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

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

6 [#]These authors contributed equally to the work

7 ¹ School of Life Sciences, Westlake University, Hangzhou, China.

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

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

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

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

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

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

20

21 Short title: Interplay between and host genetics and gut microbiome

22

23 *Correspondence to

24 Prof Ju-Sheng Zheng

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

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

27 And

28 Prof Yu-Ming Chen

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

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

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

32 And

33 Prof Jun Wang

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

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

36 Email: junwang@im.ac.cn

37

38 *Abstract*

39 There is increasing interest about the interplay between host genetics and gut
40 microbiome on human complex diseases, with prior evidence mainly derived from
41 animal models. In addition, the shared and distinct microbiome features among
42 human complex diseases remain largely unclear. We performed a microbiome
43 genome-wide association study to identify host genetic variants associated with gut
44 microbiome in a Chinese population with 1475 participants. We then conducted
45 bi-directional Mendelian randomization analyses to examine the potential causal
46 associations between gut microbiome and human complex diseases. We did not find
47 evidence supporting the causal effect of gut microbiome on human complex diseases.
48 In contrast, atrial fibrillation, chronic kidney disease and prostate cancer, as predicted
49 by the host genetics, had potential causal effect on gut microbiome. Further
50 disease-microbiome feature analysis suggested that gut microbiome features revealed
51 novel relationship among human complex diseases. These results suggest that
52 different human complex diseases share common and distinct gut microbiome
53 features, which may help re-shape our understanding about the disease etiology in
54 humans.

55

56 **Introduction**

57 Ever-increasing evidence has suggested that gut microbiome is involved in many
58 physiological processes, such as energy harvest, immune response, and neurological
59 function¹⁻³. With successes of investigation into the clinical application of fecal
60 transplants, modulation of gut microbiome has emerged as a potential treatment
61 option for some complex diseases, including inflammatory bowel disease and
62 colorectal cancer^{4,5}. However, it is still unclear whether the gut microbiome has the
63 potential to be clinically applied for the prevention or treatment of many other
64 complex diseases. Therefore, it is important to clarify the bi-directional causal
65 association between gut microbiome and human complex diseases or traits.

66

67 Mendelian randomization (MR) is a method that uses genetic variants as instrumental
68 variables to investigate the causality between an exposure and outcome in
69 observational studies⁶. Prior literature provides evidence that the composition or
70 structure of the gut microbiome can be influenced by the host genetics⁷⁻¹⁰. On the
71 other hand, host genetic variants associated with gut microbiome were rarely explored
72 in Asian populations, thus we are still lacking instrument variables to perform MR for
73 gut microbiome in Asians. This calls for novel microbiome genome-wide association
74 study (GWAS) in Asian populations.

75

76 Along with the causality issue between the gut microbiome and human complex
77 diseases, it is so far unclear whether human complex diseases had similar or unique

78 gut microbiome features. Identifying common and distinct gut microbiome features
79 across different diseases might shed light on novel relationships among the complex
80 diseases and update our understanding about the disease etiology in humans. However,
81 the composition and structure of gut microbiome are influenced by a variety of factors
82 including environment, diet and regional variation¹¹⁻¹³, which posed a key challenge
83 for the description of representative microbiome features for a specific disease.
84 Although there were several studies comparing disease-related gut microbiome¹⁴⁻¹⁶,
85 few of them has examined and compared the microbiome features across different
86 human complex diseases.

87

88 In the present study, we performed a microbiome GWAS in a Chinese cohort study:
89 the Guangzhou Nutrition and Health Study (GNHS)¹⁷, including 1475 participants.
90 Subsequently, we applied a bi-directional MR method to explore the genetically
91 predicted relationship between gut microbiome and human complex diseases. To
92 explore novel relationships among human complex diseases based on gut microbiome,
93 we investigated the shared and distinct gut microbiome features across diverse human
94 complex diseases¹⁸.

95

96 **Result**

97 **Overview of the study**

98 Our study was based on the GNHS, with 4048 participants (40-75 years old) living in
99 urban Guangzhou city recruited during 2008 and 2013¹⁷. In the GNHS, stool samples

100 were collected among 1937 participants during follow-up visits, among which 1475
101 unrelated participants without taking anti-biotics were included in our discovery
102 microbiome GWAS. We then included additional 199 participants with both genetic
103 and gut microbiome data as a replication cohort, which belonged to the control arm of
104 a case-control study of hip fracture in Guangdong Province, China ¹⁹.

105

106 For both discovery and replication cohorts, genotyping was carried out with Illumina
107 ASA-750K arrays. Quality control and relatedness filters were performed by
108 PLINK1.9 ²⁰. We conducted the genome-wide genotype imputation with 1000
109 Genomes Phase3 v5 reference panel by Minimac3²¹⁻²³. HLA region was imputed with
110 Pan-Asian reference panel and SNP2HLA v1.0.3 ²⁴⁻²⁶.

111

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

113 We performed a series of microbiome GWAS with PLINK 1.9 based on logistic
114 models for binary variables ²⁰. For continuous variables, we used GCTA with mixed
115 linear model-based association (MLMA) method ^{27,28}. We also analyzed categorical
116 variable enterotypes of the participants based on genus-level relative abundance of gut
117 microbiome, using the Jensen-Shannon Distance (JSD) and the Partitioning Around
118 Medoids (PAM) clustering algorithm ²⁹. The participants were subsequently clustered
119 into two groups according to the enterotypes (*Prevotella* vs *Bacteroides*). Thereafter,
120 we performed GWAS for enterotypes using logistic regression model to explore
121 potential associations between host genetics and enterotypes. However, we did not

122 find any genome-wide significant locus with $p < 5 \times 10^{-8}$. Furthermore, we used a
123 restricted maximum likelihood analysis (REML) with GCTA to estimate the
124 SNP-based heritability, and the estimate heritability of the enterotype was 0.055
125 (SE=0.19, Supplementary Table S2)³⁰.
126
127 To examine the association of host genetic variants with alpha diversity, we performed
128 GWAS for three indices (Shannon diversity index, Chao1 diversity indices and
129 observed OTUs index), but again no genome-wide significant signal ($p < 5 \times 10^{-8}$) was
130 found. In the discovery cohort, the heritability of alpha diversity ranged from 0.054 to
131 0.14 (SE=0.20 for all indices, Supplementary Table S2). To further investigate if there
132 is host genetic basis underlying alpha diversity, we constructed a polygenic score for
133 each alpha diversity indicator in the replication cohort, using the genetic variants
134 which showed suggestive significance ($p < 5 \times 10^{-5}$) in the discovery GWAS. The
135 polygenic score was not significantly associated with its corresponding alpha diversity
136 index in our replication cohort. Meanwhile, none of the associations with alpha
137 diversity indices reported in the literature could be replicated (Supplementary Table
138 S7)⁷.
139
140 We performed a beta diversity GWAS using a tool called MicrobiomeGWAS³¹, and
141 found that one locus at *SMARCA2* gene (rs6475456) was associated with
142 beta-diversity at a genome-wide significance level ($p = 3.96 \times 10^{-9}$). However, we could
143 not replicate the results in the replication cohort, which may be due to the limited

144 sample size of the replication cohort. In addition, prior literature had reported 73
145 genetic variants that were associated with beta diversity ^{8,13,32,33}, among which we
146 found that 3 single nucleotide polymorphisms (SNP, *UHRF2* gene-rs563779, *LHFPL3*
147 gene-rs12705241, *CTD-2135J3.4*-rs11986935) had nominal significant ($p < 0.05$)
148 association with beta-diversity in our cohort (Supplementary Table S6), although none
149 of the association survived Bonferroni correction. These studies used various methods
150 for the sequencing and calculation of beta diversity, which raised challenges to verify
151 and extrapolate results across populations.

152

153 We then took the genetic loci reported to be associated with individual taxa in prior
154 studies ^{7,8,13,33} for replication in our GNHS dataset. Although there are still some
155 signals with nominal significance ($p < 0.05$) in our study (e.g., 7 loci associated with
156 *Lachnospiraceae*, *Coprococcus* or *Bacteroides* with $p < 0.05$; Supplementary Table S5),
157 none of the associations of these genetic variants with taxa survived the Bonferroni
158 correction ($p < 1 \times 10^{-4}$). The null results may be because of various clustering
159 similarities, classifiers or reference databases to annotate taxa and different
160 sequencing methods used in these studies.

161

162 We subsequently performed GWAS discovery for individual gut microbes in our own
163 GNHS discovery dataset. For the taxa present at more than ninety percent of
164 participants and alpha diversity, we used Z-score normalization to transform the
165 distribution and carried out analysis based on a log-normal model. A MLMA test in

166 the GCTA was used to assess the association, with the first five principal components,
167 age, sex and sequencing batch fitted as fixed effects and the effects of all the SNPs
168 fitted as random effects^{27,28,34}. For other taxa present at fewer than ninety percent, we
169 transformed the absence/presence of the taxon into binary variables and used
170 PLINK1.9 to perform a logistic model, adjusted for the first five principal components,
171 age, sex and sequencing batch.

172

173 For all the gut microbiome taxa, the significant threshold was defined as 5×10^{-8} in the
174 discovery stage. As some taxa were correlated with each other, we also used an
175 eigendecomposition analysis to calculate the effective number of independent taxa on
176 each taxonomy level (phylum level: 2.3, class level: 2.9, order level: 2.9, family level:
177 5.5, genus level: 5.6, species level: 3.2)^{35,36}. We found that 6 taxa were associated
178 with host genetic variants in the discovery cohort ($p < 5 \times 10^{-8}/n$, n is the effective
179 number of independent taxa on each taxonomy level, Supplementary Table S4);
180 however, these associations were not significant ($p > 0.05$) in the replication cohort. We
181 then used a threshold of $p < 5 \times 10^{-5}$ at the GWAS discovery stage to incorporate
182 additional genetic variants which may explain a larger proportion of heritability for
183 taxa, based on which we constructed a polygenic score for each taxon in the
184 replication. We found that the polygenic scores were significantly associated with 3
185 taxa including *Coriobacteriaceae*, *Odoribacter* and *Parabacteroides_undefined* in the
186 replication set ($p < 0.05$, Methods, see also Figure 1A, 1B, 1C).

187

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

189 We used GCTA to perform a bivariate GREML (genomic-relatedness-based restricted
190 maximum-likelihood) analysis to estimate the genetic correlation between gut
191 microbiome and traits in the GNHS ^{27,37}. The traits included BMI, fasting blood sugar
192 (FBS), glycosylated hemoglobin (HbA1c), systolic blood pressure (SBP), diastolic
193 blood pressure (DBP), high density lipoprotein cholesterol (HDL-C), low density
194 lipoprotein cholesterol (LDL-C), total cholesterol (TC) and triglyceride (TG), none of
195 which could pass Bonferroni correction. Additionally, HDL-C was the only trait that
196 had nominal genetic correlation ($p < 0.05$) with gut microbes (specifically,
197 *Desulfovibrionaceae* and [*Prevotella*], Supplementary Table S3).

198

199 **Bi-directional assessment of the genetically predicted association between gut**
200 **microbiome and complex diseases/traits**

201 Using genetic variants-composed polygenic scores as genetic instruments, we
202 performed MR analysis to assess the putative causal effect of microbiome
203 (*Coriobacteriaceae*, *Odoribacter* and *Parabacteroides_undefined*) on human complex
204 diseases or traits. Inverse variance weighted (IVW) method was used for the MR
205 analysis, while other three methods (Weighted median, MR-Egger and MR-PRESSO)
206 ^{38,39} were applied to confirm the robustness of results. The horizontal pleiotropy was
207 assessed using MR-PRESSO Global test and MR-Egger Regression. For the analysis
208 of gut microbiome on complex traits, we downloaded public available GWAS
209 summary statistics of complex traits (n=58) and diseases (type 2 diabetes mellitus

210 (T2DM), atrial fibrillation (AF), colorectal cancer (CRC) and prostatic cancer (PCa))
211 reported by BioBank Japan⁴⁰⁻⁴⁴. There was no evidence that these taxa had causal
212 association with human complex diseases or traits in our MR analyses
213 (Supplementary Table S9), which may be due to the limited genetic instruments
214 discovered in our present study.

215

216 We subsequently performed a reserve MR analysis to assess the potential causal effect
217 of human complex diseases on gut microbiome features. For the reserve MR analyses,
218 the diseases of interests included T2DM, AF, coronary artery disease (CAD), chronic
219 kidney disease (CKD), Alzheimer's disease (AD), CRC and PCa, and their
220 instrumental variables for the MR analysis were based on the previous large-scale
221 GWAS in East Asians^{40,45-50}. The results suggested that AF and CKD were causally
222 associated with gut microbiome (See also Figure 2A, 2B, Supplementary Table S10).
223 Genetically predicted higher risk of AF was associated with lower abundance of
224 *Coprophilus*, *Lachnobacterium*, *Barnesiellaceae*, *Veillonellaceae* and *Mitsuokella*,
225 and higher abundance of *Alcaligenaceae*. Additionally, genetically predicted higher
226 risk of CKD could increase *Anaerostipes* abundance and higher risk of PCa could
227 decrease [*Prevotella*].

228

229 To further investigate the potential complex diseases that may be correlated with the
230 taxa affected by AF, we applied Phylogenetic Investigation of Communities by
231 Reconstruction of Unobserved States (PICRUSt) to predict the disease pathway

232 abundance⁵¹. We used Spearman's rank-order correlation to test whether 22 human
233 complex diseases were associated with the aforementioned AF-associated taxa (See
234 also Figure 2C). The heatmap indicated that cancers and neurodegenerative diseases
235 including Parkinson's disease (PD), AD, amyotrophic lateral sclerosis (ALS) as well
236 as AF were correlated with similar gut microbiome. Although the association among
237 these diseases are highly supported by previous studies⁵²⁻⁵⁴, no study has compared
238 common gut microbiome features across these different diseases.

239

240 **Microbiome features of human complex diseases**

241 To compare gut microbiome features across human diseases, we chose 22 human
242 complex diseases from predicted abundance and performed k-medoids clustering¹⁸.
243 We used an $m \times n$ matrix to perform the cluster analysis, where m is the number of
244 participants and n is the number of selected diseases. According to optimum average
245 silhouette width⁵⁵, we chose optimal number of clusters for further analysis (See also
246 Figure 3A). The plot showed that neurological diseases including ALS and AD
247 belonged to the same cluster, while PD and CRC had much similarity in gut
248 microbiome. The results also suggested that systemic lupus erythematosus (SLE) and
249 chronic myeloid leukemia (CML) shared similar gut microbiome features. Moreover,
250 we could replicate these clusters in our replication cohort, which suggested that the
251 clustering results were robust (See also Figure 3B).

252

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

254 end, we repeated the analysis among participants who took antibiotic less than two
255 weeks before stool sample collection, considering that antibiotic treatments were
256 believed to cause microbiome imbalance, and the clusters were totally different in this
257 group (See also Figure 3C). The results indicated a totally different clustering, which
258 suggested that gut microbiome indeed contributed to the correlations among diseases.
259 To further demonstrate common microbiome features across different diseases, we
260 examined the correlation of the predicted diseases with genus-level taxa. The results
261 showed that human complex diseases had shared similar gut microbiome features, as
262 well as distinct features on their own (See also Figure 4).

263

264 To validate the accuracy of the association between the predicted disease-related gut
265 microbiome features and the corresponding disease, we used T2DM as an example,
266 examining the association of predicted T2DM-related microbiome features with
267 T2DM risk in our GNHS samples. We constructed a microbiome risk score (MRS)
268 based on 16 selected taxa with predicted correlation coefficient with T2DM greater
269 than 0.2. A logistic regression was used to examine the above MRS with T2DM risk
270 in the GNHS (n=1886, with 217 T2DM cases). The results showed that higher MRS
271 was associated with lower risk of T2DM (odds ratio: 0.850, 95% confidence interval:
272 0.804 to 0.898, $p=8.75 \times 10^{-9}$).

273

274 Based on the above results, we proposed a hypothesis that related diseases might
275 share similar gut microbiome features. To test for this hypothesis, we performed

276 validation analysis by including GNHS participants who had one of the following
277 self-reported diseases: stroke (n=8), chronic hepatitis (n=19), coronary heart diseases
278 (CHD) (n=40), cataract (n=124) and insomnia (n=68). The results of k-medoids
279 clustering suggested that CHD, cataract and insomnia shared common gut
280 microbiome features, which was supported by the prior research reporting that both
281 patients suffering insomnia and women receiving cataract extraction had increased
282 risks of CHD ⁵⁶⁻⁵⁸.

283

284 **Discussion**

285 Our study is among the first to investigate the host genetics-gut microbiome
286 association in East Asian populations and reveals that several microbiome species
287 (e.g., *Coriobacteriaceae* and *Odoribacter*) are influenced by host genetics. We then
288 show that complex diseases such as atrial fibrillation, chronic kidney disease and
289 prostate cancer, have potential causal effect on gut microbiome. More interestingly,
290 our results indicate that different human complex diseases may be mechanically
291 correlated by sharing common gut microbiome features, but also maintaining their
292 own distinct microbiome features.

293

294 Previous studies and our study showed that gut microbiome had an inclination to be
295 influenced by host genetics ^{8,10,33,59}, although the successful replication tends to be
296 rare. We could not validate any of the reported genetic variants that were significantly
297 associated with gut microbiome features in prior reports, which may reflect the

298 difference in population and heterogeneity between study but also raise concerns
299 about the reproducibility. Many factors including ethnic differences,
300 gene-environment interaction and dissimilarity in sequencing methods may make it
301 hard to extrapolate results from microbiome GWAS across populations in the
302 microbiome field. Nevertheless, we successfully replicate several polygenic scores of
303 gut microbiome, and the current study represent the largest dataset in Asian
304 populations and would be a unique resource to be used in large-scale trans-ethnic
305 meta-analysis of microbiome GWAS in future.

306

307 The MR analysis in the present study did not support causal effect of gut microbiome
308 on diseases or traits, however, this result should be interpreted with caution because of
309 the limited genetic instruments derived from GWAS. In contrast, the reverse MR
310 analysis provided evidence that AF, CKD and PCa could causally influence gut
311 microbiome. As our study is among the first to investigate gene-microbiome
312 association in East Asians, we need further study in this region to confirm our results.
313 Additionally, rare and low-frequency variants may have an important impact on
314 common diseases⁶⁰, thus it will be of interest to clarify the effects of low-frequency
315 variants on gut microbiome in cohorts with large sample sizes in future.

316

317 Our results indicate that gut microbiome helps reveal novel and interesting
318 relationships among human complex diseases, suggesting that different diseases may
319 have common and distinct gut microbiome features. A prior study including

320 participants from different countries identified three microbiome clusters²⁹. Notably,
321 this study focused on classifying the individuals into distinct enterotypes regardless of
322 the individuals' health status, while in the present study we described representative
323 microbiome features for diseases of interest. The microbiome features revealed a
324 close association of AF with neurodegenerative diseases as well as cancers, which
325 was supported by prior studies showing that AF had correlation with AD and PD^{52,53},
326 and AF patients had relatively higher risks of several cancers including lung cancer
327 and CRC^{54,61}. We also observed that microbiome features of SLE and CML were
328 highly similar. Interestingly, a tyrosine kinase inhibitor of platelet-derived growth
329 factor receptor, imatinib, was widely used to treat CML and significantly ameliorated
330 survival in murine lupus autoimmune disease⁶². In addition, association between CRC
331 and PD has been reported in several observational cohorts^{63,64}. Collectively, these
332 findings strongly support our hypothesis that human complex diseases sharing similar
333 microbiome features might be mechanically correlated. Furthermore, from the
334 perspectives of risk genes of AF and neurodegenerative diseases, previous GWAS for
335 AF have identified two loci at *PITX2* gene-rs6843082 and *C9orf3* gene-rs7026071,
336 which were also associated with the risk of ALS (p=0.0138 and p=0.049, respectively)
337 ⁶⁵⁻⁶⁷.

338

339 In summary, we perform bi-directional MR analyses to examine the causal
340 relationship between gut microbiome and human complex diseases, revealing that
341 some complex diseases causally affect abundance of specific gut microbes. There is

342 no convincing evidence supporting the causal role of gut microbiome on human
343 complex diseases. The disease and gut microbiome association analysis reveals novel
344 relationships among human complex diseases, which may help re-shape our
345 understanding about the disease etiology, as well as extending clinical indications of
346 existing drugs for different diseases.

347

348 **Method**

349 **Study participants and sample collection**

350 Our study was based on the Guangzhou Nutrition and Health Study (GNHS), with
351 4048 participants (40-75 years old) living in urban Guangzhou city recruited during
352 2008 and 2013 ¹⁷. We followed up participants every three years. In the GNHS, stool
353 samples were collected among 1937 participants during follow-up visits. Among
354 those with stool samples, 1717 participants had genetic data and 1475 participants
355 with identical by decent (IBD) less than 0.185.

356

357 We included 199 participants with both genetic and gut microbiome data as a
358 replication cohort, which belonged to the control arm of a case-control study of hip
359 fracture with the participants (52-83 years old) recruited between June 2009 and
360 August 2015 in Guangdong Province, China ¹⁹.

361

362 Blood samples of all participants were collected after an overnight fasting and buffy
363 coat was separated from whole blood and stored at -80°C. Stool samples were

364 collected during the on-site visit of the participants at Sun Yat-sen University. All
365 samples were manually stirred, separated into tubes and stored at -80°C within four
366 hours.

367

368 **Genotyping data**

369 For both discovery and replication cohorts, DNA was extracted from leukocyte using
370 the TIANamp® Blood DNA Kit as per the manufacturer's instruction. DNA
371 concentrations were determined using the Qubit quantification system (Thermo
372 Scientific, Wilmington, DE, US). Extracted DNA was stored at -80°C . Genotyping
373 was carried out with Illumina ASA-750K arrays. Quality control and relatedness
374 filters were performed by PLINK1.9²⁰. Individuals with high or low proportion of
375 heterozygous genotypes (outliers defined as 3 standard deviation) were excluded⁶⁸.
376 Individuals who had different ancestries (the first two principal components ± 5
377 standard deviation from the mean) or related individuals ($\text{IBD} > 0.185$) were
378 excluded⁶⁸. Variants were mapping to the 1000 Genomes Phase3 v5 by SHAPEIT^{23,69}
379 and then we conducted the genome-wide genotype imputation with 1000 Genomes
380 Phase3 v5 reference panel by Minimac3^{21,22}. Genetic variants with imputation
381 accuracy $\text{RSQR} > 0.3$ and $\text{MAF} > 0.05$ were included in our analysis. We used
382 Pan-Asian reference panel consist of 502 participants and SNP2HLA v1.0.3 to impute
383 HLA region²⁴⁻²⁶.

384

385 **Sequencing and processing of 16S rRNA data**

386 Microbial DNA was extracted from fecal samples using the QIAamp® DNA Stool
387 Mini Kit per the manufacturer's instruction. DNA concentrations were determined
388 using the Qubit quantification system. The V3-V4 region of the 16S rRNA gene was
389 amplified from genomic DNA using primers 341F and 805R. Sequencing was
390 performed using MiSeq Reagent Kits v2 on the Illumina MiSeq System.
391
392 Fastq-files were demultiplexed by the MiSeq Controller Software. Ultra-fast sequence
393 analysis (USEARCH) was performed to trim the sequence for amplification primers,
394 diversity spacers, sequencing adapters, merge-paired and quality filter⁷⁰. Operational
395 taxonomic units (OTUs) were clustered based on 97% similarity using UPARSE⁷¹.
396 OTUs were annotated with Greengenes 13_8
397 (<https://greengenes.secondgenome.com/>). After randomly selecting 10000 reads for
398 each sample, Quantitative Insights into Microbial Ecology (QIIME) software version
399 1.9.0 was used to calculate alpha diversity (Shannon diversity index, Chao1 diversity
400 indices and observed OTUs index) based on the rarefied OTU counts⁷².

401

402 **Statistical analysis**

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

404 For each of the GNHS participants and the replication cohort, we clustered
405 participants based on genus-level relative abundance, estimating the JSD distance and
406 PAM clustering algorithm, and then defined two enterotypes according to
407 Calinski-Harabasz Index^{29,73}. PLINK 1.9 was used to perform a logistic regression

408 model for enterotypes and taxa present at fewer than ninety percent, adjusted for age,
409 sex and the first five principal components.

410

411 For beta diversity, the analyses for the genome-wide host genetic variants with beta
412 diversity was performed using MicrobiomeGWAS³¹, adjusted for covariates including
413 the first five principal components, age and sex. We filtered OTUs present at fewer
414 than ten percent of participants to calculate Bray–Curtis dissimilarity.

415

416 Alpha diversity was calculated after randomly sampling 10000 reads per sample. For
417 the taxa present at more than ninety percent of participants and alpha diversity, we
418 used Z-score normalization to transform the distribution and carried out analysis
419 based on a log-normal model. A mixed linear model based association (MLMA) test
420 in GCTA was used to assess the association, fitting the first five principal components,
421 age, sex and sequencing batch as fixed effects and the effects of all the SNPs as
422 random effects^{27,28,34}. For other taxa present at fewer than ninety percent, we
423 transformed the absence/presence of the taxon into binary variables and used
424 PLINK1.9 to perform a logistic model, adjusted for the first five principal components,
425 age, sex and sequencing batch. For all the gut microbiome features, the significant
426 threshold was defined as $5 \times 10^{-8} / n$ (n is the effective number of independent taxa on
427 each taxonomy level) in the discovery stage. We estimated genomic inflation factors
428 with LDSC v1.0.1 at local server⁷⁴.

429

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

431 We used the GREML method in GCTA to estimate the proportion of variance
432 explained by all SNPs³⁰. The taxa were divided into two groups based on whether the
433 taxa were present in the ninety percent of participants or not. For alpha diversity and
434 taxa, our model was adjusted for constrain covariates including sex and sequencing
435 batch, as well as quantitative covariates including the first five principal components
436 and age. The model was adjusted for the same covariates except for sequencing batch
437 for analysis of enterotype.

438

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

440 We used GCTA to perform a bivariate GREML analysis to estimate the genetic
441 correlation between gut microbiome and traits in the GNHS^{27,37}. The gut microbiome
442 was divided into two groups according to the previous description. We used
443 continuous variables to taxa present at more than ninety percent of participants and
444 traits. For taxa present at fewer than ninety percent of participants, we used binary
445 variables according to the absence/presence of taxa. This analysis included traits such
446 as BMI, FBS, HbA1c, SBP, DBP, HDL-C, LDL-C, TC and TG.

447

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

449 We selected lead SNPs using PLINK v1.9 with the ‘—clump’ command to clump
450 SNPs that p value $< 5 \times 10^{-5}$ and $r^2 < 0.1$ within 0.1 cM. We used beta coefficients as
451 weight to construct polygenic scores for taxa and alpha diversity. For alpha diversity

452 and taxa present at more than ninety percent participants, we constructed weighted
453 polygenic scores and performed the analysis on a general linear model with a negative
454 binomial distribution to test for association between the polygenic scores and taxa,
455 adjusted for the first five principal components, age, sex and sequencing batch. We
456 used weighted polygenic scores and logistic regression to the absence/presence taxa,
457 adjusted for the same covariates as in the above analysis. Taxa with significance
458 ($p < 0.05$) in the replication cohort were included for the further analysis.

459

460 **The effective number of independent taxa**

461 As some taxa were correlated with each other, we used an eigendecomposition
462 analysis to calculate the effective number of independent taxa on each taxonomy level
463 ^{35,36}. Matrix M is an $m \times n$ matrix, where m is the number of participants and n is the
464 number of taxa on the corresponding taxa level. Matrix A is the variance-covariance
465 matrix of matrix M . P is the matrix of eigenvectors. $\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ is the diagonal
466 matrix comprised of the ordered eigenvalues, which can be calculated as:

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

467 The effective number of independent taxa can be calculated as:

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

468

469 **Bi-directional MR analysis**

470 In the analysis of potential causal effect of gut microbiome features on diseases, we
471 used independent genetic variants (selected as part of the polygenic score analysis) as

472 the instrument variables. In the analysis of potential causal effect of diseases on gut
473 microbiome features, we selected genetic variants that were replicated in East Asian
474 populations as instrument variables. As all instrument variables were from East Asian
475 populations, we chose independent genetic variants ($r^2 < 0.1$) based on GNHS cohort.
476 We identified the best proxy ($r^2 > 0.8$) based on GNHS cohort or discarded the variant
477 if no proxy was available. We used inverse variance weighted (IVW) method to
478 estimate effect size. To confirm the robustness of results, we performed other three
479 MR methods including weighted median, MR-Egger and MR-PRESSO. To assess the
480 presence of horizontal pleiotropy, we performed MR-PRESSO Global test and
481 MR-Egger Regression. Effect sizes of gut microbiome on traits were dependent on
482 units of traits⁴³ (Supplementary table S1). Results of human complex diseases on the
483 absence/presence gut microbiome were presented as risk of presence (vs absence) of
484 the microbiome per log odds difference of the disease. Results of diseases on other
485 gut microbiome and alpha diversity were presented as changes in abundance of taxa
486 (10-SD of log transformed) per log odds difference of the respective disease.
487
488 The statistical significance of gut microbiome on traits and diseases was defined as
489 $p < 0.0008$ ($0.05/62$). In addition, the statistical significance of diseases on gut
490 microbiome features was defined as $p < 0.05/n$ (n is the effective number of
491 independent taxa on the corresponding taxonomy level). Results that could not pass
492 Bonferroni adjustment but $p < 0.05$ in all four MR methods were considered as
493 potential causal relationships. We performed MR analyses on R v3.5.3.

494

495 **Pathway analysis**

496 We used OTUs by QIIME and annotated the variation of functional genes with
497 Phylogenetic Investigation of Communities by Reconstruction of Unobserved States
498 (PICRUST) ⁵¹. The pathways and diseases were annotated using KEGG ⁷⁵⁻⁷⁷. We used
499 Spearman's rank-order correlation to investigate association of predicted pathway or
500 diseases abundance with taxa. In the heatmap, diseases were clustered with 'hcluster'
501 function on R. To test whether non-normalized pathway or disease abundance was
502 associated with each other, we used SPIEC-EASI to test the interaction relationship,
503 and then used Cytoscape v3.7.2 to visualize the interaction network ^{78,79}.

504

505 **Construction of the microbiome risk score**

506 To construct a microbiome risk score for T2DM, we used a Spearman's rank-order
507 correlation to select taxa with the absolute value of correlation coefficient higher than
508 0.2. Score for each taxon abundance <5% quantile in our study was defined as 0. For
509 those above 5%, score for each taxon showing negative association with T2DM was
510 defined as 1; score for each taxon showing positive association with T2DM was
511 defined as -1. We then summed up values from all taxa. We selected logistic
512 regression model to estimate association of the MRS with T2DM risk, and linear
513 model to estimate the association of the MRS with the continuous variables, adjusted
514 for age, sex, dietary energy intake, alcohol intake and BMI at the time of sample
515 collection.

516

517 **Clustering diseases**

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

519 We performed PAM algorithm based on predicted abundance of diseases or average

520 relative abundance after Z-score normalization⁸⁰. PAM algorithm searches k medoids

521 among the observations and then found nearest medoids to minimize the dissimilarity

522 among clusters¹⁸. Given a set of objects $x = (x_1, x_2, \dots, x_n)$, the dissimilarity between

523 objects x_i and x_j is denoted by $d(i,j)$. The assignment of object i to the

524 representative object j is denoted by z_{ij} . z_{ij} is a binary variable and is 1 if object i

525 belongs to the cluster of the representative object j. The function to minimize the

526 model is given by:

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

527

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

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

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

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

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

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

534 The silhouette width for object i can be calculated by:

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

535 Then we calculated the average of silhouette width for each object. The cluster

536 number is determined by the number of which the average silhouette width is
537 maximum.

538

539 **Acknowledgments**

540 This study was funded by National Natural Science Foundation of China (81903316,
541 81773416), Westlake University (101396021801) and the 5010 Program for Clinical
542 Researches (2007032) of the Sun Yat-sen University (Guangzhou, China). The
543 authors declare no conflict of interest. We thank the Westlake University
544 Supercomputer Center for providing computing and data analysis service for the
545 present project.

546

547 **Data availability**

548 The raw data for 16 S rRNA gene sequences are available in the CNSA
549 (<https://db.cngb.org/cnsa/>) of CNGBdb at accession number CNP0000829.

550

551

552

553 **References**

- 554 1. Awany, D. *et al.* Host and Microbiome Genome-Wide Association Studies:
555 Current State and Challenges. *Front Genet* **9**, 637 (2018).
- 556 2. Bull, M.J. & Plummer, N.T. Part 1: The Human Gut Microbiome in Health
557 and Disease. *Integr Med (Encinitas)* **13**, 17-22 (2014).
- 558 3. Lynch, J.B. & Hsiao, E.Y. Microbiomes as sources of emergent host
559 phenotypes. *Science* **365**, 1405 (2019).
- 560 4. Allegretti, J.R., Mullish, B.H., Kelly, C. & Fischer, M. The evolution of the
561 use of faecal microbiota transplantation and emerging therapeutic indications.
562 *The Lancet* **394**, 420-431 (2019).
- 563 5. Wong, S.H. & Yu, J. Gut microbiota in colorectal cancer: mechanisms of
564 action and clinical applications. *Nature Reviews Gastroenterology &*
565 *Hepatology* (2019).

-
- 566 6. Davies, N.M., Holmes, M.V. & Davey Smith, G. Reading Mendelian
567 randomisation studies: a guide, glossary, and checklist for clinicians. *BMJ* **362**,
568 k601 (2018).
- 569 7. Turpin, W. *et al.* Association of host genome with intestinal microbial
570 composition in a large healthy cohort. *Nat Genet* **48**, 1413-1417 (2016).
- 571 8. Wang, J. *et al.* Genome-wide association analysis identifies variation in
572 vitamin D receptor and other host factors influencing the gut microbiota. *Nat*
573 *Genet* **48**, 1396-1406 (2016).
- 574 9. Goodrich, J.K. *et al.* Human genetics shape the gut microbiome. *Cell* **159**,
575 789-99 (2014).
- 576 10. Goodrich, J.K. *et al.* Genetic Determinants of the Gut Microbiome in UK
577 Twins. *Cell Host Microbe* **19**, 731-43 (2016).
- 578 11. Ganesan, K., Chung, S.K., Vanamala, J. & Xu, B. Causal Relationship
579 between Diet-Induced Gut Microbiota Changes and Diabetes: A Novel
580 Strategy to Transplant *Faecalibacterium prausnitzii* in Preventing Diabetes. *Int*
581 *J Mol Sci* **19**(2018).
- 582 12. He, Y. *et al.* Regional variation limits applications of healthy gut microbiome
583 reference ranges and disease models. *Nat Med* **24**, 1532-1535 (2018).
- 584 13. Rothschild, D. *et al.* Environment dominates over host genetics in shaping
585 human gut microbiota. *Nature* **555**, 210-215 (2018).
- 586 14. Duvallet, C., Gibbons, S.M., Gurry, T., Irizarry, R.A. & Alm, E.J.
587 Meta-analysis of gut microbiome studies identifies disease-specific and shared
588 responses. *Nature communications* **8**, 1784-1784 (2017).
- 589 15. Cheng, S. *et al.* Identifying psychiatric disorder-associated gut microbiota
590 using microbiota-related gene set enrichment analysis. *Briefings in*
591 *Bioinformatics* (2019).
- 592 16. Jackson, M.A. *et al.* Gut microbiota associations with common diseases and
593 prescription medications in a population-based cohort. *Nature*
594 *Communications* **9**, 2655 (2018).
- 595 17. Cao, Y. *et al.* Association of magnesium in serum and urine with carotid
596 intima-media thickness and serum lipids in middle-aged and elderly Chinese: a
597 community-based cross-sectional study. *European journal of nutrition*
598 **55**(2015).
- 599 18. Kaufman, L. & Rousseeuw, P. Partitioning Around Medoids (Program PAM).
600 68-125 (1990).
- 601 19. Sun, L.-L. *et al.* Associations between the dietary intake of antioxidant
602 nutrients and the risk of hip fracture in elderly Chinese: A case-control study.
603 *The British journal of nutrition* **112**, 1-9 (2014).
- 604 20. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and
605 population-based linkage analyses. *American journal of human genetics* **81**,
606 559-575 (2007).
- 607 21. Das, S. *et al.* Next-generation genotype imputation service and methods.
608 *Nature Genetics* **48**, 1284 (2016).
- 609 22. Clarke, L. *et al.* The international Genome sample resource (IGSR): A

-
- 610 worldwide collection of genome variation incorporating the 1000 Genomes
611 Project data. *Nucleic Acids Research* **45**, D854-D859 (2016).
- 612 23. Delaneau, O. *et al.* Integrating sequence and array data to create an improved
613 1000 Genomes Project haplotype reference panel. *Nature Communications* **5**,
614 3934 (2014).
- 615 24. Okada, Y. *et al.* Risk for ACPA-positive rheumatoid arthritis is driven by
616 shared HLA amino acid polymorphisms in Asian and European populations.
617 *Hum Mol Genet* **23**, 6916-26 (2014).
- 618 25. Pillai, N.E. *et al.* Predicting HLA alleles from high-resolution SNP data in
619 three Southeast Asian populations. *Hum Mol Genet* **23**, 4443-51 (2014).
- 620 26. Jia, X. *et al.* Imputing Amino Acid Polymorphisms in Human Leukocyte
621 Antigens. *PLOS ONE* **8**, e64683 (2013).
- 622 27. Yang, J., Lee, S.H., Goddard, M.E. & Visscher, P.M. GCTA: a tool for
623 genome-wide complex trait analysis. *Am J Hum Genet* **88**, 76-82 (2011).
- 624 28. Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M. & Price, A.L.
625 Advantages and pitfalls in the application of mixed-model association
626 methods. *Nat Genet* **46**, 100-6 (2014).
- 627 29. Arumugam, M. *et al.* Enterotypes of the human gut microbiome. *Nature* **473**,
628 174-80 (2011).
- 629 30. Lee, S.H., Wray, N.R., Goddard, M.E. & Visscher, P.M. Estimating missing
630 heritability for disease from genome-wide association studies. *Am J Hum*
631 *Genet* **88**, 294-305 (2011).
- 632 31. Hua, X. *et al.* MicrobiomeGWAS: a tool for identifying host genetic variants
633 associated with microbiome composition. *bioRxiv*, 031187 (2015).
- 634 32. Ruhlemann, M.C. *et al.* Application of the distance-based F test in an mGWAS
635 investigating beta diversity of intestinal microbiota identifies variants in
636 SLC9A8 (NHE8) and 3 other loci. *Gut Microbes* **9**, 68-75 (2018).
- 637 33. Bonder, M.J. *et al.* The effect of host genetics on the gut microbiome. *Nat*
638 *Genet* **48**, 1407-1412 (2016).
- 639 34. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for
640 human height. *Nature Genetics* **42**, 565 (2010).
- 641 35. Bretherton, C.S., Widmann, M., Dymnikov, V.P., Wallace, J.M. & Bladé, I.
642 The Effective Number of Spatial Degrees of Freedom of a Time-Varying Field.
643 *Journal of Climate* **12**, 1990-2009 (1999).
- 644 36. Wang, H. *et al.* Genotype-by-environment interactions inferred from genetic
645 effects on phenotypic variability in the UK Biobank. *Science Advances* **5**,
646 eaaw3538 (2019).
- 647 37. Lee, S.H., Yang, J., Goddard, M.E., Visscher, P.M. & Wray, N.R. Estimation of
648 pleiotropy between complex diseases using single-nucleotide
649 polymorphism-derived genomic relationships and restricted maximum
650 likelihood. *Bioinformatics* **28**, 2540-2 (2012).
- 651 38. Sanna, S. *et al.* Causal relationships among the gut microbiome, short-chain
652 fatty acids and metabolic diseases. *Nature Genetics* **51**, 600-605 (2019).
- 653 39. Verbanck, M., Chen, C.-Y., Neale, B. & Do, R. Detection of widespread

-
- 654 horizontal pleiotropy in causal relationships inferred from Mendelian
655 randomization between complex traits and diseases. *Nature Genetics* **50**,
656 693-698 (2018).
- 657 40. Low, S.K. *et al.* Identification of six new genetic loci associated with atrial
658 fibrillation in the Japanese population. *Nat Genet* **49**, 953-958 (2017).
- 659 41. Suzuki, K. *et al.* Identification of 28 new susceptibility loci for type 2 diabetes
660 in the Japanese population. *Nat Genet* **51**, 379-386 (2019).
- 661 42. Akiyama, M. *et al.* Genome-wide association study identifies 112 new loci for
662 body mass index in the Japanese population. *Nat Genet* **49**, 1458-1467 (2017).
- 663 43. Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese
664 population links cell types to complex human diseases. *Nat Genet* **50**, 390-400
665 (2018).
- 666 44. Matoba, N. *et al.* GWAS of smoking behaviour in 165,436 Japanese people
667 reveals seven new loci and shared genetic architecture. *Nat Hum Behav* **3**,
668 471-477 (2019).
- 669 45. Lu, X.F. *et al.* Genome-wide association study in Han Chinese identifies four
670 new susceptibility loci for coronary artery disease. *Nature Genetics* **44**, 890-+
671 (2012).
- 672 46. Marzec, J. *et al.* A genetic study and meta-analysis of the genetic
673 predisposition of prostate cancer in a Chinese population. *Oncotarget* **7**,
674 21393-403 (2016).
- 675 47. Okada, Y. *et al.* Meta-analysis identifies multiple loci associated with kidney
676 function-related traits in east Asian populations. *Nat Genet* **44**, 904-9 (2012).
- 677 48. Zeng, C. *et al.* Identification of Susceptibility Loci and Genes for Colorectal
678 Cancer Risk. *Gastroenterology* **150**, 1633-1645 (2016).
- 679 49. Zhou, X. *et al.* Identification of genetic risk factors in the Chinese population
680 implicates a role of immune system in Alzheimer's disease pathogenesis. *Proc*
681 *Natl Acad Sci U S A* **115**, 1697-1706 (2018).
- 682 50. Foo, J.N. *et al.* Genome-wide association study of Parkinson's disease in East
683 Asians. *Hum Mol Genet* **26**, 226-232 (2017).
- 684 51. Langille, M.G.I. *et al.* Predictive functional profiling of microbial
685 communities using 16S rRNA marker gene sequences. *Nature Biotechnology*
686 **31**, 814 (2013).
- 687 52. Canga, Y. *et al.* Assessment of Atrial Conduction Times in Patients with
688 Newly Diagnosed Parkinson's Disease. *Parkinsons Dis* **2018**, 2916905 (2018).
- 689 53. Ihara, M. & Washida, K. Linking Atrial Fibrillation with Alzheimer's Disease:
690 Epidemiological, Pathological, and Mechanistic Evidence. *J Alzheimers Dis*
691 **62**, 61-72 (2018).
- 692 54. Conen, D. *et al.* Risk of Malignant Cancer Among Women With New-Onset
693 Atrial FibrillationAtrial Fibrillation and Risk of CancerAtrial Fibrillation and
694 Risk of Cancer. *JAMA Cardiology* **1**, 389-396 (2016).
- 695 55. Rousseeuw, P.J. Silhouettes: A graphical aid to the interpretation and
696 validation of cluster analysis. *Journal of Computational and Applied*
697 *Mathematics* **20**, 53-65 (1987).

-
- 698 56. Hu, F.B. *et al.* Prospective Study of Cataract Extraction and Risk of Coronary
699 Heart Disease in Women. *American Journal of Epidemiology* **153**, 875-881
700 (2001).
- 701 57. Javaheri, S. & Redline, S. Insomnia and Risk of Cardiovascular Disease. *Chest*
702 **152**, 435-444 (2017).
- 703 58. Strand, L.B. *et al.* Self-reported sleep duration and coronary heart disease
704 mortality: A large cohort study of 400,000 Taiwanese adults. *International*
705 *Journal of Cardiology* **207**, 246-251 (2016).
- 706 59. Blekhan, R. *et al.* Host genetic variation impacts microbiome composition
707 across human body sites. *Genome Biol* **16**, 191 (2015).
- 708 60. Cirulli, E.T. & Goldstein, D.B. Uncovering the roles of rare variants in
709 common disease through whole-genome sequencing. *Nat Rev Genet* **11**,
710 415-25 (2010).
- 711 61. Vinter, N., Christesen Amanda, M.S., Fenger-Grøn, M., Tjønneland, A. &
712 Frost, L. Atrial Fibrillation and Risk of Cancer: A Danish Population-Based
713 Cohort Study. *Journal of the American Heart Association* **7**, e009543 (2018).
- 714 62. Zoja, C. *et al.* Imatinib ameliorates renal disease and survival in murine lupus
715 autoimmune disease. *Kidney International* **70**, 97-103 (2006).
- 716 63. Boursi, B., Mamtani, R., Haynes, K. & Yang, Y.-X. Parkinson's disease and
717 colorectal cancer risk-A nested case control study. *Cancer epidemiology* **43**,
718 9-14 (2016).
- 719 64. Xie, X., Luo, X. & Xie, M. Association between Parkinson's disease and risk
720 of colorectal cancer. *Parkinsonism & Related Disorders* **35**, 42-47 (2017).
- 721 65. van Rheenen, W. *et al.* Genome-wide association analyses identify new risk
722 variants and the genetic architecture of amyotrophic lateral sclerosis. *Nat*
723 *Genet* **48**, 1043-8 (2016).
- 724 66. Lambert, J.C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new
725 susceptibility loci for Alzheimer's disease. *Nat Genet* **45**, 1452-8 (2013).
- 726 67. Pankratz, N. *et al.* Meta-analysis of Parkinson's disease: identification of a
727 novel locus, RIT2. *Annals of neurology* **71**, 370-384 (2012).
- 728 68. Anderson, C.A. *et al.* Data quality control in genetic case-control association
729 studies. *Nat Protoc* **5**, 1564-73 (2010).
- 730 69. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method
731 for thousands of genomes. *Nature Methods* **9**, 179 (2011).
- 732 70. Edgar, R.C. Search and clustering orders of magnitude faster than BLAST.
733 *Bioinformatics* **26**, 2460-2461 (2010).
- 734 71. Edgar, R.C. UPARSE: highly accurate OTU sequences from microbial
735 amplicon reads. *Nat Methods* **10**, 996-8 (2013).
- 736 72. Caporaso, J.G. *et al.* QIIME allows analysis of high-throughput community
737 sequencing data. *Nat Methods* **7**, 335-6 (2010).
- 738 73. Caliński, T. & Harabasz, J. A dendrite method for cluster analysis.
739 *Communications in Statistics* **3**, 1-27 (1974).
- 740 74. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding
741 from polygenicity in genome-wide association studies. *Nature Genetics* **47**,

-
- 742 291-295 (2015).
743 75. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes.
744 *Nucleic Acids Res* **28**, 27-30 (2000).
745 76. Kanehisa, M., Sato, Y., Furumichi, M., Morishima, K. & Tanabe, M. New
746 approach for understanding genome variations in KEGG. *Nucleic Acids Res*
747 **47**, D590-D595 (2019).
748 77. Kanehisa, M. Toward understanding the origin and evolution of cellular
749 organisms. *Protein Sci* (2019).
750 78. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of
751 biomolecular interaction networks. *Genome research* **13**, 2498-2504 (2003).
752 79. Kurtz, Z.D. *et al.* Sparse and Compositionally Robust Inference of Microbial
753 Ecological Networks. *PLOS Computational Biology* **11**, e1004226 (2015).
754 80. Reynolds, A.P., Richards, G., de la Iglesia, B. & Rayward-Smith, V.J.
755 Clustering Rules: A Comparison of Partitioning and Hierarchical Clustering
756 Algorithms. *Journal of Mathematical Modelling and Algorithms* **5**, 475-504
757 (2006).
758

759 **Figure legends**

760 **Figure 1 Association of host polygenic score with gut microbiome.** The participants
761 were divided into high and low polygenic score group according to median levels of
762 the polygenic score. The dots on the right of the box represent the distribution of
763 polygenic score. The dash line in the box is the position of median line and the solid
764 line is the position of mean line. The length of box depends on upper quartile and
765 lower quartile of datum. Sample size at the discovery stage is 1475, and that at
766 replication stage is 199. **(A).** Correlation of *Coriobacteriaceae* abundance with the
767 polygenic score (including 45 lead SNPs, Supplementary Table S8). **(B).** Correlation
768 of undefined species belonging to *Parabacteroides* genus
769 (*Parabacteroides_undefined*) with the polygenic score (including 32 lead SNPs,
770 Supplementary Table S8). **(C).** Correlation of *Odoribacter* presence with the
771 polygenic score (including 43 lead SNPs, Supplementary Table S8).

772

773 **Figure 2 Effect of host genetically predicted higher atrial fibrillation risk on gut**
774 **microbiome. (A).** Causal association of atrial fibrillation with abundance of
775 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Coprophilus*. The effect sizes
776 of atrial fibrillation on taxa are changes in abundance of bacteria (10-SD of
777 log-transformed) per genetically determined higher log odds of atrial fibrillation. **(B).**
778 Causal association of atrial fibrillation with presence of *Barnesiellaceae*,
779 *Veillonellaceae_undefined* and *Mitsuokella*. The effect size of atrial fibrillation on
780 taxa are present as odds ratio increase in log odds of atrial fibrillation. **(C).** The heat
781 map shows correlation of AF-associated taxa with predicted diseases. The grey
782 components show no significance of correlation with Bonferroni correction ($p > 0.05 /$
783 $(5.6 * 22)$, $p > 0.0004$).

784

785 **Figure 3 Association and cluster of diseases predicted by the gut microbiome. (A).**
786 Plot of clusters in Guangzhou Nutrition and Health Study (GNHS) cohort (n=1919).
787 **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5 clusters in
788 antibiotic-taking participants (n=18). The optimal cluster is 5 in GNHS cohort and 6
789 in the replication. The clusters share consistent components between two studies. In
790 contrast, components are different between antibiotic-taking participants and control
791 groups. Dimension1 (Dim1) and dimension2 (Dim2) can explain 40.1% and 13.1%
792 variance, respectively in GNHS cohort. The annotation for variables is as following.
793 AT: African trypanosomiasis, AD: Alzheimer's disease, V1: Amoebiasis, ALS:
794 Amyotrophic lateral sclerosis, BC: Bladder cancer, CD: Chagas disease, CML:
795 Chronic myeloid leukemia, CRC: Colorectal cancer, V2: Hepatitis C, HD:
796 Huntington's disease, HCM: Hypertrophic cardiomyopathy, V3: Influenza A, PD:
797 Parkinson's disease, V4: Pathways in cancer, V5: Prion disease, PCa: Prostate cancer,
798 RCC: Renal cell carcinoma, SLE: Systemic lupus erythematosus, V6: Tuberculosis,
799 T1DM: Type I diabetes mellitus, T2DM: Type II diabetes mellitus, V7: Vibrio
800 cholerae infection. **(D).** Plot of clusters in GNHS patients. Patients get only one of the
801 follow diseases: stroke (n=8), chronic hepatitis (n=19), coronary heart diseases (n=40),

802 cataract (n=124) and insomnia (n=68). **(E)**. Gut microbiome-predicted network of
803 relationship among different human complex diseases. The interaction is determined
804 by SPIEC-EASI with non-normalized predicted abundance data.

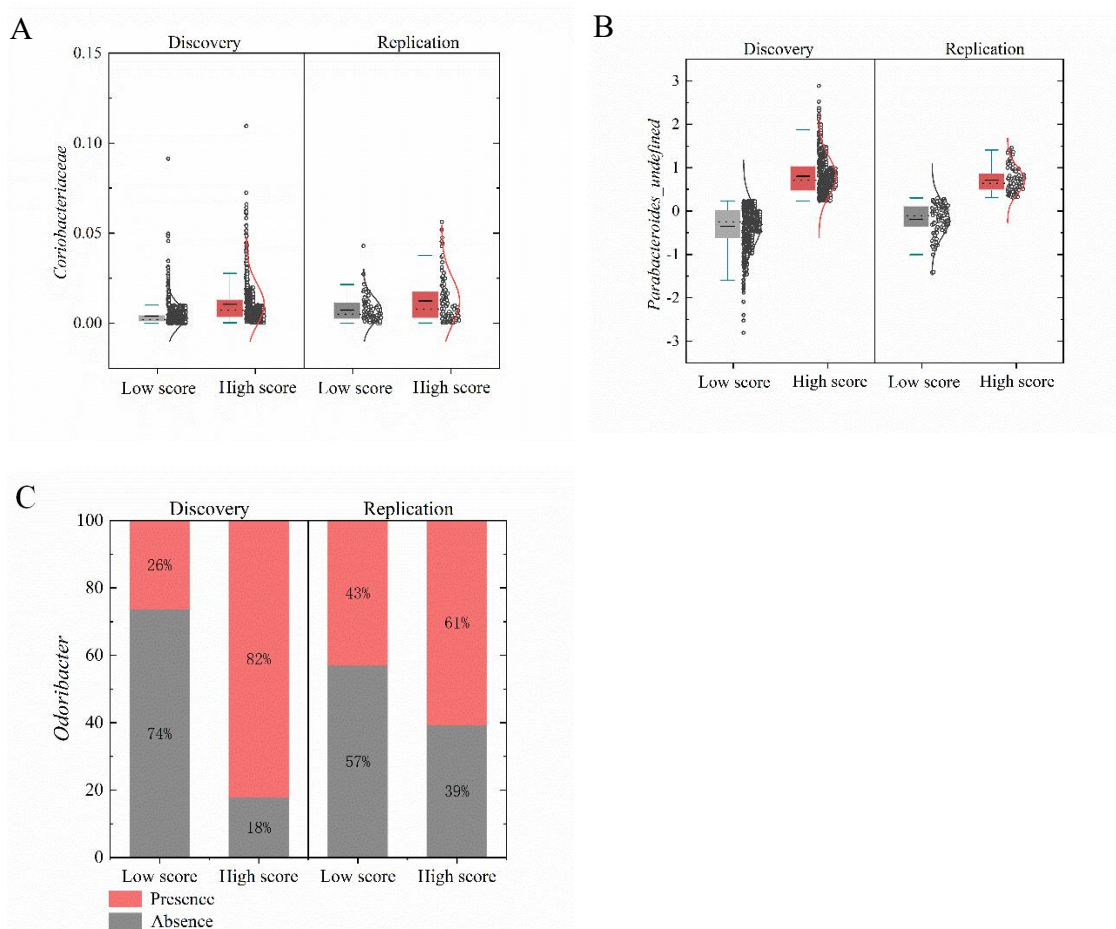
805

806 **Figure 4 Correlation of the human complex diseases with gut microbiome on**
807 **genus level.** The heat map shows correlation of predicted diseases and gut
808 microbiome on genus level. The grey components show no significance of correlation
809 with Bonferroni correction ($p > 0.05 / (5.6 * 22)$, $p > 0.0004$).

810

811

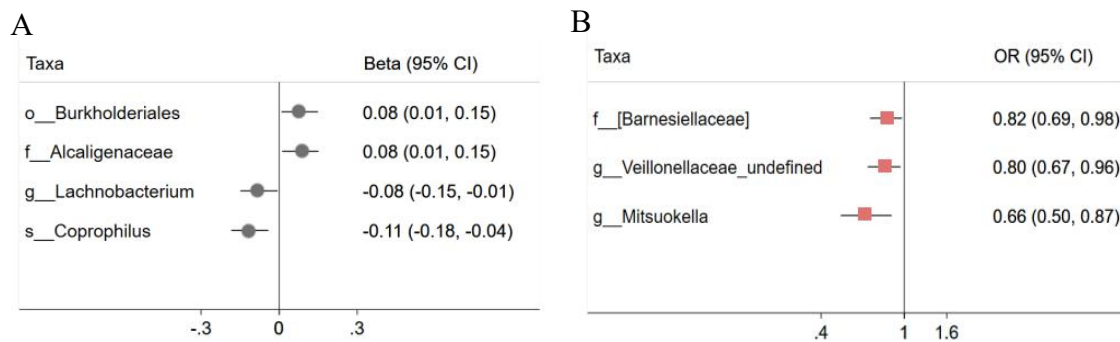
812 **Figure 1 Association of host polygenic score with gut microbiome.** The participants
813 were divided into high and low polygenic score group according to median levels of
814 the polygenic score. The dots on the right of the box represent the distribution of
815 polygenic score. The dash line in the box is the position of median line and the solid
816 line is the position of mean line. The length of box depends on upper quartile and
817 lower quartile of datum. Sample size at the discovery stage is 1475, and that at
818 replication stage is 199. **(A).** Correlation of *Coriobacteriaceae* abundance with the
819 polygenic score (including 45 lead SNPs, Supplementary Table S8). **(B).** Correlation
820 of undefined species belonging to *Parabacteroides* genus
821 (*Parabacteroides_undefined*) with the polygenic score (including 32 lead SNPs,
822 Supplementary Table S8). **(C).** Correlation of *Odoribacter* presence with the
823 polygenic score (including 43 lead SNPs, Supplementary Table S8).



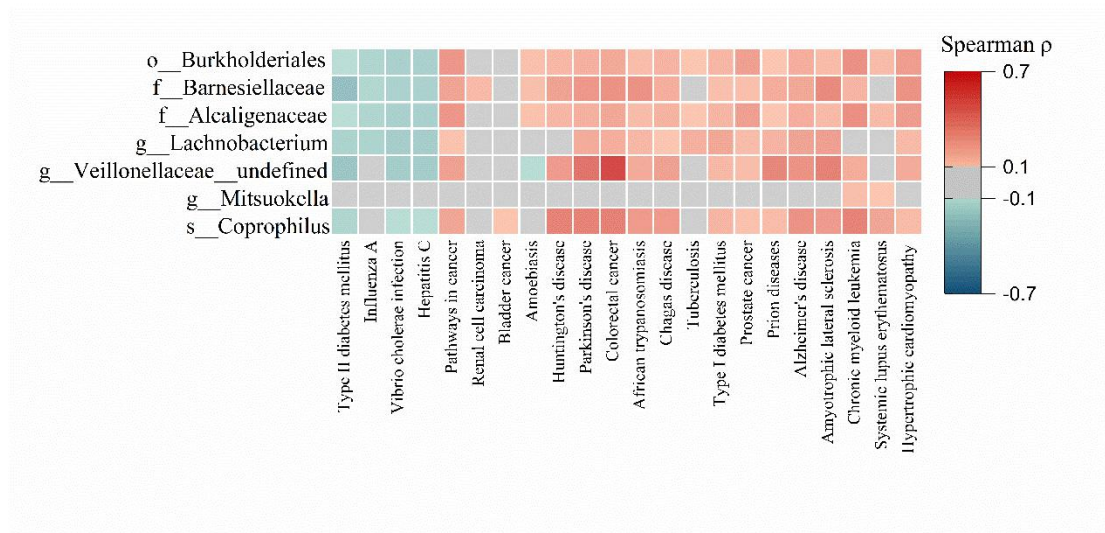
824

825

826 **Figure 2 Effect of host genetically predicted higher atrial fibrillation risk on gut**
 827 **microbiome. (A).** Causal association of atrial fibrillation with abundance of
 828 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Coprophilus*. The effect sizes
 829 of atrial fibrillation on taxa are changes in abundance of bacteria (10-SD of
 830 log-transformed) per genetically determined higher log odds of atrial fibrillation. **(B).**
 831 Causal association of atrial fibrillation with presence of *Barnesiellaceae*,
 832 *Veillonellaceae_undefined* and *Mitsuokella*. The effect size of atrial fibrillation on
 833 taxa are present as odds ratio increase in log odds of atrial fibrillation. **(C).** The heat
 834 map shows correlation of AF-associated taxa with predicted diseases. The grey
 835 components show no significance of correlation with Bonferroni correction ($p > 0.05 /$
 836 $(5.6 * 22)$, $p > 0.0004$).
 837



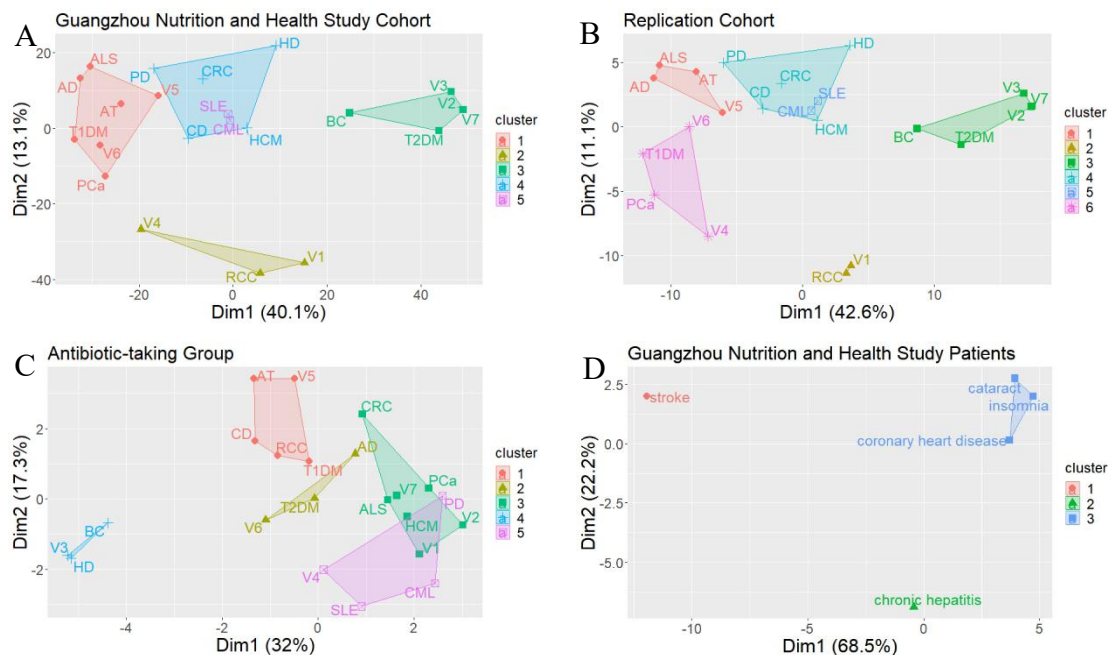
838 C



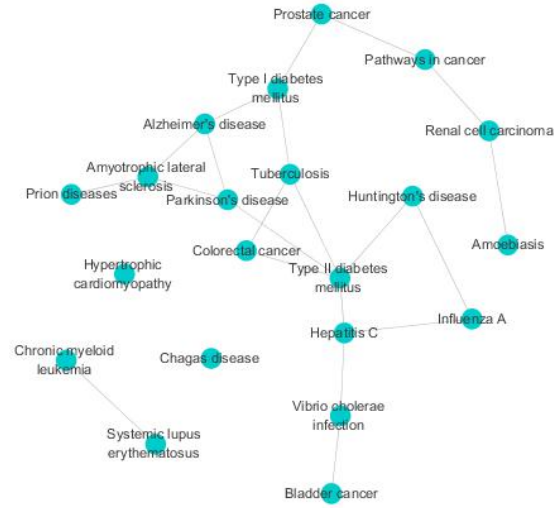
839

840 **Figure 3 Association and cluster of diseases predicted by the gut microbiome. (A).**
 841 Plot of clusters in Guangzhou Nutrition and Health Study (GNHS) cohort (n=1919).
 842 **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5 clusters in
 843 antibiotic-taking participants (n=18). The optimal cluster is 5 in GNHS cohort and 6
 844 in the replication. The clusters share consistent components between two studies. In
 845 contrast, components are different between antibiotic-taking participants and control
 846 groups. Dimension1 (Dim1) and dimension2 (Dim2) can explain 40.1% and 13.1%
 847 variance, respectively in GNHS cohort. The annotation for variables is as following.
 848 AT: African trypanosomiasis, AD: Alzheimer's disease, V1: Amoebiasis, ALS:
 849 Amyotrophic lateral sclerosis, BC: Bladder cancer, CD: Chagas disease, CML:
 850 Chronic myeloid leukemia, CRC: Colorectal cancer, V2: Hepatitis C, HD:
 851 Huntington's disease, HCM: Hypertrophic cardiomyopathy, V3: Influenza A, PD:
 852 Parkinson's disease, V4: Pathways in cancer, V5: Prion disease, PCa: Prostate cancer,
 853 RCC: Renal cell carcinoma, SLE: Systemic lupus erythematosus, V6: Tuberculosis,
 854 T1DM: Type I diabetes mellitus, T2DM: Type II diabetes mellitus, V7: Vibrio
 855 cholerae infection. **(D).** Plot of clusters in GNHS patients. Patients get only one of the
 856 follow diseases: stroke (n=8), chronic hepatitis (n=19), coronary heart diseases (n=40),
 857 cataract (n=124) and insomnia (n=68). **(E).** Gut microbiome-predicted network of
 858 relationship among different human complex diseases. The interaction is determined
 859 by SPIEC-EASI with non-normalized predicted abundance data.

860
861



E



862

863

864 **Figure 4 Correlation of the human complex diseases with gut microbiome.** The
 865 heat map shows correlation of predicted diseases and gut microbiome on genus level.
 866 The grey components show no significance of correlation with Bonferroni correction
 867 ($p > 0.05 / (5.6 * 22)$, $p > 0.0004$).

