, beta
110 estimates were highly correlated (Pearson $r>0.999$).

111
112 Three loci were strongly associated with microbial variation at study-wide significance, as
113 shown on a Manhattan plot showing the lowest resulting p-value for each SNP tested against
114 each of the 2,801 taxa (**Figure 1, Supplementary Table 1**). There was no evidence of excess
115 false positive rate in the GWAS (median $\lambda_{GC}=1.0051$) (**Figure 1B, Supplementary Table 9**).
116 After conditional analysis, the strongest association by far ($p=5.0\times 10^{-35}$) involved members of
117 class *Actinobacteria* and rs3940549, a variant in the *LCT-MCM6-ZRANB3* locus region which
118 is in high LD ($r^2=0.87$) with the well-described *LCT* variant rs4988235 causing lactase
119 persistence in adults of European ancestry (**Extended Data Figure 3**). In total, 29 taxa were
120 associated with the *LCT-MCM6* region, including 18 below study-wide significance (**Figure 1,**
121 **Supplementary Table 1**). These involved *Bifidobacterium*-related *Actinobacteriota* and three
122 taxa from the GTDB *Firmicutes_A* phylum which included 2 uncultured species defined from
123 metagenome-assembled reference genomes (*UBA3855 sp900316885* and *CAG-81*
124 *sp000435795*). The association of these three *Firmicutes_A* with *LCT* was still genome-wide
125 significant after adjusting for *Bifidobacterium* abundances, as were 11 other taxa associated with
126 the *LCT-MCM6* region (**Supplementary Table 2**). Additionally, the abundance of these *LCT-*
127 *MCM6*-associated taxa were not, or very weakly associated with the *Bifidobacterium*
128 abundances. A variant in *ABO* (rs545971), expressing the histo-blood group ABO system
129 transferase, was strongly associated ($p=1.1\times 10^{-12}$) with levels of *Faecalicatena lactaris*. There
130 was evidence for a second independent signal at *ABO* associated with the *Collinsella* genus
131 (chr9:133271182; $p=2.5\times 10^{-8}$). Rs187309577 and rs143507801 in *MED13L*, expressing the
132 Mediator complex subunit 13L, were found to be associated with genus *Enterococcus*
133 ($p=1.8\times 10^{-12}$) and the *Enterococcus faecalis* species ($p=7.26\times 10^{-11}$), respectively.

134
135 Details on the replication of previously-reported microbiome GWAS signals in our study are
136 included in **Supplementary Note**.

137 138 **Gut microbial keystone taxa associate with genetic variation**

139
140 In total, we identified 31 distinct genetic variants associated ($p<5.0\times 10^{-8}$) with 39 microbial taxa
141 related to identified keystone species as listed by Banerjee *et al.* (2018)^{28,31}, which included the
142 *Actinobacteria* class, *Helicobacter pylori*, *Bacteroides stercoris*, *Bacteroides thetaiotaomicron*,
143 *Ruminococcus bromii*, *Klebsiella pneumoniae*, *Proteus mirabilis*, *Akkermansia muciniphila*, and
144 the archaeon *Methanobrevibacter smithii* (**Figure 1c, Supplementary Table 1**). Keystone
145 species are defined as members of a microbial community exerting selective modulation and not
146 broad effects on microbiome composition variation. Only one documented keystone species
147 from Banerjee *et al.*, *Bacteroides fragilis*, was not associated with genetic variation in our
148 study²⁸. Although a lot of computationally identified keystone species remain to be
149 experimentally verified^{32,33}, this observation suggests that they would generally associate with

150 human genetic variation. This would indicate an intimate association with the human gut niche
151 in line with their reported key ecological roles in microbiome modulation and functioning. Our
152 work highlights novel human genotypes associating with keystone taxa (**Supplementary Table**
153 **1**), which could further improve our understanding of their ecology.

154

155 **Combined effect of genetics and diet on *LCT*-associated taxa**

156

157 We compared the abundances of 4 bacterial taxa strongly associated with the *LCT* locus
158 (*Bifidobacterium* genus, *Negativibacillus* genus, *UBA3855 sp900316885* and *CAG-81*
159 *sp000435795*) in individuals with different rs4988235 genotypes and dairy diets (**Figure 2a**).
160 The abundance of *Bifidobacterium* in individuals producing lactase through adulthood
161 (rs4988235:TT) was unaffected by dairy intake. However, lactose-intolerant individuals
162 (rs4988235:CC) self-reporting a regular dairy diet had a significant increase in *Bifidobacterium*
163 abundance ($p=1.75\times 10^{-13}$; Wilcoxon-rank test). An intermediate genotype (rs4988235:CT) was
164 linked to an intermediate increase in *Bifidobacterium* abundance (**Figure 2a**). This trend did not
165 seem to be affected by age³⁴ (**Extended Data Figure 4**). Additionally, we observed a moderate
166 negative correlation between *Bifidobacterium* abundances and age in rs4988235:CC individuals
167 reporting a regular dairy diet (Spearman's $\rho=-0.17$, $p=1.9\times 10^{-6}$) and in rs4988235:CC
168 individuals reporting a low or lactose free diet ($\rho=-0.19$, $p=0.002$). Furthermore, the Spearman
169 correlation between the *Bifidobacterium* residual abundance and dairy diet was still significant
170 ($\rho=-0.22$, $p=2\times 10^{-12}$) in rs4988235:CC individuals. This indicated that the associations with age
171 were consistent in individuals with and without regular dairy intake, and did not confound the
172 association between *Bifidobacterium* and dairy diet.

173

174 An inverse pattern was observed for the abundance distributions of *Negativibacillus* and
175 uncultured *CAG-81 sp000435795*, for which abundances decreased in lactose-intolerant
176 individuals reporting dairy intake, as compared to rs4988235:TT individuals consuming dairy
177 products (**Figure 2a**). Levels of *UBA3855 sp900316885* were unaffected by a dairy diet in
178 lactose-intolerant individuals but were surprisingly lower in rs4988235:TT individuals who
179 reported dairy intake ($p=8.23\times 10^{-5}$). These opposite and contrasting effects of dairy on
180 associated bacterial abundances in lactose-intolerant individuals could reflect competition for
181 lactose in the gut. *CAG-81* abundances were the most negatively correlated with those of the
182 other *LCT*-associated taxa (**Extended Data Figure 5**), which suggests that this competition
183 could be strong and prevalent enough to drive co-association at the *LCT* locus, possibly mediated
184 by lactose intake (**Figure 2b**).

185

186 **Functional profiling of CAZymes in 11 *Bifidobacterium sp.***

187

188 Of all 11 *Bifidobacterium* species prevalent enough in our study population to be included in the
189 GWAS, only *B. dentium* was not associated with the *LCT* locus ($p=1.70\times 10^{-2}$), nor was it co-
190 abundant with any other *Bifidobacterium* species (**Extended Data Figure 6a**). *B. dentium* has
191 previously been suggested to have different metabolic abilities³⁵. A clustering of carbohydrate-
192 active enzymes (CAZyme) profiles from reference genomes of all 11 *Bifidobacterium* species
193 revealed that *B. dentium* clustered apart from the 10 other species, which grouped consistently
194 with their co-abundance patterns (**Extended Data Figure 6b**). *B. dentium* harbored more genes
195 encoding CAZyme families with preferred fiber/plant-related substrates (GH94, GH26, GH53)
196 than other *Bifidobacterium* species, which seemed to harbor more milk oligosaccharide-targeting
197 CAZyme families (GH129, GH112) than *B. dentium* (**Extended Data Figure 6b**), which could
198 relate to the observed association differences. This suggests that bacterial metabolic abilities can
199 be strong drivers of co-abundance, and of association with human genetic variation.

200

201 **Impacts of genotype and fiber intake on *ABO*-associated taxa**

202
203 A variety of bacteria metabolize blood antigens, with potential applications in synthetic universal
204 donor blood production^{36,37}. Gut bacteria are particularly exposed to A- and B-antigens in the
205 gut mucosa of secretor individuals³⁸. Our associations of *Faecalicatena lactaris* ($p=1.10\times 10^{-12}$)
206 and *Collinsella* ($p=2.59\times 10^{-8}$) with *ABO* suggest a possible metabolic link with blood antigens.
207 A comparison of CAZyme profiles across a set of reference genomes revealed 3 CAZymes with
208 blood-related activities in *F. lactaris* (GH110³⁹, GH136⁴⁰, CBM32⁴¹), but none in any of 9
209 *Collinsella* species (**Figure 3**). More mucus-targeting and less fiber-degrading enzymes were
210 found in *F. lactaris* than *Collinsella*, suggesting distinct functions in the gut.

211
212 As previously reported⁴, neither *ABO* blood types, nor secretor status had an impact on alpha
213 and beta diversity (**Extended Data Figure 7a**). However, we observed that the effect of *ABO*
214 genotypes on *F. lactaris* levels, underlying the association, were largely driven by secretor status,
215 with increased abundances in secretor individuals from genotype groups rs545971:CT and
216 rs545971:TT, A and AB blood type groups, but not in rs545971:CC genotype, or B and O blood
217 types individuals (**Figure 4a**). Levels in non-secretors did not vary across *ABO* genotypes or
218 blood types. Despite a slight increase in blood type A secretors, *Collinsella* only remained
219 minimally affected by secretor status or blood group (**Extended Data Figure 7b**). Taken
220 together, this suggests that the secretion of soluble A and B-antigens strongly affects *F. lactaris*
221 in the gut, possibly through reduced opportunity to use them as substrate. Both levels of *F.*
222 *lactaris* and *Collinsella* were significantly higher when individuals were predicted to secrete A-
223 , B- and AB-antigens in their gut mucosa (**Extended Data Figure 7c**).

224
225 A high fiber diet is thought to induce a metabolic transition from mucus-degrading to fiber-
226 degrading activities in the colon, as carbohydrates from fiber are more easily metabolized⁴². The
227 increase in *F. lactaris* abundances in A/B/AB-secretors (defined as secreting A-, B- and AB-
228 antigens) compared to non- A/B/AB-secretors remained strongly significant irrespective of fiber
229 intake ($p=1.15\times 10^{-9}$ in the low-fiber diet group, and $p=4.4\times 10^{-3}$ in the high-fiber diet group),
230 suggesting that either *F. lactaris* has a strong affinity for secreted A/B/AB-antigens, does not
231 efficiently degrade dietary fiber, or will not easily switch to it as an energy source (**Figure 4b**).
232 *F. lactaris* levels were increased in non-A/B/AB-secretors with a high fiber diet compared to a
233 low fiber diet, implying a switch to fiber degradation or interaction with fiber-degrading bacteria
234 (**Figure 4b**). *Collinsella* variation in both A/B/AB-secretors and non-A/B/AB-secretors with
235 high- and low-fiber diets was similar to the compounded abundances of 13 major mucin-
236 degrading species in the human gut⁴³, suggesting a similar ecological response in stark contrast
237 with *F. lactaris* (**Figure 4b, Figure 4c**).

238 239 ***MED13L*-associated *E. faecalis* as a putative link with CRC**

240
241 The allele frequency of the *MED13L* rs143507801 variant (A>G), associated with levels of
242 *Enterococcus faecalis* ($p=7.26\times 10^{-11}$), was low (MAF=0.0111), consistent with reported allele
243 frequencies in the gnomAD database⁴⁴. In our study population, 131 individuals carried
244 rs143507801:G allele, 130 being heterozygous (GA) and only one being homozygous (GG). We
245 observed that *E. faecalis* levels were increased in heterozygous rs143507801:GA individuals
246 (**Figure 5**). *E. faecalis* is a gut commensal, but also an opportunist pathogen believed to play a
247 role in colorectal cancer (CRC) development, possibly through direct damaging of colorectal
248 cells⁴⁵⁻⁴⁷. *MED13L* and *MED13* encode for Mediator transcriptional coactivator complex
249 modules associating with RNA polymerase II⁴⁸, and as such specifically interact with cyclin-
250 dependent kinase 8 (CDK8) modules described for their oncogenic activation of transcription
251 during colon tumorigenesis⁴⁹. Consequently, we observed slightly higher levels of *E. faecalis*

252 ($p=0.014$) in 14 individuals enrolled in FR02 with a history of CRC at the time of sampling
253 (**Figure 5**). Groups of individuals segregated by allelic variant and CRC status could not be
254 compared robustly due to small sample size. Taken together, these results suggest a possible link
255 between *E. faecalis* and CRC through the MED13 activation of CDK8 in colorectal tumors,
256 which will need to be investigated further.

257 **MR highlights possible causal effect of *Morganella* on MDD**

258 Interpreting results of causal inference prediction using bacterial information entails to particular
259 caution, due to the possibility of multiple and unaccounted confounding factors¹⁰, but can be
260 useful to highlight potential focus for future research. Here, we predicted 96 causal effects in
261 both microbe to disease and disease to microbe directions using bidirectional Mendelian
262 Randomization (MR). Of these, 34 were from microbial levels as exposure to disease as
263 outcome, with a large proportion of causal effects in psychiatric and neurological diseases
264 (**Supplementary Table 5**). For example, MR suggested an increased abundance of
265 *Faecalicoccus* may have a causal effect on anorexia nervosa (OR=1.8 per SD increase in
266 bacterial abundance; CI_{95%}=1.3-2.5; $p=2.0\times 10^{-4}$, MR method IVW) (**Methods**). Other examples
267 included increasing abundances of *Morganella* and *Raoultella* predicted to have causal effects
268 on major depressive disorder (MDD) (**Supplementary Table 5**). When MR was performed in
269 the reverse direction, using disease risk as an exposure and microbial level as an outcome, most
270 predicted causal effects involved autoimmune and inflammatory diseases but the strongest
271 predicted causal effect involved type 2 diabetes (T2D) (**Supplementary Table 6**). Doubling the
272 genetic risk of T2D (possibly accompanied by external factors such as hypoglycaemic
273 medications or metformin intake) was predicted to reduce levels of the uncultured *CAG-345*
274 *sp000433315* species (*Firmicutes* phylum) by 0.14 SD (SE=0.04, $p=3.0\times 10^{-4}$, MR method
275 IVW). A few other examples included some degree of literature validation, such as the higher
276 genetic risk for primary sclerosing cholangitis (PSC) causally impacting levels of the cholesterol-
277 reducing *Eubacterium_R coprostanoligenes*⁵⁰. Furthermore, a higher genetic risk for coeliac
278 disease (CD) was predicted to increase abundances in 4 species previously reported to be more
279 abundant in CD patients than controls⁵¹ (**Supplementary Table 6**). Finally, a higher genetic risk
280 for multiple sclerosis (MS) was predicted to cause a reduction in the abundance of
281 *Lactobacillus_B ruminis*, consistent with the report that *Lactobacillus sp.* can reduce symptom
282 severity in an animal model of MS⁵².

283 The availability in our study dataset of up to 16 years of electronic health record follow-up after
284 the initial sampling of the microbiota allowed for observational validation of predicted effects
285 using MR. Of all causal predictions identified using MR, only the effect of *Morganella* on MDD
286 could be validated by a statistically significant association with incident MDD (n=181 cases;
287 HR=1.11, CI_{95%}=1.01-1.22, per SD increase of bacterial abundance), after accounting for age, sex
288 and BMI (**Figure 6**). In our GWAS, *Morganella* variation in the study population associated
289 with a variant (rs192436108; $p=6.16\times 10^{-8}$) in the *PDE1A* locus, which has previously been
290 linked to depression^{53,54} and psychiatric disorders⁵⁵. We did not find that the development of
291 MDD could be linked to an abnormal incidence of microbiome-related diseases (**Table S8**).
292 Taken together, these predicted links between *Morganella* and MDD suggest more efforts should
293 be deployed into exploring the possible roles of this bacterium as part of the brain-gut axis
294 metabolic modulation of health.

295 **Discussion**

296 Through GWAS and the subsequent investigation of functional and ecological factors
297 contributing to the most robust human-microbe associations, we present a diverse and global
298
299
300
301
302

303 picture of human-microbe interactions in a single cohort of ~6,000 European individuals. We
304 find 3 genetic loci to be strongly associated with gut microbial variation. Two of these loci, *LCT*
305 and *ABO*, are well-known and very segregated in human populations, possibly explaining why
306 our homogenous European cohort identified them as being associated so strongly. A third more
307 mysterious association with the *MED13L* locus highlights possible links with cancer while
308 causal inference highlights several diseases as being causally linked to gut microbes.

309
310 Lactase persistence, or the continued ability to digest lactose into adulthood, is the most strongly
311 selected single-gene trait over the last 10,000 years in multiple human populations⁵⁶, believed to
312 have spread amongst humans with the advent of animal domestication and the culturally
313 transmitted practice of dairying⁵⁷. In our study, as in previous work^{3,5,6,10}, the association of *LCT*
314 variants with *Actinobacteria*, more specifically *Bifidobacterium*, is by far the most statistically
315 significant, suggesting a profound interaction between *Actinobacteria* and the human gut, in line
316 with their reported keystone activities²⁹. We reported a strong increase of *Bifidobacterium* levels
317 in genetically lactose-intolerant people reporting a regular consumption of dairy products⁸. This
318 increase was not confounded by age in adults, despite *Bifidobacterium* levels generally
319 decreasing with age in our cohort. While self-reported dietary information is not entirely reliable
320 due to various reasons^{58,59}, our study population was large and the differences were significant
321 enough to consider this a robust observation, which can be explained by the evolutionary
322 adaptation of *Bifidobacterium* to specifically metabolize human and bovine milk
323 oligosaccharides⁶⁰. In lactase-deficient adults, consumed lactose is likely to become available
324 for colonic bacteria as an energy source to compete for. Hints of a possible competitive
325 relationship between *Bifidobacterium* and *Negativibacillus* were revealed, which could depend
326 upon lactose intake and should be investigated in functional studies.

327
328 Two considerations stem from our findings. First, the genetic determinants of lactose intolerance
329 are known to vary across ethnicity⁶¹ and cross-population heterogeneity in the *LCT*-
330 *Bifidobacterium* association was recently reported¹¹. As more non-European-centric genetic
331 studies are conducted worldwide^{11,62,63}, examining this combined interaction between dairy diet
332 and *Bifidobacterium* in different genetic backgrounds could bring new insights. Secondly,
333 despite recent progresses, lactose intolerance is still largely underdiagnosed, and genetic
334 prediction rates from large population studies exceed lactose intolerance prevalence rates
335 obtained using physical tests⁶¹. In our study, we lacked information on lactose malabsorption
336 symptoms in lactose intolerant individuals reporting a regular dairy diet. Lactose-free (<0.01%
337 lactose content) or low-lactose (<0.1%) dairy products have been available in Finland since 1978
338 and are popular among people experiencing symptoms of lactose malabsorption. Our data did
339 not allow us to make the distinction between lactose-intolerant individuals aware of their
340 symptoms and consuming low-lactose products as a result, from intolerant individuals unaware
341 of the cause of their symptoms while consuming dairy. The latter would either experience
342 discomfort symptoms without knowingly implicating their lactose intake, or the ability of higher
343 concentration of *Bifidobacterium* to degrade lactose in their intestines may alleviate the
344 perceived symptoms of discomfort associated with lactose intolerance, therefore encouraging
345 individuals to continue consuming indigestible lactose asymptotically⁶⁴. This possible
346 probiotic effect should be investigated in controlled studies.

347
348 The *ABO* gene expresses a glycosyltransferase in many cell types, which determines the ABO
349 blood group of an individual by modifying the oligosaccharides on cell surface glycoproteins. A
350 comparison of humans and non-human primates has identified *ABO* (along with the MHC) as
351 harboring ancient multiallelic polymorphisms that are maintained across species^{65,66}. Many
352 infectious diseases such as norovirus infection, bacterial meningitis, malaria, cholera⁶⁷, or even
353 more recently SARS-CoV-2^{68,69} are associated with host blood type and secretor status⁶⁷,

354 suggesting that infection could be a driver of a strong balancing selection that has maintained
355 *ABO* polymorphisms. Furthermore, blood type variation has been linked to various chronic
356 diseases⁶⁷, such as heart and vascular diseases, gastric cancers, diabetes, asthma or even
357 dementia⁶⁷. As many of these chronic diseases are also associated with dysbiosis of the gut
358 microbiota, this prompts an interesting but largely unexplored parallel between gut commensals,
359 blood types and disease³⁸. Our study confirms previous findings⁴ that secretor status or blood
360 types do not seem to globally affect gut microbial alpha- or beta-diversity. It also confirms
361 reports from two very recent studies: first, a meta-analysis across five German cohorts, using
362 16S rRNA sequencing to characterize the gut microbiota, linked *Bacteroides* and
363 *Faecalibacterium* to *ABO* and *FUT2*⁷⁰. The second study functionally associated bacterial
364 lactose and galactose degradation genes to *ABO* variation in a cohort of 3,432 Chinese
365 individuals⁷¹. Taken together, these findings suggest a broad association of *ABO* polymorphisms
366 with microbial variation in various human populations.

367
368 An important research effort aiming to enzymatically produce synthetic universal donor blood
369 has driven a push for screening a large diversity of CAZymes, including bacteria, revealing
370 substrate affinities for blood antigens across various microbes^{36,37}. Here we highlight *F. lactaris*
371 (formerly *Ruminococcus lactaris*), as a mucin-degrading commensal likely able to digest blood
372 antigens through its predicted GH110, GH136 and CBM32 CAZyme family genes³⁹⁻⁴¹. *F.*
373 *lactaris* is strongly associated with *ABO* genetic variation in our European cohort, and is
374 differentially abundant in people according to their predicted gut mucosal secretion of A/B/AB-
375 antigens. Interestingly, our findings are not consistent with *F. lactaris* switching to a fiber-
376 degrading activity in individuals reporting a high fiber diet, unlike other mucin-degrading
377 bacteria in our study and in the literature⁴² and *Collinsella*, another *ABO*-associated taxon. Our
378 work suggests that some gut commensals such as *F. lactaris* appear to be very efficient and
379 adapted metaboliser of A/B/AB-antigens in the gut, despite their predicted ability to degrade
380 simpler carbohydrates in fiber. This could be an example of ecological niche differentiation in
381 the gut, with impacts on associated *F. lactaris* microbial communities, of which *Collinsella*, also
382 associated with *ABO*, may belong.

383
384 Although validation of the association is inconclusive because of the low prevalence of CRC
385 cases and genetic variation in our study population, the association of *MEDI3L* rs143507801
386 variant with *Enterococcus faecalis* suggested a putative link with CRC. It has been shown that
387 *MEDI3* could directly link a cyclin-dependent kinase 8 (CDK8) module to Mediator^{72,73}, which
388 is a colorectal cancer oncogene, amplified in colorectal tumors and activating transcription
389 driving colon tumorigenesis leading to CRC⁴⁹. This could explain a long suspected link between
390 *E. faecalis* and development of CRC after having been found in higher concentrations in CRC
391 patients than healthy individuals^{46,74}. The suspected mode of action of *E. faecalis* on CRC
392 development is currently unclear, but could be linked to extracellular free radical production
393 directly leading to DNA break, point mutation and chromosomal instability in colorectal cells⁴⁷.
394 Although we saw a trend of *E. faecalis* being increased in abundance in individuals with a history
395 of CRC, and in *MEDI3L* variation, more focused work including incident CRC and a larger
396 sample size will be required to precisely pinpoint a link between this bacterium and CRC through
397 the Mediator complex, if any.

398
399 Besides suggesting a link between gut microbes and autoimmune and inflammatory diseases, in
400 line with previous studies⁷⁵, causal inference analysis highlighted a very particular and promising
401 example of interplay between a gut microbe and a complex disease. Among other suggested
402 links with psychiatric diseases, we predicted increasing abundances of *Morganella* and
403 *Klebsiella* (ex-*Raoultella*^{76,77}) to have causal effects on MDD. Members of the
404 *Enterobacteriaceae* family, such as these two genera, have previously been found in higher

405 levels in MDD patients⁷⁸. Although caution is required when interpreting predictions of
406 causality⁷⁹, increasing evidence suggests that gut microbes are likely to influence host behavior
407 via a systemic modulation of hormones and metabolites along the gut-brain axis^{80–82}. Importantly,
408 our MR-based result was consistent with observed hazards using follow-up observational data
409 up to 16 years after sampling. This observation supports previous experimental results showing
410 an increase of IgM and IgA-related immune response against *Morganella* secreted
411 lipopolysaccharide in major depression⁸³. A recent retrospective cohort study performed on 311
412 individuals including 156 MDD cases highlighted bacterial functions, metabolites and species
413 involved in the interaction between the gut microbiome and MDD⁸⁴. Although *Morganella* was
414 not specifically highlighted, levels of several other *Enterobacteriales* species were found to
415 significantly differ between MDD patients and healthy controls⁸⁴. Taken together, our findings
416 highlight the intimate influence of the gut-brain axis on humans, with more mechanistic studies
417 required to untangle and further interpret these predictions.
418

419 Our study highlights the benefits of increasing sample size to increase the statistical power for
420 discovery. Although the *LCT* locus has been reported multiple times to be associated with
421 bacterial taxa, our work is the first to report study-wide significant associations in a single cohort,
422 at the strongest significance ever reported. The association with *Bifidobacterium* in our study
423 was even stronger than the recent findings that used integrative data from 18,473 individuals in
424 28 different cohorts¹¹, emphasizing the importance of standardized methodology and
425 homogeneity in participant ethnicity (especially when studying geographically distributed traits
426 like lactose intolerance⁸⁵). *ABO* allelic variation is also notoriously affected by geography⁸⁶,
427 which could explain why some meta-analyses in non-homogenous populations could miss it.
428 Also, metagenomic sequencing with standardized, robust taxonomic definitions^{87,88} can provide
429 species-level characterization of microbial profiles in the gut of individuals, unlike 16S rRNA-
430 based studies. An example from our work is the observation that *Bifidobacterium dentium* was
431 prevalent but not associated with the *LCT* locus like all other *Bifidobacterium* species in the
432 population. Observed difference in CAZymes commonly found in other *Bifidobacterium* species
433 may explain this difference³⁵. This should be confirmed in future experiments using more
434 deeply-sequenced metagenomes unambiguously linking function to particular MAGs.
435 Furthermore, GTDB taxonomic standardization results in greater taxon granularity, i.e. smaller,
436 more discrete clades of similar phylogenetic depth than commonly known lineages or species^{87,88}.
437 In theory, this would increase overall accuracy⁸⁹, as a weak association with a poorly-defined
438 lineage may be caused by a strong association with a well-defined subset of that lineage, defined
439 as a coherent group using GTDB⁸⁸. Finally, a myriad of microbial taxa that are now solely
440 defined and represented by uncultured metagenome-assembled genomes (MAGs) in the GTDB
441 database were found to be independently associated with various loci. Along with recent reports
442 that the more gut microbiome diversity is explored, the more novel, unknown species are
443 discovered^{90,91}, this suggests that many discoveries are yet to be made in the field of human
444 microbiome studies.
445

446 **Acknowledgements**

447

448 The study protocol of FINRISK 2002 was approved by the Coordinating Ethical Committee of
449 the Helsinki and Uusimaa Hospital District (Ref. 558/E3/2001). All participants signed an
450 informed consent. The study was conducted according to the World Medical Association
451 Declaration of Helsinki on ethical principles. All necessary patient/participant consent has been
452 obtained and the appropriate institutional forms have been archived.
453

454 We thank all participants of the FINRISK 2002 survey for their contributions to this work. The
455 FINRISK surveys are mainly funded by budgetary funds from the Finnish Institute for Health

456 and Welfare with additional funding from several domestic foundations. YQ was partially
457 supported by The Albert Shimmins Fund (Faculty of Science Postgraduate Writing-Up Award
458 2020). MI was supported by the Munz Chair of Cardiovascular Prediction and Prevention. VS
459 was supported by the Finnish Foundation for Cardiovascular Research. LL was supported by
460 Academy of Finland (decision 295741). TN was supported by the Emil Aaltonen Foundation,
461 the Finnish Medical Foundation, the Paavo Nurmi Foundation, and the Academy of Finland
462 (grant 321351). ASH was supported by the Academy of Finland, grant no. 321356. RL receives
463 funding support from NIEHS (5P42ES010337), NCATS (5UL1TR001442), NIDDK
464 (U01DK061734, R01DK106419, P30DK120515, R01DK121378, R01DK124318), and DOD
465 PRCRP (W81XWH-18-2-0026). This study was supported by the Victorian Government's
466 Operational Infrastructure Support (OIS) program, and by core funding from: the UK Medical
467 Research Council (MR/L003120/1), the British Heart Foundation (RG/13/13/30194;
468 RG/18/13/33946) and the National Institute for Health Research [Cambridge Biomedical
469 Research Centre at the Cambridge University Hospitals NHS Foundation Trust] [*]. This work
470 was supported by Health Data Research UK, which is funded by the UK Medical Research
471 Council, Engineering and Physical Sciences Research Council, Economic and Social Research
472 Council, Department of Health and Social Care (England), Chief Scientist Office of the Scottish
473 Government Health and Social Care Directorates, Health and Social Care Research and
474 Development Division (Welsh Government), Public Health Agency (Northern Ireland), British
475 Heart Foundation and Wellcome.

476
477 *The views expressed are those of the authors and not necessarily those of the NHS, the NIHR
478 or the Department of Health and Social Care.

479

480 **Author contributions**

481
482 YQ, MI, VS and GM designed the work; ASH, PJ, JGS, LV, MB, QZ, ATr, YV-B, TN, LL, RK,
483 VS and GM acquired the data; YQ, YL, SCR, JGS, LL, ATo and GM analysed the data; RL, SC,
484 MJ, TN, LL, RK, VS, MI and GM supervised the work. All authors wrote the manuscript and
485 gave final approval of the version to be published.

486

487 **Competing interests**

488
489 VS has consulted for Novo Nordisk and Sanofi and received honoraria from these companies.
490 He also has ongoing research collaboration with Bayer AG, all unrelated to this study. RL serves
491 as a consultant or advisory board member for Anylam/Regeneron, Arrowhead Pharmaceuticals,
492 AstraZeneca, Bird Rock Bio, Boehringer Ingelheim, Bristol-Myer Squibb, Celgene, Cirius,
493 CohBar, Conatus, Eli Lilly, Galmed, Gemphire, Gilead, Glympse bio, GNI, GRI Bio, Inipharm,
494 Intercept, Ionis, Janssen Inc., Merck, Metacrine, Inc., NGM Biopharmaceuticals, Novartis, Novo
495 Nordisk, Pfizer, Prometheus, Promethera, Sanofi, Siemens and Viking Therapeutics. In addition,
496 his institution has received grant support from Allergan, Boehringer-Ingelheim, Bristol-Myers
497 Squibb, Cirius, Eli Lilly and Company, Galectin Therapeutics, Galmed Pharmaceuticals, GE,
498 Genfit, Gilead, Intercept, Grail, Janssen, Madrigal Pharmaceuticals, Merck, NGM
499 Biopharmaceuticals, NuSirt, Pfizer, pH Pharma, Prometheus, and Siemens. He is also co-founder
500 of Liponexus, Inc. The remaining authors declare no competing interests.

501

502 **References (for main text only)**

503

504 1. Belizário, J. E. & Napolitano, M. Human microbiomes and their roles in dysbiosis,
505 common diseases, and novel therapeutic approaches. *Front. Microbiol.* 6, (2015).

- 506 2. Levy, M., Kolodziejczyk, A. A., Thaïss, C. A. & Elinav, E. Dysbiosis and the immune
507 system. *Nat Rev Immunol* 17, 219–232 (2017).
- 508 3. Blekhnman, R. et al. Host genetic variation impacts microbiome composition across
509 human body sites. *Genome Biol* 16, 191 (2015).
- 510 4. Davenport, E. R. et al. ABO antigen and secretor statuses are not associated with gut
511 microbiota composition in 1,500 twins. *BMC Genomics* 17, 941 (2016).
- 512 5. Goodrich, J. K. et al. Genetic Determinants of the Gut Microbiome in UK Twins. *Cell*
513 *Host & Microbe* 19, 731–743 (2016).
- 514 6. Bonder, M. J. et al. The effect of host genetics on the gut microbiome. *Nat Genet* 48,
515 1407–1412 (2016).
- 516 7. Turpin, W. et al. Association of host genome with intestinal microbial composition in a
517 large healthy cohort. *Nat Genet* 48, 1413–1417 (2016).
- 518 8. Wang, J. et al. Genome-wide association analysis identifies variation in vitamin D
519 receptor and other host factors influencing the gut microbiota. *Nat Genet* 48, 1396–1406
520 (2016).
- 521 9. Rothschild, D. et al. Environment dominates over host genetics in shaping human gut
522 microbiota. *Nature* 555, 210–215 (2018).
- 523 10. Hughes, D. A. et al. Genome-wide associations of human gut microbiome variation and
524 implications for causal inference analyses. *Nat Microbiol* 5, 1079–1087 (2020).
- 525 11. Kurilshikov, A., Medina-Gomez, C., Bacigalupe, R. *et al.* Large-scale association
526 analyses identify host factors influencing human gut microbiome composition. *Nat Genet* **53**,
527 156–165 (2021).
- 528 12. Kolde, R. et al. Host genetic variation and its microbiome interactions within the
529 Human Microbiome Project. *Genome Med* 10, 6 (2018).
- 530 13. Rühlemann, M. C. et al. Application of the distance-based F test in an mGWAS
531 investigating β diversity of intestinal microbiota identifies variants in SLC9A8 (NHE8) and 3
532 other loci. *Gut Microbes* 9, 68–75 (2018).
- 533 14. Goodrich, J. K. et al. Human Genetics Shape the Gut Microbiome. *Cell* 159, 789–799
534 (2014).
- 535 15. Xie, H. et al. Shotgun Metagenomics of 250 Adult Twins Reveals Genetic and
536 Environmental Impacts on the Gut Microbiome. *Cell Systems* 3, 572–584.e3 (2016).
- 537 16. Lim, M. Y. et al. The effect of heritability and host genetics on the gut microbiota and
538 metabolic syndrome. *Gut* 66, 1031–1038 (2017).
- 539 17. Le Roy, C. I. et al. Heritable components of the human fecal microbiome are associated
540 with visceral fat. *Gut Microbes* 9, 61–67 (2018).
- 541 18. Goodrich, J. K., Davenport, E. R., Clark, A. G. & Ley, R. E. The Relationship Between
542 the Human Genome and Microbiome Comes into View. *Annu. Rev. Genet.* 51, 413–433
543 (2017).
- 544 19. Kurilshikov, A., Wijmenga, C., Fu, J. & Zhernakova, A. Host Genetics and Gut
545 Microbiome: Challenges and Perspectives. *Trends in Immunology* 38, 633–647 (2017).
- 546 20. David, L. A. et al. Diet rapidly and reproducibly alters the human gut microbiome.
547 *Nature* 505, 559–563 (2014).
- 548 21. Falony, G. et al. Population-level analysis of gut microbiome variation. *Science* 352,
549 560–564 (2016).
- 550 22. Zhernakova, A. et al. Population-based metagenomics analysis reveals markers for gut
551 microbiome composition and diversity. *Science* 352, 565–569 (2016).
- 552 23. Eng, A. & Borenstein, E. Taxa-function robustness in microbial communities.
553 *Microbiome* 6, 45 (2018).
- 554 24. Ferrer, M. et al. Microbiota from the distal guts of lean and obese adolescents exhibit
555 partial functional redundancy besides clear differences in community structure: Metaproteomic
556 insights associated to human obesity. *Environ Microbiol* 15, 211–226 (2013).

- 557 25. Moya, A. & Ferrer, M. Functional Redundancy-Induced Stability of Gut Microbiota
558 Subjected to Disturbance. *Trends in Microbiology* 24, 402–413 (2016).
- 559 26. Louca, S. et al. Function and functional redundancy in microbial systems. *Nat Ecol*
560 *Evol* 2, 936–943 (2018).
- 561 27. Louca, S. et al. High taxonomic variability despite stable functional structure across
562 microbial communities. *Nat Ecol Evol* 1, 0015 (2017).
- 563 28. Banerjee, S., Schlaeppi, K. & van der Heijden, M. G. A. Keystone taxa as drivers of
564 microbiome structure and functioning. *Nat Rev Microbiol* 16, 567–576 (2018).
- 565 29. Trosvik, P. & de Muinck, E. J. Ecology of bacteria in the human gastrointestinal tract—
566 identification of keystone and foundation taxa. *Microbiome* 3, 44 (2015).
- 567 30. Shetty, S. A., Hugenholtz, F., Lahti, L., Smidt, H. & de Vos, W. M. Intestinal
568 microbiome landscaping: insight in community assemblage and implications for microbial
569 modulation strategies. *FEMS Microbiol. Rev.* 41, 182–199 (2017).
- 570 31. Chia, L. W. et al. Deciphering the trophic interaction between *Akkermansia*
571 *muciniphila* and the butyrogenic gut commensal *Anaerostipes caccae* using a
572 metatranscriptomic approach. *Antonie van Leeuwenhoek* 111, 859–873 (2018).
- 573 32. Banerjee, S., Schlaeppi, K. & van der Heijden, M. G. A. Reply to ‘Can we predict
574 microbial keystones?’ *Nat Rev Microbiol* 17, 194 (2019).
- 575 33. Röttjers, L. & Faust, K. Can we predict keystones? *Nat Rev Microbiol* 17, 193 (2019).
- 576 34. Kato, K. et al. Age-Related Changes in the Composition of Gut *Bifidobacterium*
577 *Species*. *Curr Microbiol* 74, 987–995 (2017).
- 578 35. Engevik, M. A. et al. *Bifidobacterium dentium* Fortifies the Intestinal Mucus Layer via
579 Autophagy and Calcium Signaling Pathways. *mBio* 10, e01087-19, /mbio/10/3/mBio.01087-
580 19.atom (2019).
- 581 36. Rahfeld, P. & Withers, S. G. Toward universal donor blood: Enzymatic conversion of
582 A and B to O type. *J. Biol. Chem.* 295, 325–334 (2020).
- 583 37. Liu, Q. P. et al. Bacterial glycosidases for the production of universal red blood cells.
584 *Nat Biotechnol* 25, 454–464 (2007).
- 585 38. Arnolds, K. L., Martin, C. G. & Lozupone, C. A. Blood type and the microbiome-
586 untangling a complex relationship with lessons from pathogens. *Current Opinion in*
587 *Microbiology* 56, 59–66 (2020).
- 588 39. Liu, Q. P. et al. Identification of a GH110 Subfamily of α 1,3-Galactosidases: NOVEL
589 ENZYMES FOR REMOVAL OF THE α 3GAL XENOTRANSPLANTATION ANTIGEN. *J.*
590 *Biol. Chem.* 283, 8545–8554 (2008).
- 591 40. Pichler, M. J. et al. Butyrate producing colonic Clostridiales metabolise human milk
592 oligosaccharides and cross feed on mucin via conserved pathways. *Nat Commun* 11, 3285
593 (2020).
- 594 41. Ficko-Blean, E. & Boraston, A. B. The Interaction of a Carbohydrate-binding Module
595 from a *Clostridium perfringens* N -Acetyl- β -hexosaminidase with Its Carbohydrate Receptor. *J.*
596 *Biol. Chem.* 281, 37748–37757 (2006).
- 597 42. Desai, M. S. et al. A Dietary Fiber-Deprived Gut Microbiota Degrades the Colonic
598 Mucus Barrier and Enhances Pathogen Susceptibility. *Cell* 167, 1339-1353.e21 (2016).
- 599 43. Tailford, L. E., Crost, E. H., Kavanaugh, D. & Juge, N. Mucin glycan foraging in the
600 human gut microbiome. *Front. Genet.* 6, (2015).
- 601 44. Genome Aggregation Database Consortium et al. The mutational constraint spectrum
602 quantified from variation in 141,456 humans. *Nature* 581, 434–443 (2020).
- 603 45. Amarnani, R. & Rapose, A. Colon cancer and enterococcus bacteremia co-affection: A
604 dangerous alliance. *Journal of Infection and Public Health* 10, 681–684 (2017).
- 605 46. Khan, Z., Siddiqui, N. & Saif, M. W. *Enterococcus Faecalis* Infective Endocarditis and
606 *Colorectal Carcinoma*: Case of New Association Gaining Ground. *Gastroenterol Res* 11, 238–
607 240 (2018).

- 608 47. Huycke, M. M., Abrams, V. & Moore, D. R. Enterococcus faecalis produces
609 extracellular superoxide and hydrogen peroxide that damages colonic epithelial cell DNA.
610 Carcinogenesis 23, 529–536 (2002).
- 611 48. Allen, B. L. & Taatjes, D. J. The Mediator complex: a central integrator of
612 transcription. Nat Rev Mol Cell Biol 16, 155–166 (2015).
- 613 49. Firestein, R. et al. CDK8 is a colorectal cancer oncogene that regulates β -catenin
614 activity. Nature 455, 547–551 (2008).
- 615 50. Li, L., Batt, S. M., Wannemuehler, M., Dispirito, A. & Beitz, D. C. Effect of feeding of
616 a cholesterol-reducing bacterium, Eubacterium coprostanoligenes, to germ-free mice. Lab.
617 Anim. Sci. 48, 253–255 (1998).
- 618 51. Marasco, G. et al. Gut Microbiota and Celiac Disease. Dig Dis Sci 61, 1461–1472
619 (2016).
- 620 52. Lavasani, S. et al. A Novel Probiotic Mixture Exerts a Therapeutic Effect on
621 Experimental Autoimmune Encephalomyelitis Mediated by IL-10 Producing Regulatory T
622 Cells. PLoS ONE 5, e9009 (2010).
- 623 53. Tomita, H. et al. G protein-linked signaling pathways in bipolar and major depressive
624 disorders. Front. Genet. 4, (2013).
- 625 54. Wong, M.-L. et al. Phosphodiesterase genes are associated with susceptibility to major
626 depression and antidepressant treatment response. Proceedings of the National Academy of
627 Sciences 103, 15124–15129 (2006).
- 628 55. Schork, A. J. et al. A genome-wide association study of shared risk across psychiatric
629 disorders implicates gene regulation during fetal neurodevelopment. Nat Neurosci 22, 353–361
630 (2019).
- 631 56. Burger, J. et al. Low Prevalence of Lactase Persistence in Bronze Age Europe Indicates
632 Ongoing Strong Selection over the Last 3,000 Years. Current Biology S0960982220311878
633 (2020) doi:10.1016/j.cub.2020.08.033.
- 634 57. Gerbault, P. et al. Evolution of lactase persistence: an example of human niche
635 construction. Phil. Trans. R. Soc. B 366, 863–877 (2011).
- 636 58. Hebert, J. R. et al. Social Desirability Trait Influences on Self-Reported Dietary
637 Measures among Diverse Participants in a Multicenter Multiple Risk Factor Trial. The Journal
638 of Nutrition 138, 226S–234S (2008).
- 639 59. Schoeller, D. A. How Accurate Is Self-Reported Dietary Energy Intake? Nutrition
640 Reviews 48, 373–379 (2009).
- 641 60. Sakanaka, M. et al. Evolutionary adaptation in fucosyllactose uptake systems supports
642 bifidobacteria-infant symbiosis. Sci. Adv. 5, eaaw7696 (2019).
- 643 61. Storhaug, C. L., Fosse, S. K. & Fadnes, L. T. Country, regional, and global estimates
644 for lactose malabsorption in adults: a systematic review and meta-analysis. The Lancet
645 Gastroenterology & Hepatology 2, 738–746 (2017).
- 646 62. Liu, X., Tang, S., Zhong, H. et al. A genome-wide association study for gut
647 metagenome in Chinese adults illuminates complex diseases. Cell Discov 7, 9 (2021).
- 648 63. Sirugo, G., Williams, S. M. & Tishkoff, S. A. The Missing Diversity in Human Genetic
649 Studies. Cell 177, 26–31 (2019).
- 650 64. Szilagy, A. Adaptation to Lactose in Lactase Non Persistent People: Effects on
651 Intolerance and the Relationship between Dairy Food Consumption and Evaluation of Diseases.
652 Nutrients 7, 6751–6779 (2015).
- 653 65. Ségurel, L., Gao, Z. & Przeworski, M. Ancestry runs deeper than blood: The
654 evolutionary history of ABO points to cryptic variation of functional importance: Insights &
655 Perspective. BioEssays n/a-n/a (2013) doi:10.1002/bies.201300030.
- 656 66. Segurel, L. et al. The ABO blood group is a trans-species polymorphism in primates.
657 Proceedings of the National Academy of Sciences 109, 18493–18498 (2012).

658 67. Ewald, D. R. & Sumner, S. C. J. Blood type biochemistry and human disease: Blood
659 type biochemistry and human disease. *WIREs Syst Biol Med* 8, 517–535 (2016).

660 68. Ellinghaus, D. et al. Genomewide Association Study of Severe Covid-19 with
661 Respiratory Failure. *N Engl J Med* NEJMoa2020283 (2020) doi:10.1056/NEJMoa2020283.

662 69. Shelton, J.F., Shastri, A.J., Ye, C. et al. Trans-ancestry analysis reveals genetic and
663 nongenetic associations with COVID-19 susceptibility and severity. *Nat Genet* 53, 801–808
664 (2021). doi:10.1101/2020.09.04.20188318.

665 70. Rühlemann, M.C., Hermes, B.M., Bang, C. et al. Genome-wide association study in
666 8,956 German individuals identifies influence of ABO histo-blood groups on gut microbiome.
667 *Nat Genet* 53, 147–155 (2021).

668 71. Liu, X. et al. Inter-determination of blood metabolite levels and gut microbiome
669 supported by Mendelian randomization.
670 <http://biorxiv.org/lookup/doi/10.1101/2020.06.30.181438> (2020)
671 doi:10.1101/2020.06.30.181438.

672 72. Knuesel, M. T., Meyer, K. D., Bernecky, C. & Taatjes, D. J. The human CDK8
673 subcomplex is a molecular switch that controls Mediator coactivator function. *Genes &
674 Development* 23, 439–451 (2009).

675 73. Tsai, K.-L. et al. A conserved Mediator–CDK8 kinase module association regulates
676 Mediator–RNA polymerase II interaction. *Nat Struct Mol Biol* 20, 611–619 (2013).

677 74. De Almeida, C. et al. Differential Responses of Colorectal Cancer Cell Lines to
678 *Enterococcus faecalis*’ Strains Isolated from Healthy Donors and Colorectal Cancer Patients.
679 *JCM* 8, 388 (2019).

680 75. Marchesi, J. R. et al. The gut microbiota and host health: a new clinical frontier. *Gut*
681 65, 330–339 (2016).

682 76. Ma, Y. et al. Proposal for reunification of the genus *Raoultella* with the genus
683 *Klebsiella* and reclassification of *Raoultella electrica* as *Klebsiella electrica* comb. nov. *Res
684 Microbiol* 103851 (2021) doi:10.1016/j.resmic.2021.103851.

685 77. Wyres, K. L., Lam, M. M. C. & Holt, K. E. Population genomics of *Klebsiella*
686 *pneumoniae*. *Nat Rev Microbiol* 18, 344–359 (2020).

687 78. Jiang, H. et al. Altered fecal microbiota composition in patients with major depressive
688 disorder. *Brain, Behavior, and Immunity* 48, 186–194 (2015).

689 79. Wade, K. H. & Hall, L. J. Improving causality in microbiome research: can human
690 genetic epidemiology help? *Wellcome Open Res* 4, 199 (2020).

691 80. Foster, J. A. & McVey Neufeld, K.-A. Gut–brain axis: how the microbiome influences
692 anxiety and depression. *Trends in Neurosciences* 36, 305–312 (2013).

693 81. Fung, T. C., Olson, C. A. & Hsiao, E. Y. Interactions between the microbiota, immune
694 and nervous systems in health and disease. *Nat Neurosci* 20, 145–155 (2017).

695 82. Valles-Colomer, M. et al. The neuroactive potential of the human gut microbiota in
696 quality of life and depression. *Nat Microbiol* 4, 623–632 (2019).

697 83. Maes, M., Kubera, M. & Leunis, J.-C. The gut-brain barrier in major depression:
698 intestinal mucosal dysfunction with an increased translocation of LPS from gram negative
699 enterobacteria (leaky gut) plays a role in the inflammatory pathophysiology of depression.
700 *Neuro Endocrinol. Lett.* 29, 117–124 (2008).

701 84. Yang, J. et al. Landscapes of bacterial and metabolic signatures and their interaction in
702 major depressive disorders. *Sci Adv* 6, (2020).

703 85. Mattar, R., Mazo, & Carrilho. Lactose intolerance: diagnosis, genetic, and clinical
704 factors. *CEG* 113 (2012) doi:10.2147/CEG.S32368.

705 86. Bodmer, W. Genetic Characterization of Human Populations: From ABO to a Genetic
706 Map of the British People. *Genetics* 199, 267–279 (2015).

707 87. Parks, D. H. et al. A standardized bacterial taxonomy based on genome phylogeny
708 substantially revises the tree of life. *Nat Biotechnol* 36, 996–1004 (2018).

- 709 88. Parks, D. H. et al. A complete domain-to-species taxonomy for Bacteria and Archaea.
710 Nat Biotechnol (2020) doi:10.1038/s41587-020-0501-8.
- 711 89. Méric, G., Wick, R. R., Watts, S. C., Holt, K. E. & Inouye, M. Correcting index
712 databases improves metagenomic studies. <http://biorxiv.org/lookup/doi/10.1101/712166> (2019)
713 doi:10.1101/712166.
- 714 90. Pasolli, E. et al. Extensive Unexplored Human Microbiome Diversity Revealed by
715 Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. Cell
716 176, 649-662.e20 (2019).
- 717 91. Almeida, A. et al. A unified catalog of 204,938 reference genomes from the human gut
718 microbiome. Nat Biotechnol (2020) doi:10.1038/s41587-020-0603-3.

719

720 **Figure Legends (for main text only)**

721

722 **Figure 1. Genome-wide association of human genetic and gut microbial variations.** (a)
723 Manhattan plot aggregating the top associations with microbial variation. Each SNP was tested
724 against each of the 2,801 taxa and the Manhattan plot shows the lowest resulting p-value for each
725 SNP. Loci with associations above study-wide significance level ($p < 3.8 \times 10^{-11}$; red dashed line)
726 are annotated with the human locus name and the corresponding associated microbial taxa. The
727 blue dashed line denotes genome-wide significance level ($p < 5 \times 10^{-8}$). Of all genome-wide
728 significant associations shown on the Manhattan plot, 320/567 (56.4%) involved 265 lead SNPs
729 with MAF between 1% and 5%, and 247/567 (43.6%) involved 185 lead SNPs with MAF >5%.
730 P-values denote significance of joint analysis model using GTCA-COJO. (b) The distribution of
731 genomic inflation factor (λ_{GC}) in 2,801 tested taxa [median(λ_{GC})=1.0051; mean(λ_{GC})=1.0059].
732 (c) Tree-based visualization of the taxonomic diversity of genome-wide associated microbial
733 taxa. The central root of the tree represents the Bacteria domain, the first connected node
734 represents phylum, the second connected node class, the third order and the fourth family. Every
735 node represents at least one associated taxa in the GWAS at genome-wide significance level.
736 The three smaller trees on the right highlight all taxonomic groups containing at least one taxon
737 identified as associated with the *LCT-MCM6*, *ABO*, and *MED13L* loci (blue edges and nodes
738 denote taxa associated at study-wide significance level and purple edges and nodes denote taxa
739 associated at genome-wide significance level). The main tree is annotated to indicate phyla
740 harbouring >10 distinct genome-wide associated taxa, as well as previously described keystone
741 taxa.

742

743 **Figure 2. Interaction of human genotype, dairy diet and gut bacterial variation with the**
744 ***LCT* locus.** (a) The 4 panels present variation in microbial relative abundances (not CLR-
745 transformed) for the 4 taxa associated at study-wide significance level with the *LCT* locus at
746 $p < 3.8 \times 10^{-11}$: *Bifidobacterium*, *Negativibacillus*, *UBA3855 sp900316885* and *CAG-81*
747 *sp000435795*. Abundances are compared across stratified groups of individuals from the FR02
748 cohort according to *LCT-MCM6*:rs4988235 genotype and self-reported dietary lactose intake
749 (red: regular dairy diet; blue: lactose-free diet). Sample sizes for groups of individuals self-
750 reporting a regular dairy diet: rs4988235:TT (n=1,786), CT (n=2,413), CC (n=736); self-
751 reporting a non-regular dairy diet or lactose-free diet: TT (n=150), CT (n=198), CC (n=245). All
752 statistical comparisons denote the p-values of Wilcoxon rank test on the distributions of
753 untransformed relative abundances. Only significantly different comparisons ($p < 0.05$) are
754 indicated. For all box plots, the central line, box and whiskers represent the median, interquartile
755 range (IQR) and 1.5 times the IQR, respectively. Violin plots represent the distribution density
756 of the data points; (b) Host genetics and gut microbes interact in the context of dairy intake and
757 lactose intolerance.

758

759 **Figure 3. Functional profiling of reference genomes from two bacterial taxa associated with**
760 **the *ABO* locus** CAZyme distribution patterns in *F. lactaris* and *Collinsella* reference genomes
761 (from the GTDB release 89 index used to classify metagenomes in this study). The heatmap
762 indicates species abundance in corresponding CAZyme families, corresponding to the total count
763 of detected families for each species divided by the number of reference genomes examined for
764 the same species. Values <1 (white to light blue) indicate that less than one copy per genome of
765 the corresponding CAZyme family was detected for each species, values >1 (light blue to dark
766 blue) indicate that more than one copy per genome was detected. Preferred substrate groups are
767 based on literature search and descriptions on CAZypedia.org.
768

769 **Figure 4. Effect of host genetics and dietary fiber intake on gut abundance variation of two**
770 **bacterial taxa associated with the *ABO* locus.** (a) *ABO*-associated *F. lactaris* relative
771 abundances (not CLR-transformed) are compared across stratified groups of individuals from
772 the FR02 cohort according to (left panel): *ABO*:rs4988235 genotype and predicted secretor status
773 (blue: secretor status conferred by *FUT2* rs601338:AG/AA genotype; red: non-secretor status
774 conferred by *FUT2* rs601338:GG genotype) and (right panel) according to predicted A, AB, B
775 and O blood types, and predicted secretor status. Sample sizes for compared groups: secretor
776 status with rs545971:C/C (n=1,538), C/T (n=2,493), T/T (n=1,050) and blood group A
777 (n=2,178), AB (n=460), B (n=900), O (n=1,543); non-secretor status with rs545971:C/C
778 (n=266), C/T (n=437), T/T (n=175) and blood group A (n=383), AB (n=80), B (n=148), O
779 (n=267). (b) *ABO*-associated *F. lactaris* and *Collinsella sp.* relative abundances, as well as
780 compounded abundances from 13 mucin-degrading species from Tailford *et al.* (2015), are
781 compared across stratified groups of individuals from the FR02 cohort according to the predicted
782 A/B/AB-antigen secretion status and dietary fiber intake. Secretion status was defined to
783 segregate individuals. A/B/AB-antigen secretors were defined as secretor individuals from blood
784 types A, AB and B. Non-A/B/AB-antigen secretors were defined as non-secretor individuals and
785 O-antigen secretors. Fiber intake was compared in individual groups from the top and bottom
786 quartiles of total fiber score (see **Methods**). Sample sizes for compared groups of individuals:
787 A/B/AB-antigen secretors (n=1393) following a low-fiber diet (n=723) or a fiber-rich diet
788 (n=670), or non- A/B/AB-antigen secretors (n=952) following a low-fiber diet (n=490) or a
789 fiber-rich diet (n=462). All statistical comparisons denote the p-values of Wilcoxon rank test on
790 the distributions of untransformed relative abundances. For all box plots (B and C), the central
791 line, box and whiskers represent the median, interquartile range (IQR) and 1.5 times the IQR,
792 respectively. Violin plots represent the distribution density of the data points. (c) Host genetics
793 and gut microbes interact in the context of fiber intake, secretor status and blood types.
794

795 **Figure 5. Effect of host genetics and prevalent colorectal cancer on gut levels of**
796 ***Enterococcus faecalis* associated with *MEDI3L* variation across participants of the FR02**
797 **cohort.** Abundances are compared across individuals grouped according to (left panel):
798 *MEDI3L*:rs143507801 genotype, (right panel): colorectal cancer (CRC) prevalence according
799 to the Finnish Cancer Registry. The comparison between *E. faecalis* variation and
800 *MEDI3L*:rs143507801 reflects the GWAS results (**Supplementary Table 1**). The comparison
801 of *E. faecalis* abundances in individuals with or without a history of CRC at the time of sampling
802 was performed using a Wilcoxon rank test. Sample sizes for compared groups of individuals:
803 rs143507801:A/A (n=5,825), G/A (n=130) (Note: only 1/5959 individual in our cohort was
804 G/G); with CRC (n=14), without a history of CRC at baseline (n=5,941). For all box plots, the
805 central line, box and whiskers represent the median, interquartile range (IQR) and 1.5 times the
806 IQR, respectively. Violin plots represent the distribution density of the data points.
807

808 **Figure 6. MR-based causal effects and incident depression analysis link *Morganella* with**
809 **major depressive disorder.** Forest plot (in blue) representing the magnitude of the effect on

810 MDD risk per 1-s.d. increase in bacterial abundance. MR analysis was carried out with 28 genetic
811 instruments and their effect sizes from FR02 (5,959 samples) and MR-base summary statistics
812 (173,005 samples). In red is shown the hazard ratio for incident MDD in the FR02 cohort up to
813 16 years after baseline sampling, using Cox model (see Methods). Error bars represent the 95%
814 confidence intervals.

815 **Methods**

816 **Study population**

817
818
819
820 The FINRISK study has been extensively described elsewhere⁹². FINRISK population surveys
821 have been performed every 5 years since 1972 to monitor trends of cardiovascular and other non-
822 communicable disease risk factors in the Finnish population^{92,93}. The study population of this
823 study consists in the participants of FINRISK 2002 (FR02) study, including men and women
824 aged between 25 and 74 years from six geographical areas of Finland^{92,94,95}. The sampling was
825 stratified by sex, region and 10-year age group so that each stratum had 250 participants. The
826 overall participation rate was 65.5% (n = 8,798). Participants filled out a questionnaire, then
827 participated in a clinical examination carried out by specifically trained nurses and gave a blood
828 sample from which various laboratory measurements were performed. They also received a
829 sampling kit and instructions to donate a stool sample at home and mailed it to the Finnish
830 Institute for Health and Welfare in an overnight mail. The survey was done in Finland during
831 winter months (January to March 2002), with average temperatures well below 0°C. Special care
832 was additionally taken to ensure that samples did not remain sitting in a post office more than
833 >1 day, or over the weekend. Upon reception in THL, samples were immediately frozen to -
834 20°C and kept unthawed until shipped to the University of California San Diego (USA), where
835 they were processed and sequenced. The use of antibiotics was recorded from participants in the
836 FINRISK 2002 questionnaire and by linking with prescription registry. In addition, participants
837 in each study site were asked whether they have an acute infection and were, as a general rule,
838 asked to reschedule their examinations and stool sampling if they had.

839
840 The follow-up of the cohort took place by record linkage of the study data with the Finnish
841 national electronic health registers (Hospital Discharge Register and Causes of Death Register),
842 which provide in practice 100% coverage of relevant health events in Finnish residents. For
843 present analyses involving follow-up data, we used a follow-up which extended until
844 31/12/2018.

845
846 The study protocol of FR02 was approved by the Coordinating Ethical Committee of the Helsinki
847 and Uusimaa Hospital District (Ref. 558/E3/2001). All participants signed an informed consent.
848 The study was conducted according to the World Medical Association's Declaration of Helsinki
849 on ethical principles.

850 **Cohort phenotype metadata and specific dietary information**

851
852
853 The phenotype data in this study comprised of demographic characteristics, life habits, disease
854 history, clinical measurements, laboratory test results and follow-up electronic health records
855 (EHRs). More specifically, baseline dietary factors were collected. Details of the method have
856 been described previously⁹³. To broadly assess diet information within the cohort participants, a
857 binary variable was used to indicate whether individuals were self-reporting to follow various
858 possible dietary restrictions. Dietary consumption of specific food product categories was also
859 reported. Habitual diet was assessed using a food propensity questionnaire (FPQ) which
860 contained 42 food items or groups and had choices ranging from 1-6 for consumption frequency

861 ranging from “Less than once a month” to “Once a day or more often”. The consumption
862 frequencies were converted to frequencies per month ranging from 0.5 times per month to 30,
863 45 or 60 times per month. Food items that are rarely eaten more than once a day were given the
864 value of 30 times per month. Food items that are often eaten multiple times a day such as fresh
865 vegetables, breads, etc. were given a value of 60 times per month. Food items that fall in between
866 these two groups were given 45 points.

867

868 **Self-reporting of lactose-free diet and dietary fiber consumption**

869

870 Allelic distribution at the *LCT-MCM6*:rs4988235 variant responsible for lactase persistence in
871 Europeans was as following in our study population: 1,936 (35%) individuals had the T/T allele
872 conferring a lactase persistence phenotype through adulthood, allowing them to digest lactose,
873 while 981 (18%) individuals had the C/C allele conferring lactose intolerance. Most individuals
874 (n=2,611, 47%) had the intermediate allele C/T making them likely to be able to digest lactose.
875 Most individuals reported a regular dairy intake in their diet (n=5,002, 89%), while 706 (12.5%)
876 individuals reported a regular lactose-free diet.

877

878 A total fiber consumption score was calculated from the questionnaires, reflecting the overall
879 consumption of a combination of various fiber sources such as high-fiber bread, vegetables
880 (vegetarian dishes, fresh vegetables and boiled vegetables and legumes) and fruits, berries and
881 natural juices. The resulting total fiber index values ranged from 9 (low dietary fiber intake) to
882 48 (high dietary fiber intake), with a median of 33. Comparisons of the effects of low- vs. high-
883 fiber diets were made between the 1st (n=1,213) and 4th (n=1,132) quartiles of the total fiber
884 index.

885

886 **Genotyping, imputation and quality control**

887

888 The genotyping was performed on Illumina genome-wide SNP arrays (the HumanCoreExome
889 BeadChip, the Human610-Quad BeadChip and the HumanOmniExpress) and has been described
890 previously⁹⁶. Stringent criteria were applied to remove samples and variants of low quality.
891 Samples with call rate <95%, sex discrepancies, excess heterozygosity and non-European
892 ancestry were excluded. Variants with call rate <98%, deviation from Hardy-Weinberg
893 Equilibrium ($p < 1 \times 10^{-6}$), and minor allele count < 3 were filtered. Data was pre-phased by using
894 Eagle2 v2.3⁹⁷. Imputation was performed using IMPUTE2 v2.3.0⁹⁸ with two Finnish-population-
895 specific reference panels: 2,690 high-coverage whole-genome sequencing and 5,092 whole-
896 exome sequencing samples. To evaluate the imputation quality, we compared the sample allele
897 frequencies with reference populations and examined imputation quality (INFO scores)
898 distributions. Imputed SNPs with INFO >0.7 were kept for analysis. Post imputation quality
899 control was carried out by using plink v2.0⁹⁹. Samples with >10% missing rate were removed.
900 Individuals with extreme height or BMI values were further excluded (31 individuals with
901 height < 1.47m; 5 with BMI > 50 were removed). Both genotyped and imputed SNPs were kept
902 for analysis if they met the following criteria: call rate >90%, no significant deviation from
903 Hardy-Weinberg Equilibrium ($p > 1.0 \times 10^{-6}$), and minor allele frequency >1%. SNP filtering was
904 based on all individuals for which genotype information was available (n=7280), not on the 5,959
905 individuals selected subsequently for GWAS after QC. The post-QC dataset comprised
906 7,967,866 SNPs.

907

908 **Metagenomic sequencing from stool samples**

909

910 Stool samples were collected by participants and mailed overnight to Finnish Institute for Health
911 and Welfare for storing at -20°C; the samples were sequenced at the University of California

912 San Diego in 2017. No special arrangements were taken regarding the temperature of the samples
913 when they were shipped from the field clinics to the lab in THL but, as the survey was done in
914 Finland during winter months (January to March 2002) with average temperatures well below
915 0°C. Special care was anyway additionally taken to ensure that samples did not remain sitting in
916 a post office over the weekend. The gut microbiome was characterized by shallow shotgun
917 metagenomics sequencing with Illumina HiSeq 4000 Systems. We successfully performed stool
918 shotgun sequencing in n=7,231 individuals. The detailed procedures for DNA extraction, library
919 preparation and sequence processing have been previously described⁹⁵. Adapter and host
920 sequences were removed. To preserve the quality of data while retaining most of the disease
921 cases, samples with a total number of sequenced reads lower than 400,000 were removed.

922

923 **Taxonomic profiling, quality filtering and data transformation**

924

925 Taxonomic profiling of FR02 metagenomes was performed as follows: briefly, raw shotgun
926 metagenomic sequencing reads were mapped using the *k*-mer-based metagenomic classification
927 tool Centrifuge v1.0.4¹⁰⁰ to an index database custom-built to encompass reference genomes that
928 followed the taxonomic nomenclature introduced and updated in the GTDB release 89⁸⁷⁻⁸⁹. This
929 implies that unless specified otherwise, all taxonomic names in our study refer to their
930 nomenclature in GTDB, which can be related to the original NCBI nomenclature using the
931 GTDB database server: https://gtdb.ecogenomic.org/taxon_history/. The same profiling
932 approach has also been used and described in recent studies from our consortium^{94,95,101}. Our
933 study present results involving *Faecalicatena lactaris*, which is called differently in NCBI and
934 subsequent GTDB releases. A particular note on the evolution of this nomenclature can be found
935 in the **Supplementary Note**.

936

937 Gut microbial composition was represented as the relative abundance of taxa. For each
938 metagenome at phylum, class, order, family, genus and species levels, the relative abundance of
939 a taxon was computed as the proportion of reads assigned to the clade rooted at this taxon among
940 total classified reads. The relative abundance of a taxon with no reads assigned in a metagenome
941 was considered as zero in the corresponding profile. For the purpose of this association study
942 and because of reduced accuracy and power when considering rare taxa, we focused on common
943 and relatively abundant microbial taxa, defined as prevalent in >25% studied individuals, and
944 defined with at least 10 mapped reads per individual. For the purpose of association, and as
945 previous studies have reported that only some microbial taxa are inheritable¹⁰², we also removed
946 taxa with zero SNP-heritability. This filtering resulted in a microbial dataset composed of a total
947 of 2,801 taxa, including 59 phyla, 95 classes, 187 orders, 415 families, 922 genera and 1,123
948 species.

949

950 Taxonomic profiles derived from sequencing data are by nature compositional because of an
951 arbitrary total imposed by the instrument¹⁰³. The compositional data of microbial taxa is not
952 independent and can lead to inappropriate use of linear regression. To overcome this artificial
953 bias, all relative abundance values were transformed by centre-log-ratio (CLR)¹⁰⁴. More
954 information about data transformation can be found in the **Supplementary Note**.

955

956 When visually comparing relative abundances in groups of individuals throughout the
957 manuscript, we used untransformed relative abundances, for better interpretability. Alpha
958 (Shannon index) and beta (Bray-Curtis distance) diversity were calculated at genus level used
959 functions in the R package *vegan* v2.5-6. We did not find a correlation between sequencing depth
960 and Shannon diversity index (Spearman's $\rho = -0.001598$, $p = 0.90$) in n=5,959 samples (**Extended
961 Data Figure 8**). Additionally, to define CLR-transformed abundances of higher taxonomic
962 levels than species, we summed the raw abundances of all taxa (e.g. species) belonging to a

963 specific higher taxonomic taxon (e.g. genus), and then applied a CLR transformation.
964 Additionally, we observed that Eastern and Western Finnish populations did not have different
965 microbiome diversity, despite having overall slightly different lifestyles, and mortality rates. To
966 further investigate this, we visualized potential geographical effects using a PCoA plot on beta-
967 diversity (Bray-Curtis dissimilarity) from metagenomic profiles of samples used in the GWAS
968 from our study (n=5,959; **Extended Data Figure 9**).

970 **Genome-wide association analysis**

971
972 The protocol followed in this study was described elsewhere¹⁰⁵. Briefly, linear mixed model
973 (LMM) implemented in BOLT-LMM v2.3.2¹⁰⁶ was used to search for genome-wide associations
974 accounting for the individual similarity. Since BOLT-LMM only accepts <1 million SNPs in
975 modelling the genetic relationship matrix, SNPs were pruned at the threshold of $r^2 < 0.1$ (plink2⁹⁹,
976 command *--indep-pairwise 1000 80 0.1*), resulting in 106,201 independent SNPs. This list of
977 independent SNPs was used to estimate heritability using BOLT-LMM. Additionally, BOLT-
978 LMM automatically performs leave-one-chromosome-out (LOCO) analysis to avoid proximal
979 contamination. Although LMM accounts for the cryptic relatedness in individuals, there are still
980 large population structure cannot be addressed. Thus, the top 10 genetic principal components
981 (calculated by FlashPCA v2.0¹⁰⁷ based on the pruned SNPs mentioned above) were included as
982 covariates. Age, gender, and genotyping batch were adjusted. We did not adjust for microbiome
983 sequencing batch, as we observed that it had a no effect on microbiome composition variation
984 (**Extended Data Figure 9**). As no genetic variant was reported to have large effect size on gut
985 microbiota, statistic estimates were based on infinitesimal model which assumes small non-zero
986 effect for large number of genetic variants. To identify independent associations, GCTA-COJO
987 v1.91.3¹⁰⁸ was used to conduct approximate conditional and joint analysis using individual
988 genetic data. Window size was set to 10 Mb, assuming SNPs on different chromosomes or more
989 than 10 Mb distance are uncorrelated. The resulting effect size (beta coefficient) indicated the
990 number of standard deviation changes of a taxon's CLR transformed abundance corresponding
991 to one effective allele increase of SNP. Additionally, for all but 2 reported SNPs (rs146740485
992 and rs2797225), the effect allele was the reference allele in the GWAS cohort.

993
994 As microbes interact non-independently with each other in the gut, as part of larger ecological
995 and functional communities, matSpDlite v1.0^{109,110} was used to estimate the number of
996 independent tests based on eigenvalue variance, the larger the eigenvalue variance the smaller
997 the number of effective tests. The number of independent tests was 1,328 for 2,801 tested taxa.
998 We used this information to calculate a Bonferroni-adjusted study-wide significant level for
999 significant associations, which was set to $5 \times 10^{-8} / 1328 = 3.8 \times 10^{-11}$. A genome-wide significant
1000 threshold was set as 5×10^{-8} . The identified SNPs were annotated using ANNOVAR
1001 v2018Apr16¹¹¹ and grouped into genetic loci using 200kb window flanking the top SNPs.

1002
1003 We also examined whether antibiotic prescription prior to baseline sampling could be an
1004 important confounder of results. We obtained individual information on the prescription of any
1005 antibiotic up to 1 month prior to baseline fecal sampling, corresponding to 250 individuals out
1006 of 5,959 (4.2%). We examined whether individual microbial profiles (via beta-diversity
1007 estimates using Bray-Curtis dissimilarity) were broadly affected by recent antibiotic prescription
1008 and observed a slight effect along PCoAs with significant variance explained (**Extended Data**
1009 **Figure 9C**). After repeating the GWAS for all microbial taxa for which we initially had found
1010 at least one significantly associated locus, this time adjusting for prior antibiotic prescription
1011 status (yes vs. no) (**Supplementary Table 9**), we found that recent antibiotic prescriptions had
1012 very minor effects on the GWAS association results. Adjusting for antibiotic prescription did not
1013 change any study-wide significant associations and only 32 out of 567 genome-wide associations

1014 moved slightly above $p=5\times 10^{-8}$ (the largest p-value was 3.2×10^{-7}), which is likely by chance
1015 given inclusion of any additional covariate (**Supplementary Table 9**). In addition, the beta
1016 estimates with and without the adjustment of antibiotics usage were highly consistent (Pearson
1017 $r=0.9999487$).

1018
1019 One important association in our study involved *F. lactaris* abundance and variants in the ABO
1020 locus. We observed the distribution of *F. lactaris* abundance in our GWAS cohort ($n=5,959$) was
1021 slightly bimodal (**Extended Data Figure 10**). To investigate whether a logistic model gives the
1022 same result for this taxon, we arbitrarily coded *F. lactaris* abundance as “1” if the relative
1023 abundance was higher than 5×10^{-4} ($n=2866$), and “0” if smaller ($n=3093$). Akaike Information
1024 Criterion (AIC) value was smaller for logistic than for linear model (AIC=8196 vs AIC=12463,
1025 respectively), and the strongest association was also observed in the same top SNP (rs545971,
1026 $p=5.5\times 10^{-18}$) as when using linear regression (rs545971, $p=1.1\times 10^{-12}$).

1027 1028 **Replication of previously reported associations**

1029
1030 To evaluate the reproducibility of our results with previously reported associations, we collected
1031 GWAS summary results from 8 studies published in peer-reviewed journals at the time of this
1032 work^{3,6-10,102,112}. These studies reported associations between 548 SNPs and microbial features.
1033 ANNOVAR was used to annotate the reported SNPs to the hg38 human reference genome¹¹¹
1034 and we used plink2⁹⁹ to identify a further 15,427 SNPs in high LD ($r^2>0.8$, within 5 Mbp) with
1035 any of these 548 SNPs. To assess replication, we first examined whether previously reported
1036 associations could be matched in our results to identical or linked SNPs, with an association
1037 below the Bonferroni-corrected suggestive significance threshold, which was set to $0.05/548 =$
1038 9.124×10^{-5} . More details about the replication methods and the use of GTDB taxonomic system
1039 can be found in **Supplementary Note**.

1040 1041 **Prediction of ABO blood groups and secretor status**

1042
1043 SNP-based typing of ABO histo-blood group was performed. A combination of four SNPs¹¹⁴
1044 was used for the prediction, and a 98% concordance with phenotypically typed ABO histo-blood
1045 group has been reported for this method⁴. For blood group allele A, the two different types A1
1046 and A2 were predicted by rs507666 and rs8176704 respectively. Blood group allele B was
1047 inferred from rs8176746 and blood group allele O was predicted by rs687289. As the
1048 combination of these SNPs are exclusive, no haplotype information was needed. To validate the
1049 accuracy of prediction, we compared it with the prediction using a different combination of
1050 SNPs⁶⁸. The two predictions were highly consistent, with over 99.9% concordance. In addition,
1051 the distribution of ABO groups was consistent with the population distribution found in public
1052 database. Secretor status was predicted by the genotype of *FUT2* variant rs601338, where AA
1053 or AG genotypes are secretors and GG genotypes are non-secretors. An 100% concordance
1054 between the variation in rs601338 and secretor status was reported in a study on Finnish
1055 individuals¹¹⁵.

1056 1057 **Bidirectional two-sample Mendelian randomization (MR) analysis**

1058
1059 Causal relationships between diseases and gut microbiota were investigated at genus and species
1060 levels only to maximize interpretability. In total, 213 species and 148 genera associated with at
1061 least one variant at genome-wide significant level ($p<1\times 10^{-8}$) were included. GWAS summary
1062 results were collected for 46 diseases from MR-Base¹¹⁶ (**Supplementary Table 4**). These
1063 included 12 autoimmune or inflammatory diseases, 9 cardiometabolic diseases, 13 psychiatric

1064 or neurological diseases, cardiovascular diseases, 4 bone diseases and 8 cancers. For disease with
1065 more than one GWAS records, the record with the largest sample size was kept.

1066
1067 Bi-directional causal inference was performed to infer causal effects of microbial abundance
1068 variation (exposure) on disease risk (outcome), and of disease (exposure) on microbial
1069 abundance levels (outcome). To select the SNP instruments for microbial exposures in our study
1070 (**Supplementary Table 7**), we followed recommendations from a previous study showing that
1071 associated SNPs below a significance threshold of $p < 1 \times 10^{-5}$ had the largest explained variance
1072 on microbial features¹¹⁷. For each taxon, GCTA-COJO was used to perform a conditional
1073 analysis to select independently associated SNPs at $p < 1 \times 10^{-5}$. F-statistics were calculated to
1074 estimate the strength of instruments for each bacterial exposures, and were found to be > 10 for
1075 all exposures (**Supplementary Table 5**). SNP instruments for disease exposures were selected
1076 at genome-wide significant threshold ($p < 5 \times 10^{-8}$). Subsequently LD-clumping with a strict
1077 threshold ($r^2 < 0.001$ in 1000G EUR within 10 Mb windows) was conducted to select independent
1078 instruments with the lowest p values for taxa and diseases, respectively.

1079
1080 Details about the precise methods used for MR inference can be found in **Supplementary Note**.

1081

1082 **Cox proportional hazards regression**

1083

1084 Cox proportional hazards regression was conducted to test the association between baseline
1085 abundance of gut microbe and incident major depression (16 years follow-up, $n=181$ incident
1086 events). Microbial abundances were CLR-transformed and standardized to zero-mean and unit-
1087 variance. The Cox models were stratified by sex and adjusted for age and log-transformed BMI,
1088 with time-on-study as the time scale. Participants with prevalent major depression at baseline
1089 were excluded. R function *coxph()* in the R package *survival* v3.1-8 was used for this analysis.

1090

1091 **Profiling of carbohydrate-active enzymes (CAZymes) in bacterial genomes**

1092

1093 The standalone run_dbCAN2 v2.0.11 tool¹¹⁸ (https://github.com/linnabrown/run_dbcan) was
1094 used to scan for the presence of CAZyme genes from public assembled bacterial genomes taken
1095 from the GTDB release 89 reference. We used a CAZyme reference database taken from the
1096 CAZy database¹¹⁹ (31st July 2019 update). In total, we scanned 327 *Bifidobacterium sp.*, 2
1097 *Faecalicatena lactaris* and 15 *Collinsella sp.* reference genomes included in GTDB release 89.
1098 Three methods were compared as part of the run_dbCAN2 procedure (HMMER, DIAMOND,
1099 and Hotpep). We considered a positive detection result when all three methods agreed on a
1100 CAZyme family identification. Identification of preferred reported substrates for the various
1101 CAZyme families was done manually from key publications^{42,120}, from literature searches and
1102 from the CAZypedia website¹²¹. Certain CAZyme families have a broad range of substrates,
1103 many of which are still unknown, which results in our reported preferred substrates to be as
1104 accurate as possible, but non-exhaustive.

1105

1106 **Carbon impact and offsetting**

1107

1108 We used GreenAlgorithms v1.0¹²² to estimate that the main computational work in this study
1109 had a carbon impact of at least 2,660 kg CO₂e, corresponding to 233 tree-years. As a commitment
1110 to the reduction of carbon emissions associated with computation in research, we consequently
1111 funded planting of 30 trees through a local Australian charity, which across their lifetime will
1112 sequester a combined estimated 8,040 kg CO₂e, or 3 times the amount of CO₂e generated by this
1113 study.

1114
1115 **Data availability**
1116
1117 Complete summary statistics of microbial taxa with genome-wide significant hits are available
1118 in the NHGRI-EBI GWAS Catalog (<https://www.ebi.ac.uk/gwas/>), GCP ID: GCP000228; study
1119 accession numbers GCST90032172-GCST90032644. The metagenomic data from FINRISK
1120 2002 samples are available from the European Genome-Phenome Archive (study ID:
1121 EGAS00001005020). The phenotype data contain sensitive information from healthcare
1122 registers and are not publicly available because to avoid compromising research participant
1123 privacy/consent. They are available through the THL biobank upon submission of a research
1124 plan and signing a data transfer agreement (<https://thl.fi/en/web/thl-biobank/for-researchers/application-process>). Additional databases used in this work include GTDB release
1125 89 (<https://gtdb.ecogenomic.org/>) and CAZy (last accessed 31/07/2019) (<http://www.cazy.org/>).
1126
1127

1128 **Code availability**

1129
1130 Scripts used to analyze non-identifiable data in this study have been made available on Zenodo
1131 (doi: 10.5281/zenodo.5641303).
1132

1133 **Methods-specific references**

- 1134
1135 92. Borodulin, K. et al. Cohort Profile: The National FINRISK Study. *International Journal*
1136 *of Epidemiology* 47, 696–696i (2018).
1137 93. Borodulin, K. et al. Forty-year trends in cardiovascular risk factors in Finland. *The*
1138 *European Journal of Public Health* 25, 539–546 (2015).
1139 94. Liu, Y. et al. Early prediction of liver disease using conventional risk factors and gut
1140 microbiome-augmented gradient boosting.
1141 <http://medrxiv.org/lookup/doi/10.1101/2020.06.24.20138933> (2020)
1142 doi:10.1101/2020.06.24.20138933.
1143 95. Salosensaari, A., Laitinen, V., Havulinna, A.S. *et al.* Taxonomic signatures of cause-
1144 specific mortality risk in human gut microbiome. *Nat Commun* 12, 2671 (2021).
1145 96. FinnGen et al. Polygenic and clinical risk scores and their impact on age at onset and
1146 prediction of cardiometabolic diseases and common cancers. *Nat Med* 26, 549–557 (2020).
1147 97. Loh, P.-R. et al. Reference-based phasing using the Haplotype Reference Consortium
1148 panel. *Nat Genet* 48, 1443–1448 (2016).
1149 98. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and
1150 accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat*
1151 *Genet* 44, 955–959 (2012).
1152 99. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and
1153 richer datasets. *GigaSci* 4, 7 (2015).
1154 100. Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: rapid and sensitive
1155 classification of metagenomic sequences. *Genome Res.* 26, 1721–1729 (2016).
1156 101. Ruuskanen, M. O. et al. Links between gut microbiome composition and fatty liver
1157 disease in a large population sample. *Gut Microbes* 13, 1–22 (2021).
1158 102. Goodrich, J. K., Davenport, E. R., Waters, J. L., Clark, A. G. & Ley, R. E. Cross-
1159 species comparisons of host genetic associations with the microbiome. *Science* 352, 532–535
1160 (2016).
1161 103. Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome
1162 Datasets Are Compositional: And This Is Not Optional. *Front. Microbiol.* 8, 2224 (2017).
1163 104. Aitchison, J., Barceló-Vidal, C., Martín-Fernández, J. A. & Pawlowsky-Glahn, V. [No
1164 title found]. *Mathematical Geology* 32, 271–275 (2000).

1165 105. Qin, Y. et al. Genome-wide association and Mendelian randomization analysis
1166 prioritizes bioactive metabolites with putative causal effects on common diseases.
1167 <http://medrxiv.org/lookup/doi/10.1101/2020.08.01.20166413> (2020)
1168 doi:10.1101/2020.08.01.20166413.

1169 106. Loh, P.-R. et al. Efficient Bayesian mixed-model analysis increases association power
1170 in large cohorts. *Nat Genet* 47, 284–290 (2015).

1171 107. Abraham, G., Qiu, Y. & Inouye, M. FlashPCA2: principal component analysis of
1172 Biobank-scale genotype datasets. *Bioinformatics* 33, 2776–2778 (2017).

1173 108. Genetic Investigation of ANthropometric Traits (GIANT) Consortium et al.
1174 Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional
1175 variants influencing complex traits. *Nat Genet* 44, 369–375 (2012).

1176 109. Li, J. & Ji, L. Adjusting multiple testing in multilocus analyses using the eigenvalues of
1177 a correlation matrix. *Heredity* 95, 221–227 (2005).

1178 110. Nyholt, D. R. A Simple Correction for Multiple Testing for Single-Nucleotide
1179 Polymorphisms in Linkage Disequilibrium with Each Other. *The American Journal of Human*
1180 *Genetics* 74, 765–769 (2004).

1181 111. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
1182 variants from high-throughput sequencing data. *Nucleic Acids Res* 38, e164 (2010).

1183 112. Davenport, E. R. et al. Seasonal variation in human gut microbiome composition. *PLoS*
1184 *One* 9, e90731 (2014).

1185 113. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide
1186 complex trait analysis. *Am J Hum Genet* 88, 76–82 (2011).

1187 114. Paré, G. et al. Novel Association of ABO Histo-Blood Group Antigen with Soluble
1188 ICAM-1: Results of a Genome-Wide Association Study of 6,578 Women. *PLoS Genet* 4,
1189 e1000118 (2008).

1190 115. Wacklin, P. et al. Secretor Genotype (FUT2 gene) Is Strongly Associated with the
1191 Composition of Bifidobacteria in the Human Intestine. *PLoS ONE* 6, e20113 (2011).

1192 116. Hemani, G. et al. The MR-Base platform supports systematic causal inference across
1193 the human phenome. *eLife* 7, e34408 (2018).

1194 117. Sanna, S. et al. Causal relationships among the gut microbiome, short-chain fatty acids
1195 and metabolic diseases. *Nat Genet* 51, 600–605 (2019).

1196 118. Zhang, H. et al. dbCAN2: a meta server for automated carbohydrate-active enzyme
1197 annotation. *Nucleic Acids Research* 46, W95–W101 (2018).

1198 119. Cantarel, B. L. et al. The Carbohydrate-Active EnZymes database (CAZy): an expert
1199 resource for Glycogenomics. *Nucleic Acids Research* 37, D233–D238 (2009).

1200 120. Cantarel, B. L., Lombard, V. & Henrissat, B. Complex Carbohydrate Utilization by the
1201 Healthy Human Microbiome. *PLoS ONE* 7, e28742 (2012).

1202 121. The CAZypedia Consortium. Ten years of CAZypedia: a living encyclopedia of
1203 carbohydrate-active enzymes. *Glycobiology* 28, 3–8 (2018).

1204 122. Lannelongue, L., Grealey, J., Inouye, M., Green Algorithms: Quantifying the Carbon
1205 Footprint of Computation. *Adv. Sci.* 2021, 8, 2100707.

1206
1207