

The interplay between host genetics and the gut microbiome reveals common and distinct microbiome features for complex human diseases

Fengzhe Xu

Westlake University

Yuanqing Fu

Westlake University

Tingyu Sun

Sun Yat-Sen University

Zengliang Jiang

Westlake University

Zelei Miao

Westlake University

Menglei Shuai

Westlake University

Wanglong Gou

Westlake University

Chuwen Ling

Sun Yat-Sen University

Jian Yang

University of Queensland

Jun Wang

Chinese Academy of Sciences

Yu-ming Chen

Sun Yat-Sen University

Ju-Sheng Zheng (✉ zhengjusheng@westlake.edu.cn)

Westlake University <https://orcid.org/0000-0001-6560-4890>

Research

Keywords: Gut Microbiome, Host Genetics, Bi-directional Mendelian Randomization Analyses, Disease-Microbiome Features

Posted Date: September 14th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-40335/v3>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on October 8th, 2020. See the published version at <https://doi.org/10.1186/s40168-020-00923-9>.

1 **The interplay between host genetics and the gut microbiome reveals common**
2 **and distinct microbiome features for complex human diseases**

3 Fengzhe Xu^{1#}, Yuanqing Fu^{1,3#}, Ting-yu Sun², Zengliang Jiang^{1,3}, Zelei Miao¹,
4 Menglei Shuai¹, Wanglong Gou¹, Chu-wen Ling², Jian Yang^{4,5}, Jun Wang^{6*}, Yu-ming
5 Chen^{2*}, Ju-Sheng Zheng^{1,3,7*}

6 [#]These authors contributed equally to the work

7 ¹ Zhejiang Provincial Laboratory of Life Sciences and Biomedicine, Key Laboratory
8 of Growth Regulation and Translational Research of Zhejiang Province, School of
9 Life Sciences, Westlake University, Hangzhou, China.

10 ² Guangdong Provincial Key Laboratory of Food, Nutrition and Health; Department of
11 Epidemiology, School of Public Health, Sun Yat-sen University, Guangzhou, China.

12 ³ Institute of Basic Medical Sciences, Westlake Institute for Advanced Study,
13 Hangzhou, China.

14 ⁴ Institute for Molecular Bioscience, The University of Queensland, Brisbane, QLD,
15 Australia.

16 ⁵ Institute for Advanced Research, Wenzhou Medical University, Wenzhou, Zhejiang
17 325027, China

18 ⁶ CAS Key Laboratory for Pathogenic Microbiology and Immunology, Institute of
19 Microbiology, Chinese Academy of Sciences, Beijing, China.

20 ⁷ MRC Epidemiology Unit, University of Cambridge, Cambridge, UK.

21

22 Short title: Interplay between host genetics and gut microbiome

23

24 *Correspondence to

25 Prof Ju-Sheng Zheng

26 School of Life Sciences, Westlake University, 18 Shilongshan Rd, Cloud Town,

27 Hangzhou, China. Tel: +86 (0)57186915303. Email: zhengjusheng@westlake.edu.cn

28 And

29 Prof Yu-Ming Chen

30 Guangdong Provincial Key Laboratory of Food, Nutrition and Health; Department of

31 Epidemiology, School of Public Health, Sun Yat-sen University, Guangzhou, China.

32 Email: chenyum@mail.sysu.edu.cn

33 And

34 Prof Jun Wang

35 CAS Key Laboratory for Pathogenic Microbiology and Immunology, Institute of

36 Microbiology, Chinese Academy of Sciences, Beijing, China.

37 Email: junwang@im.ac.cn

39 *Abstract*

40 **Background** Interest in the interplay between host genetics and the gut microbiome
41 in complex human diseases is increasing, with prior evidence mainly being derived
42 from animal models. In addition, the shared and distinct microbiome features among
43 complex human diseases remain largely unclear.

44 **Results** This analysis was based on a Chinese population with 1,475 participants. We
45 estimated the SNP-based heritability, which suggested that *Desulfovibrionaceae* and
46 *Odoribacter* had significant heritability estimates (0.456 and 0.476, respectively). We
47 performed a microbiome genome-wide association study to identify host genetic
48 variants associated with the gut microbiome. We then conducted bidirectional
49 Mendelian randomization analyses to examine the potential causal associations
50 between the gut microbiome and complex human diseases. We found that
51 *Saccharibacteria* could potentially decrease the concentration of serum creatinine and
52 increase the estimated glomerular filtration rate. On the other hand, atrial fibrillation,
53 chronic kidney disease and prostate cancer, as predicted by host genetics, had
54 potential causal effects on the abundance of some specific gut microbiota. For
55 example, atrial fibrillation increased the abundance of *Burkholderiales* and
56 *Alcaligenaceae* and decreased the abundance of *Lachnobacterium*, *Bacteroides*
57 *coprophilus*, *Barnesiellaceae*, ~~an undefined genus in the family undefined genus in~~
58 ~~family~~ *Veillonellaceae* and *Mitsuokella*. Further disease-microbiome feature analysis
59 suggested that systemic lupus erythematosus and chronic myeloid leukaemia shared

60 common gut microbiome features.

61 **Conclusions** These results suggest that different complex human diseases share
62 common and distinct gut microbiome features, which may help reshape our
63 understanding of disease aetiology in humans.

64 **Keywords** Gut Microbiome; Host Genetics; Bidirectional Mendelian Randomization
65 Analyses; Disease-Microbiome Features

66

67 **Background**

68 Ever increasing evidence has suggested that the gut microbiome is involved in many
69 physiological processes, such as energy harvesting, the immune response, and
70 neurological function [1-3]. With successes of investigation into the clinical
71 application of faecal transplants, the modulation of the gut microbiome has emerged
72 as a potential treatment option for some complex diseases, including inflammatory
73 bowel disease and colorectal cancer [4, 5]. However, it is still unclear whether the gut
74 microbiome has the potential to be clinically applied for the prevention or treatment
75 of many other complex diseases. Therefore, it is important to clarify the bidirectional
76 causal association between the gut microbiome and complex human diseases or traits.

77

78 Mendelian randomization (MR) is a method that uses genetic variants as instrumental
79 variables to investigate the causality between an exposure and an outcome in
80 observational studies [6]. Prior studies provide evidence that the composition or
81 structure of the gut microbiome can be influenced by host genetics [7-10]. On the
82 other hand, host genetic variants associated with the gut microbiome are rarely
83 explored in Asian populations; thus, we still lack instrumental variables to perform
84 MR for the gut microbiome in Asians. This calls for a novel microbiome genome-
85 wide association study (GWAS) in Asian populations.

86

87 Along with the causality issue between the gut microbiome and complex human
88 diseases, it is unclear whether complex human diseases have similar or unique gut
89 microbiome features. The identification of common and distinct gut microbiome
90 features across different diseases may shed light on novel relationships among the
91 complex diseases and update our understanding of the disease aetiology in humans.
92 However, the composition and structure of the gut microbiome are influenced by a
93 variety of factors, including the environment, diet and regional variation [11-13],
94 which poses a key challenge for the description of representative microbiome features
95 for a specific disease. Although there were several studies comparing disease-related
96 gut microbiome features [14-16], few of them examined and compared the
97 microbiome features across different ~~human~~-complex human diseases.

98

99 In the present study, we performed a microbiome GWAS in a Chinese cohort, the
100 Guangzhou Nutrition and Health Study (GNHS) [17], including 1475 participants.
101 Subsequently, we applied a bidirectional MR method to explore the genetically
102 predicted relationship between the gut microbiome and complex human diseases. To
103 explore novel relationships among complex human diseases based on the gut
104 microbiome, we investigated the shared and distinct gut microbiome features across
105 diverse complex human diseases.

106

107 **Results**

108 **Overview of the study**

109 Our study was based on the GNHS, with 4048 participants (40-75 years old) living in
110 the urban area of Guangzhou city recruited during 2008 and 2013 [17]. In the GNHS,
111 stool samples were collected among 1937 participants during follow-up visits, among
112 whom 1475 unrelated participants not taking antibiotics were included in our
113 discovery microbiome GWAS. We then included an additional 199 participants with
114 both genetic data and gut microbiome data as a replication cohort, which belonged to
115 the control arm of a case-control study of hip fracture in Guangdong Province, China
116 [18] (see also Figure 1).

117

118 **SNP-based heritability of the gut microbiome**

119 The heritability of alpha diversity ranged from 0.035 to 0.103 (SE: from 0.174 to
120 0.193, Supplementary Table S3). Significant heritability estimates were observed for
121 several taxa (see also Figure 2, Supplementary Table S3), with crude p values <0.05 .
122 To further correct the multiple testing, we calculated the effective number of
123 independent taxa in each taxonomic level (phylum level: 2.3, class level: 2.9, order
124 level: 2.9, family level: 5.5, genus level: 5.6, species level: 3.2), as some taxa were
125 highly correlated with each other. The results suggested that *Desulfovibrionaceae* and
126 *Odoribacter* were heritable ($p < 0.05/n$, where n is the effective number of independent

127 taxa). Notably, among the suggestively heritable taxa in our cohort
128 [*Paraprevotellaceae*], *Veillonellaceae*, *Desulfovibrionaceae*, *Pasteurellaceae*,
129 *Odoribacter*, *Paraprevotella*, *Veillonella* and *Bifidobacterium* had nominally
130 significant heritability estimates in prior literature [7, 19-21].

131

132 **Association of host genetics with gut microbiome features**

133 We generated categorical variable enterotypes (*Prevotella* vs *Bacteroides*) of the
134 participants based on the genus-level relative abundance of the gut microbiome [22].
135 Thereafter, we performed a GWAS for enterotypes using a logistic regression model
136 to explore potential associations between host genetics and enterotypes. However, we
137 did not find any genome-wide significant loci ($p < 5 \times 10^{-8}$).

138

139 To examine the association of host genetic variants with alpha diversity, we performed
140 a GWAS for four indices (Shannon diversity index, Chao1 diversity indices, observed
141 OTU index and phylogenetic diversity), but again, no genome-wide significant signal
142 ($p < 5 \times 10^{-8}$) was found. To further investigate whether there is a host genetic basis
143 underlying alpha diversity, we constructed a polygenic score for each alpha diversity
144 indicator in the replication cohort using the genetic variants that showed suggestive
145 significance ($p < 5 \times 10^{-5}$) in the discovery GWAS. The polygenic score was not
146 significantly associated with its corresponding alpha diversity index in our replication

147 cohort. Furthermore, none of the associations with alpha diversity indices reported in
148 the literature could be replicated (Supplementary Table S8) [7].

149

150 The beta diversity GWAS was performed with MicrobiomeGWAS based on Bray–
151 Curtis dissimilarity [23]. We found that one locus at the *SMARCA2* gene (rs6475456)
152 was associated with beta diversity at a genome-wide significant level ($p=3.96\times 10^{-9}$).
153 However, we could not replicate the results in the replication cohort, which might be
154 due to the limited sample size of the replication cohort. In addition, prior literature
155 had reported 73 genetic variants that were associated with beta diversity [8, 13, 24,
156 25], among which we found that 3 single-nucleotide polymorphisms (SNPs, *UHRF2*
157 gene-rs563779, *LHFPL3* gene-rs12705241, *CTD-2135J3.4*-rs11986935) had
158 nominally significant ($p<0.05$) associations with beta diversity in our cohort
159 (Supplementary Table S7), although none of the associations survived Bonferroni
160 correction. These studies used various methods for the sequencing and calculation of
161 beta diversity, which raised challenges to verify and extrapolate results across
162 populations.

163

164 We subsequently performed a discovery GWAS for individual gut microbes in our
165 own GNHS discovery dataset. For the taxa ($n=114$) present in not fewer than ninety
166 percent of participants, we carried out an analysis based on a log-normal model. For

167 other taxa (n=88) present in fewer than ninety percent, we transformed the
168 absence/presence of the taxon into binary variables and used a logistic model to
169 prevent zero inflation (Supplementary Table S1). For all the gut microbiome taxa, the
170 significance threshold was defined as 5×10^{-8} in the discovery stage. We found that 6
171 taxa were associated with host genetic variants in the discovery cohort ($p < 5 \times 10^{-8}/n$,
172 where n is the effective number of independent taxa in each taxonomic level,
173 Supplementary Table S5); however, these associations were not significant ($p > 0.05$)
174 in the replication cohort. We then took the genetic loci reported to be associated with
175 individual taxa in prior studies [7, 8, 13, 25] for replication in our GNHS dataset.
176 Although none of the associations of these genetic variants with taxa survived the
177 Bonferroni correction ($p < 1 \times 10^{-4}$), we found that *STPG2*-rs4699323 had a nominally
178 significant association ($p < 0.05$) with *Clostridiales* (Beta: -0.131 [95% CI: -0.233, -
179 0.029], $p = 0.012$; Supplementary Table S6). We then used a threshold of $p < 5 \times 10^{-5}$ at
180 the discovery GWAS stage to incorporate additional genetic variants that might
181 explain a larger proportion of heritability for taxa, and based on this, we constructed a
182 polygenic score for each taxon in the replication. We found that the polygenic scores
183 were significantly associated with 5 taxa, including *Saccharibacteria* (also known as
184 *TM7* phylum), *Clostridiaceae*, *Comamonadaceae*, *Klebsiella* and *Desulfovibrio d168*,
185 in the replication set ($p < 0.05$, Methods, see also Supplementary Figure 1,
186 Supplementary Table S9).

187

188 Genetic correlation of gut microbiome and traits

189 WhileAs the associations of the microbiome with complex diseases and traits have
190 been widely reported [26], the genetic correlation between the gut microbiome and
191 traits of interest is less clear. Therefore, we applied bivariate GREML analysis to
192 address this question. The traits included BMI, fasting blood sugar (FBS),
193 glycosylated haemoglobin (HbA1c), systolic blood pressure (SBP), diastolic blood
194 pressure (DBP), high density lipoprotein cholesterol (HDL-C), low density
195 lipoprotein cholesterol (LDL-C), total cholesterol (TC) and triglyceride (TG), none of
196 which could pass Bonferroni correction. HDL-C was the only trait that had nominal
197 genetic correlation ($p < 0.05$) with gut microbes (specifically, *Desulfovibrionaceae* and
198 [*Prevotella*], Supplementary Table S4).

199

200 Bidirectional assessment of the genetically predicted association between the gut 201 microbiome and complex diseases/traits

202 Using genetic-variant-composed polygenic scores as genetic instruments, we
203 performed MR analysis to assess the putative causal effect of the microbiome
204 (*Saccharibacteria*, *Clostridiaceae*, *Comamonadaceae*, *Klebsiella* and *Desulfovibrio*
205 *d168*) on complex human diseases or traits. The inverse variance weighted (IVW)
206 method was used for the MR analysis, while the other three methods (weighted
207 median, MR-Egger and MR-PRESSO) were applied to confirm the robustness of the

208 results. Horizontal pleiotropy was assessed using the MR-PRESSO global test and
209 MR-Egger regression. For the analysis of the gut microbiome on complex traits, we
210 downloaded publicly available GWAS summary statistics of complex traits (n=58)
211 and diseases (type 2 diabetes mellitus (T2DM), atrial fibrillation (AF), colorectal
212 cancer (CRC) and rheumatoid arthritis (RA)) reported by BioBank Japan [27-32]. The
213 results suggested that *Saccharibacteria* (per 1-SD higher in the log-transformed
214 abundance) could potentially decrease the concentration of serum creatinine (-0.011
215 [95% CI: -0.019, -0.003], $p=0.007$) and increase the estimated glomerular filtration
216 rate (eGFR) (0.012 [95% CI: 0.004, 0.020], $p=0.003$, Supplementary Table S10),
217 which might help improve renal function. We did not find evidence of pleiotropic
218 effects: genetic variants associated with *Saccharibacteria* were not associated with
219 any of the above traits (58 complex traits and 4 disease outcomes, $p<0.05/62$). These
220 taxa were not causally associated with other complex diseases or traits in our MR
221 analyses, which might be due to the limited genetic instruments discovered in our
222 present study.

223

224 We subsequently performed a ~~reverse~~ MR analysis to assess the potential
225 causal effect of complex human diseases on gut microbiome features. For the
226 ~~reverse~~ MR analyses, the diseases of interest included T2DM, AF, coronary
227 artery disease (CAD), chronic kidney disease (CKD), Alzheimer's disease (AD), CRC
228 and prostatic cancer (PCa), and their instrumental variables for the MR analysis were

229 based on previous large-scale GWASs in East Asians [27, 33-38]. The results
230 suggested that AF and CKD were causally associated with the gut microbiome (see
231 also Figure 3A, 3B, Supplementary Table S11). Specifically, genetically predicted
232 higher risk of AF (per log odds) was associated with a lower abundance of
233 *Lachnobacterium* (Beta: -0.078 [95% CI: -0.148, -0.006], $p=0.034$), *Bacteroides*
234 *coprophilus* (Beta: -0.113 [95% CI: -0.184, -0.041], $p=0.002$), *Barnesiellaceae* (odds
235 ratio: 0.818 [95% CI: 0.686, 0.976], $p=0.026$), an undefined genus in the family
236 ~~undefined-genus-in-family~~ *Veillonellaceae* (odds ratio: 0.801 [95% CI: 0.669, 0.960],
237 $p=0.017$) and *Mitsuokella* (odds ratio: 0.657 [95% CI: 0.496, 0.870], $p=0.003$), and
238 higher abundance of *Burkholderiales* (Beta: 0.079 [95% CI: 0.009, 0.150], $p=0.027$)
239 and *Alcaligenaceae* (Beta: 0.082 [95% CI: 0.012, 0.152], $p=0.022$). Additionally,
240 genetically predicted higher risk of CKD could increase *Anaerostipes* (Beta: 0.291
241 [95% CI: 0.057, 0.524], $p=0.015$) abundance, and a higher risk of PCa could decrease
242 [*Prevotella*] (odds ratio: -0.758 [95% CI: -1.354, -0.162], $p=0.013$).

243

244 **Microbiome features of ~~human~~-complex human diseases**

245 To further investigate the potential complex diseases that may be correlated with the
246 taxa affected by AF, we applied Phylogenetic Investigation of Communities by
247 Reconstruction of Unobserved States (PICRUSt) to predict the disease pathway
248 abundance [39]. We used Spearman's rank-order correlation to test whether the

249 relative abundances of predicted diseases based on PICRUSt were associated with the
250 aforementioned AF-associated taxa (see also Supplementary Figure 2, Supplementary
251 Table S12). The heatmap indicated that cancers and neurodegenerative diseases,
252 including Parkinson's disease (PD), AD, amyotrophic lateral sclerosis (ALS) and AF,
253 were correlated with similar gut microbiomes. Although the association among these
254 diseases is highly supported by previous studies [40-42], no study has compared
255 common gut microbiome features across these different diseases.

256

257 To compare gut microbiome features across human diseases, we used the predicted
258 disease abundance based on PICRUSt and performed k-medoid clustering. According
259 to the optimum average silhouette width [43], we chose the optimal number of
260 clusters for further analysis. The plot showed that neurological diseases, including
261 ALS and AD, belonged to the same cluster, while PD and CRC had much similarity in
262 the gut microbiome. The results also suggested that systemic lupus erythematosus
263 (SLE) and chronic myeloid leukaemia (CML) shared similar gut microbiome features
264 (see also Figure 4A, 4B). Moreover, we could replicate these clusters in our
265 replication cohort, which suggested that the clustering results were robust (see also
266 Figure 4C).

267

268 We further asked whether the gut microbiome contributed to the novel clustering. To

269 this end, we repeated the analysis among participants who took antibiotics less than
270 two weeks before stool sample collection, considering that antibiotic treatments were
271 believed to cause microbiome imbalance. We used the Jaccard similarity coefficient to
272 estimate the cluster difference among the GNHS cohort, the replication cohort and the
273 antibiotic group. The similarity between the GNHS cohort and the replication cohort
274 was higher than that between the GNHS cohort and the antibiotic group (Jaccard
275 similarity coefficient: 0.61 versus 0.11). The results indicated a different clustering,
276 which suggested that the gut microbiome indeed contributed to the correlations
277 among diseases (see also Figure 4D). To further demonstrate common microbiome
278 features across different diseases, we examined the correlation of the predicted
279 diseases with genus-level taxa. The results showed that ~~human~~-complex human
280 diseases shared similar gut microbiome features, as well as distinct features on their
281 own (see also Figure 5, Supplementary Table S13).

282

283 To validate whether the disease-related gut microbiome features annotated by KEGG
284 would be associated with the risk of the disease in a real-world community-based
285 cohort, we used T2DM as an example, examining the association of predicted T2DM-
286 related microbiome features with T2DM risk in our GNHS cohort. We constructed a
287 microbiome risk score (MRS) based on 16 selected taxa with predicted correlation
288 coefficients with T2DM greater than 0.2. A logistic regression model was used to
289 examine the association between MRS and T2DM risk in GNHS (n=1886, with 217

290 T2DM cases). The results showed that MRS was positively associated with the risk of
291 T2DM (odds ratio: 1.176 [95% CI: 1.114, 1.244], $p=8.75\times 10^{-9}$).

292

293 **Discussion**

294 Our study is among the first to investigate host genetics-gut microbiome associations
295 in East Asian populations and reveals that several microbiome species (e.g.,
296 *Saccharibacteria* and *Klebsiella*) are influenced by host genetics. We found that
297 *Saccharibacteria* might causally improve renal function by affecting renal function
298 biomarkers (i.e., creatinine and eGFR). On the other hand, complex diseases such as
299 atrial fibrillation, chronic kidney disease and prostate cancer have potential causal
300 effects on the gut microbiome. More interestingly, our results indicated that different
301 complex diseases may be ~~mechanistically~~mechanically correlated by sharing common
302 gut microbiome features but also maintaining their own distinct microbiome features.

303

304 Previous studies and our study showed that the gut microbiome had an inclination to
305 be influenced by host genetics [8, 10, 25, 44, 45]. The results suggested that
306 *Desulfovibrionaceae* and *Odoribacter* had nominally significant heritability estimates,
307 which were consistent with prior results [7]. We also identified several suggestively
308 heritable taxa that were nominally significant in previous studies [19-21]. In addition,
309 we successfully constructed polygenic scores for *Clostridiaceae* and

310 *Comamonadaceae*, both of which have been identified to be heritable or suggested to
311 be heritable [7, 45].

312

313 We could not replicate any of the reported genetic variants that were significantly
314 associated with gut microbiome features in prior reports, which might be due to
315 multiple reasons. One of the major reasons may be that the massive multiple testing in
316 insufficiently large samples in prior microbiome GWASs may potentially lead to
317 false-positive findings. In addition, other factors, including ethnic differences,
318 heterogeneity between studies, gene-environment interactions and dissimilarity in
319 sequencing methods, might also make it difficult to extrapolate results from
320 microbiome GWASs across populations in the microbiome field. Nevertheless, we
321 successfully replicated several polygenic scores of the gut microbiome. The current
322 study represents the largest dataset, to the best of our knowledge, in Asian populations
323 and may serve as a unique resource for large-scale trans-ethnic meta-analyses of
324 microbiome GWASs in the future.

325

326 MR analysis showed that *Saccharibacteria* might decrease the concentration of serum
327 creatinine and increase eGFR. Little is known about *Saccharibacteria* as one of the
328 uncultivated phyla, and previous studies have shown that it might be essential for the
329 immune response, oral inflammation and inflammatory bowel disease [46-48]. Our

330 results also provided a genetic instrument of *Saccharibacteria* for further causal
331 analysis with other complex diseases. The ~~reverse~~ MR analysis provided
332 evidence that AF, CKD and PCa could causally influence the gut microbiome. The
333 rare and low-frequency variants may have an important impact on common diseases
334 [49]; thus, it will be of interest to clarify the effects of low-frequency variants on the
335 gut microbiome in cohorts with large sample sizes in the future.

336

337 Our results indicate that the gut microbiome helps reveal novel and interesting
338 relationships among complex human diseases, and different diseases may have
339 common and distinct gut microbiome features. A prior study including participants
340 from different countries identified three microbiome clusters [22]. Notably, this study
341 focused on classifying the individuals into distinct enterotypes regardless of the
342 individuals' health status, while in the present study, we described representative
343 microbiome features for diseases of interest. We provided an approach to interpret the
344 data from mechanistic studies based on the microbiome. The microbiome features
345 revealed a close association of AF with neurodegenerative diseases as well as cancers,
346 which was supported by prior studies showing that AF had a correlation with AD and
347 PD [40, 41], and AF patients had relatively higher risks of several cancers, including
348 lung cancer and CRC [42, 50]. We also observed that the microbiome features of SLE
349 and CML were highly similar. Interestingly, a tyrosine kinase inhibitor of platelet-
350 derived growth factor receptor, imatinib, was widely used to treat CML and

351 significantly ameliorated survival in murine models of SLE [51]. In addition, a close
352 association between CRC and PD has also been reported in several observational
353 cohorts [52, 53]. Collectively, these findings strongly supported our hypothesis that
354 complex human diseases sharing similar microbiome features might be
355 ~~mechanistically~~mechanically correlated. Furthermore, from the perspectives of risk
356 genes of AF and neurodegenerative diseases, previous GWASs for AF identified two
357 loci at *PITX2* gene-rs6843082 and *C9orf3* gene-rs7026071, which were also
358 associated with a risk of ALS ($p=0.0138$ and $p=0.049$, respectively) [54-56].

359

360 We acknowledge several limitations of our study. First, the participants were of East
361 Asian ancestry; thus, factors such as ethnic differences and gene-environment
362 interactions might make it difficult to generalize the prior results to our study and
363 extrapolate our results to different ethnic populations. Second, although our analysis
364 included participants with the identical by descent (IBD) <0.185 , the vertical
365 transmission of the microbiome from parent to offspring might still partially affect the
366 SNP-based heritability estimates and polygenic scores [20, 57]. Third, genetic factors
367 could explain only a small proportion of the variance in gut microbiome features;
368 thus, the power to detect the causal relationship was limited. Therefore, large-scale
369 studies are warranted to reveal potential relationships between the gut microbiome
370 and complex traits.

371

372 Conclusions

373 In summary, we reveal some causal relationships between ~~the abundance of~~ the gut
374 microbiome and complex human diseases or traits. The disease and gut microbiome
375 feature analysis revealed novel relationships among complex human diseases, which
376 may help reshape our understanding of disease aetiology and provide some clues for
377 extending the clinical indications of existing drugs for different diseases.

378

379 Method**380 Study participants and sample collection**

381 Our study was based on the Guangzhou Nutrition and Health Study (GNHS), with
382 4048 participants (40-75 years old) living in the urban area of Guangzhou city
383 recruited during 2008 and 2013 [17]. We followed up with participants every three
384 years. In the GNHS, stool samples were collected from 1937 participants during
385 follow-up visits. Among those with stool samples, 1717 participants had genetic data,
386 and IBD for 1475 participants was less than 0.185.

387

388 We included 199 participants with both genetic data and gut microbiome data as a
389 replication cohort, which belonged to the control arm of a case-control study of hip

390 fracture with the participants (52-83 years old) recruited between June 2009 and
391 August 2015 in Guangdong Province, China [18].

392

393 Blood samples of all participants were collected after overnight fasting, and the buffy
394 coat was separated from whole blood and stored at -80°C . Stool samples were
395 collected during the on-site visit of the participants at Sun Yat-sen University. All
396 samples were manually stirred, separated into tubes and stored at -80°C within four
397 hours.

398

399 **Genotyping data**

400 For both the discovery and replication cohorts, DNA was extracted from leukocytes
401 using the TIANamp® Blood DNA Kit (DP348, TianGen Biotech Co, Ltd, China)
402 according to the manufacturer's instructions. DNA concentrations were determined
403 using the Qubit quantification system (Thermo Scientific, Wilmington, DE, US).
404 Extracted DNA was stored at -80°C . Genotyping was carried out with Illumina ASA-
405 750K arrays. Quality control and relatedness filters were performed by PLINK1.9
406 [58]. Individuals with a high or low proportion of heterozygous genotypes (outliers
407 defined as 3 standard deviations) were excluded [59]. Individuals who had different
408 ancestries (the first two principal components ± 5 standard deviations from the mean)
409 or related individuals ($\text{IBD} > 0.185$) were excluded [59]. Variants were mapped to the

410 1000 Genomes Phase 3 v5 by SHAPEIT [60, 61], and then we conducted genome-
411 wide genotype imputation with the 1000 Genomes Phase 3 v5 reference panel by
412 Minimac3 [62, 63]. Genetic variants with imputation accuracy RSQR > 0.3 and
413 MAF > 0.05 were included in our analysis. We used the Pan-Asian reference panel,
414 consisting of 502 participants, and SNP2HLA v1.0.3 to impute the HLA region [64-
415 66].

416

417 **Sequencing and processing of 16S rRNA gene data**

418 Microbial DNA was extracted from faecal samples using the QIAamp® DNA Stool
419 Mini Kit per the manufacturer's instructions. DNA concentrations were determined
420 using the Qubit quantification system. The V3-V4 region of the 16S rRNA gene was
421 amplified from genomic DNA using primers 341F (CCTACGGGNGGCWGCAG)
422 and 805R (GACTACHVGGGTATCTAATCC). At the step of amplicon generation, 2
423 µL sterile water was used as negative controls in the PCR reaction system. At the
424 subsequent step of sequencing, no sequencing negative controls were included, since
425 no contamination of PCR products was observed. The pooled amplicons were
426 sequenced using MiSeq Reagent Kits v2 on the Illumina MiSeq System with 2x250
427 bp pair-end sequencing.

428

429 Fastq files were demultiplexed by MiSeq Controller Software. Ultra-fast sequence

430 analysis (USEARCH) was performed to trim the sequence for amplification primers,
431 diversity spacers, sequencing adapters, and merged paired-end reads [67]. The low-
432 quality reads (Phred quality scores \leq 30) were removed. Operational taxonomic units
433 (OTUs) were clustered based on 97% similarity using UPARSE [68]. We removed the
434 OTUs present only in one sample. OTUs were annotated with Greengenes 13_8
435 (<https://greengenes.secondgenome.com/>) [69]. After randomly selecting 10000 reads
436 for each sample, Quantitative Insights into Microbial Ecology (QIIME) software
437 version 1.9.0 was used to calculate alpha diversity (Shannon diversity index, Chao1
438 diversity indices and the observed OTU index and phylogenetic diversity) based on
439 the rarefied OTU counts [70].

440

441 **Statistical analysis**

442 **Proportion of variance explained by all SNPs**

443 We used the GREML method in GCTA to estimate the proportion of variance
444 explained by all SNPs [71]. The taxa were divided into two groups based on whether
445 the taxa were present in ninety percent of participants. Our model was adjusted for
446 age and sex. The power of GREML analysis was calculated with the GCTA power
447 calculator [72].

448

449 **Genome-wide association analysis of gut microbiome features**

450 For each of the GNHS participants and the replication cohort, we clustered
451 participants based on genus-level relative abundance, estimating the JSD distance and
452 PAM clustering algorithm, and then we defined two enterotypes according to the
453 Calinski-Harabasz index [22, 73]. We calculated the genetic principal components of
454 ancestry from genome-wide genetic variants to estimate the population structure.
455 PLINK 1.9 was used to perform a logistic regression model for enterotypes and taxa
456 present in fewer than ninety percent, adjusted for age, sex, sequencing batch and the
457 first five genetic principal components of ancestry.

458

459 For beta diversity, the analysis for the genome-wide host genetic variants with beta
460 diversity was performed using MicrobiomeGWAS [23], adjusted for covariates
461 including the first five genetic principal components of ancestry, age and sex.

462

463 Alpha diversity was calculated after randomly sampling 10000 reads per sample. For
464 the taxa present in no fewer than ninety percent of participants and alpha diversity, we
465 used Z-score normalization to transform the distribution and carried out analysis
466 based on a log-normal model. A mixed linear model-based association (MLMA) test
467 in GCTA was used to assess the association, fitting the first five genetic principal
468 components of ancestry, age, sex and sequencing batch as fixed effects and the effects
469 of all the SNPs as random effects [74-76]. For other taxa present in fewer than ninety

470 percent of participants, we transformed the absence/presence of the taxon into binary
471 variables and used PLINK1.9 to perform a logistic model, adjusted for the first five
472 genetic principal components of ancestry, age, sex and sequencing batch. For all the
473 gut microbiome features, the significance threshold was defined as $5 \times 10^{-8}/n$ (n is the
474 effective number of independent taxa in each taxonomic level) in the discovery stage.
475 QUANTO software was used for power calculations
476 (<http://biostats.usc.edu/Quanto.html>). We estimated genomic inflation factors with
477 LDSC v1.0.1 at the local server [77].

478

479 **Genetic correlation of gut microbiome and traits**

480 We used GCTA to perform a bivariate GREML analysis to estimate the genetic
481 correlation between the gut microbiome and traits in GNHS participants [74, 78]. The
482 gut microbiome was divided into two groups according to the previous description.
483 We used continuous variables for taxa present in no fewer than ninety percent of
484 participants. For taxa present in fewer than ninety percent of participants, we used
485 binary variables according to the absence/presence of taxa. This analysis included
486 traits such as BMI, FBS, HbA1c, SBP, DBP, HDL-C, LDL-C, TC and TG. The power
487 of bivariate GREML analysis was calculated with the GCTA power calculator [72].

488

489

490 **Constructing polygenic scores for taxa and alpha diversity**

491 We selected lead SNPs using PLINK v1.9 with the ‘—clump’ command to clump
492 SNPs with a p value $< 5 \times 10^{-5}$ and $r^2 < 0.1$ within 0.1 cM. We used beta coefficients as
493 the weight to construct polygenic scores for taxa and alpha diversity. For alpha
494 diversity and taxa present in no fewer than ninety percent of participants, we
495 constructed weighted polygenic scores and performed the analysis on a general linear
496 model with a negative binomial distribution to test for association between the
497 polygenic scores and taxa, adjusted for the first five genetic principal components of
498 ancestry, age, sex and sequencing batch. We used weighted polygenic scores and
499 logistic regression to the absence/presence taxa, adjusted for the same covariates as in
500 the above analysis. Taxa with significance ($p < 0.05$) in the replication cohort were
501 included for further analysis.

502

503 **The effective number of independent taxa**

504 As some taxa were correlated with each other, we used an eigendecomposition
505 analysis to calculate the effective number of independent taxa for each taxonomic
506 level [79, 80]. Matrix M is an $m \times n$ matrix, where m is the number of participants and
507 n is the number of total taxa in the corresponding taxonomic level. Matrix A is the
508 variance-covariance matrix of matrix M . P is the matrix of eigenvectors.
509 $\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ is the diagonal matrix composed of the ordered eigenvalues,

510 which can be calculated as

511
$$\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\} = P^{-1}AP$$

512 The effective number of independent taxa can be calculated as

513
$$\frac{(\sum_{i=1}^n \lambda_i)^2}{\sum_{i=1}^n \lambda_i^2}$$

514

515 **Bidirectional MR analysis**

516 In the analysis of the potential causal effect of gut microbiome features on diseases,
517 we used independent genetic variants (selected as part of the polygenic score analysis)
518 as the instrumental variables. For each trait, we excluded instrumental variables that
519 showed a significant association with the trait ($p < 0.05/n$, where n is number of
520 independent genetic variants). In the analysis of the potential causal effect of diseases
521 on gut microbiome features, we selected genetic variants that were replicated in East
522 Asian populations as instrumental variables. As all instrumental variables were from
523 East Asian populations, we chose independent genetic variants ($r^2 < 0.1$) based on the
524 GNHS cohort. We identified the best proxy ($r^2 > 0.9$) based on the GNHS cohort or
525 discarded the variant if no proxy was available. We used the inverse variance
526 weighted (IVW) method to estimate the effect size. To confirm the robustness of the
527 results, we performed three other MR methods, including weighted median, MR-
528 Egger and MR-PRESSO [81-83]. To assess the presence of horizontal pleiotropy, we

529 performed the MR-PRESSO **g**Global test and MR-Egger Regression. The magnitude
530 of the effect of the gut microbiome on traits was dependent on the units of traits
531 (Supplementary Table S1). The results of the effects of complex human diseases on
532 the absence/presence of specific gut microbes are presented as the risk of the presence
533 (vs absence) of the microbe per the log odds difference of the disease. The results of
534 the effects of diseases on other gut microbes were presented as changes in the
535 abundance of taxa (1-SD of log transformed) per the log odds difference of the
536 respective disease.

537

538 The statistical significance of the effects of the gut microbiome on traits and diseases
539 was defined as $p < 0.0008$ ($0.05/62$). In addition, the statistical significance of the
540 effects of diseases on gut microbiome features was defined as $p < 0.05/n$ (where n is
541 the effective number of independent taxa on the corresponding taxonomic level). The
542 results that could not pass Bonferroni adjustment but $p < 0.05$ in all four MR methods
543 were considered potential causal relationships. We performed MR analyses with R
544 v3.5.3.

545

546 **Pathway analysis**

547 We used OTUs by QIIME and annotated the variation of functional genes with
548 Phylogenetic Investigation of Communities by Reconstruction of Unobserved States

549 (PICRUSt) [39]. The pathways and diseases were annotated using KEGG [84-86]. We
550 used Spearman's rank-order correlation to investigate the association of the predicted
551 pathway or disease abundance with AF-associated taxa and genus-level taxa. In the
552 heatmap, diseases were clustered with the 'hcluster' function in R. To test whether the
553 non-normalized pathway or disease abundance were associated with each other, we
554 used SPIEC-EASI to test the interaction relationship and then used Cytoscape v3.7.2
555 to visualize the interaction network [87, 88].

556

557 **Construction of the microbiome risk score**

558 The microbiome risk score was constructed to validate the accuracy of the association
559 between the predicted disease-related gut microbiome features and the corresponding
560 disease. As we have a large sample size for T2DM cases (n=217 cases) in our cohort,
561 we constructed a microbiome risk score of T2DM as an example. We used Spearman's
562 rank-order correlation to select taxa with an absolute value of correlation coefficient
563 higher than 0.2. The score for each taxon abundance in the <5% quantile in our study
564 was defined as 0. For those above 5%, the score for each taxon showing an inverse
565 association with T2DM was defined as -1; the score for each taxon showing a positive
566 association with T2DM was defined as 1. We then summed values from all taxa. We
567 selected a logistic regression model to estimate the association of the MRS with
568 T2DM risk and a linear model to estimate the association of the MRS with the

569 continuous variables, adjusted for age, sex, dietary energy intake, alcohol intake and
570 BMI at the time of sample collection.

571

572 **Clustering diseases**

573 The clustering analysis was carried out with ‘cluster’ and ‘factoextra’ for plot in R.

574 We performed the PAM algorithm based on the predicted abundance of diseases or the

575 average relative abundance after Z-score normalization [89]. The PAM algorithm

576 searches k medoids among the observations and then finds the nearest medoids to

577 minimize the dissimilarity among clusters [90]. Given a set of objects $x = (x_1,$

578 $x_2, \dots, x_n)$, the dissimilarity between objects x_i and x_j is denoted by $d(i, j)$. The

579 assignment of object i to the representative object j is denoted by z_{ij} . z_{ij} is a binary

580 variable and is 1 if object i belongs to the cluster of the representative object j . The

581 function to minimize the model is given by

582
$$\sum_{i=1}^n \sum_{j=1}^n d(i, j) z_{ij}$$

583 To identify the optimal cluster number, we used the ‘pamk’ function in R to determine

584 the optimum average silhouette width. For each object i , we defined N_i as the

585 average dissimilarity between object i and all other objects within its cluster. For the

586 remaining clusters, $b(i, w)$ represents the average dissimilarity between i and all

587 objects in cluster w . The minimum dissimilarity M_i can be calculated by

$$M_i = \min \forall w (b(i, w)).$$

589 The silhouette width for object i can be calculated by

$$sw_i = \frac{M_i - N_i}{\max(M_i, N_i)}$$

591 Then, we calculated the average silhouette width for each object. The cluster number
 592 is determined by the number at which the average silhouette width is maximum. We
 593 estimated the Jaccard similarity coefficient to quantify the cluster difference between
 594 groups. The Jaccard similarity coefficient is positively associated with the similarity
 595 of clusters. Given objects i and j , as well as groups A and B, there are four situations,
 596 as follows:

- 597 (1) S1: in both groups A and B, objects i and j belong to the same cluster;
- 598 (2) S2: in group A, objects i and j belong to the same cluster; in group B, they belong
 599 to different clusters;
- 600 (3) S3: in group A, objects i and j belong to different clusters; in group A, they belong
 601 to the same cluster; and
- 602 (4) S4: In both groups A and B, objects i and j belong to different clusters.

603 a , b , c and d represent the numbers of S1, S2, S3 and S4, respectively. The Jaccard
 604 similarity coefficient can be calculated by the following formula:

$$J = \frac{a}{a + b + c}$$

606

607

608 Availability of data and materials

609 The raw data for 16S rRNA gene sequences are available in the CNSA

610 (<https://db.cngb.org/cnsa/>) of CNGBdb at accession number CNP0000829. Original R

611 scripts are available in GitHub

612 (<https://github.com/hsufengzhe/microbiome/tree/master>). Requests for the metadata

613 from this study can be submitted via email to zhengjusheng@westlake.edu.cn. A

614 proposal is also required for approval.

615

616 **Acknowledgements**

617 We thank the Westlake University Supercomputer Center for providing computing
618 and data analysis services for the present project.

619 **Ethics approval and consent to participate**

620 This study was approved by the Ethics Committee of the School of Public Health at
621 Sun Yat-sen University and Ethics Committee of Westlake University, and all
622 participants provided written informed consent.

623 **Consent for publication**

624 Not applicable.

625 **Competing interests**

626 The authors declare no conflicts of interest.

627 **Authors' contributions**

628 JSZ, YMC and JW initiated and led the study. JY assisted with the data analyses.
629 FZX, YQF and JSZ analysed the data and wrote the manuscript. TYS and CWL
630 collected the data. ZLJ, ZLM, MLS and WLG analysed the data. All authors read and
631 approved the final manuscript.

632 **Funding**

633 This study was funded by the National Natural Science Foundation of China (81903316,
634 81773416), the Zhejiang Ten-thousand Talents Program (101396522001) and the 5010

635 Program for Clinical Research (2007032) of Sun Yat-sen University (Guangzhou,
636 China). Dr Jun Wang was supported by the National Key Research and Development
637 Program of China (grant number 2018YFC2000500), the Strategic Priority Research
638 Program of the Chinese Academy of Sciences (grant number XDB29020000), and the
639 National Natural Science Foundation of China (grant number 31771481 and 91857101).
640 _____

 641 **References**

- 642 1. Awany D, Allali I, Dalvie S, Hemmings S, Mwaikono KS, Thomford NE, et
 643 al. Host and Microbiome Genome-Wide Association Studies: Current State
 644 and Challenges. *Frontiers in genetics*. 2019;9:637-;
 645 doi:10.3389/fgene.2018.00637.
- 646 2. Bull MJ, Plummer NT. Part 1: The Human Gut Microbiome in Health and
 647 Disease. *Integrative medicine (Encinitas, Calif)*. 2014;13(6):17-22.
- 648 3. Lynch JB, Hsiao EY. Microbiomes as sources of emergent host phenotypes.
 649 *Science*. 2019;365(6460):1405-9; doi:10.1126/science.aay0240.
- 650 4. Allegretti JR, Mullish BH, Kelly C, Fischer M. The evolution of the use of
 651 faecal microbiota transplantation and emerging therapeutic indications. *The*
 652 *Lancet*. 2019;394(10196):420-31; doi:10.1016/S0140-6736(19)31266-8.
- 653 5. Wong SH, Yu J. Gut microbiota in colorectal cancer: mechanisms of action
 654 and clinical applications. *Nature Reviews Gastroenterology & Hepatology*.
 655 2019; doi:10.1038/s41575-019-0209-8.
- 656 6. Davies NM, Holmes MV, Davey Smith G. Reading Mendelian randomisation
 657 studies: a guide, glossary, and checklist for clinicians. *BMJ*. 2018;362:k601;
 658 doi:10.1136/bmj.k601.
- 659 7. Turpin W, Espin-Garcia O, Xu W, Silverberg MS, Kevans D, Smith MI, et al.
 660 Association of host genome with intestinal microbial composition in a large
 661 healthy cohort. *Nat Genet*. 2016;48(11):1413-7; doi:10.1038/ng.3693.
- 662 8. Wang J, Thingholm LB, Skieceviciene J, Rausch P, Kummén M, Hov JR, et al.
 663 Genome-wide association analysis identifies variation in vitamin D receptor
 664 and other host factors influencing the gut microbiota. *Nat Genet*.
 665 2016;48(11):1396-406; doi:10.1038/ng.3695.
- 666 9. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, et al.
 667 Human genetics shape the gut microbiome. *Cell*. 2014;159(4):789-99;
 668 doi:10.1016/j.cell.2014.09.053.
- 669 10. Goodrich JK, Davenport ER, Beaumont M, Jackson MA, Knight R, Ober C, et
 670 al. Genetic Determinants of the Gut Microbiome in UK Twins. *Cell Host*
 671 *Microbe*. 2016;19(5):731-43; doi:10.1016/j.chom.2016.04.017.
- 672 11. Ganesan K, Chung SK, Vanamala J, Xu B. Causal Relationship between Diet-
 673 Induced Gut Microbiota Changes and Diabetes: A Novel Strategy to
 674 Transplant *Faecalibacterium prausnitzii* in Preventing Diabetes. *Int J Mol Sci*.
 675 2018;19(12); doi:10.3390/ijms19123720.
- 676 12. He Y, Wu W, Zheng HM, Li P, McDonald D, Sheng HF, et al. Regional
 677 variation limits applications of healthy gut microbiome reference ranges and
 678 disease models. *Nat Med*. 2018;24(10):1532-5; doi:10.1038/s41591-018-0164-
 679 x.
- 680 13. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, et
 681 al. Environment dominates over host genetics in shaping human gut
 682 microbiota. *Nature*. 2018;555(7695):210-5; doi:10.1038/nature25973.

-
- 683 14. Duvallet C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut
684 microbiome studies identifies disease-specific and shared responses. *Nature*
685 *communications*. 2017;8(1):1784-; doi:10.1038/s41467-017-01973-8.
- 686 15. Cheng S, Han B, Ding M, Wen Y, Ma M, Zhang L, et al. Identifying
687 psychiatric disorder-associated gut microbiota using microbiota-related gene
688 set enrichment analysis. *Briefings in Bioinformatics*. 2019;
689 doi:10.1093/bib/bbz034.
- 690 16. Jackson MA, Verdi S, Maxan M-E, Shin CM, Zierer J, Bowyer RCE, et al.
691 Gut microbiota associations with common diseases and prescription
692 medications in a population-based cohort. *Nature Communications*.
693 2018;9(1):2655; doi:10.1038/s41467-018-05184-7.
- 694 17. Cao Y, Wang C, Guan K, Xu Y, Su Y-X, Chen YM. Association of magnesium
695 in serum and urine with carotid intima-media thickness and serum lipids in
696 middle-aged and elderly Chinese: a community-based cross-sectional study.
697 *European journal of nutrition*. 2015;55; doi:10.1007/s00394-015-0839-8.
- 698 18. Sun L-L, Li B-L, Xie H-L, Fan F, Yu W-Z, Wu B-H, et al. Associations
699 between the dietary intake of antioxidant nutrients and the risk of hip fracture
700 in elderly Chinese: A case-control study. *The British journal of nutrition*.
701 2014;112:1-9; doi:10.1017/S0007114514002773.
- 702 19. Lim MY, You HJ, Yoon HS, Kwon B, Lee JY, Lee S, et al. The effect of
703 heritability and host genetics on the gut microbiota and metabolic syndrome.
704 *Gut*. 2017;66(6):1031-8; doi:10.1136/gutjnl-2015-311326.
- 705 20. Davenport ER. Elucidating the role of the host genome in shaping microbiome
706 composition. *Gut microbes*. 2016;7(2):178-84;
707 doi:10.1080/19490976.2016.1155022.
- 708 21. Davenport ER, Cusanovich DA, Michelini K, Barreiro LB, Ober C, Gilad Y.
709 Genome-Wide Association Studies of the Human Gut Microbiota. *PloS one*.
710 2015;10(11):e0140301-e; doi:10.1371/journal.pone.0140301.
- 711 22. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, et al.
712 Enterotypes of the human gut microbiome. *Nature*. 2011;473(7346):174-80;
713 doi:10.1038/nature09944.
- 714 23. Hua X, Song L, Yu G, Goedert JJ, Abnet CC, Landi MT, et al.
715 MicrobiomeGWAS: a tool for identifying host genetic variants associated with
716 microbiome composition. *bioRxiv*. 2015:031187; doi:10.1101/031187.
- 717 24. Ruhlemann MC, Degenhardt F, Thingholm LB, Wang J, Skieceviciene J,
718 Rausch P, et al. Application of the distance-based F test in an mGWAS
719 investigating beta diversity of intestinal microbiota identifies variants in
720 SLC9A8 (NHE8) and 3 other loci. *Gut Microbes*. 2018;9(1):68-75;
721 doi:10.1080/19490976.2017.1356979.
- 722 25. Bonder MJ, Kurilshikov A, Tigchelaar EF, Mujagic Z, Imhann F, Vila AV, et al.
723 The effect of host genetics on the gut microbiome. *Nat Genet*.
724 2016;48(11):1407-12; doi:10.1038/ng.3663.

-
- 725 26. Tang WHW, Kitai T, Hazen SL. Gut Microbiota in Cardiovascular Health and
726 Disease. *Circulation research*. 2017;120(7):1183-96;
727 doi:10.1161/CIRCRESAHA.117.309715.
- 728 27. Low SK, Takahashi A, Ebana Y, Ozaki K, Christophersen IE, Ellinor PT, et al.
729 Identification of six new genetic loci associated with atrial fibrillation in the
730 Japanese population. *Nat Genet*. 2017;49(6):953-8; doi:10.1038/ng.3842.
- 731 28. Suzuki K, Akiyama M, Ishigaki K, Kanai M, Hosoe J, Shojima N, et al.
732 Identification of 28 new susceptibility loci for type 2 diabetes in the Japanese
733 population. *Nat Genet*. 2019;51(3):379-86; doi:10.1038/s41588-018-0332-4.
- 734 29. Akiyama M, Okada Y, Kanai M, Takahashi A, Momozawa Y, Ikeda M, et al.
735 Genome-wide association study identifies 112 new loci for body mass index in
736 the Japanese population. *Nat Genet*. 2017;49(10):1458-67;
737 doi:10.1038/ng.3951.
- 738 30. Kanai M, Akiyama M, Takahashi A, Matoba N, Momozawa Y, Ikeda M, et al.
739 Genetic analysis of quantitative traits in the Japanese population links cell
740 types to complex human diseases. *Nat Genet*. 2018;50(3):390-400;
741 doi:10.1038/s41588-018-0047-6.
- 742 31. Matoba N, Akiyama M, Ishigaki K, Kanai M, Takahashi A, Momozawa Y, et
743 al. GWAS of smoking behaviour in 165,436 Japanese people reveals seven
744 new loci and shared genetic architecture. *Nat Hum Behav*. 2019;3(5):471-7;
745 doi:10.1038/s41562-019-0557-y.
- 746 32. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, et al. Genetics of
747 rheumatoid arthritis contributes to biology and drug discovery. *Nature*.
748 2014;506(7488):376-81; doi:10.1038/nature12873.
- 749 33. Lu XF, Wang LY, Chen SF, He L, Yang XL, Shi YY, et al. Genome-wide
750 association study in Han Chinese identifies four new susceptibility loci for
751 coronary artery disease. *Nature Genetics*. 2012;44(8):890-+;
752 doi:10.1038/ng.2337.
- 753 34. Marzec J, Mao X, Li M, Wang M, Feng N, Gou X, et al. A genetic study and
754 meta-analysis of the genetic predisposition of prostate cancer in a Chinese
755 population. *Oncotarget*. 2016;7(16):21393-403; doi:10.18632/oncotarget.7250.
- 756 35. Okada Y, Sim X, Go MJ, Wu JY, Gu D, Takeuchi F, et al. Meta-analysis
757 identifies multiple loci associated with kidney function-related traits in east
758 Asian populations. *Nat Genet*. 2012;44(8):904-9; doi:10.1038/ng.2352.
- 759 36. Zeng C, Matsuda K, Jia WH, Chang J, Kweon SS, Xiang YB, et al.
760 Identification of Susceptibility Loci and Genes for Colorectal Cancer Risk.
761 *Gastroenterology*. 2016;150(7):1633-45; doi:10.1053/j.gastro.2016.02.076.
- 762 37. Zhou X, Chen Y, Mok KY, Zhao Q, Chen K, Chen Y, et al. Identification of
763 genetic risk factors in the Chinese population implicates a role of immune
764 system in Alzheimer's disease pathogenesis. *Proceedings of the National
765 Academy of Sciences*. 2018;115(8):1697; doi:10.1073/pnas.1715554115.
- 766 38. Gan W, Walters RG, Holmes MV, Bragg F, Millwood IY, Banasik K, Chen Y,

-
- 767 Du H, Iona A, Mahajan A, et al: Evaluation of type 2 diabetes genetic risk
768 variants in Chinese adults: findings from 93,000 individuals from the China
769 Kadoorie Biobank. *Diabetologia*. 2016;59(7):1446-1457.
- 770 39. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA,
771 et al. Predictive functional profiling of microbial communities using 16S
772 rRNA marker gene sequences. *Nature Biotechnology*. 2013;31:814;
773 doi:10.1038/nbt.2676.
- 774 40. Canga Y, Emre A, Yuksel GA, Karatas MB, Yelgec NS, Gurkan U, et al.
775 Assessment of Atrial Conduction Times in Patients with Newly Diagnosed
776 Parkinson's Disease. *Parkinsons Dis*. 2018;2018:2916905;
777 doi:10.1155/2018/2916905.
- 778 41. Ihara M, Washida K. Linking Atrial Fibrillation with Alzheimer's Disease:
779 Epidemiological, Pathological, and Mechanistic Evidence. *J Alzheimers Dis*.
780 2018;62(1):61-72; doi:10.3233/JAD-170970.
- 781 42. Conen D, Wong JA, Sandhu RK, Cook NR, Lee I-M, Buring JE, et al. Risk of
782 Malignant Cancer Among Women With New-Onset Atrial Fibrillation. *JAMA*
783 *Cardiology*. 2016;1(4):389-96; doi:10.1001/jamacardio.2016.0280.
- 784 43. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation
785 of cluster analysis. *Journal of Computational and Applied Mathematics*.
786 1987;20:53-65.
- 787 44. Blekhman R, Goodrich JK, Huang K, Sun Q, Bukowski R, Bell JT, et al. Host
788 genetic variation impacts microbiome composition across human body sites.
789 *Genome Biol*. 2015;16:191; doi:10.1186/s13059-015-0759-1.
- 790 45. Goodrich JK, Davenport ER, Clark AG, Ley RE. The Relationship Between
791 the Human Genome and Microbiome Comes into View. *Annu Rev Genet*.
792 2017;51:413-33; doi:10.1146/annurev-genet-110711-155532.
- 793 46. Kuehbachner T, Rehman A, Lepage P, Hellmig S, Fölsch UR, Schreiber S, et al.
794 Intestinal TM7 bacterial phylogenies in active inflammatory bowel disease.
795 *Journal of Medical Microbiology*. 2008;57(12):1569-76.
- 796 47. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, et al. Cultivation
797 of a human-associated TM7 phylotype reveals a reduced genome and epibiotic
798 parasitic lifestyle. *Proceedings of the National Academy of Sciences of the*
799 *United States of America*. 2015;112(1):244-9; doi:10.1073/pnas.1419038112.
- 800 48. Bor B, Bedree JK, Shi W, McLean JS, He X. Saccharibacteria (TM7) in the
801 Human Oral Microbiome. *Journal of Dental Research*. 2019;98(5):500-9;
802 doi:10.1177/0022034519831671.
- 803 49. Cirulli ET, Goldstein DB. Uncovering the roles of rare variants in common
804 disease through whole-genome sequencing. *Nat Rev Genet*. 2010;11(6):415-
805 25; doi:10.1038/nrg2779.
- 806 50. Vinter N, Christesen Amanda MS, Fenger-Grøn M, Tjønneland A, Frost L.
807 Atrial Fibrillation and Risk of Cancer: A Danish Population-Based Cohort
808 Study. *Journal of the American Heart Association*. 2018;

-
- 809 doi:10.1161/JAHA.118.009543.
- 810 51. Zoja C, Corna D, Rottoli D, Zanchi C, Abbate M, Remuzzi G. Imatinib
811 ameliorates renal disease and survival in murine lupus autoimmune disease.
812 *Kidney International*. 2006;70(1):97-103.
- 813 52. Boursi B, Mamtani R, Haynes K, Yang Y-X. Parkinson's disease and colorectal
814 cancer risk-A nested case control study. *Cancer Epidemiol*. 2016;43:9-14;
815 doi:10.1016/j.canep.2016.05.007.
- 816 53. Xie X, Luo X, Xie M. Association between Parkinson's disease and risk of
817 colorectal cancer. *Parkinsonism & Related Disorders*. 2017;35:42-7;
818 doi:10.1016/j.parkreldis.2016.11.011.
- 819 54. van Rheenen W, Shatunov A, Dekker AM, McLaughlin RL, Diekstra FP, Pulit
820 SL, et al. Genome-wide association analyses identify new risk variants and the
821 genetic architecture of amyotrophic lateral sclerosis. *Nat Genet*.
822 2016;48(9):1043-8; doi:10.1038/ng.3622.
- 823 55. Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C,
824 et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci
825 for Alzheimer's disease. *Nat Genet*. 2013;45(12):1452-8; doi:10.1038/ng.2802.
- 826 56. Pankratz N, Beecham GW, DeStefano AL, Dawson TM, Doheny KF, Factor
827 SA, et al. Meta-analysis of Parkinson's disease: identification of a novel locus,
828 RIT2. *Ann Neurol*. 2012;71(3):370-84; doi:10.1002/ana.22687.
- 829 57. Zhao L, Wang G, Siegel P, He C, Wang H, Zhao W, et al. Quantitative genetic
830 background of the host influences gut microbiomes in chickens. *Sci Rep*.
831 2013;3:1163-; doi:10.1038/srep01163.
- 832 58. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al.
833 PLINK: a tool set for whole-genome association and population-based linkage
834 analyses. *American journal of human genetics*. 2007;81(3):559-75;
835 doi:10.1086/519795.
- 836 59. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan
837 KT. Data quality control in genetic case-control association studies. *Nat*
838 *Protoc*. 2010;5(9):1564-73; doi:10.1038/nprot.2010.116.
- 839 60. Delaneau O, Marchini J, Zagury J-F. A linear complexity phasing method for
840 thousands of genomes. *Nature Methods*. 2011;9:179; doi:10.1038/nmeth.1785.
- 841 61. Delaneau O, Marchini J, The Genomes Project C, McVean GA, Donnelly P,
842 Lunter G, et al. Integrating sequence and array data to create an improved
843 1000 Genomes Project haplotype reference panel. *Nature Communications*.
844 2014;5:3934; doi:10.1038/ncomms4934.
- 845 62. Das S, Forer L, Schön herr S, Sidore C, Locke AE, Kwong A, et al. Next-
846 generation genotype imputation service and methods. *Nature Genetics*.
847 2016;48:1284; doi:10.1038/ng.3656.
- 848 63. Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, et al. The
849 international Genome sample resource (IGSR): A worldwide collection of
850 genome variation incorporating the 1000 Genomes Project data. *Nucleic Acids*

-
- 851 Research. 2016; doi:10.1093/nar/gkw829.
- 852 64. Okada Y, Kim K, Han B, Pillai NE, Ong RT, Saw WY, et al. Risk for ACPA-
853 positive rheumatoid arthritis is driven by shared HLA amino acid
854 polymorphisms in Asian and European populations. *Hum Mol Genet.*
855 2014;23(25):6916-26; doi:10.1093/hmg/ddu387.
- 856 65. Pillai NE, Okada Y, Saw WY, Ong RT, Wang X, Tantoso E, et al. Predicting
857 HLA alleles from high-resolution SNP data in three Southeast Asian
858 populations. *Hum Mol Genet.* 2014;23(16):4443-51;
859 doi:10.1093/hmg/ddu149.
- 860 66. Jia X, Han B, Onengut-Gumuscu S, Chen W-M, Concannon PJ, Rich SS, et al.
861 Imputing Amino Acid Polymorphisms in Human Leukocyte Antigens. *PLOS*
862 *ONE.* 2013;8(6):e64683; doi:10.1371/journal.pone.0064683.
- 863 67. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
864 *Bioinformatics.* 2010;26(19):2460-1; doi:10.1093/bioinformatics/btq461.
- 865 68. Edgar RC. UPARSE: highly accurate OTU sequences from microbial
866 amplicon reads. *Nat Methods.* 2013;10(10):996-8; doi:10.1038/nmeth.2604.
- 867 69. Second Genome, Inc: the Greengenes
868 Databases.<http://greengenes.secondgenome.com/>. Accessed 12 Mar 2019.
- 869 70. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello
870 EK, et al. QIIME allows analysis of high-throughput community sequencing
871 data. *Nat Methods.* 2010;7(5):335-6; doi:10.1038/nmeth.f.303.
- 872 71. Lee SH, Wray NR, Goddard ME, Visscher PM. Estimating missing heritability
873 for disease from genome-wide association studies. *Am J Hum Genet.*
874 2011;88(3):294-305; doi:10.1016/j.ajhg.2011.02.002.
- 875 72. Visscher PM, Hemani G, Vinkhuyzen AAE, Chen G-B, Lee SH, Wray NR, et
876 al. Statistical Power to Detect Genetic (Co)Variance of Complex Traits Using
877 SNP Data in Unrelated Samples. *PLOS Genetics.* 2014;10(4):e1004269;
878 doi:10.1371/journal.pgen.1004269.
- 879 73. Caliński T, Harabasz J. A dendrite method for cluster analysis.
880 *Communications in Statistics.* 1974;3(1):1-27;
881 doi:10.1080/03610927408827101.
- 882 74. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide
883 complex trait analysis. *Am J Hum Genet.* 2011;88(1):76-82;
884 doi:10.1016/j.ajhg.2010.11.011.
- 885 75. Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and
886 pitfalls in the application of mixed-model association methods. *Nat Genet.*
887 2014;46(2):100-6; doi:10.1038/ng.2876.
- 888 76. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al.
889 Common SNPs explain a large proportion of the heritability for human height.
890 *Nature Genetics.* 2010;42:565; doi:10.1038/ng.608.
- 891 77. Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Patterson N, et
892 al. LD Score regression distinguishes confounding from polygenicity in

-
- 893 genome-wide association studies. *Nature Genetics*. 2015;47(3):291-5;
894 doi:10.1038/ng.3211.
- 895 78. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR. Estimation of
896 pleiotropy between complex diseases using single-nucleotide polymorphism-
897 derived genomic relationships and restricted maximum likelihood.
898 *Bioinformatics*. 2012;28(19):2540-2; doi:10.1093/bioinformatics/bts474.
- 899 79. Wang H, Zhang F, Zeng J, Wu Y, Kemper KE, Xue A, et al. Genotype-by-
900 environment interactions inferred from genetic effects on phenotypic
901 variability in the UK Biobank. *Science Advances*. 2019;5(8):eaaw3538;
902 doi:10.1126/sciadv.aaw3538.
- 903 80. Bretherton CS, Widmann M, Dymnikov VP, Wallace JM, Bladé I. The
904 Effective Number of Spatial Degrees of Freedom of a Time-Varying Field.
905 *Journal of Climate*. 1999;12(7):1990-2009; doi:10.1175/1520-
906 0442(1999)012<1990:Tenosd>2.0.Co;2.
- 907 81. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in
908 Mendelian Randomization with Some Invalid Instruments Using a Weighted
909 Median Estimator. *Genet Epidemiol*. 2016;40(4):304-14;
910 doi:10.1002/gepi.21965.
- 911 82. Burgess S, Thompson SG. Interpreting findings from Mendelian
912 randomization using the MR-Egger method. *Eur J Epidemiol*. 2017;32(5):377-
913 89; doi:10.1007/s10654-017-0255-x.
- 914 83. Verbanck M, Chen C-Y, Neale B, Do R. Detection of widespread horizontal
915 pleiotropy in causal relationships inferred from Mendelian randomization
916 between complex traits and diseases. *Nature Genetics*. 2018;50(5):693-8;
917 doi:10.1038/s41588-018-0099-7.
- 918 84. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes.
919 *Nucleic Acids Res*. 2000;28(1):27-30; doi:10.1093/nar/28.1.27.
- 920 85. Kanehisa M, Sato Y, Furumichi M, Morishima K, Tanabe M. New approach
921 for understanding genome variations in KEGG. *Nucleic Acids Res*. 2019;
922 doi:10.1093/nar/gky962.
- 923 86. Kanehisa M. Toward understanding the origin and evolution of cellular
924 organisms. *Protein Sci*. 2019; doi:10.1002/pro.3715.
- 925 87. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al.
926 Cytoscape: a software environment for integrated models of biomolecular
927 interaction networks. *Genome Res*. 2003;13(11):2498-504;
928 doi:10.1101/gr.1239303.
- 929 88. Kurtz ZD, Müller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA.
930 Sparse and Compositionally Robust Inference of Microbial Ecological
931 Networks. *PLOS Computational Biology*. 2015;
932 doi:10.1371/journal.pcbi.1004226.
- 933 89. Reynolds AP, Richards G, de la Iglesia B, Rayward-Smith VJ. Clustering
934 Rules: A Comparison of Partitioning and Hierarchical Clustering Algorithms.

-
- 935 Journal of Mathematical Modelling and Algorithms. 2006;5(4):475-504;
936 doi:10.1007/s10852-005-9022-1.
937 90. Kaufman L, Rousseeuw P. Partitioning Around Medoids (Program PAM).
938 John Wiley & Sons, Inc; 1990. p. 68-125.

939 **Figure legends**

940 **Figure 1 Study overview.** The figure shows the highlights of our study. First, we
941 performed a microbiome genome-wide association study in a Chinese population
942 (Step A). We validated significant genetic variants reported in previous studies and
943 replicated our results in an independent cohort. Second, we investigated the causal
944 relationship between the gut microbiome and complex human diseases using host
945 genetics as instrumental variables for bidirectional Mendelian randomization (MR)
946 analysis (Step B). For the analysis of the effects of the gut microbiome on complex
947 traits, we used publicly available GWAS summary statistics of complex traits (n=58)
948 and diseases (type 2 diabetes mellitus (T2DM), atrial fibrillation (AF), colorectal
949 cancer (CRC) and rheumatoid arthritis) reported by BioBank Japan [27-32]. For the
950 ~~reverse~~ MR analyses, the diseases of interest included T2DM (cases: 7,109;
951 non-cases: 86,022), AF (cases: 8,180; non-cases: 28,612), coronary artery disease
952 (cases: 1,515; non-cases: 5,019), chronic kidney disease (n=71,149), Alzheimer's
953 disease (cases: 477; non-cases: 442), CRC (cases: 8,027; non-cases: 22,577) and
954 prostatic cancer (cases: 495; non-cases: 640) reported in the previous large-scale
955 GWASs in East Asians [27,33-38]. Finally, we identified common and distinct gut
956 microbiome features across different diseases (Step C).

957

958 **Figure 2 The SNP-based heritability of the gut microbiome.** The plot shows the
959 taxa with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, where n
960 is the effective number of independent taxa in each taxonomic level.

961

962 **Figure 3 Effect of host genetically predicted higher atrial fibrillation risk on the**
963 **gut microbiome. (A).** Causal association of atrial fibrillation with the abundance of
964 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Bacteroides coprophilus*. The
965 magnitude of the effect of atrial fibrillation on taxa is dependent on changes in the
966 abundance of bacteria (1-SD of the log-transformed abundance) per genetically
967 determined higher log odds of atrial fibrillation. **(B).** Causal association of atrial
968 fibrillation with the presence of *Barnesiellaceae*, ~~an undefined genus in the~~
969 ~~family undefined genus in family~~ *Veillonellaceae* and *Mitsuokella*. The magnitude of
970 the effect of atrial fibrillation on taxa is presented as an odds ratio increase in the log
971 odds of atrial fibrillation.

972

973 **Figure 4 Association and cluster of diseases predicted by the gut microbiome.**

974 **(A)**. Plot of clusters in the Guangzhou Nutrition and Health Study (GNHS) cohort
975 (n=1919). **(B)**. Plot of cluster results in the replication cohort (n=217). **(C)**. Plot of 5
976 clusters in antibiotic-taking participants (n=18). The optimal cluster was 5 in the
977 GNHS cohort and 6 in the replication cohort. The clusters share consistent
978 components between the two studies. In contrast, components are different between
979 antibiotic-taking participants and control groups. Dimension1 (Dim1) and dimension2
980 (Dim2) explained 40.1% and 13.1% of the variance, respectively, in the GNHS
981 cohort. The annotation for variables is as follows. AT: African trypanosomiasis, AD:
982 Alzheimer's disease, V1: Amoebiasis, ALS: Amyotrophic lateral sclerosis, BC:
983 Bladder cancer, CD: Chagas disease, CML: Chronic myeloid leukaemia, CRC:
984 Colorectal cancer, V2: Hepatitis C, HD: Huntington's disease, HCM: Hypertrophic
985 cardiomyopathy, V3: Influenza A, PD: Parkinson's disease, V4: Pathways in cancer,
986 V5: Prion disease, PCa: Prostate cancer, RCC: Renal cell carcinoma, SLE: Systemic
987 lupus erythematosus, V6: Tuberculosis, T1DM: Type I diabetes mellitus, T2DM:
988 Type II diabetes mellitus, V7: *Vibrio cholerae* infection. **(D)**. Gut microbiome-
989 predicted network of relationships among different ~~human~~-complex human diseases.
990 The relationship between diseases is determined by SPIEC-EASI with non-
991 normalized predicted abundance data. The diseases that shared the same edge had the
992 gut microbiome-predicted correlation.

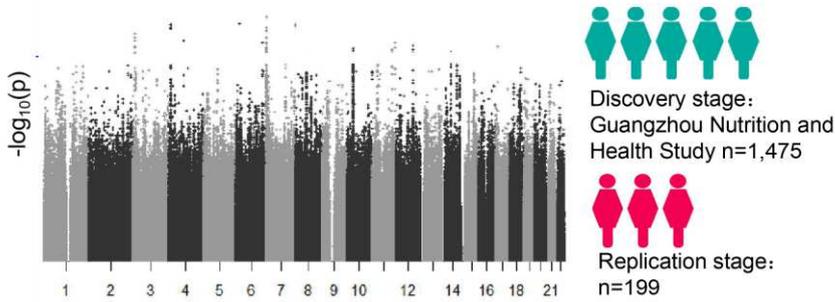
993

994 **Figure 5 Correlation of ~~human~~-complex human diseases with the gut**
995 **microbiome.** The heatmap shows Spearman's correlation of predicted diseases and
996 the gut microbiome at the genus level. The grey components show no significant
997 correlation with Bonferroni correction ($p > 0.05 / (5.6 * 22)$, $p > 0.0004$).

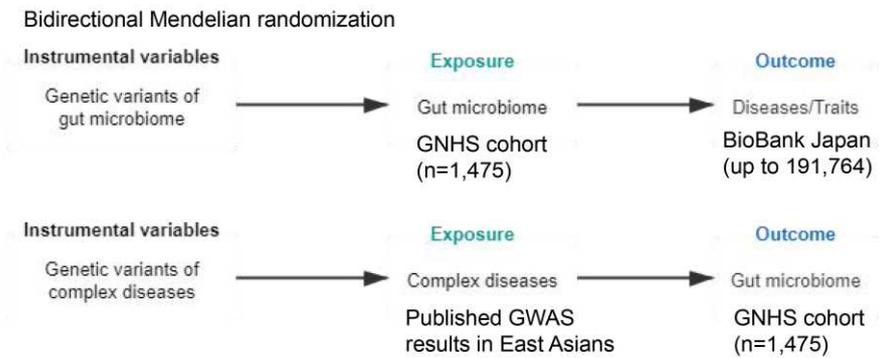
998

999 **Figure 1 Study overview.** The figure shows the highlights of our study. First, we
1000 performed a microbiome genome-wide association study in a Chinese population
1001 (Step A). We validated significant genetic variants reported in previous studies and
1002 replicated our results in an independent cohort. Second, we investigated the causal
1003 relationship between the gut microbiome and complex human diseases using host
1004 genetics as instrumental variables for bidirectional Mendelian randomization (MR)
1005 analysis (Step B). For the analysis of the effects of the gut microbiome on complex
1006 traits, we used publicly available GWAS summary statistics of complex traits (n=58)
1007 and diseases (type 2 diabetes mellitus (T2DM), atrial fibrillation (AF), colorectal
1008 cancer (CRC) and rheumatoid arthritis) reported by BioBank Japan [27-32]. For the
1009 ~~reverse~~ MR analyses, the diseases of interest included T2DM (cases: 7,109;
1010 non-cases: 86,022), AF (cases: 8,180; non-cases: 28,612), coronary artery disease
1011 (cases: 1,515; non-cases: 5,019), chronic kidney disease (n=71,149), Alzheimer's
1012 disease (cases: 477; non-cases: 442), CRC (cases: 8,027; non-cases: 22,577) and
1013 prostatic cancer (cases: 495; non-cases: 640) reported in the previous large-scale
1014 GWASs in East Asians [27, 33-38]. Finally, we identified common and distinct gut
1015 microbiome features across different diseases (Step C).

A. Association of host genetics with gut microbiome in a Chinese population.

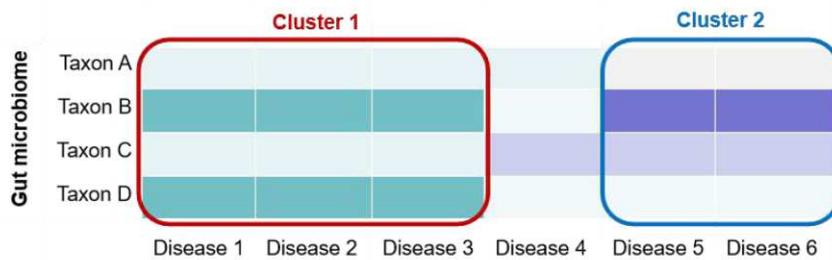


B. The causal relationships between gut microbiome and complex human diseases.

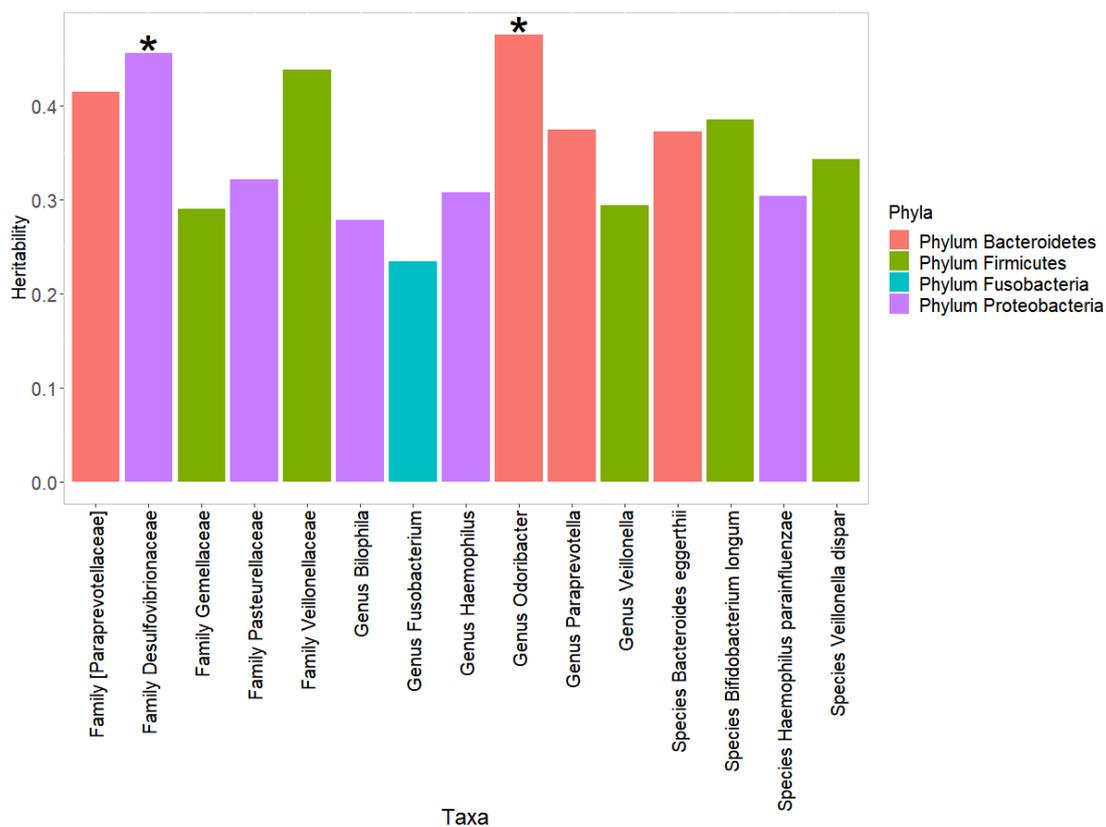


C. The shared and distinct microbiome features among complex human diseases.

1,919 participants from GNHS cohort.



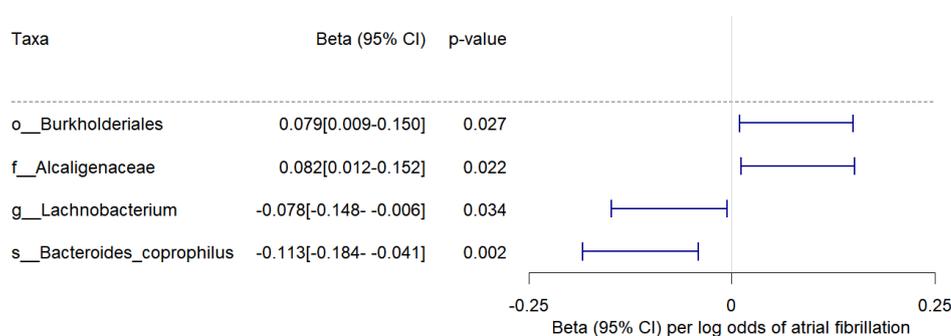
1017 **Figure 2 The SNP-based heritability of the gut microbiome.** The plot shows the
 1018 taxa with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, where n
 1019 is the effective number of independent taxa in each taxonomic level.



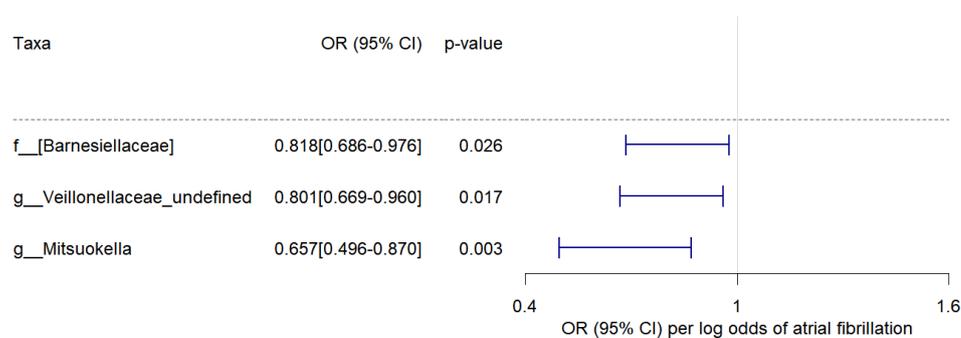
1020

1021 **Figure 3 Effect of host genetically predicted higher atrial fibrillation risk on the**
 1022 **gut microbiome. (A).** Causal association of atrial fibrillation with the abundance of
 1023 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Bacteroides coprophilus*. The
 1024 magnitude of the effect of atrial fibrillation on taxa is dependent on changes in the
 1025 abundance of bacteria (1-SD of the log-transformed abundance) per genetically
 1026 determined higher log odds of atrial fibrillation. **(B).** Causal association of atrial
 1027 fibrillation with the presence of *Barnesiellaceae*, an undefined genus in the family
 1028 ~~undefined genus in family~~ *Veillonellaceae* and *Mitsuokella*. The magnitude of the
 1029 effect of atrial fibrillation on taxa is presented as an odds ratio increase in the log odds
 1030 of atrial fibrillation.

A



B

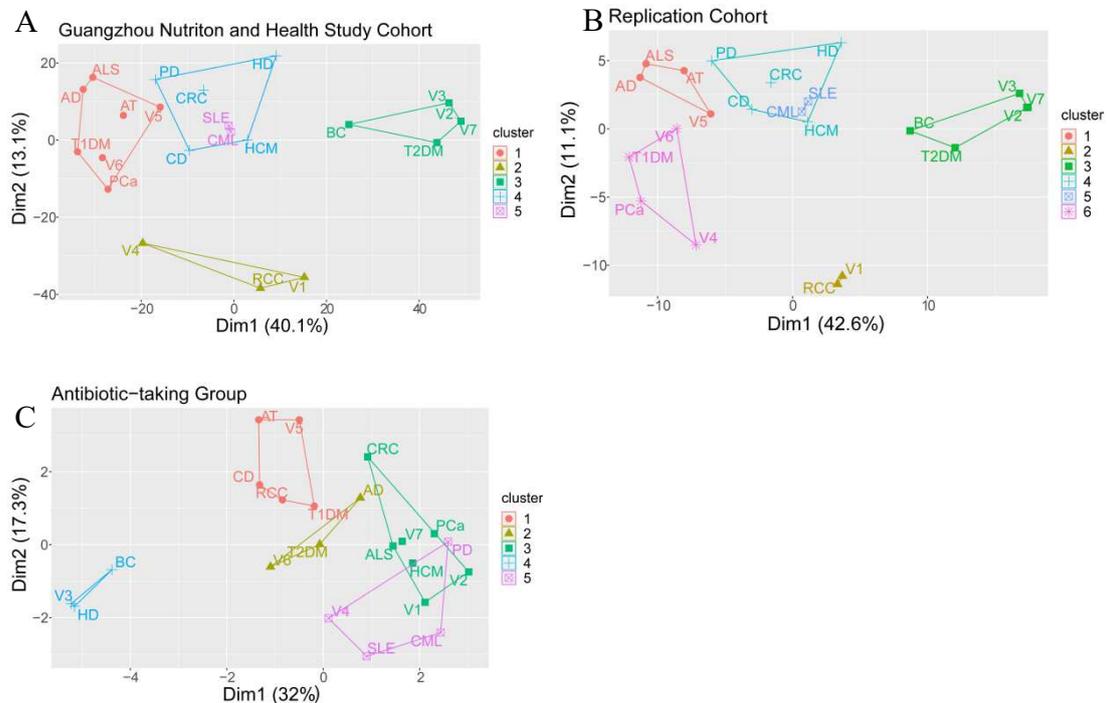


1031

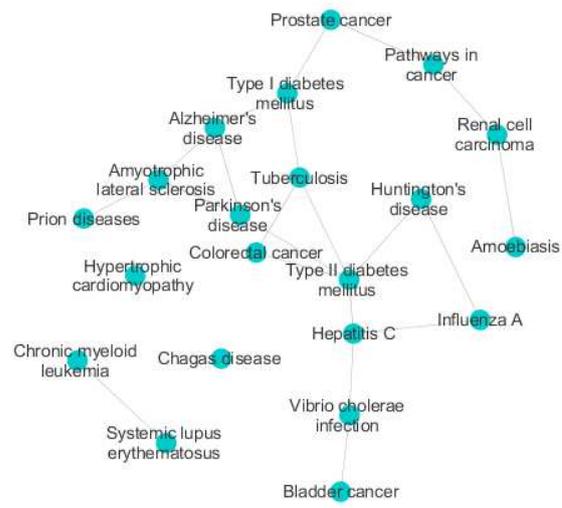
1032

1033

1034 **Figure 4 Association and cluster of diseases predicted by the gut microbiome.**
 1035 **(A).** Plot of clusters in the Guangzhou Nutrition and Health Study (GNHS) cohort
 1036 (n=1919). **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5
 1037 clusters in antibiotic-taking participants (n=18). The optimal cluster was 5 in the
 1038 GNHS cohort and 6 in the replication cohort. The clusters share consistent
 1039 components between the two studies. In contrast, components are different between
 1040 antibiotic-taking participants and control groups. Dimension1 (Dim1) and dimension2
 1041 (Dim2) explained 40.1% and 13.1% of the variance, respectively, in the GNHS
 1042 cohort. The annotation for variables is as follows. AT: African trypanosomiasis, AD:
 1043 Alzheimer's disease, V1: Amoebiasis, ALS: Amyotrophic lateral sclerosis, BC:
 1044 Bladder cancer, CD: Chagas disease, CML: Chronic myeloid leukaemia, CRC:
 1045 Colorectal cancer, V2: Hepatitis C, HD: Huntington's disease, HCM: Hypertrophic
 1046 cardiomyopathy, V3: Influenza A, PD: Parkinson's disease, V4: Pathways in cancer,
 1047 V5: Prion disease, PCa: Prostate cancer, RCC: Renal cell carcinoma, SLE: Systemic
 1048 lupus erythematosus, V6: Tuberculosis, T1DM: Type I diabetes mellitus, T2DM:
 1049 Type II diabetes mellitus, V7: *Vibrio cholerae* infection. **(D).** Gut microbiome-
 1050 predicted network of relationships among different human-complex human diseases.
 1051 The relationship between diseases is determined by SPIEC-EASI with non-
 1052 normalized predicted abundance data. The diseases that shared the same edge had the
 1053 gut microbiome-predicted correlation.



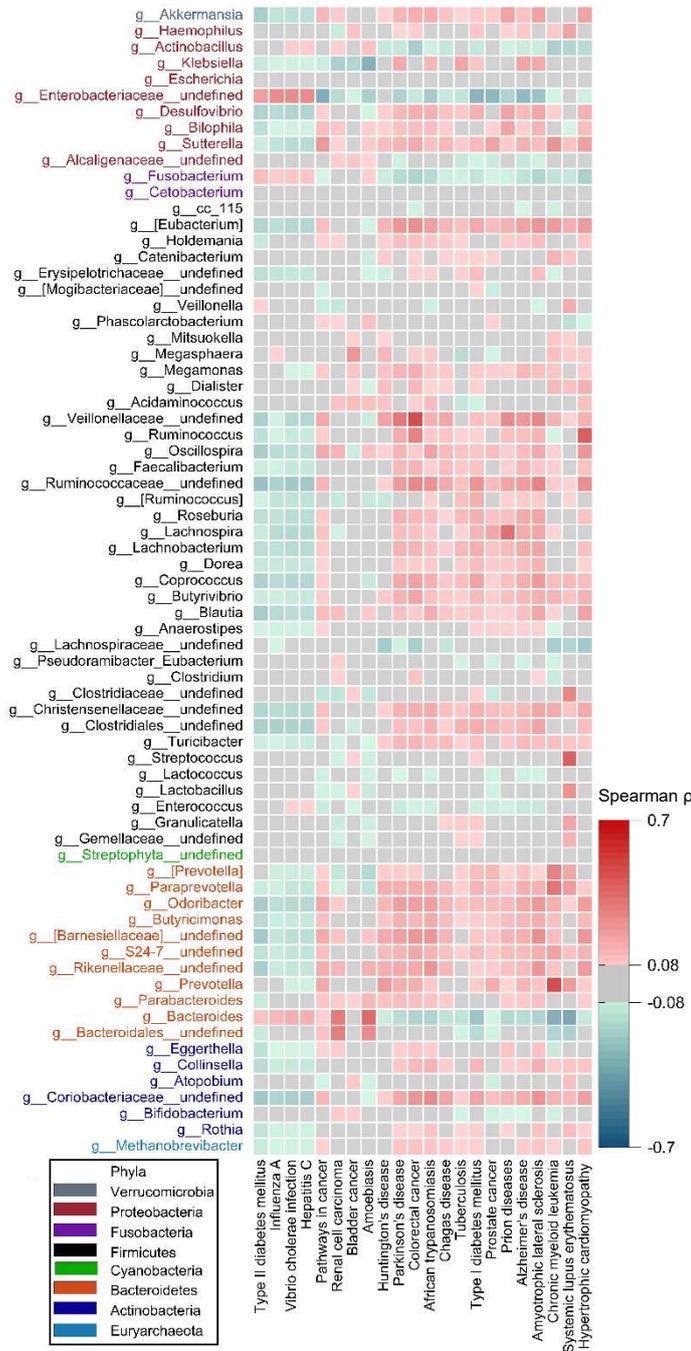
D



1054

1055

1056 Figure 5 Correlation of ~~the human~~-complex human diseases with the gut microbiome.
 1057 The ~~heat map~~ heatmap shows Spearman's correlation of predicted diseases and the gut
 1058 microbiome ~~on~~at the genus level. The grey components show no ~~significance of~~
 1059 significant correlation with Bonferroni correction ($p > 0.05 / ((5.6 * 22), p > 0.0004)$).



1060

1061

1062 **Supplementary Tables [Supplementary Tables.xls]**

1063

1064 **Supplementary Table S1 Transformation of traits in BioBank Japan and taxa in GNHS**1065 **Supplementary Table S2 Required effect size (beta) to reach 80% of power in GNHS cohort**1066 **Supplementary Table S3 Heritability of taxa, enterotype and alpha diversity**1067 **Supplementary Table S4 Significant genetic correlations of gut microbiome and metabolic**
1068 **traits**1069 **Supplementary Table S5 Significant associations of all taxa with SNPs identified in the**
1070 **discovery stage before adjustment (p<5e-8)**1071 **Supplementary Table S6 Replication of genetic variants associated with taxa**1072 **Supplementary Table S7 Replication of genetic variants associated with beta diversity**1073 **Supplementary Table S8 Replication of genetic variants associated with alpha diversity**1074 **Supplementary Table S9 Lead SNPs used to construct polygenic scores**1075 **Supplementary Table S10 MR analysis of gut microbiota on traits and diseases**1076 **Supplementary Table S11 MR analysis of diseases on gut microbiota features**1077 **Supplementary Table S12 Spearman's correlation of certain taxa and complex diseases**1078 **Supplementary Table S13 Spearman's correlation of gut microbiota on genus level and**
1079 **characteristics**

1080

1081 **Supplementary Figures [Supplementary Figures.pdf]**

1082

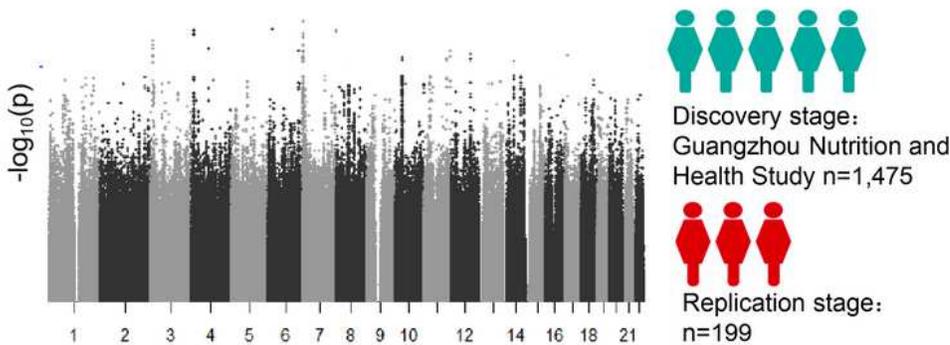
1083 **Supplementary Figure 1 Genome-wide analysis results of taxa.**1084 **Supplementary Figure 2 Spearman's correlation of the relative abundance of AF-associated**
1085 **taxa with the relative level of diseases predicted by PICRUSt.**

1086

1087

Figures

A. Association of host genetics with gut microbiome in a Chinese population.



B. The causal relationships between gut microbiome and human complex diseases.

Bi-directional Mendelian randomization

Instrumental variables

Genetic variants of gut microbiome

Exposure

Gut microbiome
GNHS cohort
(n=1,475)

Outcome

Diseases/Traits
BioBank Japan
(up to 191,764)

Instrumental variables

Genetic variants of complex diseases

Exposure

Complex diseases
Published GWAS
results in East Asians

Outcome

Gut microbiome
GNHS cohort
(n=1,475)

C. The shared and distinct microbiome features among human complex diseases.

1,919 participants from GNHS cohort.

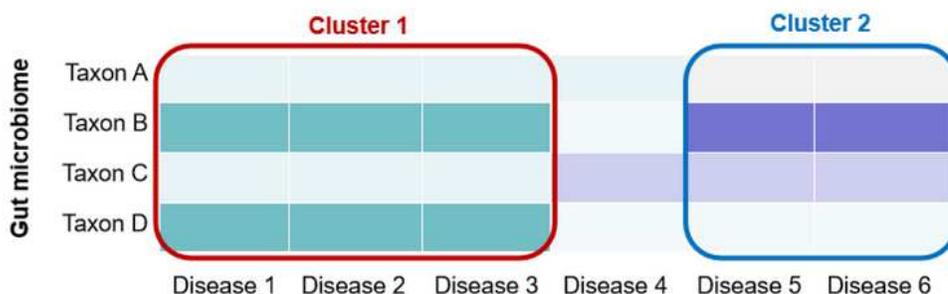


Figure 1

Study overview. The figure shows the highlights of our study. First, we performed a microbiome genome-wide association study in a Chinese population (Step A). We validated significant genetic variants reported in previous studies and replicated our results in an independent cohort. Second, we investigated

the causal relationship between the gut microbiome and complex human diseases using host genetics as instrumental variables for bidirectional Mendelian randomization (MR) analysis (Step B). For the analysis of the effects of the gut microbiome on complex traits, we used publicly available GWAS summary statistics of complex traits (n=58) and diseases (type 2 diabetes mellitus (T2DM), atrial fibrillation (AF), colorectal cancer (CRC) and rheumatoid arthritis) reported by BioBank Japan [27-32]. For the reserve MR analyses, the diseases of interest included T2DM (cases: 7,109; non-cases: 86,022), AF (cases: 8,180; non-cases: 28,612), coronary artery disease (cases: 1,515; non-cases: 5,019), chronic kidney disease (n=71,149), Alzheimer's disease (cases: 477; non-cases: 442), CRC (cases: 8,027; non-cases: 22,577) and prostatic cancer (cases: 495; non-cases: 640) reported in the previous large-scale GWASs in East Asians [27,33-38]. Finally, we identified common and distinct gut microbiome features across different diseases (Step C).

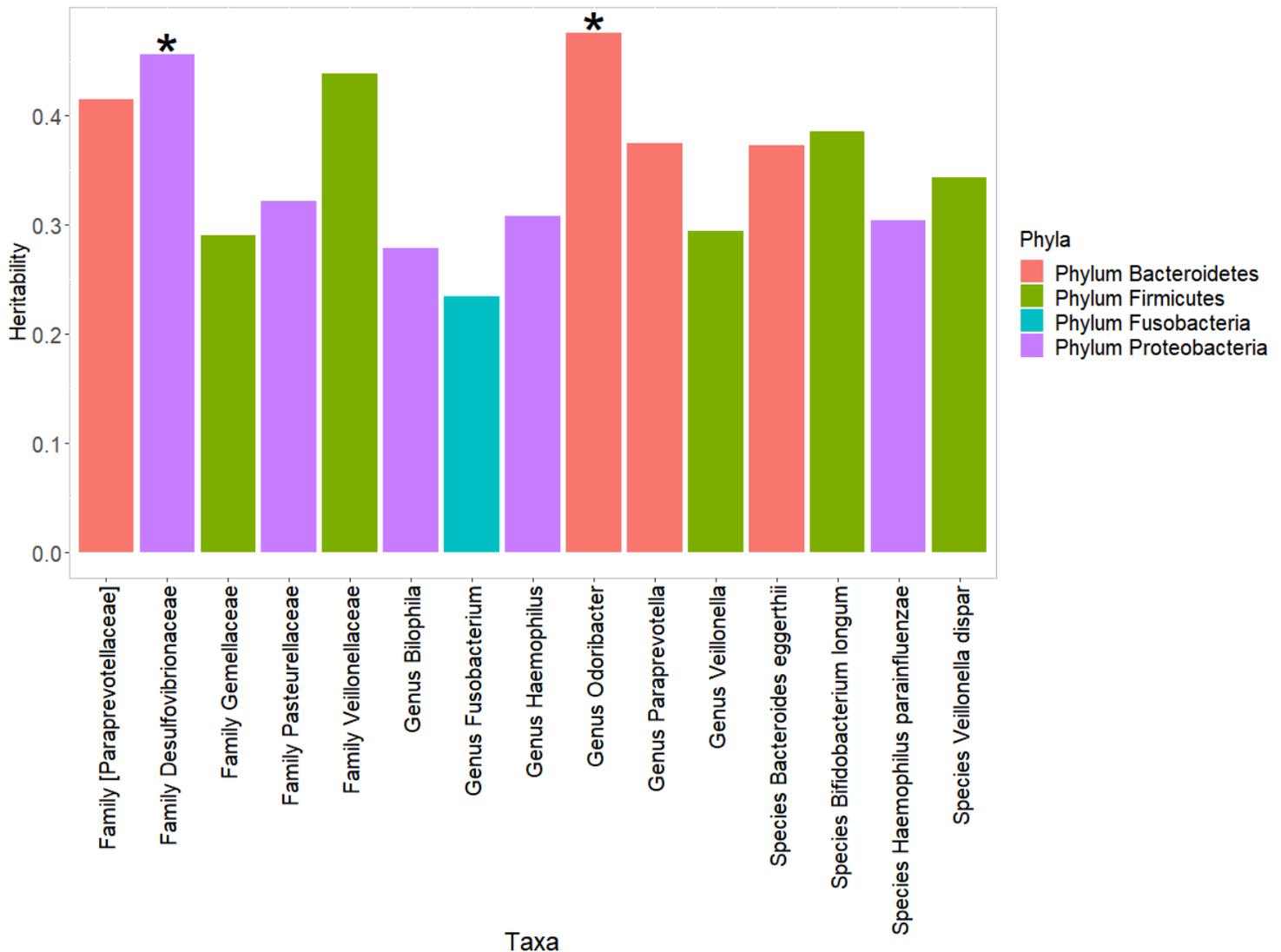
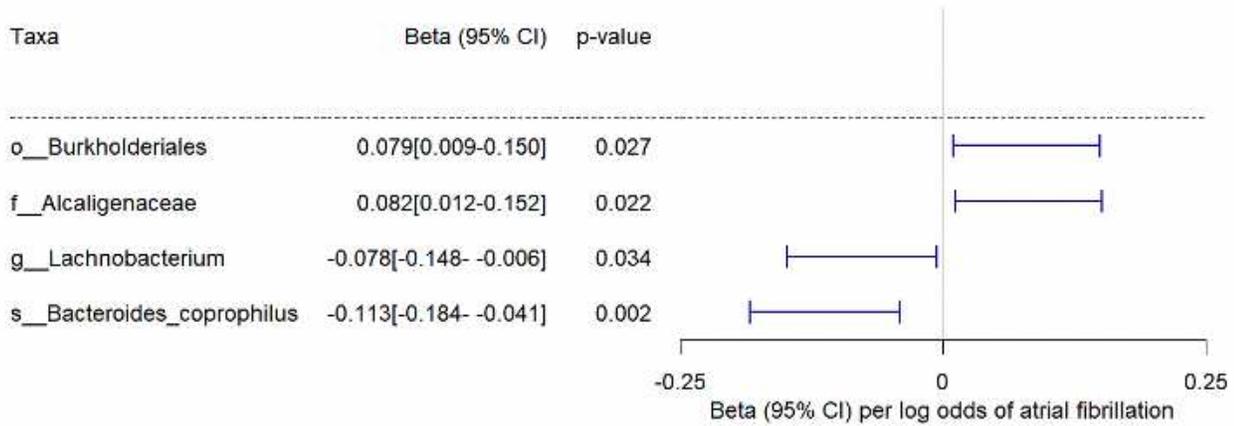
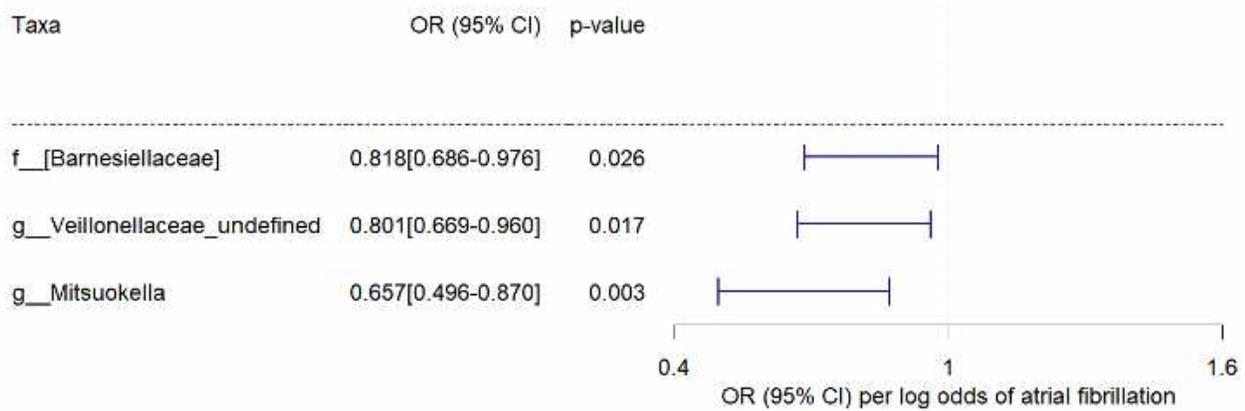


Figure 2

The SNP-based heritability of the gut microbiome. The plot shows the taxa with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, where n is the effective number of independent taxa in each taxonomic level.

A**B****Figure 3**

Effect of host genetically predicted higher atrial fibrillation risk on the gut microbiome. (A). Causal association of atrial fibrillation with the abundance of Burkholderiales, Alcaligenaceae, Lachnobacterium and Bacteroides coprophilus. The magnitude of the effect of atrial fibrillation on taxa is dependent on changes in the abundance of bacteria (1-SD of the log-transformed abundance) per genetically determined higher log odds of atrial fibrillation. (B). Causal association of atrial fibrillation with the presence of Barnesiellaceae, undefined genus in family Veillonellaceae and Mitsuokella. The magnitude of the effect of atrial fibrillation on taxa is presented as an odds ratio increase in the log odds of atrial fibrillation.

participants and control groups. Dimension1 (Dim1) and dimension2 (Dim2) explained 40.1% and 13.1% of the variance, respectively, in the GNHS cohort. The annotation for variables is as follows. AT: African trypanosomiasis, AD: Alzheimer's disease, V1: Amoebiasis, ALS: Amyotrophic lateral sclerosis, BC: Bladder cancer, CD: Chagas disease, CML: Chronic myeloid leukaemia, CRC: Colorectal cancer, V2: Hepatitis C, HD: Huntington's disease, HCM: Hypertrophic cardiomyopathy, V3: Influenza A, PD: Parkinson's disease, V4: Pathways in cancer, V5: Prion disease, PCa: Prostate cancer, RCC: Renal cell carcinoma, SLE: Systemic lupus erythematosus, V6: Tuberculosis, T1DM: Type I diabetes mellitus, T2DM: Type II diabetes mellitus, V7: Vibrio cholerae infection. (D). Gut microbiome-predicted network of relationships among different human complex diseases. The relationship between diseases is determined by SPIEC-EASI with non-normalized predicted abundance data. The diseases that shared the same edge had the gut microbiome-predicted correlation.

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTables.xls](#)
- [SupplementaryFigures.pdf](#)