
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: chenym@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*, undefined genus in family *Veillonellaceae* and
58 *Mitsuokella*. Further disease-microbiome feature analysis suggested that systemic
59 lupus erythematosus and chronic myeloid leukaemia shared common gut microbiome

60 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 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 As the associations of the microbiome with complex diseases and traits have been
190 widely reported [26], the genetic correlation between the gut microbiome and traits of
191 interest is less clear. Therefore, we applied bivariate GREML analysis to address this
192 question. The traits included BMI, fasting blood sugar (FBS), glycosylated
193 haemoglobin (HbA1c), systolic blood pressure (SBP), diastolic blood pressure (DBP),
194 high density lipoprotein cholesterol (HDL-C), low density lipoprotein cholesterol
195 (LDL-C), total cholesterol (TC) and triglyceride (TG), none of which could pass
196 Bonferroni correction. HDL-C was the only trait that had nominal genetic correlation
197 ($p < 0.05$) with gut microbes (specifically, *Desulfovibrionaceae* and [*Prevotella*],
198 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 reserve MR analysis to assess the potential causal effect
225 of complex human diseases on gut microbiome features. For the reserve MR analyses,
226 the diseases of interest included T2DM, AF, coronary artery disease (CAD), chronic
227 kidney disease (CKD), Alzheimer's disease (AD), CRC and prostatic cancer (PCa),
228 and their instrumental variables for the MR analysis were based on previous large-

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

242

243 **Microbiome features of human complex diseases**

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

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

255

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

266

267 We further asked whether the gut microbiome contributed to the novel clustering. To
268 this end, we repeated the analysis among participants who took antibiotics less than

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

281

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

290 T2DM (odds ratio: 1.176 [95% CI: 1.114, 1.244], $p=8.75\times 10^{-9}$).

291

292 **Discussion**

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

302

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

310 be heritable [7, 45].

311

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

324

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

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

335

336 Our results indicate that the gut microbiome helps reveal novel and interesting
337 relationships among complex human diseases, and different diseases may have
338 common and distinct gut microbiome features. A prior study including participants
339 from different countries identified three microbiome clusters [22]. Notably, this study
340 focused on classifying the individuals into distinct enterotypes regardless of the
341 individuals' health status, while in the present study, we described representative
342 microbiome features for diseases of interest. We provided an approach to interpret the
343 data from mechanistic studies based on the microbiome. The microbiome features
344 revealed a close association of AF with neurodegenerative diseases as well as cancers,
345 which was supported by prior studies showing that AF had a correlation with AD and
346 PD [40, 41], and AF patients had relatively higher risks of several cancers, including
347 lung cancer and CRC [42, 50]. We also observed that the microbiome features of SLE
348 and CML were highly similar. Interestingly, a tyrosine kinase inhibitor of platelet-
349 derived growth factor receptor, imatinib, was widely used to treat CML and
350 significantly ameliorated survival in murine models of SLE [51]. In addition, a close

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

358

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

370

371 **Conclusions**

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

377

378 **Method**

379 **Study participants and sample collection**

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

386

387 We included 199 participants with both genetic data and gut microbiome data as a
388 replication cohort, which belonged to the control arm of a case-control study of hip
389 fracture with the participants (52-83 years old) recruited between June 2009 and

390 August 2015 in Guangdong Province, China [18].

391

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

397

398 **Genotyping data**

399 For both the discovery and replication cohorts, DNA was extracted from leukocytes
400 using the TIANamp® Blood DNA Kit (DP348, TianGen Biotech Co, Ltd, China)
401 according to the manufacturer's instructions. DNA concentrations were determined
402 using the Qubit quantification system (Thermo Scientific, Wilmington, DE, US).
403 Extracted DNA was stored at -80°C . Genotyping was carried out with Illumina ASA-
404 750K arrays. Quality control and relatedness filters were performed by PLINK1.9
405 [58]. Individuals with a high or low proportion of heterozygous genotypes (outliers
406 defined as 3 standard deviations) were excluded [59]. Individuals who had different
407 ancestries (the first two principal components ± 5 standard deviations from the mean)
408 or related individuals ($\text{IBD} > 0.185$) were excluded [59]. Variants were mapped to the
409 1000 Genomes Phase 3 v5 by SHAPEIT [60, 61], and then we conducted genome-

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

415

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

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

427

428 Fastq files were demultiplexed by MiSeq Controller Software. Ultra-fast sequence
429 analysis (USEARCH) was performed to trim the sequence for amplification primers,

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

439

440 **Statistical analysis**

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

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

447

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

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

457

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

461

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

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

477

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

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

487

488

489 Constructing polygenic scores for taxa and alpha diversity

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

501

502 The effective number of independent taxa

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

509 which can be calculated as

510
$$\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\} = P^{-1}AP$$

511 The effective number of independent taxa can be calculated as

512
$$\frac{(\sum_{i=1}^n \lambda_i)^2}{\sum_{i=1}^n \lambda_i^2}$$

513

514 **Bidirectional MR analysis**

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

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

536

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

544

545 **Pathway analysis**

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

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

555

556 **Construction of the microbiome risk score**

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

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

570

571 **Clustering diseases**

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

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

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

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

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

577 $x_2, \dots, x_n)$, the dissimilarity between objects x_i and x_j is denoted by $d(i, j)$. The

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

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

580 function to minimize the model is given by

581
$$\sum_{i=1}^n \sum_{j=1}^n d(i, j) z_{ij}$$

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

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

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

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

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

587
$$M_i = \min \forall w (b(i, w)).$$

588 The silhouette width for object i can be calculated by

589
$$sw_i = \frac{M_i - N_i}{\max(M_i, N_i)}$$

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

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

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

604
$$J = \frac{a}{a + b + c}$$

605

606 **Availability of data and materials**

607 The raw data for 16S rRNA gene sequences are available in the CNSA
608 (<https://db.cngb.org/cnsa/>) of CNGBdb at accession number CNP0000829. Original R
609 scripts are available in GitHub
610 (<https://github.com/hsufengzhe/microbiome/tree/master>). Requests for the metadata
611 from this study can be submitted via email to zhengjusheng@westlake.edu.cn. A
612 proposal is also required for approval.

613

614 **Acknowledgements**

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

617 **Ethics approval and consent to participate**

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

621 **Consent for publication**

622 Not applicable.

623 **Competing interests**

624 The authors declare no conflicts of interest.

625 **Authors' contributions**

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

630 **Funding**

631 This study was funded by the National Natural Science Foundation of China (81903316,
632 81773416), the Zhejiang Ten-thousand Talents Program (101396522001) and the 5010

633 Program for Clinical Research (2007032) of Sun Yat-sen University (Guangzhou,
634 China).

References

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

-
- 675 microbiota. *Nature*. 2018;555(7695):210-5; doi:10.1038/nature25973.
- 676 14. Duvallat C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut
677 microbiome studies identifies disease-specific and shared responses. *Nature*
678 *communications*. 2017;8(1):1784-; doi:10.1038/s41467-017-01973-8.
- 679 15. Cheng S, Han B, Ding M, Wen Y, Ma M, Zhang L, et al. Identifying
680 psychiatric disorder-associated gut microbiota using microbiota-related gene
681 set enrichment analysis. *Briefings in Bioinformatics*. 2019;
682 doi:10.1093/bib/bbz034.
- 683 16. Jackson MA, Verdi S, Maxan M-E, Shin CM, Zierer J, Bowyer RCE, et al.
684 Gut microbiota associations with common diseases and prescription
685 medications in a population-based cohort. *Nature Communications*.
686 2018;9(1):2655; doi:10.1038/s41467-018-05184-7.
- 687 17. Cao Y, Wang C, Guan K, Xu Y, Su Y-X, Chen YM. Association of magnesium
688 in serum and urine with carotid intima-media thickness and serum lipids in
689 middle-aged and elderly Chinese: a community-based cross-sectional study.
690 *European journal of nutrition*. 2015;55; doi:10.1007/s00394-015-0839-8.
- 691 18. Sun L-L, Li B-L, Xie H-L, Fan F, Yu W-Z, Wu B-H, et al. Associations
692 between the dietary intake of antioxidant nutrients and the risk of hip fracture
693 in elderly Chinese: A case-control study. *The British journal of nutrition*.
694 2014;112:1-9; doi:10.1017/S0007114514002773.
- 695 19. Lim MY, You HJ, Yoon HS, Kwon B, Lee JY, Lee S, et al. The effect of
696 heritability and host genetics on the gut microbiota and metabolic syndrome.
697 *Gut*. 2017;66(6):1031-8; doi:10.1136/gutjnl-2015-311326.
- 698 20. Davenport ER. Elucidating the role of the host genome in shaping microbiome
699 composition. *Gut microbes*. 2016;7(2):178-84;
700 doi:10.1080/19490976.2016.1155022.
- 701 21. Davenport ER, Cusanovich DA, Michelini K, Barreiro LB, Ober C, Gilad Y.
702 Genome-Wide Association Studies of the Human Gut Microbiota. *PloS one*.
703 2015;10(11):e0140301-e; doi:10.1371/journal.pone.0140301.
- 704 22. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, et al.
705 Enterotypes of the human gut microbiome. *Nature*. 2011;473(7346):174-80;
706 doi:10.1038/nature09944.
- 707 23. Hua X, Song L, Yu G, Goedert JJ, Abnet CC, Landi MT, et al.
708 MicrobiomeGWAS: a tool for identifying host genetic variants associated with
709 microbiome composition. *bioRxiv*. 2015:031187; doi:10.1101/031187.
- 710 24. Ruhlemann MC, Degenhardt F, Thingholm LB, Wang J, Skieceviciene J,
711 Rausch P, et al. Application of the distance-based F test in an mGWAS
712 investigating beta diversity of intestinal microbiota identifies variants in
713 SLC9A8 (NHE8) and 3 other loci. *Gut Microbes*. 2018;9(1):68-75;
714 doi:10.1080/19490976.2017.1356979.
- 715 25. Bonder MJ, Kurilshikov A, Tigchelaar EF, Mujagic Z, Imhann F, Vila AV, et
716 al. The effect of host genetics on the gut microbiome. *Nat Genet*.

-
- 717 2016;48(11):1407-12; doi:10.1038/ng.3663.
- 718 26. Tang WHW, Kitai T, Hazen SL. Gut Microbiota in Cardiovascular Health and
719 Disease. *Circulation research*. 2017;120(7):1183-96;
720 doi:10.1161/CIRCRESAHA.117.309715.
- 721 27. Low SK, Takahashi A, Ebana Y, Ozaki K, Christophersen IE, Ellinor PT, et al.
722 Identification of six new genetic loci associated with atrial fibrillation in the
723 Japanese population. *Nat Genet*. 2017;49(6):953-8; doi:10.1038/ng.3842.
- 724 28. Suzuki K, Akiyama M, Ishigaki K, Kanai M, Hosoe J, Shojima N, et al.
725 Identification of 28 new susceptibility loci for type 2 diabetes in the Japanese
726 population. *Nat Genet*. 2019;51(3):379-86; doi:10.1038/s41588-018-0332-4.
- 727 29. Akiyama M, Okada Y, Kanai M, Takahashi A, Momozawa Y, Ikeda M, et al.
728 Genome-wide association study identifies 112 new loci for body mass index in
729 the Japanese population. *Nat Genet*. 2017;49(10):1458-67;
730 doi:10.1038/ng.3951.
- 731 30. Kanai M, Akiyama M, Takahashi A, Matoba N, Momozawa Y, Ikeda M, et al.
732 Genetic analysis of quantitative traits in the Japanese population links cell
733 types to complex human diseases. *Nat Genet*. 2018;50(3):390-400;
734 doi:10.1038/s41588-018-0047-6.
- 735 31. Matoba N, Akiyama M, Ishigaki K, Kanai M, Takahashi A, Momozawa Y, et
736 al. GWAS of smoking behaviour in 165,436 Japanese people reveals seven
737 new loci and shared genetic architecture. *Nat Hum Behav*. 2019;3(5):471-7;
738 doi:10.1038/s41562-019-0557-y.
- 739 32. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, et al. Genetics of
740 rheumatoid arthritis contributes to biology and drug discovery. *Nature*.
741 2014;506(7488):376-81; doi:10.1038/nature12873.
- 742 33. Lu XF, Wang LY, Chen SF, He L, Yang XL, Shi YY, et al. Genome-wide
743 association study in Han Chinese identifies four new susceptibility loci for
744 coronary artery disease. *Nature Genetics*. 2012;44(8):890-+;
745 doi:10.1038/ng.2337.
- 746 34. Marzec J, Mao X, Li M, Wang M, Feng N, Gou X, et al. A genetic study and
747 meta-analysis of the genetic predisposition of prostate cancer in a Chinese
748 population. *Oncotarget*. 2016;7(16):21393-403; doi:10.18632/oncotarget.7250.
- 749 35. Okada Y, Sim X, Go MJ, Wu JY, Gu D, Takeuchi F, et al. Meta-analysis
750 identifies multiple loci associated with kidney function-related traits in east
751 Asian populations. *Nat Genet*. 2012;44(8):904-9; doi:10.1038/ng.2352.
- 752 36. Zeng C, Matsuda K, Jia WH, Chang J, Kweon SS, Xiang YB, et al.
753 Identification of Susceptibility Loci and Genes for Colorectal Cancer Risk.
754 *Gastroenterology*. 2016;150(7):1633-45; doi:10.1053/j.gastro.2016.02.076.
- 755 37. Zhou X, Chen Y, Mok KY, Zhao Q, Chen K, Chen Y, et al. Identification of
756 genetic risk factors in the Chinese population implicates a role of immune
757 system in Alzheimer's disease pathogenesis. *Proceedings of the National
758 Academy of Sciences*. 2018;115(8):1697; doi:10.1073/pnas.1715554115.

-
- 759 38. Gan W, Walters RG, Holmes MV, Bragg F, Millwood IY, Banasik K, Chen Y,
760 Du H, Iona A, Mahajan A, et al: Evaluation of type 2 diabetes genetic risk
761 variants in Chinese adults: findings from 93,000 individuals from the China
762 Kadoorie Biobank. *Diabetologia*. 2016;59(7):1446-1457. doi:10.1007/s00125-
763 016-3920-9
- 764 39. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA,
765 et al. Predictive functional profiling of microbial communities using 16S
766 rRNA marker gene sequences. *Nature Biotechnology*. 2013;31:814;
767 doi:10.1038/nbt.2676.
- 768 40. Canga Y, Emre A, Yuksel GA, Karatas MB, Yelgec NS, Gurkan U, et al.
769 Assessment of Atrial Conduction Times in Patients with Newly Diagnosed
770 Parkinson's Disease. *Parkinsons Dis*. 2018;2018:2916905;
771 doi:10.1155/2018/2916905.
- 772 41. Ihara M, Washida K. Linking Atrial Fibrillation with Alzheimer's Disease:
773 Epidemiological, Pathological, and Mechanistic Evidence. *J Alzheimers Dis*.
774 2018;62(1):61-72; doi:10.3233/JAD-170970.
- 775 42. Conen D, Wong JA, Sandhu RK, Cook NR, Lee I-M, Buring JE, et al. Risk of
776 Malignant Cancer Among Women With New-Onset Atrial Fibrillation. *JAMA*
777 *Cardiology*. 2016;1(4):389-96; doi:10.1001/jamacardio.2016.0280.
- 778 43. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation
779 of cluster analysis. *Journal of Computational and Applied Mathematics*.
780 1987;20:53-65; doi:[https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
- 781 44. Blekhman R, Goodrich JK, Huang K, Sun Q, Bukowski R, Bell JT, et al. Host
782 genetic variation impacts microbiome composition across human body sites.
783 *Genome Biol*. 2015;16:191; doi:10.1186/s13059-015-0759-1.
- 784 45. Goodrich JK, Davenport ER, Clark AG, Ley RE. The Relationship Between
785 the Human Genome and Microbiome Comes into View. *Annu Rev Genet*.
786 2017;51:413-33; doi:10.1146/annurev-genet-110711-155532.
- 787 46. Kuehbach T, Rehman A, Lepage P, Hellmig S, Fölsch UR, Schreiber S, et al.
788 Intestinal TM7 bacterial phylogenies in active inflammatory bowel disease.
789 *Journal of Medical Microbiology*. 2008;57(12):1569-76;
790 doi:<https://doi.org/10.1099/jmm.0.47719-0>.
- 791 47. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, et al. Cultivation
792 of a human-associated TM7 phylotype reveals a reduced genome and epibiotic
793 parasitic lifestyle. *Proceedings of the National Academy of Sciences of the*
794 *United States of America*. 2015;112(1):244-9; doi:10.1073/pnas.1419038112.
- 795 48. Bor B, Bedree JK, Shi W, McLean JS, He X. Saccharibacteria (TM7) in the
796 Human Oral Microbiome. *Journal of Dental Research*. 2019;98(5):500-9;
797 doi:10.1177/0022034519831671.
- 798 49. Cirulli ET, Goldstein DB. Uncovering the roles of rare variants in common
799 disease through whole-genome sequencing. *Nat Rev Genet*. 2010;11(6):415-
800 25; doi:10.1038/nrg2779.

-
- 801 50. Vinter N, Christesen Amanda MS, Fenger-Grøn M, Tjønneland A, Frost L.
802 Atrial Fibrillation and Risk of Cancer: A Danish Population-Based Cohort
803 Study. *Journal of the American Heart Association*. 2018;
804 doi:10.1161/JAHA.118.009543.
- 805 51. Zoja C, Corna D, Rottoli D, Zanchi C, Abbate M, Remuzzi G. Imatinib
806 ameliorates renal disease and survival in murine lupus autoimmune disease.
807 *Kidney International*. 2006;70(1):97-103;
808 doi:<https://doi.org/10.1038/sj.ki.5001528>.
- 809 52. Boursi B, Mamtani R, Haynes K, Yang Y-X. Parkinson's disease and colorectal
810 cancer risk-A nested case control study. *Cancer Epidemiol*. 2016;43:9-14;
811 doi:10.1016/j.canep.2016.05.007.
- 812 53. Xie X, Luo X, Xie M. Association between Parkinson's disease and risk of
813 colorectal cancer. *Parkinsonism & Related Disorders*. 2017;35:42-7;
814 doi:10.1016/j.parkreldis.2016.11.011.
- 815 54. van Rheenen W, Shatunov A, Dekker AM, McLaughlin RL, Diekstra FP, Pulit
816 SL, et al. Genome-wide association analyses identify new risk variants and the
817 genetic architecture of amyotrophic lateral sclerosis. *Nat Genet*.
818 2016;48(9):1043-8; doi:10.1038/ng.3622.
- 819 55. Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C,
820 et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci
821 for Alzheimer's disease. *Nat Genet*. 2013;45(12):1452-8; doi:10.1038/ng.2802.
- 822 56. Pankratz N, Beecham GW, DeStefano AL, Dawson TM, Doheny KF, Factor
823 SA, et al. Meta-analysis of Parkinson's disease: identification of a novel locus,
824 RIT2. *Ann Neurol*. 2012;71(3):370-84; doi:10.1002/ana.22687.
- 825 57. Zhao L, Wang G, Siegel P, He C, Wang H, Zhao W, et al. Quantitative genetic
826 background of the host influences gut microbiomes in chickens. *Sci Rep*.
827 2013;3:1163-; doi:10.1038/srep01163.
- 828 58. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al.
829 PLINK: a tool set for whole-genome association and population-based linkage
830 analyses. *American journal of human genetics*. 2007;81(3):559-75;
831 doi:10.1086/519795.
- 832 59. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan
833 KT. Data quality control in genetic case-control association studies. *Nat*
834 *Protoc*. 2010;5(9):1564-73; doi:10.1038/nprot.2010.116.
- 835 60. Delaneau O, Marchini J, Zagury J-F. A linear complexity phasing method for
836 thousands of genomes. *Nature Methods*. 2011;9:179; doi:10.1038/nmeth.1785.
- 837 61. Delaneau O, Marchini J, The Genomes Project C, McVean GA, Donnelly P,
838 Lunter G, et al. Integrating sequence and array data to create an improved
839 1000 Genomes Project haplotype reference panel. *Nature Communications*.
840 2014;5:3934; doi:10.1038/ncomms4934.
- 841 62. Das S, Forer L, Schönherr S, Sidore C, Locke AE, Kwong A, et al. Next-
842 generation genotype imputation service and methods. *Nature Genetics*.

-
- 843 2016;48:1284; doi:10.1038/ng.3656.
- 844 63. Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, et al. The
845 international Genome sample resource (IGSR): A worldwide collection of
846 genome variation incorporating the 1000 Genomes Project data. *Nucleic Acids
847 Research*. 2016; doi:10.1093/nar/gkw829.
- 848 64. Okada Y, Kim K, Han B, Pillai NE, Ong RT, Saw WY, et al. Risk for ACPA-
849 positive rheumatoid arthritis is driven by shared HLA amino acid
850 polymorphisms in Asian and European populations. *Hum Mol Genet*.
851 2014;23(25):6916-26; doi:10.1093/hmg/ddu387.
- 852 65. Pillai NE, Okada Y, Saw WY, Ong RT, Wang X, Tantoso E, et al. Predicting
853 HLA alleles from high-resolution SNP data in three Southeast Asian
854 populations. *Hum Mol Genet*. 2014;23(16):4443-51;
855 doi:10.1093/hmg/ddu149.
- 856 66. Jia X, Han B, Onengut-Gumuscu S, Chen W-M, Concannon PJ, Rich SS, et al.
857 Imputing Amino Acid Polymorphisms in Human Leukocyte Antigens. *PLOS
858 ONE*. 2013;8(6):e64683; doi:10.1371/journal.pone.0064683.
- 859 67. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
860 *Bioinformatics*. 2010;26(19):2460-1; doi:10.1093/bioinformatics/btq461.
- 861 68. Edgar RC. UPARSE: highly accurate OTU sequences from microbial
862 amplicon reads. *Nat Methods*. 2013;10(10):996-8; doi:10.1038/nmeth.2604.
- 863 69. Second Genome, Inc: the Greengenes
864 Databases.<http://greengenes.secondgenome.com/>. Accessed 12 Mar 2019.
- 865 70. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello
866 EK, et al. QIIME allows analysis of high-throughput community sequencing
867 data. *Nat Methods*. 2010;7(5):335-6; doi:10.1038/nmeth.f.303.
- 868 71. Lee SH, Wray NR, Goddard ME, Visscher PM. Estimating missing heritability
869 for disease from genome-wide association studies. *Am J Hum Genet*.
870 2011;88(3):294-305; doi:10.1016/j.ajhg.2011.02.002.
- 871 72. Visscher PM, Hemani G, Vinkhuyzen AAE, Chen G-B, Lee SH, Wray NR, et
872 al. Statistical Power to Detect Genetic (Co)Variance of Complex Traits Using
873 SNP Data in Unrelated Samples. *PLOS Genetics*. 2014;10(4):e1004269;
874 doi:10.1371/journal.pgen.1004269.
- 875 73. Caliński T, Harabasz J. A dendrite method for cluster analysis.
876 *Communications in Statistics*. 1974;3(1):1-27;
877 doi:10.1080/03610927408827101.
- 878 74. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide
879 complex trait analysis. *Am J Hum Genet*. 2011;88(1):76-82;
880 doi:10.1016/j.ajhg.2010.11.011.
- 881 75. Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and
882 pitfalls in the application of mixed-model association methods. *Nat Genet*.
883 2014;46(2):100-6; doi:10.1038/ng.2876.
- 884 76. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al.

-
- 885 Common SNPs explain a large proportion of the heritability for human height.
886 Nature Genetics. 2010;42:565; doi:10.1038/ng.608.
- 887 77. Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Patterson N, et
888 al. LD Score regression distinguishes confounding from polygenicity in
889 genome-wide association studies. Nature Genetics. 2015;47(3):291-5;
890 doi:10.1038/ng.3211.
- 891 78. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR. Estimation of
892 pleiotropy between complex diseases using single-nucleotide polymorphism-
893 derived genomic relationships and restricted maximum likelihood.
894 Bioinformatics. 2012;28(19):2540-2; doi:10.1093/bioinformatics/bts474.
- 895 79. Wang H, Zhang F, Zeng J, Wu Y, Kemper KE, Xue A, et al. Genotype-by-
896 environment interactions inferred from genetic effects on phenotypic
897 variability in the UK Biobank. Science Advances. 2019;5(8):eaaw3538;
898 doi:10.1126/sciadv.aaw3538.
- 899 80. Bretherton CS, Widmann M, Dymnikov VP, Wallace JM, Bladé I. The
900 Effective Number of Spatial Degrees of Freedom of a Time-Varying Field.
901 Journal of Climate. 1999;12(7):1990-2009; doi:10.1175/1520-
902 0442(1999)012<1990:Tenosd>2.0.Co;2.
- 903 81. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in
904 Mendelian Randomization with Some Invalid Instruments Using a Weighted
905 Median Estimator. Genet Epidemiol. 2016;40(4):304-14;
906 doi:10.1002/gepi.21965.
- 907 82. Burgess S, Thompson SG. Interpreting findings from Mendelian
908 randomization using the MR-Egger method. Eur J Epidemiol. 2017;32(5):377-
909 89; doi:10.1007/s10654-017-0255-x.
- 910 83. Verbanck M, Chen C-Y, Neale B, Do R. Detection of widespread horizontal
911 pleiotropy in causal relationships inferred from Mendelian randomization
912 between complex traits and diseases. Nature Genetics. 2018;50(5):693-8;
913 doi:10.1038/s41588-018-0099-7.
- 914 84. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes.
915 Nucleic Acids Res. 2000;28(1):27-30; doi:10.1093/nar/28.1.27.
- 916 85. Kanehisa M, Sato Y, Furumichi M, Morishima K, Tanabe M. New approach
917 for understanding genome variations in KEGG. Nucleic Acids Res. 2019;
918 doi:10.1093/nar/gky962.
- 919 86. Kanehisa M. Toward understanding the origin and evolution of cellular
920 organisms. Protein Sci. 2019; doi:10.1002/pro.3715.
- 921 87. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al.
922 Cytoscape: a software environment for integrated models of biomolecular
923 interaction networks. Genome Res. 2003;13(11):2498-504;
924 doi:10.1101/gr.1239303.
- 925 88. Kurtz ZD, Müller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA.
926 Sparse and Compositionally Robust Inference of Microbial Ecological

-
- 927 Networks. PLOS Computational Biology. 2015;
928 doi:10.1371/journal.pcbi.1004226.
- 929 89. Reynolds AP, Richards G, de la Iglesia B, Rayward-Smith VJ. Clustering
930 Rules: A Comparison of Partitioning and Hierarchical Clustering Algorithms.
931 Journal of Mathematical Modelling and Algorithms. 2006;5(4):475-504;
932 doi:10.1007/s10852-005-9022-1.
- 933 90. Kaufman L, Rousseeuw P. Partitioning Around Medoids (Program PAM).
934 John Wiley & Sons, Inc; 1990. p. 68-125.

935 **Figure legends**

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

953

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

957

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

967

968 **Figure 4 Association and cluster of diseases predicted by the gut microbiome.**
969 **(A).** Plot of clusters in the Guangzhou Nutrition and Health Study (GNHS) cohort

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

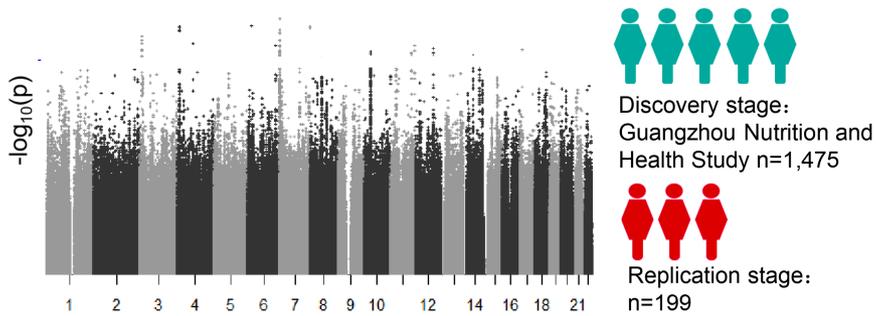
988

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

993

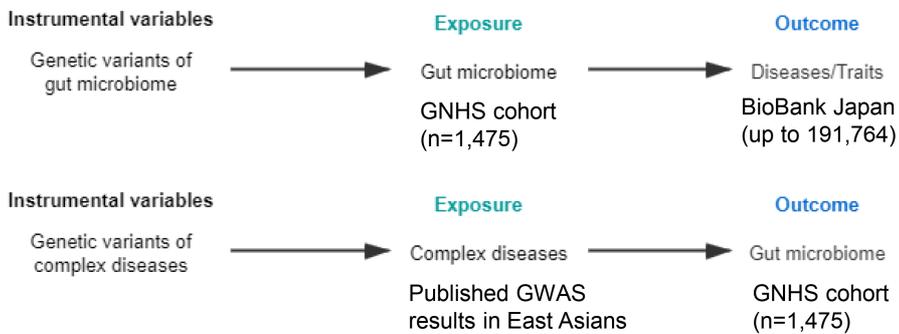
994 **Figure 1 Study overview.** The figure shows the highlights of our study. First, we
995 performed a microbiome genome-wide association study in a Chinese population
996 (Step A). We validated significant genetic variants reported in previous studies and
997 replicated our results in an independent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

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



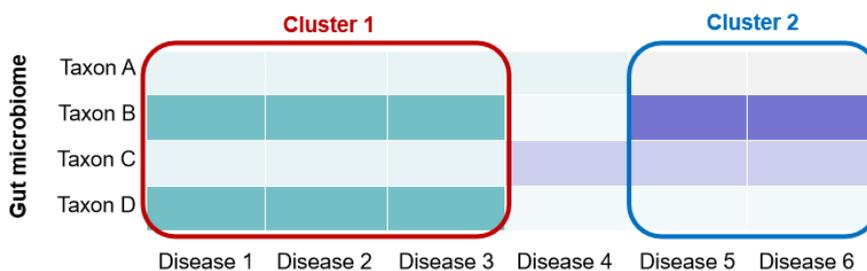
B. The causal relationships between gut microbiome and human complex diseases.

Bi-directional Mendelian randomization

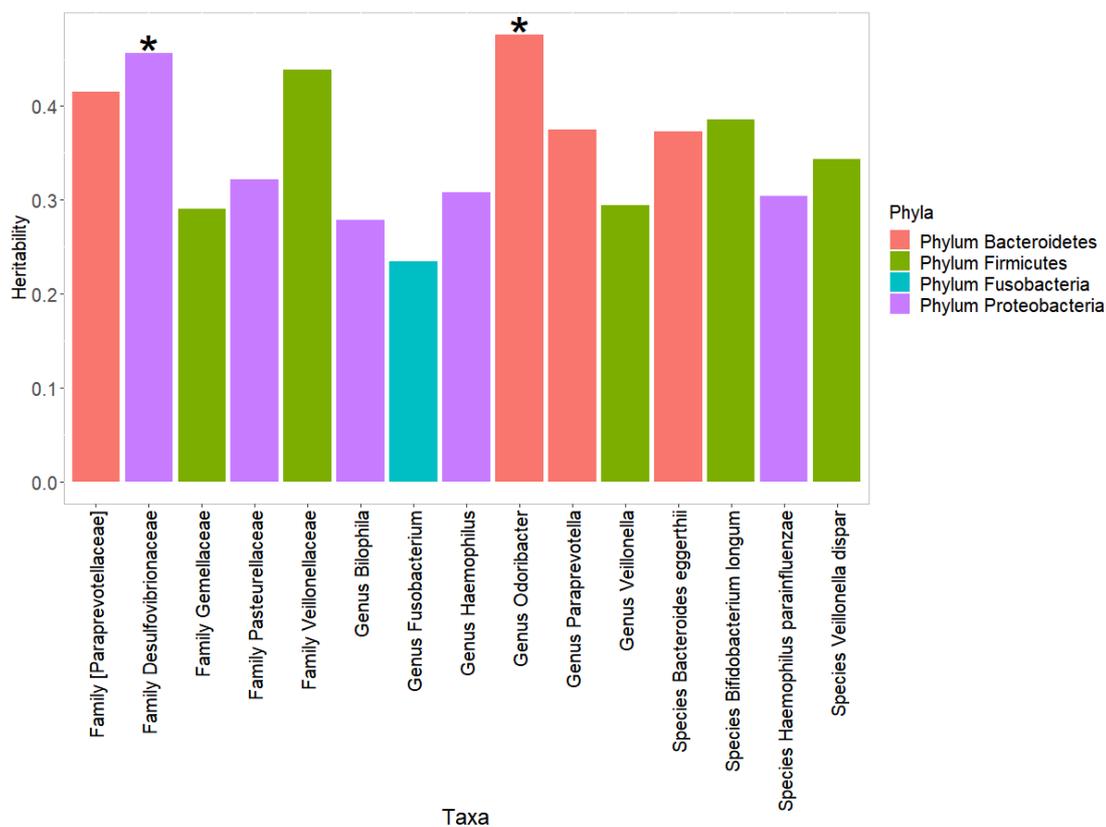


C. The shared and distinct microbiome features among human complex diseases.

1,919 participants from GNHS cohort.



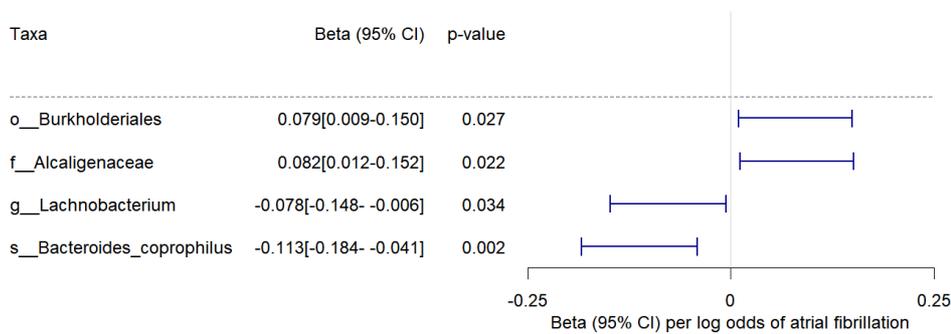
1012 **Figure 2 The SNP-based heritability of the gut microbiome.** The plot shows the
 1013 taxa with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, where n
 1014 is the effective number of independent taxa in each taxonomic level.



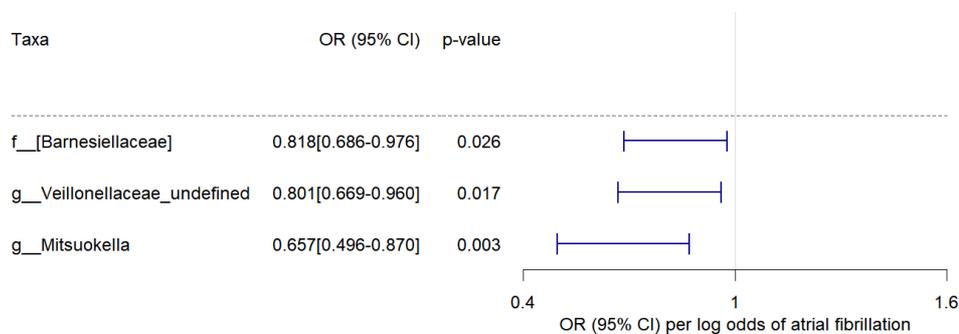
1015

1016 **Figure 3 Effect of host genetically predicted higher atrial fibrillation risk on the**
 1017 **gut microbiome. (A).** Causal association of atrial fibrillation with the abundance of
 1018 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Bacteroides coprophilus*. The
 1019 magnitude of the effect of atrial fibrillation on taxa is dependent on changes in the
 1020 abundance of bacteria (1-SD of the log-transformed abundance) per genetically
 1021 determined higher log odds of atrial fibrillation. **(B).** Causal association of atrial
 1022 fibrillation with the presence of *Barnesiellaceae*, undefined genus in family
 1023 *Veillonellaceae* and *Mitsuokella*. The magnitude of the effect of atrial fibrillation on
 1024 taxa is presented as an odds ratio increase in the log odds of atrial fibrillation.

A



B



1025

1026

1027

1028 **Figure 4 Association and cluster of diseases predicted by the gut microbiome.**
1029 **(A).** Plot of clusters in the Guangzhou Nutrition and Health Study (GNHS) cohort
1030 (n=1919). **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5
1031 clusters in antibiotic-taking participants (n=18). The optimal cluster was 5 in the
1032 GNHS cohort and 6 in the replication cohort. The clusters share consistent
1033 components between the two studies. In contrast, components are different between
1034 antibiotic-taking participants and control groups. Dimension1 (Dim1) and dimension2
1035 (Dim2) explained 40.1% and 13.1% of the variance, respectively, in the GNHS
1036 cohort. The annotation for variables is as follows. AT: African trypanosomiasis, AD:
1037 Alzheimer's disease, V1: Amoebiasis, ALS: Amyotrophic lateral sclerosis, BC:
1038 Bladder cancer, CD: Chagas disease, CML: Chronic myeloid leukaemia, CRC:
1039 Colorectal cancer, V2: Hepatitis C, HD: Huntington's disease, HCM: Hypertrophic
1040 cardiomyopathy, V3: Influenza A, PD: Parkinson's disease, V4: Pathways in cancer,
1041 V5: Prion disease, PCa: Prostate cancer, RCC: Renal cell carcinoma, SLE: Systemic
1042 lupus erythematosus, V6: Tuberculosis, T1DM: Type I diabetes mellitus, T2DM:
1043 Type II diabetes mellitus, V7: *Vibrio cholerae* infection. **(D).** Gut microbiome-
1044 predicted network of relationships among different human complex diseases. The
1045 relationship between diseases is determined by SPIEC-EASI with non-normalized
1046 predicted abundance data. The diseases that shared the same edge had the gut
1047 microbiome-predicted correlation.

1048

1049

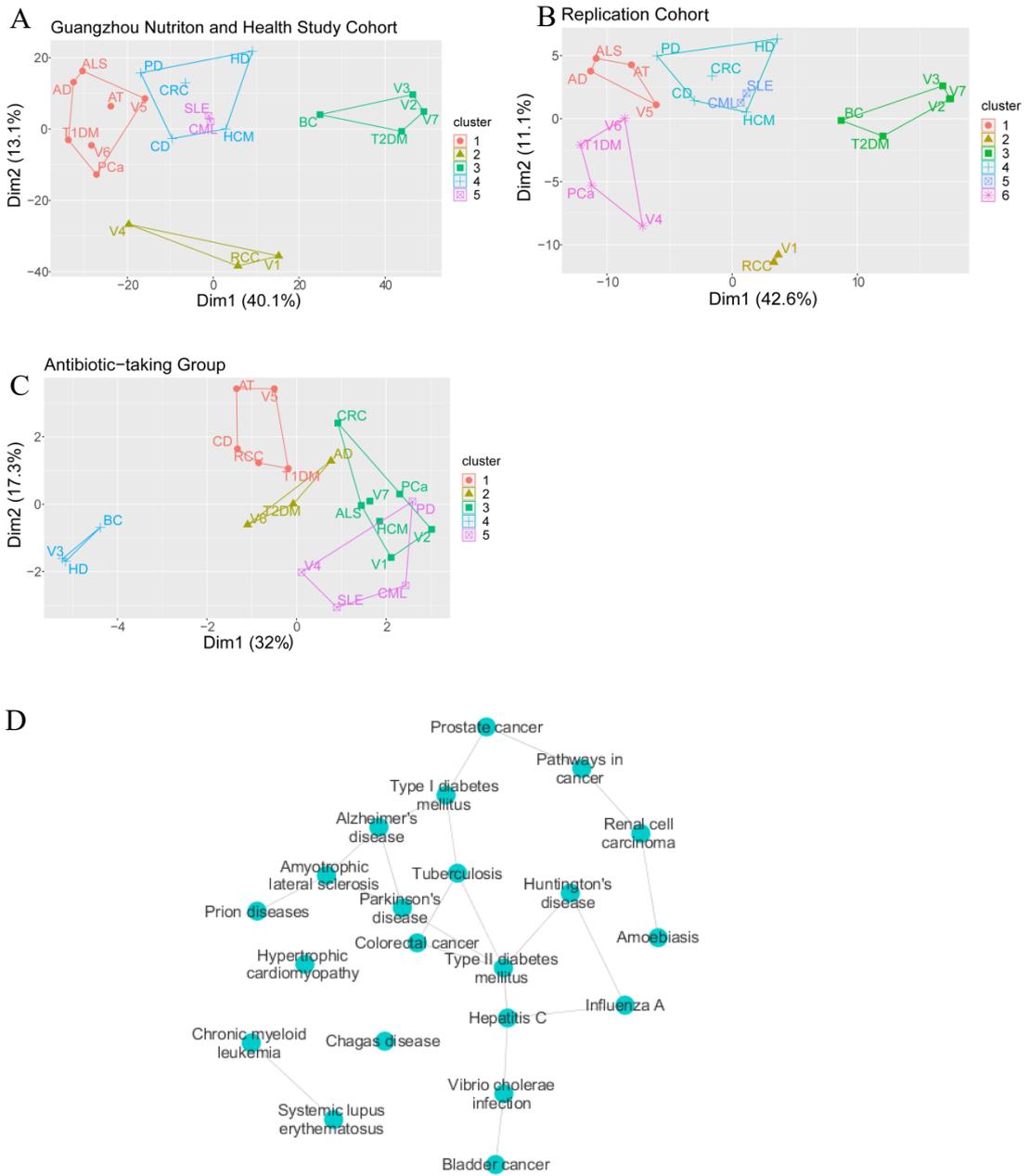
1050

1051

1052

1053

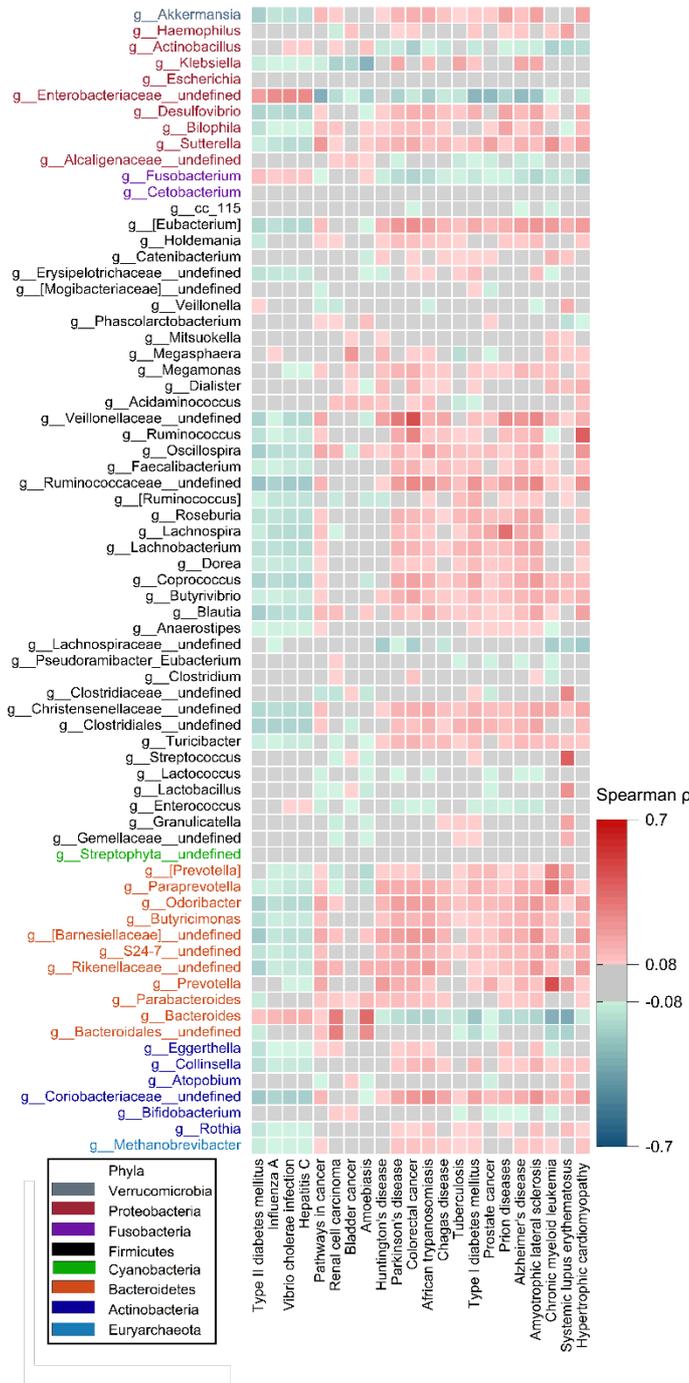
1054



1055

1056

1057 Figure 5 Correlation of the human complex diseases with the gut microbiome. The
 1058 heat map heatmap shows Spearman's correlation of predicted diseases and the gut
 1059 microbiome onat the genus level. The grey components show no significance of
 1060 significant correlation with Bonferroni correction ($p > 0.05 / ((5.6 * 22), p > 0.0004)$).



1061

1062

1063 **Supplementary Tables [Supplementary Tables.xls]**

1064

1065 **Supplementary Table S1 Transformation of traits in BioBank Japan and taxa in GNHS**1066 **Supplementary Table S2 Required effect size (beta) to reach 80% of power in GNHS cohort**1067 **Supplementary Table S3 Heritability of taxa, enterotype and alpha diversity**1068 **Supplementary Table S4 Significant genetic correlations of gut microbiome and metabolic**
1069 **traits**1070 **Supplementary Table S5 Significant associations of all taxa with SNPs identified in the**
1071 **discovery stage before adjustment($p < 5e-8$)**1072 **Supplementary Table S6 Replication of genetic variants associated with taxa**1073 **Supplementary Table S7 Replication of genetic variants associated with beta diversity**1074 **Supplementary Table S8 Replication of genetic variants associated with alpha diversity**1075 **Supplementary Table S9 Lead SNPs used to construct polygenic scores**1076 **Supplementary Table S10 MR analysis of gut microbiota on traits and diseases**1077 **Supplementary Table S11 MR analysis of diseases on gut microbiota features**1078 **Supplementary Table S12 Spearman' s correlation of certain taxa and complex diseases**1079 **Supplementary Table S13 Spearman' s correlation of gut microbiota on genus level and**
1080 **characteristics**

1081

1082 **Supplementary Figures [Supplementary Figures.pdf]**

1083

1084 **Supplementary Figure 1 Genome-wide analysis results of taxa.**1085 **Supplementary Figure 2 Spearman's correlation of the relative abundance of AF-associated**
1086 **taxa with the relative level of diseases predicted by PICRUSt.**

1087

1088

1089