
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^{1,3#}, 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 Translation Research of Zhejiang Province, School of Life
9 Sciences, Westlake University, Hangzhou, China.

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

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

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

17 ⁵ Institute for Advanced Research, Wenzhou Medical University, Wenzhou, Zhejiang
18 325027, China

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

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

22

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

24

25 *Correspondence to

26 Prof Ju-Sheng Zheng

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

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

29 And

30 Prof Yu-Ming Chen

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

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

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

34 And

35 Prof Jun Wang

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

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

38 Email: junwang@im.ac.cn

39

40 *Abstract*

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

45 **Results** The analysis was based on a Chinese population with 1,475 participants. We
46 estimated the SNP-based heritability, which suggested that *Desulfovibrionaceae* and
47 *Odoribacter* had significant heritability estimates (0.456 and 0.476, respectively). We
48 performed a microbiome genome-wide association study to identify host genetic
49 variants associated with gut microbiome. We then conducted bi-directional Mendelian
50 randomization analyses to examine the potential causal associations between gut
51 microbiome and human complex diseases. We found that *Saccharibacteria* (per 1-SD
52 higher in the log-transformed abundance) could potentially decrease the concentration
53 of serum creatinine (Beta: -0.011 [95%CI: -0.019, -0.003], $p=0.007$) and increase
54 estimated glomerular filtration rate (Beta: 0.012 [95%CI: 0.004, 0.020], $p=0.003$). On
55 the other hand, atrial fibrillation, chronic kidney disease and prostate cancer, as
56 predicted by the host genetics, had potential causal effect on the abundance of some
57 specific gut microbiota. For example, atrial fibrillation (per log odds) could increase
58 the abundance of *Burkholderiales* (Beta: 0.079[95%CI: 0.009, 0.150], $p=0.027$) and
59 *Alcaligenaceae* (Beta: 0.082[95%CI: 0.012, 0.152], $p=0.022$), and decrease the
60 abundance of *Lachnobacterium* (Beta:-0.078[95%CI: -0.148, -0.006], $p=0.034$),
61 *Bacteroides coprophilus* (Beta: -0.113[95%CI: -0.184, -0.041], $p=0.002$),

62 *Barnesiellaceae* (odds ratio: 0.818[95%CI: 0.686, 0.976], $p=0.026$), *Veillonellaceae*
63 undefined (odds ratio: 0.801[95%CI: 0.669, 0.960], $p=0.017$) as well as *Mitsuokella*
64 (odds ratio: 0.657[95%CI: 0.496, 0.870], $p=0.003$). Further disease-microbiome
65 feature analysis suggested that systemic lupus erythematosus and chronic myeloid
66 leukemia shared common gut microbiome features.

67 **Conclusions** These results suggest that different human complex diseases share
68 common and distinct gut microbiome features, which may help re-shape our
69 understanding about the disease etiology in humans.

70 **Keywords** Gut Microbiome; Host Genetics; Bi-directional Mendelian Randomization
71 Analyses; Disease-Microbiome Features

72

73 **Background**

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

83

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

92

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

95 gut microbiome features. Identifying common and distinct gut microbiome features
96 across different diseases may shed light on novel relationships among the complex
97 diseases and update our understanding about the disease etiology in humans. However,
98 the composition and structure of gut microbiome are influenced by a variety of factors
99 including environment, diet and regional variation [11-13], which poses a key
100 challenge for the description of representative microbiome features for a specific
101 disease. Although there were several studies comparing disease-related gut
102 microbiome [14-16], few of them examined and compared the microbiome features
103 across different human complex diseases.

104

105 In the present study, we performed a microbiome GWAS in a Chinese cohort study:
106 the Guangzhou Nutrition and Health Study (GNHS) [17], including 1475 participants.
107 Subsequently, we applied a bi-directional MR method to explore the genetically
108 predicted relationship between gut microbiome and human complex diseases. To
109 explore novel relationships among human complex diseases based on gut microbiome,
110 we investigated the shared and distinct gut microbiome features across diverse human
111 complex diseases.

112

113 **Results**

114 **Overview of the study**

115 Our study was based on the GNHS, with 4048 participants (40-75 years old) living in
116 urban Guangzhou city recruited during 2008 and 2013 [17]. In the GNHS, stool

117 samples were collected among 1937 participants during follow-up visits, among
118 which 1475 unrelated participants without taking anti-biotics were included in our
119 discovery microbiome GWAS. We then included additional 199 participants with both
120 genetic and gut microbiome data as a replication cohort, which belonged to the
121 control arm of a case-control study of hip fracture in Guangdong Province, China
122 [18](See also Figure 1).

123

124 **SNP-based heritability of gut microbiome**

125 The heritability of alpha diversity ranged from 0.035 to 0.103 (SE: from 0.174 to
126 0.193, Supplementary Table S3). Significant heritability estimates were observed for
127 several taxa (See also Figure 2, supplementary table S3) with a crude $p < 0.05$. To
128 further correct the multiple testing, we calculated the effective number of independent
129 taxa on each taxonomy level (phylum level: 2.3, class level: 2.9, order level: 2.9,
130 family level: 5.5, genus level: 5.6, species level: 3.2), as some taxa were highly
131 correlated with each other. The results suggested that *Desulfovibrionaceae* and
132 *Odoribacter* were heritable ($p < 0.05/n$, n is the effective number of independent taxa).
133 Notably, among the suggestively heritable taxa in our cohort, [*Paraprevotellaceae*],
134 *Veillonellaceae*, *Desulfovibrionaceae*, *Pasteurellaceae*, *Odoribacter*, *Paraprevotella*,
135 *Veillonella* and *Bifidobacterium* had nominally significant heritability estimates in
136 prior literatures [7, 19-21].

137

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

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

144

145 To examine the association of host genetic variants with alpha diversity, we performed
146 GWAS for four indices (Shannon diversity index, Chao1 diversity indices, observed
147 OTUs index and Phylogenetic diversity), but again no genome-wide significant signal
148 ($p < 5 \times 10^{-8}$) was found. To further investigate if there is host genetic basis underlying
149 alpha diversity, we constructed a polygenic score for each alpha diversity indicator in
150 the replication cohort, using the genetic variants which showed suggestive
151 significance ($p < 5 \times 10^{-5}$) in the discovery GWAS. The polygenic score was not
152 significantly associated with its corresponding alpha diversity index in our replication
153 cohort. Meanwhile, none of the associations with alpha diversity indices reported in
154 the literature could be replicated (Supplementary Table S8) [7].

155

156 The beta diversity GWAS was performed with MicrobiomeGWAS based on
157 Bray–Curtis dissimilarity [23]. We found that one locus at *SMARCA2* gene
158 (rs6475456) was associated with beta-diversity at a genome-wide significance level
159 ($p = 3.96 \times 10^{-9}$). However, we could not replicate the results in the replication cohort,
160 which might be due to the limited sample size of the replication cohort. In addition,

161 prior literature had reported 73 genetic variants that were associated with beta
162 diversity [8, 13, 24, 25], among which we found that 3 single nucleotide
163 polymorphisms (SNP, *UHRF2* gene-rs563779, *LHFPL3* gene-rs12705241,
164 *CTD-2135J3.4*-rs11986935) had nominal significant ($p < 0.05$) association with
165 beta-diversity in our cohort (Supplementary Table S7), although none of the
166 association survived Bonferroni correction. These studies used various methods for
167 the sequencing and calculation of beta diversity, which raised challenges to verify and
168 extrapolate results across populations.

169

170 We subsequently performed GWAS discovery for individual gut microbes in our own
171 GNHS discovery dataset. For the taxa (n=114) present in not fewer than ninety
172 percent of participants, we carried out analysis based on a log-normal model. For
173 other taxa (n=88) present in fewer than ninety percent, we transformed the
174 absence/presence of the taxon into binary variables and used a logistic model to
175 prevent zero inflation (Supplementary table S1). For all the gut microbiome taxa, the
176 significant threshold was defined as 5×10^{-8} in the discovery stage. We found that 6
177 taxa were associated with host genetic variants in the discovery cohort ($p < 5 \times 10^{-8}/n$, n
178 is the effective number of independent taxa on each taxonomy level, Supplementary
179 Table S5); however, these associations were not significant ($p > 0.05$) in the replication
180 cohort. We then took the genetic loci reported to be associated with individual taxa in
181 prior studies [7, 8, 13, 25] for replication in our GNHS dataset. Although none of the
182 associations of these genetic variants with taxa survived the Bonferroni correction

183 ($p < 1 \times 10^{-4}$), we found that *STPG2*-rs4699323 had a nominally significant association
184 ($p < 0.05$) with *Clostridiales* (Beta: -0.131 [-0.233 - -0.029], $p = 0.012$; Supplementary
185 Table S6). We then used a threshold of $p < 5 \times 10^{-5}$ at the GWAS discovery stage to
186 incorporate additional genetic variants which might explain a larger proportion of
187 heritability for taxa, based on which we constructed a polygenic score for each taxon
188 in the replication. We found that the polygenic scores were significantly associated
189 with 5 taxa including *Saccharibacteria* (also known as *TM7* phylum), *Clostridiaceae*,
190 *Comamonadaceae*, *Klebsiella* and *Desulfovibrio d168* in the replication set ($p < 0.05$,
191 Methods, see also Supplementary Figure 1, Supplementary Table S9).

192

193 Genetic correlation of gut microbiome and traits

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

204

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

207 Using genetic variants-composed polygenic scores as genetic instruments, we
208 performed MR analysis to assess the putative causal effect of microbiome
209 (*Saccharibacteria*, *Clostridiaceae*, *Comamonadaceae*, *Klebsiella* and *Desulfovibrio*
210 *d168*) on human complex diseases or traits. Inverse variance weighted (IVW) method
211 was used for the MR analysis, while other three methods (Weighted median,
212 MR-Egger and MR-PRESSO) were applied to confirm the robustness of results. The
213 horizontal pleiotropy was assessed using MR-PRESSO Global test and MR-Egger
214 Regression. For the analysis of gut microbiome on complex traits, we downloaded
215 public available GWAS summary statistics of complex traits (n=58) and diseases
216 (type 2 diabetes mellitus (T2DM), atrial fibrillation (AF), colorectal cancer (CRC)
217 and rheumatoid arthritis (RA)) reported by BioBank Japan [27-32]. The result
218 suggested that *Saccharibacteria* (per 1-SD higher in the log-transformed abundance)
219 could potentially decrease the concentration of serum creatinine (-0.011[95%CI:
220 -0.019, -0.003], $p=0.007$) and increase estimated glomerular filtration
221 rate(eGFR)(0.012[95%CI: 0.004, 0.020], $p=0.003$, Supplementary Table S10), which
222 might help improve renal function. We did not find evidence of pleiotropic effect:
223 genetic variants associated with *Saccharibacteria* were not associated with any of the
224 above traits (58 complex traits and 4 disease outcomes, $p<0.05/62$). These taxa were
225 not causally associated with other human complex diseases or traits in our MR
226 analyses, which might be due to the limited genetic instruments discovered in our

227 present study.

228

229 We subsequently performed a reserve MR analysis to assess the potential causal effect
230 of human complex diseases on gut microbiome features. For the reserve MR analyses,
231 the diseases of interests included T2DM, AF, coronary artery disease (CAD), chronic
232 kidney disease (CKD), Alzheimer's disease (AD), CRC and prostatic cancer (PCa),
233 and their instrumental variables for the MR analysis were based on previous
234 large-scale GWAS in East Asians [27, 33-38]. The results suggested that AF and CKD
235 were causally associated with gut microbiome (See also Figure 3A, 3B,
236 Supplementary Table S11). Specifically, genetically predicted higher risk of AF (per
237 log odds) was associated with lower abundance of *Lachnobacterium* (Beta:-
238 0.078[95%CI: -0.148, -0.006], $p=0.034$), *Bacteroides coprophilus* (Beta:
239 -0.113[95%CI: -0.184, -0.041], $p=0.002$), *Barnesiellaceae* (odds ratio: 0.818[95%CI:
240 0.686, 0.976], $p=0.026$), *Veillonellaceae* undefined (odds ratio: 0.801[95%CI: 0.669,
241 0.960], $p=0.017$) as well as *Mitsuokella* (odds ratio: 0.657[95%CI: 0.496, 0.870],
242 $p=0.003$), and higher abundance of *Burkholderiales* (Beta: 0.079[95%CI: 0.009,
243 0.150], $p=0.027$) and *Alcaligenaceae* (Beta: 0.082[95%CI: 0.012, 0.152], $p=0.022$).
244 Additionally, genetically predicted higher risk of CKD could increase *Anaerostipes*
245 (Beta: 0.291[95%CI: 0.057, 0.524], $p=0.015$) abundance and higher risk of PCa could
246 decrease [*Prevotella*] (odds ratio: -0.758[95%CI: -1.354, -0.162], $p=0.013$).

247

248 **Microbiome features of human complex diseases**

249 To further investigate the potential complex diseases that may be correlated with the
250 taxa affected by AF, we applied Phylogenetic Investigation of Communities by
251 Reconstruction of Unobserved States (PICRUSt) to predict the disease pathway
252 abundance [39]. We used Spearman's rank-order correlation to test whether the
253 relative abundances of predicted diseases based on PICRUSt were associated with the
254 aforementioned AF-associated taxa (See also Supplementary Figure 2, Supplementary
255 Table S12). The heatmap indicated that cancers and neurodegenerative diseases
256 including Parkinson's disease (PD), AD, amyotrophic lateral sclerosis (ALS) as well
257 as AF were correlated with similar gut microbiome. Although the association among
258 these diseases are highly supported by previous studies [40-42], no study has
259 compared common gut microbiome features across these different diseases.

260
261 To compare gut microbiome features across human diseases, we used the predicted
262 disease abundance based on PICRUSt and performed k-medoids clustering.

263 According to optimum average silhouette width [43], we chose optimal number of
264 clusters for further analysis. The plot showed that neurological diseases including
265 ALS and AD belonged to the same cluster, while PD and CRC had much similarity in
266 gut microbiome. The results also suggested that systemic lupus erythematosus (SLE)
267 and chronic myeloid leukemia (CML) shared similar gut microbiome features (See
268 also Figure 4A, 4B). Moreover, we could replicate these clusters in our replication
269 cohort, which suggested that the clustering results were robust (See also Figure 4C).

270

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

285

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

293 with 217 T2DM cases). The results showed that MRS was positively associated with
294 the risk of T2DM (odds ratio: 1.176[95%CI: 1.114, 1.244], $p=8.75\times 10^{-9}$).

295

296 Discussion

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

306

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

315

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

327

328 The MR analysis showed that *Saccharibacteria* might decrease the concentration of
329 serum creatinine and increase eGFR. Little is known about the *Saccharibacteria* as
330 one of the uncultivated phyla, and previous studies have shown that it might be
331 essential for the immune response, oral inflammation and inflammatory bowel disease
332 [46-48]. Our results also provided genetic instrument of *Saccharibacteria* for further
333 causal analysis with other complex diseases. The reserve MR analysis provided
334 evidence that AF, CKD and PCa could causally influence gut microbiome. The rare
335 and low-frequency variants may have an important impact on common diseases [49],
336 thus it will be of interest to clarify the effects of low-frequency variants on gut

337 microbiome in cohorts with large sample sizes in future.

338

339 Our results indicate that gut microbiome help reveal novel and interesting

340 relationships among human complex diseases, and different diseases may have

341 common and distinct gut microbiome features. A prior study including participants

342 from different countries has identified three microbiome clusters [22]. Notably, this

343 study focused on classifying the individuals into distinct enterotypes regardless of the

344 individuals' health status, while in the present study we described representative

345 microbiome features for diseases of interest. We provided an approach to interpret the

346 data from mechanistic studies based on microbiome. The microbiome features

347 revealed a close association of AF with neurodegenerative diseases as well as cancers,

348 which was supported by prior studies showing that AF had correlation with AD and

349 PD [40, 41], and AF patients had relatively higher risks of several cancers including

350 lung cancer and CRC [42, 50]. We also observed that microbiome features of SLE and

351 CML were highly similar. Interestingly, a tyrosine kinase inhibitor of platelet-derived

352 growth factor receptor, imatinib, was widely used to treat CML and significantly

353 ameliorated survival in murine models of SLE [51]. In addition, close association

354 between CRC and PD has also been reported in several observational cohorts [52, 53].

355 Collectively, these findings strongly supported our hypothesis that human complex

356 diseases sharing similar microbiome features might be mechanically correlated.

357 Furthermore, from the perspectives of risk genes of AF and neurodegenerative

358 diseases, previous GWAS for AF had identified two loci at *PITX2* gene-rs6843082

359 and *C9orf3* gene-rs7026071, which were also associated with the risk of ALS
360 ($p=0.0138$ and $p=0.049$, respectively) [54-56].

361

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

372

373 **Conclusions**

374 In summary, we reveal some causal relationships between abundance of gut
375 microbiome and human complex diseases or traits. The disease and gut microbiome
376 feature analysis reveals novel relationships among human complex diseases, which
377 may help re-shape our understanding about the disease etiology, as well as provide
378 some clues for extending clinical indications of existing drugs for different diseases.

379

380 **Method**

381 **Study participants and sample collection**

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

388

389 We included 199 participants with both genetic and gut microbiome data as a
390 replication cohort, which belonged to the control arm of a case-control study of hip
391 fracture with the participants (52-83 years old) recruited between June 2009 and
392 August 2015 in Guangdong Province, China [18].

393

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

399

400 **Genotyping data**

401 For both discovery and replication cohorts, DNA was extracted from leukocyte using
402 the TIANamp® Blood DNA Kit as per the manufacturer's instruction. DNA

403 concentrations were determined using the Qubit quantification system (Thermo
404 Scientific, Wilmington, DE, US). Extracted DNA was stored at -80°C . Genotyping
405 was carried out with Illumina ASA-750K arrays. Quality control and relatedness
406 filters were performed by PLINK1.9 [58]. Individuals with high or low proportion of
407 heterozygous genotypes (outliers defined as 3 standard deviation) were excluded [59].
408 Individuals who had different ancestries (the first two principal components ± 5
409 standard deviation from the mean) or related individuals ($\text{IBD} > 0.185$) were excluded
410 [59]. Variants were mapping to the 1000 Genomes Phase3 v5 by SHAPEIT [60, 61]
411 and then we conducted the genome-wide genotype imputation with 1000 Genomes
412 Phase3 v5 reference panel by Minimac3 [62, 63]. Genetic variants with imputation
413 accuracy $\text{RSQR} > 0.3$ and $\text{MAF} > 0.05$ were included in our analysis. We used
414 Pan-Asian reference panel consist of 502 participants and SNP2HLA v1.0.3 to impute
415 HLA region [64-66].

416

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

418 Microbial DNA was extracted from fecal samples using the QIAamp® DNA Stool
419 Mini Kit per the manufacturer's instruction. DNA concentrations were determined
420 using the Qubit quantification system. The V3-V4 region of the 16S rRNA gene was
421 amplified from genomic DNA using primers 341F (CCTACGGGNGGCWGCAG)
422 and 805R (GACTACHVGGGTATCTAATCC). The pooled amplicons were
423 sequenced using MiSeq Reagent Kits v2 on the Illumina MiSeq System with 2 ×
424 250bp pair-end sequencing.

425

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

437

438 **Statistical analysis**

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

440 We used the GREML method in GCTA to estimate the proportion of variance
441 explained by all SNPs [71]. The taxa were divided into two groups based on whether
442 the taxa were present in the ninety percent of participants or not. Our model was
443 adjusted for age and sex. The power of GREML analysis was calculated with GCTA
444 power calculator [72].

445

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

447 For each of the GNHS participants and the replication cohort, we clustered
448 participants based on genus-level relative abundance, estimating the JSD distance and
449 PAM clustering algorithm, and then defined two enterotypes according to
450 Calinski-Harabasz Index [22, 73]. We calculated the genetic principal components of
451 ancestry from genome-wide genetic variants to estimate the population structure.

452 PLINK 1.9 was used to perform a logistic regression model for enterotypes and taxa
453 present in fewer than ninety percent, adjusted for age, sex, sequencing batch and the
454 first five genetic principal components of ancestry.

455
456 For beta diversity, the analysis for the genome-wide host genetic variants with beta
457 diversity was performed using MicrobiomeGWAS [23], adjusted for covariates
458 including the first five genetic principal components of ancestry, age and sex.

459
460 Alpha diversity was calculated after randomly sampling 10000 reads per sample. For
461 the taxa present in not fewer than ninety percent of participants and alpha diversity,
462 we used Z-score normalization to transform the distribution and carried out analysis
463 based on a log-normal model. A mixed linear model based association (MLMA) test
464 in GCTA was used to assess the association, fitting the first five genetic principal
465 components of ancestry, age, sex and sequencing batch as fixed effects and the effects
466 of all the SNPs as random effects [74-76]. For other taxa present in fewer than ninety
467 percent, we transformed the absence/presence of the taxon into binary variables and
468 used PLINK1.9 to perform a logistic model, adjusted for the first five genetic

469 principal components of ancestry, age, sex and sequencing batch. For all the gut
470 microbiome features, the significant threshold was defined as $5 \times 10^{-8} / n$ (n is the
471 effective number of independent taxa on each taxonomy level) in the discovery stage.

472 The QUANTO software was used for power calculations

473 (<http://biostats.usc.edu/Quanto.html>). We estimated genomic inflation factors with
474 LDSC v1.0.1 at local server [77].

475

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

477 We used GCTA to perform a bivariate GREML analysis to estimate the genetic
478 correlation between gut microbiome and traits in the GNHS [74, 78]. The gut
479 microbiome was divided into two groups according to the previous description. We
480 used continuous variables to taxa present in not fewer than ninety percent of
481 participants. For taxa present in few than ninety percent of participants, we used
482 binary variables according to the absence/presence of taxa. This analysis included
483 traits such as BMI, FBS, HbA1c, SBP, DBP, HDL-C, LDL-C, TC and TG. The power
484 of bivariate GREML analysis was calculated with GCTA power calculator [72].

485

486

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

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

491 and taxa present in not fewer than ninety percent participants, we constructed
 492 weighted polygenic scores and performed the analysis on a general linear model with
 493 a negative binomial distribution to test for association between the polygenic scores
 494 and taxa, adjusted for the first five genetic principal components of ancestry, age, sex
 495 and sequencing batch. We used weighted polygenic scores and logistic regression to
 496 the absence/presence taxa, adjusted for the same covariates as in the above analysis.
 497 Taxa with significance ($p < 0.05$) in the replication cohort were included for further
 498 analysis.

499

500 **The effective number of independent taxa**

501 As some taxa were correlated with each other, we used an eigendecomposition
 502 analysis to calculate the effective number of independent taxa on each taxonomy level
 503 [79, 80]. Matrix M is an $m \times n$ matrix, where m is the number of participants and n is
 504 the number of total taxa on the corresponding taxa level. Matrix A is the
 505 variance-covariance matrix of matrix M . P is the matrix of eigenvectors.
 506 $\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ is the diagonal matrix comprised of the ordered eigenvalues, which
 507 can be calculated as:

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

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

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

509

510 **Bi-directional MR analysis**

511 In the analysis of potential causal effect of gut microbiome features on diseases, we
512 used independent genetic variants (selected as part of the polygenic score analysis) as
513 the instrumental variables. For each trait, we excluded instrumental variables that
514 showed the significant association with the trait ($p < 0.05/n$, n is number of
515 independent genetic variants). In the analysis of potential causal effect of diseases on
516 gut microbiome features, we selected genetic variants that were replicated in East
517 Asian populations as instrumental variables. As all instrumental variables were from
518 East Asian populations, we chose independent genetic variants ($r^2 < 0.1$) based on
519 GNHS cohort. We identified the best proxy ($r^2 > 0.9$) based on GNHS cohort or
520 discarded the variant if no proxy was available. We used inverse variance weighted
521 (IVW) method to estimate effect size. To confirm the robustness of results, we
522 performed other three MR methods including weighted median, MR-Egger and
523 MR-PRESSO [81-83]. To assess the presence of horizontal pleiotropy, we performed
524 MR-PRESSO Global test and MR-Egger Regression. Effect sizes of gut microbiome
525 on traits were dependent on units of traits (Supplementary table S1). Results of human
526 complex diseases on the absence/presence gut microbiome were presented as risk of
527 presence (vs absence) of the microbiome per log odds difference of the disease.
528 Results of diseases on other gut microbiome were presented as changes in abundance
529 of taxa (1-SD of log transformed) per log odds difference of the respective disease.
530
531 The statistical significance of gut microbiome on traits and diseases was defined as
532 $p < 0.0008$ ($0.05/62$). In addition, the statistical significance of diseases on gut

533 microbiome features was defined as $p < 0.05/n$ (n is the effective number of
534 independent taxa on the corresponding taxonomy level). Results that could not pass
535 Bonferroni adjustment but $p < 0.05$ in all four MR methods were considered as
536 potential causal relationships. We performed MR analyses on R v3.5.3.

537

538 **Pathway analysis**

539 We used OTUs by QIIME and annotated the variation of functional genes with
540 Phylogenetic Investigation of Communities by Reconstruction of Unobserved States
541 (PICRUSt) [39]. The pathways and diseases were annotated using KEGG [84-86]. We
542 used Spearman's rank-order correlation to investigate the association of predicted
543 pathway or diseases abundance with the AF-associated taxa and genus-level taxa. In
544 the heatmap, diseases were clustered with 'hcluster' function on R. To test whether
545 non-normalized pathway or disease abundance was associated with each other, we
546 used SPIEC-EASI to test the interaction relationship, and then used Cytoscape v3.7.2
547 to visualize the interaction network [87, 88].

548

549 **Construction of the microbiome risk score**

550 The microbiome risk score was constructed to validate the accuracy of the association
551 between the predicted disease-related gut microbiome features and the corresponding
552 disease. As we have a large sample size for T2DM cases ($n=217$ cases) in our cohort,
553 we constructed a microbiome risk score of T2DM as an exemplar. We used a
554 Spearman's rank-order correlation to select taxa with the absolute value of correlation

555 coefficient higher than 0.2. Score for each taxon abundance <5% quantile in our study
 556 was defined as 0. For those above 5%, score for each taxon showing inverse
 557 association with T2DM was defined as -1; score for each taxon showing positive
 558 association with T2DM was defined as 1. We then summed up values from all taxa.
 559 We selected logistic regression model to estimate association of the MRS with T2DM
 560 risk, and linear model to estimate the association of the MRS with the continuous
 561 variables, adjusted for age, sex, dietary energy intake, alcohol intake and BMI at the
 562 time of sample collection.

563

564 **Clustering diseases**

565 The clustering analysis was carried out with ‘cluster’ and ‘factoextra’ for plot on R.
 566 We performed PAM algorithm based on predicted abundance of diseases or average
 567 relative abundance after Z-score normalization [89]. PAM algorithm searches k
 568 medoids among the observations and then found nearest medoids to minimize the
 569 dissimilarity among clusters [90]. Given a set of objects $x = (x_1, x_2, \dots, x_n)$, the
 570 dissimilarity between objects x_i and x_j is denoted by $d(i,j)$. The assignment of
 571 object i to the representative object j is denoted by z_{ij} . z_{ij} is a binary variable and is
 572 1 if object i belongs to the cluster of the representative object j. The function to
 573 minimize the model is given by:

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

574 To identify the optimal cluster number, we used ‘pamk’ function in R to determine the
 575 optimum average silhouette width. For each object i, we defined N_i as the average

576 dissimilarity between object i and all other objects within its cluster. For the
 577 remaining clusters, $b(i,w)$ represents the average dissimilarity between i and all
 578 objects in cluster w . The minimum dissimilarity M_i can be calculated by:

$$579 \quad M_i = \min \forall w (b(i,w)).$$

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

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

581 Then we calculated the average of silhouette width for each object. The cluster
 582 number is determined by the number of which the average silhouette width is
 583 maximum. We estimated the Jaccard similarity coefficient to quantify the cluster
 584 difference between groups. The Jaccard similarity coefficient is positively associated
 585 with the similarity of clusters. Given object i and j , as well as group A and B, there are
 586 four situations as follows:

- 587 (1) S1: in both group A and B, object i and j belong to the same cluster;
 588 (2) S2: in group A, object i and j belong to the same cluster; in group B, they belong
 589 to different clusters;
 590 (3) S3: in group A, object i and j belong to different clusters; in group A, they belong
 591 to the same cluster;
 592 (4) S4: in both group A and B, object i and j belong to different clusters.

593 a , b , c and d represents the number of S1, S2, S3 and S4, respectively. The Jaccard
 594 similarity coefficient can be calculated by the following formula:

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

595

596 **Availability of data and materials**

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

603

604 Acknowledgments

605 We thank the Westlake University Supercomputer Center for providing computing
606 and data analysis service for the present project.

607 Ethics approval and consent to participate

608 This study was approved by Westlake University Ethics Committee
609 (20190114ZJS003).

610 Consent for publication

611 Not applicable.

612 Competing interests

613 The authors declare no conflict of interest.

614 Authors' contributions

615 JSZ, YMC and JW initiated and led the study. JY assisted with the data analyses. FZX,
616 YQF and JSZ analyzed the data and wrote the manuscript. TYS and CWL collected
617 the data. ZLJ, ZLM, MLS and WLG analyzed the data. All authors read and approved
618 the final manuscript.

619 Funding

620 This study was funded by National Natural Science Foundation of China (81903316,
621 81773416), Zhejiang Ten-thousand Talents Program (101396522001) and the 5010
622 Program for Clinical Researches (2007032) of the Sun Yat-sen University
623 (Guangzhou, China).

References

- 624
625 1. Awany D, Allali I, Dalvie S, Hemmings S, Mwaikono KS, Thomford NE, et al.
626 Host and Microbiome Genome-Wide Association Studies: Current State and
627 Challenges. *Frontiers in genetics*. 2019;9:637-; doi:10.3389/fgene.2018.00637.
- 628 2. Bull MJ, Plummer NT. Part 1: The Human Gut Microbiome in Health and
629 Disease. *Integrative medicine (Encinitas, Calif)*. 2014;13(6):17-22.
- 630 3. Lynch JB, Hsiao EY. Microbiomes as sources of emergent host phenotypes.
631 *Science*. 2019;365(6460):1405-9; doi:10.1126/science.aay0240.
- 632 4. Allegretti JR, Mullish BH, Kelly C, Fischer M. The evolution of the use of
633 faecal microbiota transplantation and emerging therapeutic indications. *The*
634 *Lancet*. 2019;394(10196):420-31; doi:10.1016/S0140-6736(19)31266-8.
- 635 5. Wong SH, Yu J. Gut microbiota in colorectal cancer: mechanisms of action
636 and clinical applications. *Nature Reviews Gastroenterology & Hepatology*.
637 2019; doi:10.1038/s41575-019-0209-8.
- 638 6. Davies NM, Holmes MV, Davey Smith G. Reading Mendelian randomisation
639 studies: a guide, glossary, and checklist for clinicians. *BMJ*. 2018;362:k601;
640 doi:10.1136/bmj.k601.
- 641 7. Turpin W, Espin-Garcia O, Xu W, Silverberg MS, Kevans D, Smith MI, et al.
642 Association of host genome with intestinal microbial composition in a large
643 healthy cohort. *Nat Genet*. 2016;48(11):1413-7; doi:10.1038/ng.3693.
- 644 8. Wang J, Thingholm LB, Skieceviciene J, Rausch P, Kummén M, Hov JR, et al.
645 Genome-wide association analysis identifies variation in vitamin D receptor
646 and other host factors influencing the gut microbiota. *Nat Genet*.
647 2016;48(11):1396-406; doi:10.1038/ng.3695.
- 648 9. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, et al.
649 Human genetics shape the gut microbiome. *Cell*. 2014;159(4):789-99;
650 doi:10.1016/j.cell.2014.09.053.
- 651 10. Goodrich JK, Davenport ER, Beaumont M, Jackson MA, Knight R, Ober C, et
652 al. Genetic Determinants of the Gut Microbiome in UK Twins. *Cell Host*
653 *Microbe*. 2016;19(5):731-43; doi:10.1016/j.chom.2016.04.017.
- 654 11. Ganesan K, Chung SK, Vanamala J, Xu B. Causal Relationship between
655 Diet-Induced Gut Microbiota Changes and Diabetes: A Novel Strategy to
656 Transplant *Faecalibacterium prausnitzii* in Preventing Diabetes. *Int J Mol Sci*.
657 2018;19(12); doi:10.3390/ijms19123720.
- 658 12. He Y, Wu W, Zheng HM, Li P, McDonald D, Sheng HF, et al. Regional
659 variation limits applications of healthy gut microbiome reference ranges and
660 disease models. *Nat Med*. 2018;24(10):1532-5;
661 doi:10.1038/s41591-018-0164-x.
- 662 13. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, et al.
663 Environment dominates over host genetics in shaping human gut microbiota.
664 *Nature*. 2018;555(7695):210-5; doi:10.1038/nature25973.
- 665 14. Duvallé C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut
666 microbiome studies identifies disease-specific and shared responses. *Nature*
667 *communications*. 2017;8(1):1784-; doi:10.1038/s41467-017-01973-8.

-
- 668 15. Cheng S, Han B, Ding M, Wen Y, Ma M, Zhang L, et al. Identifying
669 psychiatric disorder-associated gut microbiota using microbiota-related gene
670 set enrichment analysis. *Briefings in Bioinformatics*. 2019;
671 doi:10.1093/bib/bbz034.
- 672 16. Jackson MA, Verdi S, Maxan M-E, Shin CM, Zierer J, Bowyer RCE, et al.
673 Gut microbiota associations with common diseases and prescription
674 medications in a population-based cohort. *Nature Communications*.
675 2018;9(1):2655; doi:10.1038/s41467-018-05184-7.
- 676 17. Cao Y, Wang C, Guan K, Xu Y, Su Y-X, Chen YM. Association of magnesium
677 in serum and urine with carotid intima-media thickness and serum lipids in
678 middle-aged and elderly Chinese: a community-based cross-sectional study.
679 *European journal of nutrition*. 2015;55; doi:10.1007/s00394-015-0839-8.
- 680 18. Sun L-L, Li B-L, Xie H-L, Fan F, Yu W-Z, Wu B-H, et al. Associations
681 between the dietary intake of antioxidant nutrients and the risk of hip fracture
682 in elderly Chinese: A case-control study. *The British journal of nutrition*.
683 2014;112:1-9; doi:10.1017/S0007114514002773.
- 684 19. Lim MY, You HJ, Yoon HS, Kwon B, Lee JY, Lee S, et al. The effect of
685 heritability and host genetics on the gut microbiota and metabolic syndrome.
686 *Gut*. 2017;66(6):1031-8; doi:10.1136/gutjnl-2015-311326.
- 687 20. Davenport ER. Elucidating the role of the host genome in shaping microbiome
688 composition. *Gut microbes*. 2016;7(2):178-84;
689 doi:10.1080/19490976.2016.1155022.
- 690 21. Davenport ER, Cusanovich DA, Michelini K, Barreiro LB, Ober C, Gilad Y.
691 Genome-Wide Association Studies of the Human Gut Microbiota. *PLoS one*.
692 2015;10(11):e0140301-e; doi:10.1371/journal.pone.0140301.
- 693 22. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, et al.
694 Enterotypes of the human gut microbiome. *Nature*. 2011;473(7346):174-80;
695 doi:10.1038/nature09944.
- 696 23. Hua X, Song L, Yu G, Goedert JJ, Abnet CC, Landi MT, et al.
697 MicrobiomeGWAS: a tool for identifying host genetic variants associated with
698 microbiome composition. *bioRxiv*. 2015:031187; doi:10.1101/031187.
- 699 24. Ruhlemann MC, Degenhardt F, Thingholm LB, Wang J, Skieceviciene J,
700 Rausch P, et al. Application of the distance-based F test in an mGWAS
701 investigating beta diversity of intestinal microbiota identifies variants in
702 SLC9A8 (NHE8) and 3 other loci. *Gut Microbes*. 2018;9(1):68-75;
703 doi:10.1080/19490976.2017.1356979.
- 704 25. Bonder MJ, Kurilshikov A, Tigchelaar EF, Mujagic Z, Imhann F, Vila AV, et al.
705 The effect of host genetics on the gut microbiome. *Nat Genet*.
706 2016;48(11):1407-12; doi:10.1038/ng.3663.
- 707 26. Tang WHW, Kitai T, Hazen SL. Gut Microbiota in Cardiovascular Health and
708 Disease. *Circulation research*. 2017;120(7):1183-96;
709 doi:10.1161/CIRCRESAHA.117.309715.
- 710 27. Low SK, Takahashi A, Ebana Y, Ozaki K, Christophersen IE, Ellinor PT, et al.
711 Identification of six new genetic loci associated with atrial fibrillation in the

- 712 Japanese population. *Nat Genet.* 2017;49(6):953-8; doi:10.1038/ng.3842.
- 713 28. Suzuki K, Akiyama M, Ishigaki K, Kanai M, Hosoe J, Shojima N, et al.
- 714 Identification of 28 new susceptibility loci for type 2 diabetes in the Japanese
- 715 population. *Nat Genet.* 2019;51(3):379-86; doi:10.1038/s41588-018-0332-4.
- 716 29. Akiyama M, Okada Y, Kanai M, Takahashi A, Momozawa Y, Ikeda M, et al.
- 717 Genome-wide association study identifies 112 new loci for body mass index in
- 718 the Japanese population. *Nat Genet.* 2017;49(10):1458-67;
- 719 doi:10.1038/ng.3951.
- 720 30. Kanai M, Akiyama M, Takahashi A, Matoba N, Momozawa Y, Ikeda M, et al.
- 721 Genetic analysis of quantitative traits in the Japanese population links cell
- 722 types to complex human diseases. *Nat Genet.* 2018;50(3):390-400;
- 723 doi:10.1038/s41588-018-0047-6.
- 724 31. Matoba N, Akiyama M, Ishigaki K, Kanai M, Takahashi A, Momozawa Y, et
- 725 al. GWAS of smoking behaviour in 165,436 Japanese people reveals seven
- 726 new loci and shared genetic architecture. *Nat Hum Behav.* 2019;3(5):471-7;
- 727 doi:10.1038/s41562-019-0557-y.
- 728 32. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, et al. Genetics of
- 729 rheumatoid arthritis contributes to biology and drug discovery. *Nature.*
- 730 2014;506(7488):376-81; doi:10.1038/nature12873.
- 731 33. Lu XF, Wang LY, Chen SF, He L, Yang XL, Shi YY, et al. Genome-wide
- 732 association study in Han Chinese identifies four new susceptibility loci for
- 733 coronary artery disease. *Nature Genetics.* 2012;44(8):890-+;
- 734 doi:10.1038/ng.2337.
- 735 34. Marzec J, Mao X, Li M, Wang M, Feng N, Gou X, et al. A genetic study and
- 736 meta-analysis of the genetic predisposition of prostate cancer in a Chinese
- 737 population. *Oncotarget.* 2016;7(16):21393-403; doi:10.18632/oncotarget.7250.
- 738 35. Okada Y, Sim X, Go MJ, Wu JY, Gu D, Takeuchi F, et al. Meta-analysis
- 739 identifies multiple loci associated with kidney function-related traits in east
- 740 Asian populations. *Nat Genet.* 2012;44(8):904-9; doi:10.1038/ng.2352.
- 741 36. Zeng C, Matsuda K, Jia WH, Chang J, Kweon SS, Xiang YB, et al.
- 742 Identification of Susceptibility Loci and Genes for Colorectal Cancer Risk.
- 743 *Gastroenterology.* 2016;150(7):1633-45; doi:10.1053/j.gastro.2016.02.076.
- 744 37. Zhou X, Chen Y, Mok KY, Zhao Q, Chen K, Chen Y, et al. Identification of
- 745 genetic risk factors in the Chinese population implicates a role of immune
- 746 system in Alzheimer's disease pathogenesis. *Proceedings of the National*
- 747 *Academy of Sciences.* 2018;115(8):1697; doi:10.1073/pnas.1715554115.
- 748 38. Gan W, Walters RG, Holmes MV, Bragg F, Millwood IY, Banasik K, Chen Y,
- 749 Du H, Iona A, Mahajan A, et al: Evaluation of type 2 diabetes genetic risk
- 750 variants in Chinese adults: findings from 93,000 individuals from the China
- 751 Kadoorie Biobank. *Diabetologia.* 2016;59(7):1446-1457.
- 752 doi:10.1007/s00125-016-3920-9
- 753 39. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA,
- 754 et al. Predictive functional profiling of microbial communities using 16S
- 755 rRNA marker gene sequences. *Nature Biotechnology.* 2013;31:814;

- 756 doi:10.1038/nbt.2676.
- 757 40. Canga Y, Emre A, Yuksel GA, Karatas MB, Yelgec NS, Gurkan U, et al.
758 Assessment of Atrial Conduction Times in Patients with Newly Diagnosed
759 Parkinson's Disease. *Parkinsons Dis.* 2018;2018:2916905;
760 doi:10.1155/2018/2916905.
- 761 41. Ihara M, Washida K. Linking Atrial Fibrillation with Alzheimer's Disease:
762 Epidemiological, Pathological, and Mechanistic Evidence. *J Alzheimers Dis.*
763 2018;62(1):61-72; doi:10.3233/JAD-170970.
- 764 42. Conen D, Wong JA, Sandhu RK, Cook NR, Lee I-M, Buring JE, et al. Risk of
765 Malignant Cancer Among Women With New-Onset Atrial Fibrillation. *JAMA*
766 *Cardiology.* 2016;1(4):389-96; doi:10.1001/jamacardio.2016.0280.
- 767 43. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation
768 of cluster analysis. *Journal of Computational and Applied Mathematics.*
769 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).
- 770 44. Blekhman R, Goodrich JK, Huang K, Sun Q, Bukowski R, Bell JT, et al. Host
771 genetic variation impacts microbiome composition across human body sites.
772 *Genome Biol.* 2015;16:191; doi:10.1186/s13059-015-0759-1.
- 773 45. Goodrich JK, Davenport ER, Clark AG, Ley RE. The Relationship Between
774 the Human Genome and Microbiome Comes into View. *Annu Rev Genet.*
775 2017;51:413-33; doi:10.1146/annurev-genet-110711-155532.
- 776 46. Kuehbacher T, Rehman A, Lepage P, Hellmig S, Fölsch UR, Schreiber S, et al.
777 Intestinal TM7 bacterial phylogenies in active inflammatory bowel disease.
778 *Journal of Medical Microbiology.* 2008;57(12):1569-76;
779 doi:<https://doi.org/10.1099/jmm.0.47719-0>.
- 780 47. He X, McLean JS, Edlund A, Yooseph S, Hall AP, Liu S-Y, et al. Cultivation
781 of a human-associated TM7 phylotype reveals a reduced genome and epibiotic
782 parasitic lifestyle. *Proceedings of the National Academy of Sciences of the*
783 *United States of America.* 2015;112(1):244-9; doi:10.1073/pnas.1419038112.
- 784 48. Bor B, Bedree JK, Shi W, McLean JS, He X. Saccharibacteria (TM7) in the
785 Human Oral Microbiome. *Journal of Dental Research.* 2019;98(5):500-9;
786 doi:10.1177/0022034519831671.
- 787 49. Cirulli ET, Goldstein DB. Uncovering the roles of rare variants in common
788 disease through whole-genome sequencing. *Nat Rev Genet.*
789 2010;11(6):415-25; doi:10.1038/nrg2779.
- 790 50. Vinter N, Christesen Amanda MS, Fenger-Grøn M, Tjønneland A, Frost L.
791 Atrial Fibrillation and Risk of Cancer: A Danish Population-Based Cohort
792 Study. *Journal of the American Heart Association.* 2018;
793 doi:10.1161/JAHA.118.009543.
- 794 51. Zoja C, Corna D, Rottoli D, Zanchi C, Abbate M, Remuzzi G. Imatinib
795 ameliorates renal disease and survival in murine lupus autoimmune disease.
796 *Kidney International.* 2006;70(1):97-103;
797 doi:<https://doi.org/10.1038/sj.ki.5001528>.
- 798 52. Boursi B, Mamtani R, Haynes K, Yang Y-X. Parkinson's disease and colorectal
799 cancer risk-A nested case control study. *Cancer Epidemiol.* 2016;43:9-14;

-
- 800 doi:10.1016/j.canep.2016.05.007.
- 801 53. Xie X, Luo X, Xie M. Association between Parkinson's disease and risk of
802 colorectal cancer. *Parkinsonism & Related Disorders*. 2017;35:42-7;
803 doi:10.1016/j.parkreldis.2016.11.011.
- 804 54. van Rheenen W, Shatunov A, Dekker AM, McLaughlin RL, Diekstra FP, Pulit
805 SL, et al. Genome-wide association analyses identify new risk variants and the
806 genetic architecture of amyotrophic lateral sclerosis. *Nat Genet*.
807 2016;48(9):1043-8; doi:10.1038/ng.3622.
- 808 55. Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C,
809 et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci
810 for Alzheimer's disease. *Nat Genet*. 2013;45(12):1452-8; doi:10.1038/ng.2802.
- 811 56. Pankratz N, Beecham GW, DeStefano AL, Dawson TM, Doheny KF, Factor
812 SA, et al. Meta-analysis of Parkinson's disease: identification of a novel locus,
813 RIT2. *Ann Neurol*. 2012;71(3):370-84; doi:10.1002/ana.22687.
- 814 57. Zhao L, Wang G, Siegel P, He C, Wang H, Zhao W, et al. Quantitative genetic
815 background of the host influences gut microbiomes in chickens. *Sci Rep*.
816 2013;3:1163-; doi:10.1038/srep01163.
- 817 58. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al.
818 PLINK: a tool set for whole-genome association and population-based linkage
819 analyses. *American journal of human genetics*. 2007;81(3):559-75;
820 doi:10.1086/519795.
- 821 59. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan
822 KT. Data quality control in genetic case-control association studies. *Nat Protoc*.
823 2010;5(9):1564-73; doi:10.1038/nprot.2010.116.
- 824 60. Delaneau O, Marchini J, Zagury J-F. A linear complexity phasing method for
825 thousands of genomes. *Nature Methods*. 2011;9:179; doi:10.1038/nmeth.1785.
- 826 61. Delaneau O, Marchini J, The Genomes Project C, McVean GA, Donnelly P,
827 Lunter G, et al. Integrating sequence and array data to create an improved
828 1000 Genomes Project haplotype reference panel. *Nature Communications*.
829 2014;5:3934; doi:10.1038/ncomms4934.
- 830 62. Das S, Forer L, Schön herr S, Sidore C, Locke AE, Kwong A, et al.
831 Next-generation genotype imputation service and methods. *Nature Genetics*.
832 2016;48:1284; doi:10.1038/ng.3656.
- 833 63. Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, et al. The
834 international Genome sample resource (IGSR): A worldwide collection of
835 genome variation incorporating the 1000 Genomes Project data. *Nucleic Acids
836 Research*. 2016; doi:10.1093/nar/gkw829.
- 837 64. Okada Y, Kim K, Han B, Pillai NE, Ong RT, Saw WY, et al. Risk for
838 ACPA-positive rheumatoid arthritis is driven by shared HLA amino acid
839 polymorphisms in Asian and European populations. *Hum Mol Genet*.
840 2014;23(25):6916-26; doi:10.1093/hmg/ddu387.
- 841 65. Pillai NE, Okada Y, Saw WY, Ong RT, Wang X, Tantoso E, et al. Predicting
842 HLA alleles from high-resolution SNP data in three Southeast Asian
843 populations. *Hum Mol Genet*. 2014;23(16):4443-51;

- 844 doi:10.1093/hmg/ddu149.
- 845 66. Jia X, Han B, Onengut-Gumuscu S, Chen W-M, Concannon PJ, Rich SS, et al.
846 Imputing Amino Acid Polymorphisms in Human Leukocyte Antigens. *PLOS*
847 *ONE*. 2013;8(6):e64683; doi:10.1371/journal.pone.0064683.
- 848 67. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
849 *Bioinformatics*. 2010;26(19):2460-1; doi:10.1093/bioinformatics/btq461.
- 850 68. Edgar RC. UPARSE: highly accurate OTU sequences from microbial
851 amplicon reads. *Nat Methods*. 2013;10(10):996-8; doi:10.1038/nmeth.2604.
- 852 69. Second Genome, Inc: the Greengenes
853 Databases.<http://greengenes.secondgenome.com/>. Accessed 12 Mar 2019.
- 854 70. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello
855 EK, et al. QIIME allows analysis of high-throughput community sequencing
856 data. *Nat Methods*. 2010;7(5):335-6; doi:10.1038/nmeth.f.303.
- 857 71. Lee SH, Wray NR, Goddard ME, Visscher PM. Estimating missing heritability
858 for disease from genome-wide association studies. *Am J Hum Genet*.
859 2011;88(3):294-305; doi:10.1016/j.ajhg.2011.02.002.
- 860 72. Visscher PM, Hemani G, Vinkhuyzen AAE, Chen G-B, Lee SH, Wray NR, et
861 al. Statistical Power to Detect Genetic (Co)Variance of Complex Traits Using
862 SNP Data in Unrelated Samples. *PLOS Genetics*. 2014;10(4):e1004269;
863 doi:10.1371/journal.pgen.1004269.
- 864 73. Caliński T, Harabasz J. A dendrite method for cluster analysis.
865 *Communications in Statistics*. 1974;3(1):1-27;
866 doi:10.1080/03610927408827101.
- 867 74. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide
868 complex trait analysis. *Am J Hum Genet*. 2011;88(1):76-82;
869 doi:10.1016/j.ajhg.2010.11.011.
- 870 75. Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and
871 pitfalls in the application of mixed-model association methods. *Nat Genet*.
872 2014;46(2):100-6; doi:10.1038/ng.2876.
- 873 76. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al.
874 Common SNPs explain a large proportion of the heritability for human height.
875 *Nature Genetics*. 2010;42:565; doi:10.1038/ng.608.
- 876 77. Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Patterson N, et al.
877 LD Score regression distinguishes confounding from polygenicity in
878 genome-wide association studies. *Nature Genetics*. 2015;47(3):291-5;
879 doi:10.1038/ng.3211.
- 880 78. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR. Estimation of
881 pleiotropy between complex diseases using single-nucleotide
882 polymorphism-derived genomic relationships and restricted maximum
883 likelihood. *Bioinformatics*. 2012;28(19):2540-2;
884 doi:10.1093/bioinformatics/bts474.
- 885 79. Wang H, Zhang F, Zeng J, Wu Y, Kemper KE, Xue A, et al.
886 Genotype-by-environment interactions inferred from genetic effects on
887 phenotypic variability in the UK Biobank. *Science Advances*.

-
- 888 2019;5(8):eaaw3538; doi:10.1126/sciadv.aaw3538.
- 889 80. Bretherton CS, Widmann M, Dymnikov VP, Wallace JM, Bladé I. The
890 Effective Number of Spatial Degrees of Freedom of a Time-Varying Field.
891 *Journal of Climate*. 1999;12(7):1990-2009;
892 doi:10.1175/1520-0442(1999)012<1990:Tenosd>2.0.Co;2.
- 893 81. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in
894 Mendelian Randomization with Some Invalid Instruments Using a Weighted
895 Median Estimator. *Genet Epidemiol*. 2016;40(4):304-14;
896 doi:10.1002/gepi.21965.
- 897 82. Burgess S, Thompson SG. Interpreting findings from Mendelian
898 randomization using the MR-Egger method. *Eur J Epidemiol*.
899 2017;32(5):377-89; doi:10.1007/s10654-017-0255-x.
- 900 83. Verbanck M, Chen C-Y, Neale B, Do R. Detection of widespread horizontal
901 pleiotropy in causal relationships inferred from Mendelian randomization
902 between complex traits and diseases. *Nature Genetics*. 2018;50(5):693-8;
903 doi:10.1038/s41588-018-0099-7.
- 904 84. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes.
905 *Nucleic Acids Res*. 2000;28(1):27-30; doi:10.1093/nar/28.1.27.
- 906 85. Kanehisa M, Sato Y, Furumichi M, Morishima K, Tanabe M. New approach
907 for understanding genome variations in KEGG. *Nucleic Acids Res*. 2019;
908 doi:10.1093/nar/gky962.
- 909 86. Kanehisa M. Toward understanding the origin and evolution of cellular
910 organisms. *Protein Sci*. 2019; doi:10.1002/pro.3715.
- 911 87. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al.
912 Cytoscape: a software environment for integrated models of biomolecular
913 interaction networks. *Genome Res*. 2003;13(11):2498-504;
914 doi:10.1101/gr.1239303.
- 915 88. Kurtz ZD, Müller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA.
916 Sparse and Compositionally Robust Inference of Microbial Ecological
917 Networks. *PLOS Computational Biology*. 2015;
918 doi:10.1371/journal.pcbi.1004226.
- 919 89. Reynolds AP, Richards G, de la Iglesia B, Rayward-Smith VJ. Clustering
920 Rules: A Comparison of Partitioning and Hierarchical Clustering Algorithms.
921 *Journal of Mathematical Modelling and Algorithms*. 2006;5(4):475-504;
922 doi:10.1007/s10852-005-9022-1.
- 923 90. Kaufman L, Rousseeuw P. *Partitioning Around Medoids (Program PAM)*.
924 John Wiley & Sons, Inc; 1990. p. 68-125.
925

926 **Figure legends**

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

944

945 **Figure 2 The SNP-based heritability of gut microbiome.** The plot shows the taxa
 946 with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, n is the
 947 effective number of independent taxa on each taxonomy level.

948

949 **Figure 3 Effect of host genetically predicted higher atrial fibrillation risk on gut**
 950 **microbiome. (A).** Causal association of atrial fibrillation with abundance of
 951 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Bacteroides coprophilus*. The
 952 effect sizes of atrial fibrillation on taxa are changes in abundance of bacteria (1-SD of
 953 log-transformed) per genetically determined higher log odds of atrial fibrillation. **(B).**
 954 Causal association of atrial fibrillation with presence of *Barnesiellaceae*,
 955 *Veillonellaceae_undefined* and *Mitsuokella*. The effect size of atrial fibrillation on
 956 taxa are present as odds ratio increase in log odds of atrial fibrillation.

957

958 **Figure 4 Association and cluster of diseases predicted by the gut microbiome. (A).**
 959 Plot of clusters in Guangzhou Nutrition and Health Study (GNHS) cohort (n=1919).
 960 **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5 clusters in
 961 antibiotic-taking participants (n=18). The optimal cluster is 5 in GNHS cohort and 6
 962 in the replication. The clusters share consistent components between two studies. In
 963 contrast, components are different between antibiotic-taking participants and control
 964 groups. Dimension1 (Dim1) and dimension2 (Dim2) can explain 40.1% and 13.1%
 965 variance, respectively in GNHS cohort. The annotation for variables is as following.
 966 AT: African trypanosomiasis, AD: Alzheimer's disease, V1: Amoebiasis, ALS:
 967 Amyotrophic lateral sclerosis, BC: Bladder cancer, CD: Chagas disease, CML:
 968 Chronic myeloid leukemia, CRC: Colorectal cancer, V2: Hepatitis C, HD:

969 Huntington's disease, HCM: Hypertrophic cardiomyopathy, V3: Influenza A, PD:
970 Parkinson's disease, V4: Pathways in cancer, V5: Prion disease, PCa: Prostate cancer,
971 RCC: Renal cell carcinoma, SLE: Systemic lupus erythematosus, V6: Tuberculosis,
972 T1DM: Type I diabetes mellitus, T2DM: Type II diabetes mellitus, V7: Vibrio
973 cholerae infection. **(D). Gut microbiome-predicted network of relationship among**
974 **different human complex diseases. The relationship between diseases is determined**
975 **by SPIEC-EASI with non-normalized predicted abundance data. The diseases that**
976 **shared the same edge meant they had the gut microbiome-predicted correlation.**

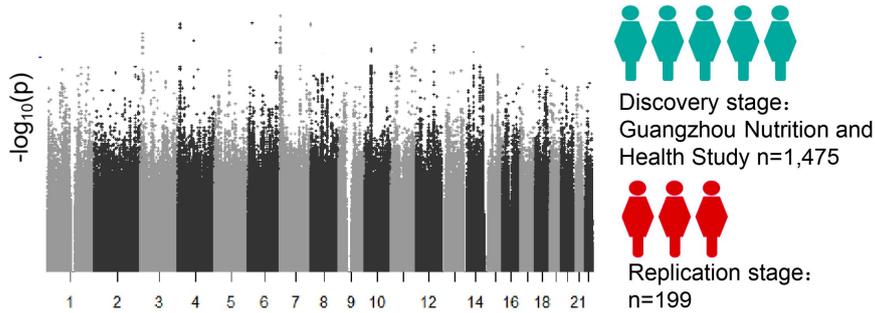
977

978 **Figure 5 Correlation of the human complex diseases with gut microbiome.** The
979 heat map shows **Spearman's correlation** of predicted diseases and gut microbiome on
980 genus level. The grey components show no significance of correlation with
981 Bonferroni correction ($p > 0.05 / (5.6 * 22)$, $p > 0.0004$).

982

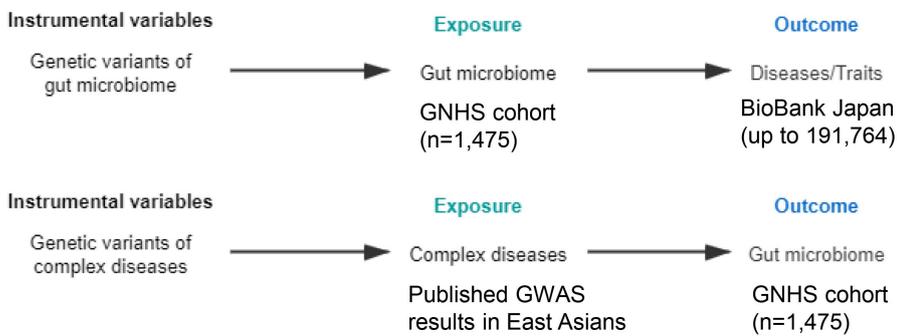
983 **Figure 1 Study overview.** The figure shows the highlights of our study. First, we
984 performed a microbiome genome-wide association study in a Chinese population
985 (Step A). We validated significant genetic variants reported in the previous studies and
986 replicated our results in an independent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

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



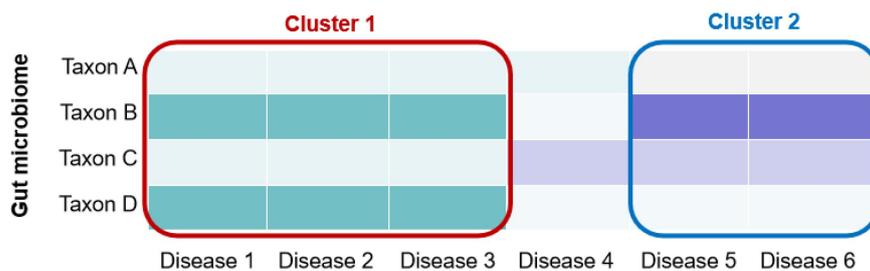
B. The causal relationships between gut microbiome and human complex diseases.

Bi-directional Mendelian randomization



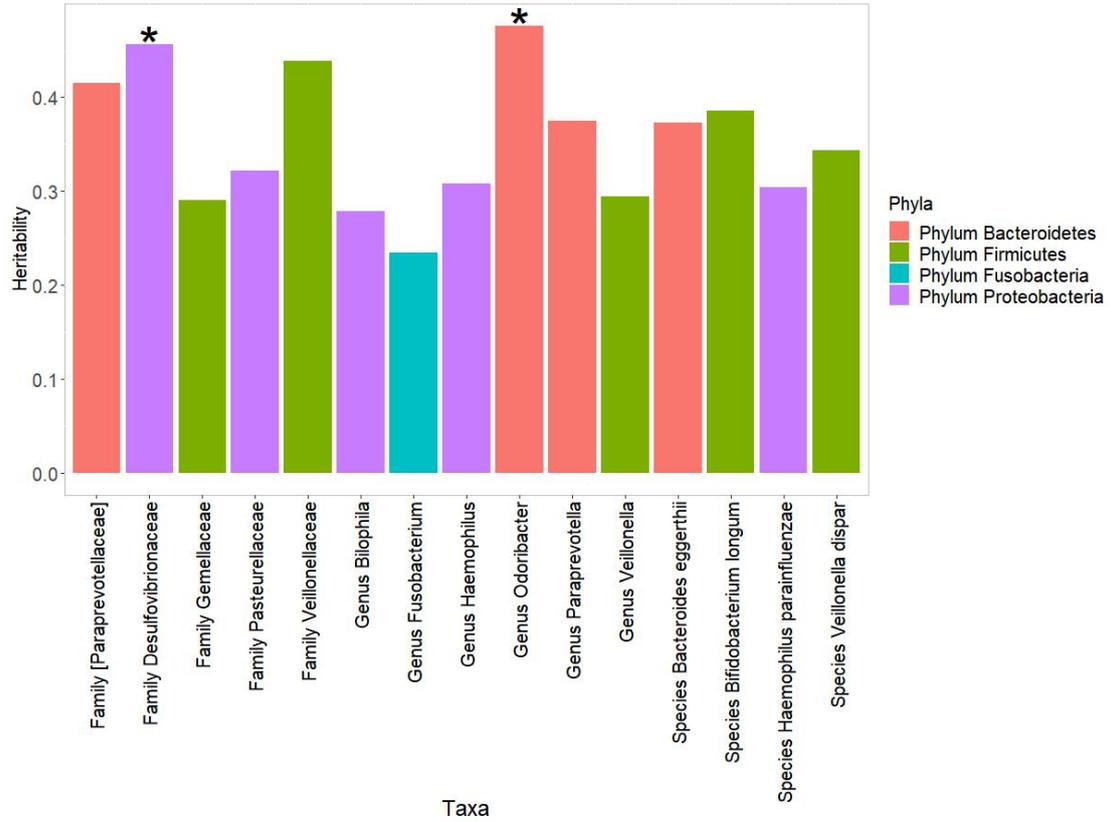
C. The shared and distinct microbiome features among human complex diseases.

1,919 participants from GNHS cohort.



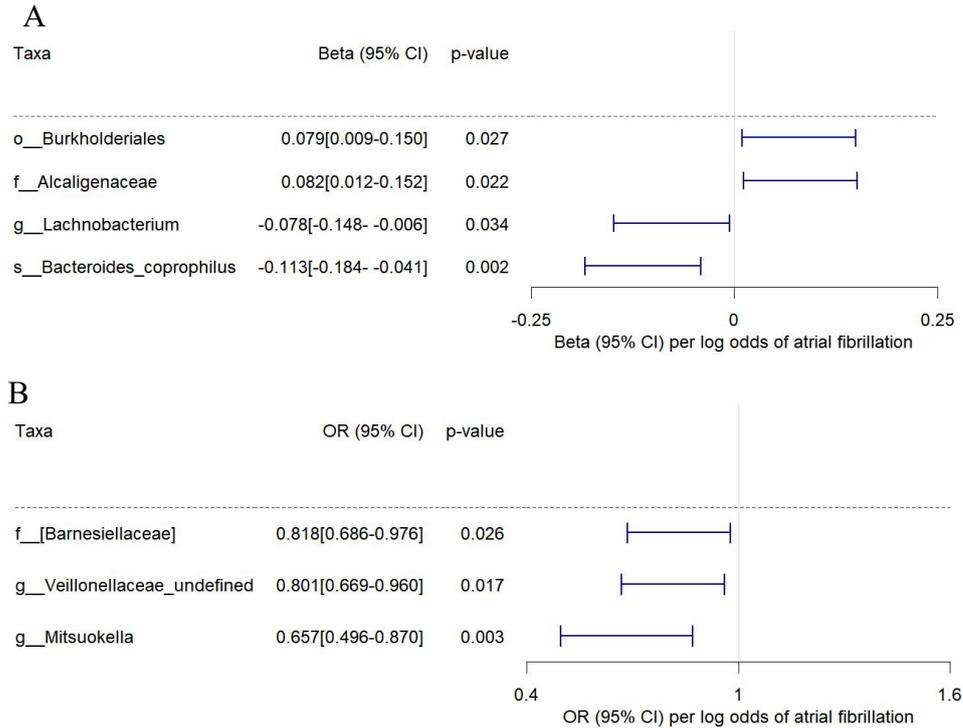
1000

1001 **Figure 2 The SNP-based heritability of gut microbiome.** The plot shows the taxa
 1002 with nominally significant heritability estimates ($p < 0.05$). * $p < 0.05/n$, n is the
 1003 effective number of independent taxa on each taxonomy level.



1004

1005 **Figure 3 Effect of host genetically predicted higher atrial fibrillation risk on gut**
 1006 **microbiome. (A).** Causal association of atrial fibrillation with abundance of
 1007 *Burkholderiales*, *Alcaligenaceae*, *Lachnobacterium* and *Bacteroides coprophilus*. The
 1008 effect sizes of atrial fibrillation on taxa are changes in abundance of bacteria (1-SD of
 1009 log-transformed) per genetically determined higher log odds of atrial fibrillation. **(B).**
 1010 Causal association of atrial fibrillation with presence of *Barnesiellaceae*,
 1011 *Veillonellaceae_undefined* and *Mitsuokella*. The effect size of atrial fibrillation on
 1012 taxa are present as odds ratio increase in log odds of atrial fibrillation.

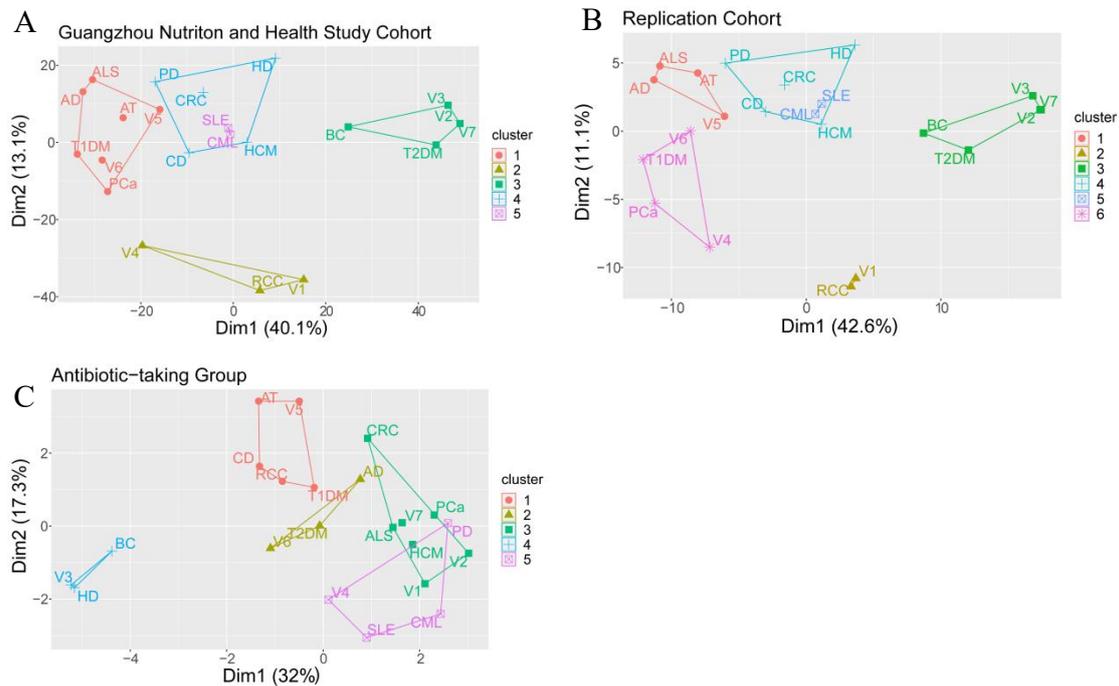


1013

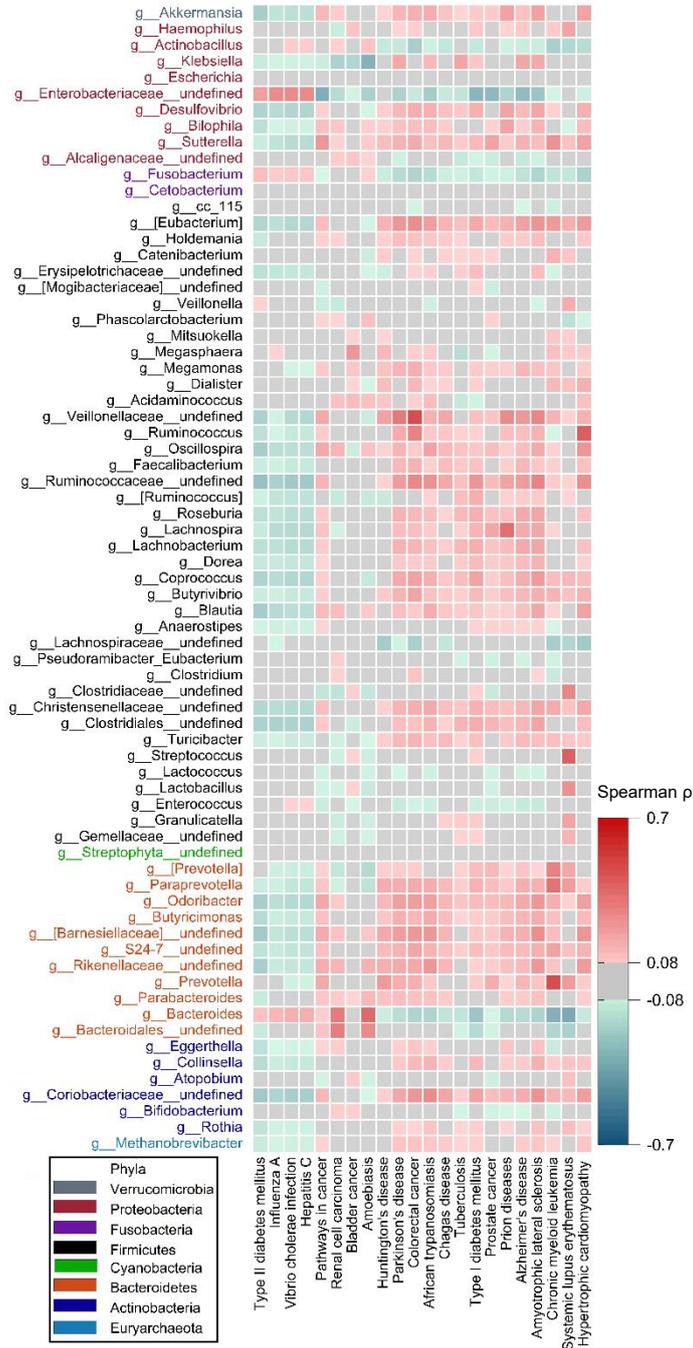
1014

1015

1016 **Figure 4 Association and cluster of diseases predicted by the gut microbiome. (A).**
 1017 Plot of clusters in Guangzhou Nutrition and Health Study (GNHS) cohort (n=1919).
 1018 **(B).** Plot of cluster results in the replication cohort (n=217). **(C).** Plot of 5 clusters in
 1019 antibiotic-taking participants (n=18). The optimal cluster is 5 in GNHS cohort and 6
 1020 in the replication. The clusters share consistent components between two studies. In
 1021 contrast, components are different between antibiotic-taking participants and control
 1022 groups. Dimension1 (Dim1) and dimension2 (Dim2) can explain 40.1% and 13.1%
 1023 variance, respectively in GNHS cohort. The annotation for variables is as following.
 1024 AT: African trypanosomiasis, AD: Alzheimer's disease, V1: Amoebiasis, ALS:
 1025 Amyotrophic lateral sclerosis, BC: Bladder cancer, CD: Chagas disease, CML:
 1026 Chronic myeloid leukemia, CRC: Colorectal cancer, V2: Hepatitis C, HD:
 1027 Huntington's disease, HCM: Hypertrophic cardiomyopathy, V3: Influenza A, PD:
 1028 Parkinson's disease, V4: Pathways in cancer, V5: Prion disease, PCa: Prostate cancer,
 1029 RCC: Renal cell carcinoma, SLE: Systemic lupus erythematosus, V6: Tuberculosis,
 1030 T1DM: Type I diabetes mellitus, T2DM: Type II diabetes mellitus, V7: Vibrio
 1031 cholerae infection. **(D).** Gut microbiome-predicted network of relationship among
 1032 different human complex diseases. The relationship between diseases is determined
 1033 by SPIEC-EASI with non-normalized predicted abundance data. The diseases that
 1034 shared the same edge meant they had the gut microbiome-predicted correlation.
 1035



1038 **Figure 5 Correlation of the human complex diseases with gut microbiome.** The
 1039 heat map shows Spearman's correlation of predicted diseases and gut microbiome on
 1040 genus level. The grey components show no significance of correlation with
 1041 Bonferroni correction ($p > 0.05 / (5.6 * 22)$, $p > 0.0004$).



1042

| | |
|------|---|
| 1043 | Supplementary Tables [<i>Supplementary Tables.xls</i>] |
| 1044 | |
| 1045 | Supplementary Table S1 Transformation of traits in BioBank Japan and taxa in GNHS |
| 1046 | Supplementary Table S2 Required effect size (beta) to reach 80% of power in GNHS cohort |
| 1047 | Supplementary Table S3 Heritability of taxa, enterotype and alpha diversity |
| 1048 | Supplementary Table S4 Significant genetic correlations of gut microbiome and metabolic |
| 1049 | traits |
| 1050 | Supplementary Table S5 Significant associations of all taxa with SNPs identified in the |
| 1051 | discovery stage before adjustment($p < 5e-8$) |
| 1052 | Supplementary Table S6 Replication of genetic variants associated with taxa |
| 1053 | Supplementary Table S7 Replication of genetic variants associated with beta diversity |
| 1054 | Supplementary Table S8 Replication of genetic variants associated with alpha diversity |
| 1055 | Supplementary Table S9 Lead SNPs used to construct polygenic scores |
| 1056 | Supplementary Table S10 MR analysis of gut microbiota on traits and diseases |
| 1057 | Supplementary Table S11 MR analysis of diseases on gut microbiota features |
| 1058 | Supplementary Table S12 Spearman's correlation of certain taxa and complex diseases |
| 1059 | Supplementary Table S13 Spearman's correlation of gut microbiota on genus level and |
| 1060 | characteristics |
| 1061 | |
| 1062 | Supplementary Figures [<i>Supplementary Figures.pdf</i>] |
| 1063 | |
| 1064 | Supplementary Figure 1 Genome-wide analysis results of taxa. |
| 1065 | Supplementary Figure 2 Spearman's correlation of the relative abundance of AF-associated |
| 1066 | taxa with the relative level of diseases predicted by PICRUSt. |