

A microbial causal mediation analytic tool for health disparity and applications in body mass index

Chan Wang

New York University Grossman School of Medicine

Jiyoung Ahn

New York University Grossman School of Medicine

Thaddeus Tarpey

New York University Grossman School of Medicine

Stella S. Yi

New York University Grossman School of Medicine

Richard B. Hayes

New York University Grossman School of Medicine

Huilin Li (✉ Huilin.Li@nyulangone.org)

New York University Grossman School of Medicine

Method Article

Keywords: Casual mediation model, Health disparity, Manipulable disparity measure, Microbiome mediator, Non-manipulable exposure

Posted Date: January 13th, 2023

DOI: <https://doi.org/10.21203/rs.3.rs-2463503/v1>

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

[Read Full License](#)

Additional Declarations: No competing interests reported.

1 **A microbial causal mediation analytic tool for health disparity and applications in body**
2 **mass index**

3
4 Chan Wang¹, Jiyoung Ahn², Thaddeus Tarpey¹, Stella S. Yi³, Richard B. Hayes², Huilin Li^{1*}

5 ¹Division of Biostatistics, Department of Population Health, New York University Grossman School of
6 Medicine, New York, 10016, NY, USA

7 ²Division of Epidemiology, Department of Population Health, New York University Grossman School of
8 Medicine, New York, 10016, NY, USA

9 ³Department of Population Health Section for Health Equity, New York University Grossman School of
10 Medicine, New York, 10016, USA.

11 *Correspondence: Huilin.Li@nyulangone.org

12
13 Emails: Chan Wang: Chan.Wang@nyulangone.org, Jiyoung Ahn: Jiyoung.Ahn@nyulangone.org,
14 Thaddeus Tarpey: Thaddeus.Tarpey@nyulangone.org, Stella S. Yi: Stella.Yi@nyulangone.org, Richard
15 B. Hayes: Richard.B.Hayes@nyulangone.org, Huilin Li: Huilin.Li@nyulangone.org

16

17

18

19

20

21 **Abstract**

22 **Background:** Emerging evidence suggests the potential mediating role of microbiome in health
23 disparities. However, no analytic framework is available to analyze microbiome as a mediator between
24 health disparity and clinical outcome, due to the unique structure of microbiome data, including high
25 dimensionality, sparsity, and compositionality.

26 **Methods:** Considering the modifiable and quantitative features of microbiome, we propose a microbial
27 causal mediation model framework, SparseMCMM_HD, to uncover the mediating role of microbiome in
28 health disparities, by depicting a plausible path from a non-manipulable exposure (e.g. race or region) to a
29 continuous outcome through microbiome. The proposed SparseMCMM_HD rigorously defines and
30 quantifies the manipulable disparity measure that would be eliminated by equalizing microbiome profiles
31 between comparison and reference groups. Moreover, two tests checking the impact of microbiome on
32 health disparity are proposed.

33 **Results:** Through three body mass index (BMI) studies selected from the curatedMetagenomicData 3.4.2
34 package and the American gut project: China vs. USA, China vs. UK, and Asian or Pacific Islander (API)
35 vs. Caucasian, we exhibit the utility of the proposed SparseMCMM_HD framework for investigating
36 microbiome's contributions in health disparities. Specifically, BMI exhibits disparities and microbial
37 community diversities are significantly distinctive between the reference and comparison groups in all
38 three applications. By employing SparseMCMM_HD, we illustrate that microbiome plays a crucial role
39 in explaining the disparities in BMI between races or regions. 11.99%, 12.90%, and 7.4% of the overall
40 disparity in BMI in China-USA, China-UK, and API-Caucasian comparisons, respectively, would be
41 eliminated if the between-group microbiome profiles were equalized; and 15, 21, and 12 species are
42 identified to play the mediating role respectively.

43 **Conclusions:** The proposed SparseMCMM_HD is an effective and validated tool to elucidate the
44 mediating role of microbiome in health disparity. Three BMI applications shed light on the utility of
45 microbiome in reducing BMI disparity by manipulating microbial profiles.

46 **Keywords:** Casual mediation model; Health disparity; Manipulable disparity measure; Microbiome
47 mediator; Non-manipulable exposure

48

49 **Background**

50 Health disparities refer to the inequalities in the quality of health, health care, and health outcomes
51 experienced by groups that are usually classified by race, ethnicity, and region. Many factors, including
52 genetics, social-economic status, culture, dietary habits, and geographical conditions, contribute to health
53 disparities between groups. Researchers have long been interested in identifying the modifiable
54 environmental determinants of health disparity to pave the way to improve health equity. However,
55 environmental exposures are often numerous, ubiquitous, descriptive, or hard to measure, which makes
56 this task difficult.

57 Gut microbiome is the aggregate of all genomes harbored by gut microbiota, which is the collection of all
58 microbes that reside in human gut. Benefiting from the advent of high throughput sequencing
59 technologies, a great number of microbiome studies have been conducted to quantitatively characterize
60 the microbiome profiling and understand its role in human health [1-4]. Gut microbiome has been closely
61 linked with host metabolic, immune, and neuroendocrine functions [5-12]. On the other hand, many
62 environmental and social factors, such as diet, drugs, lifestyle, psychological state and behavior, aid in
63 shaping gut microbial profiles [13-16]. Recently, the mediating role of microbiome between these
64 environmental exposures and various human diseases, including obesity, type 2 diabetes, inflammatory
65 bowel disease, depression, and different cancers, has been investigated and recognized [17-22]. Given the
66 modifiable and quantitative features of microbiome, we here aim to disentangle health disparities by

67 exploring the extent of the observed disparity in the outcome of interest that could be reduced if the gut
68 microbial profile was modified. In Figure 1, we propose a mediation framework to answer such questions.
69 Here, the disparity group, e.g., race or region, is the exposure denoted by R ; the gut microbial profile is
70 the mediator denoted by M ; and the continuous study outcome, e.g., body mass index (BMI), is denoted
71 by Y .

72 There are several existing mediation analysis frameworks tailored for non-manipulable exposures, such as
73 race, region, sex or socioeconomic position [27], however, due to the unique structure of microbiome
74 data, including high dimensionality, sparsity and compositionality, these approaches are not immediately
75 applicable for analyzing microbiome as a mediator for health disparity. Recently, we developed a rigorous
76 Sparse Microbial Causal Mediation Model (SparseMCMM) [12] for interrogating the mediating role of
77 microbiome in a typical three-factor (randomized treatments, microbiome as mediator, and outcome)
78 clinical trial causal study design. SparseMCMM quantifies the overall mediation effect of microbiome
79 community and the component-wise mediation effect for each individual microbe under the
80 counterfactual framework, identifies the signature causal microbes with regularization strategies, and tests
81 the mediation effects while fully acknowledging the unique structure of microbiome data. In this paper,
82 by extending SparseMCMM to a non-manipulable exposure setting, we propose a microbial causal
83 mediation framework for health disparity study and denote it as SparseMCMM_HD (SparseMCMM for
84 Health Disparity). As VanderWeele and Robinson [23] discussed, causal interpretation of a non-
85 manipulable exposure, i.e., ethnicity or region, is not definable in the traditional counterfactual
86 framework, because a hypothetical intervention on a non-manipulable exposure is not possible. Instead,
87 one can interpret the causality of health inequality by the hypothesized intervention effect on the
88 manipulable mediating variable. Thus, in SparseMCMM_HD, we aim to quantify the overall health
89 inequality on the outcome (called overall disparity), the health inequality effect that would be eliminated
90 by equalizing microbiome profiles across racial or regional groups (called manipulable disparity), and the
91 healthy inequality effect that would remain even after microbiome profiles across racial or regional

92 groups were equalized (called residual disparity). In addition, we equip two hypothesis tests to examine
93 the mediating role of microbiome in health disparity and statistically identify which specific microbes
94 contribute to it.

95 Obesity (defined via BMI) is a global epidemic and a persistent public health problem [24]. It is well
96 documented that the prevalence of adult obesity is distributed unevenly across racial groups and regions.
97 Partial effect of manipulable exposures such as diet, medication, and antibiotics use [17-19] on obesity
98 has been shown to be mediated through microbiome. In addition, accumulating evidence indicates that gut
99 microbial profile varies across ethnicities as well as geographically [25-27]. Together, these studies
100 suggest that microbiome may play a mediating role in the ethnic or regional disparity of obesity. It is
101 crucial to investigate rigorously how much health inequalities in BMI can be reduced by manipulating
102 microbiome profiles. Utilizing SparseMCMM_HD, we investigate the role of microbiome in the regional
103 and racial disparity of BMI in curated microbiome data from the curatedMetagenomicData 3.4.2 package
104 [28] and the American Gut Project (AGP) (www.americangut.org) respectively. Through these real data
105 analyses, we illustrate a clear and plausible causal path analysis to understand the current racial or
106 regional disparity in BMI and identify a comprehensive set of mediating microbial taxa. The proposed
107 analytic pipeline is available through an interactive web app at
108 <https://chanw0.shinyapps.io/sparsemcmhd/>. We believe that this novel pipeline will be useful for
109 investigating the manipulable disparity through gut microbiome and understanding the causes of health
110 disparity.

111 **Methods**

112 **SparseMCMM_HD framework**

113 **Casual mediation model.** Suppose there are I subjects from two categories of a non-manipulable
114 exposure group (e.g. race or region), J taxa, and K covariates. Subscripts i, j, k , indicate a subject, a
115 taxon, and a covariate respectively. For the i th subject, let $R_i = 1$ or 0 indicate the reference or

116 comparison group, let $\mathbf{M}_i = (M_{i1}, \dots, M_{iJ})^T$ be the microbiome relative abundance vector with the
 117 constraint $\sum_{j=1}^J M_{ij} = 1$, and let $\mathbf{X}_i = (X_{i1}, \dots, X_{iK})^T$ represent the covariates, and let Y_i be a continuous
 118 outcome of interest.

119 To statistically describe the causal relationships shown in Figure 1, following our previous work [12], we
 120 use the linear log-contrast model to regress the continuous outcome on the non-manipulable exposure,
 121 microbiome compositions, interactions between the non-manipulable exposure and microbiome
 122 compositions, while adjusting the confounding covariates:

$$123 \quad Y_i = \alpha_0 + \boldsymbol{\alpha}_X^T \mathbf{X}_i + \alpha_R R_i + \boldsymbol{\alpha}_M^T [\log(\mathbf{M}_i)] + \boldsymbol{\alpha}_C^T [\log(\mathbf{M}_i)] R_i + \epsilon_i, \quad (1)$$

subject to $\boldsymbol{\alpha}_M^T \mathbf{1} = 0$, and $\boldsymbol{\alpha}_C^T \mathbf{1} = 0$,

124 where α_0 is the intercept, α_R is the coefficient of the non-manipulable exposure, $\boldsymbol{\alpha}_X = (\alpha_{X1}, \dots, \alpha_{XK})^T$,
 125 $\boldsymbol{\alpha}_M = (\alpha_{M1}, \dots, \alpha_{MJ})^T$, and $\boldsymbol{\alpha}_C = (\alpha_{C1}, \dots, \alpha_{CJ})^T$ are the vectors of coefficients of covariates, microbiome
 126 compositions, interactions between the non-manipulable exposure and microbiome compositions,
 127 respectively. Due to the compositionality of microbiome data as $\sum_{j=1}^J M_{ij} = 1$, $\boldsymbol{\alpha}_M$ and $\boldsymbol{\alpha}_C$ are subject to
 128 $\boldsymbol{\alpha}_M^T \mathbf{1} = 0$ and $\boldsymbol{\alpha}_C^T \mathbf{1} = 0$. $\epsilon_i \sim N(0, \sigma^2)$ is the error term. On the other hand, the Dirichlet regression [29]
 129 is used to model the microbial relative abundance as a function of the non-manipulable exposure and
 130 covariates:

$$131 \quad E[M_{ij}] = \frac{\gamma_j(R_i, \mathbf{X}_i)}{\sum_{m=1}^J \gamma_m(R_i, \mathbf{X}_i)}, \quad (2)$$

$$\log\{\gamma_j(R_i, \mathbf{X}_i)\} = \beta_{0j} + \beta_{Rj} R_i + \boldsymbol{\beta}_{Xj}^T \mathbf{X}_i.$$

132 Specifically, we assume that $\mathbf{M}_i | (R_i, \mathbf{X}_i) \sim \text{Dirichlet}(\gamma_1(R_i, \mathbf{X}_i), \dots, \gamma_J(R_i, \mathbf{X}_i))$, and their microbial
 133 relative means are linked with the non-manipulable exposure and covariates (R_i, \mathbf{X}_i) in the generalized
 134 linear model fashion with a log link. β_{0j} is the intercept and β_{Rj} and $\boldsymbol{\beta}_{Xj}$ are the coefficients of the non-
 135 manipulable exposure and covariates for the j th taxon, respectively.

136 **Definition of disparity measures in the counterfactual framework.** As discussed in the Background,
 137 we propose to conceptualize an overall disparity measure (ODM) on the outcome that can be decomposed
 138 into manipulable disparity measure (MDM) and residual disparity measure (RDM). MDM represents the
 139 portion of disparity that would be eliminated by equalizing microbiome profiles between comparison and
 140 reference groups, and RDM represents the portion that would remain even after microbiome profiles
 141 between comparison and reference groups were equalized. With the counterfactual notation,
 142 mathematically we have:

$$143 \quad \text{ODM} = \text{MDM} + \text{RDM},$$

$$144 \quad \text{MDM} = E[E[Y_{\mathbf{M}_x(1)}|R = 1, \mathbf{x}]] - E[E[Y_{\mathbf{M}_x(0)}|R = 1, \mathbf{x}]], \text{ and}$$

$$145 \quad \text{RDM} = E[E[Y_{\mathbf{M}_x(0)}|R = 1, \mathbf{x}] - E[Y_{\mathbf{M}_x(0)}|R = 0, \mathbf{x}]].$$

146 Here, $\mathbf{M}_x(0)$ ($\mathbf{M}_x(1)$) is a random value from the microbiome distribution of the reference (comparison)
 147 population with given covariates \mathbf{x} . Y_m denotes an individual's potential counterfactual outcome if his or
 148 her microbial mediators were set to \mathbf{m} , where \mathbf{m} can be $\mathbf{M}_x(0)$ or $\mathbf{M}_x(1)$. $E[Y_{\mathbf{M}_x(0)}|R = 0, \mathbf{x}]$
 149 ($E[Y_{\mathbf{M}_x(1)}|R = 1, \mathbf{x}]$) denotes the expected outcome for a reference (comparison) individual with given
 150 covariates \mathbf{x} , $E[Y_{\mathbf{M}_x(0)}|R = 1, \mathbf{x}]$ denotes the expected outcome for a comparison individual with given
 151 covariates \mathbf{x} if their microbial mediators were set to a random value from that of the reference population
 152 with the same covariates \mathbf{x} .

153 **MDM, RDM, and ODM expressions.** Two assumptions must be satisfied for the identification of MDM,
 154 RDM, and ODM [23, 30]. The effect of non-manipulable exposure R on outcome Y are unconfounded
 155 conditional on all covariates \mathbf{X} , i.e., $Y \perp\!\!\!\perp R \mid \mathbf{X}$ and the effects of mediator \mathbf{M} on outcome Y are
 156 unconfounded conditional on the non-manipulable exposure R and all covariates \mathbf{X} , i.e., $Y \perp\!\!\!\perp \mathbf{M} \mid R, \mathbf{X}$.

157 With these sufficient identifiability assumptions and the models (1)-(2) proposed in the
 158 SparseMCMM_HD framework, disparity measures MDM, RDM, and ODM can be further expressed,
 159 respectively, as follows (see Section S1 for the detailed derivations):

160
$$\text{MDM} = \sum_{j=1}^J (\alpha_{Mj} + \alpha_{Cj}) \{E[\log(M_j)|R = 1, \mathbf{x}] - E[\log(M_j)|R = 0, \mathbf{x}]\},$$

161
$$\text{RDM} = \alpha_R + \boldsymbol{\alpha}_C^T E[\log(\mathbf{M})|R = 0, \mathbf{x}] = \alpha_R + \sum_{j=1}^J \alpha_{Cj} E[\log(M_j)|R = 0, \mathbf{x}],$$

162 and

163
$$\begin{aligned} \text{ODM} &= \text{MDM} + \text{RDM} \\ &= \alpha_R + \sum_{j=1}^J (\alpha_{Mj} + \alpha_{Cj}) E[\log(M_j)|R = 1, \mathbf{x}] - \sum_{j=1}^J \alpha_{Mj} E[\log(M_j)|R = 0, \mathbf{x}], \end{aligned}$$

164 where $E[\log(M_j)|R = r, \mathbf{x}] = \psi[\gamma_j(R = r, \mathbf{x})] - \psi[\sum_{m=1}^J \gamma_m(R = r, \mathbf{x})]$, $\gamma_j(R = r, \mathbf{x}) =$

165 $\exp(\beta_{0j} + \beta_{Rj}r + \boldsymbol{\beta}_{Xj}^T \mathbf{x})$, $r = 0$ or 1 , and $\psi(\cdot) = \frac{d}{dx} \ln(\Gamma(x))$ is the digamma function, with given

166 covariates \mathbf{x} .

167 Note that these mathematical expressions of RDM and MDM are the same as the formulas of causal
 168 direct effect of treatment and mediation effect through microbiome correspondingly on the outcome in the
 169 typical three-factor causal design based on the traditional causal mediation inference, developed in our
 170 SparseMCMM [12]. Analogous to ME in SparseMCMM, MDM is the summation of individual mediation
 171 effects from each taxon MDM_j : $\text{MDM} := \sum_{j=1}^J MDM_j$ and $MDM_j = (\alpha_{Mj} + \alpha_{Cj})\{E[\log(M_j)|R = 1, \mathbf{x}] -$
 172 $E[\log(M_j)|R = 0, \mathbf{x}]\}$. MDM_j thus is non-zero only when both the j th microbial effect on the outcome and
 173 the exposure effect on the j th taxon are not zero. Therefore, SparseMCMM_HD illuminates the mediating
 174 role of microbiome in the health disparity of outcome, and quantifies the manipulable disparity for overall
 175 microbiome community and for each specific taxon, respectively.

176 **Parameter estimation.** Note that in [12], we have demonstrated the excellent performance of
 177 SparseMCMM in terms of estimation by extensive simulations and real data analysis in various scenarios.
 178 Thus for SparseMCMM_HD, we directly employ the same two-step procedure to estimate the regression

179 parameters in models (1)-(2) to obtain the estimated RDM, MDM, MDM_j for each taxon, and ODM.
180 Furthermore, SparseMCMM_HD has the full capability to perform variable selection to select the
181 signature causal microbes that play mediating roles in the disparity of the continuous outcome with
182 regularization strategies. Specifically, L_1 norm and group-lasso penalties are incorporated for variable
183 selection meanwhile addressing the heredity condition.

184 **Hypothesis tests for manipulable disparity.** Similarly, we employ the hypothesis tests for mediation
185 effects in SparseMCMM to examine whether microbiome has any mediation effect on the disparity in an
186 outcome, at both community and taxon levels. Specifically, regarding the null hypothesis of no
187 manipulable disparity $H_0: MDM = 0$, the first test statistic is defined as $OMD = \widehat{MDM}$, the estimator of the
188 manipulable disparity. Meanwhile, we consider another null hypothesis, $H_0: MDM_j = 0, \forall j \in \{1, \dots, J\}$
189 and define the second test statistic as $CMD = \sum_{j=1}^J \widehat{MDM}_j^2$, the summation of the squared estimators of
190 individual mediation effects across all taxa. Permutation procedure is employed to assess the significance
191 of these two test statistics. This provides a mechanism to check whether microbiome has any impact on
192 health disparity that could be potentially eliminated through microbiome.

193 **Implementation.** The simulation evaluation results regarding the estimation and testing of
194 SparseMCMM [12] are applicable to SparseMCMM_HD framework. Therefore, the proposed
195 SparseMCMM_HD is a validated analytic tool to illuminate the mediating role of microbiome in the
196 disparity of outcome, and quantifies the manipulable disparity for overall microbiome community and for
197 each specific taxon, respectively. In practice, we perform both parameter estimation and hypothesis
198 testing using the analytical procedures in the SparseMCMM package and illustrate the proposed
199 SparseMCMM_HD pipeline through an interactive web app
200 (https://chanw0.shinyapps.io/sparsemcm_hd/).

201 **Control for confounding covariates**

202 Due to the non-manipulable nature of the exposure in health disparity research, in principle, it is
203 impossible to design a randomized trial on the exposure of interest to eliminate the potential confounding
204 effect on the interested causal pathway. Many studies on health disparity are observational and usually
205 include significant degrees of confounding, due to factors such as lifestyle, health status, and disease
206 history. We want to emphasize that it is a necessary step to control for confounding covariates while
207 utilizing the proposed SparseMCMM_HD to estimate RDM, MDM, and ODM in a typical observational
208 study. Specifically, we propose to perform propensity score matching (PSM) [31], which is a commonly
209 used method in biomedical research to create a balanced covariate distribution between two groups, to
210 control confounding covariates in our applications (see Section S2). Standardized mean difference (SMD)
211 is used to evaluate the balance of the covariate distributions between groups. A SMD that is less than 0.1
212 indicates a balanced distribution [32]. The matched data will then be used to quantify RDM, MDM, and
213 ODM, and examine whether the microbiome could reduce the health disparity between two non-
214 manipulable exposure groups. The control for confounding covariates procedure has been included as a
215 preprocessing step in the proposed SparseMCMM_HD analytic pipeline.

216 **curatedMetagenomicDataV3.4.2**

217 The curatedMetagenomicData 3.4.2 package [28] provides a curated human microbiome meta dataset
218 aggregated from 86 shotgun sequencing cohorts in 6 body sites. The raw sequencing data were processed
219 using the same bioinformatics protocol and pipelines. Each sample has 6 types of data available including
220 gene family, marker abundance, marker presence, pathway abundance, pathway coverage, and taxonomic
221 (relative) abundance. The taxonomic abundance was calculated with MetaPhlan3, and metabolic
222 functional potential was calculated with HUMAnN3. The manually curated clinical and phenotypic
223 metadata are available as well. More details can be found in the curatedMetagenomicData package
224 document [28]. Here we focus on healthy subjects to explore the relationship among region, microbiome,
225 and BMI. Specifically, we chose subjects from all cohorts based on the following inclusion criteria: 1)

226 healthy status; 2) no missing values in BMI, gender, and age; 3) age ≥ 18 ; 4) no pregnant; 5) currently no
227 antibiotic use; 6) currently no alcohol consumption; 7) no smoking; and 8) fecal sample with more than
228 1,250 sample reads. In addition, when multiple samples available for a subject, we randomly selected one
229 sample. Overall, we identified 4,868 healthy adults from different regions. Here we further focus on three
230 regional groups which have large sample sizes: China (n=570), United States (USA; n=350), and United
231 Kingdom (UK; n=1019) for the analysis in the main text. Specifically, we conducted two comparison
232 studies: China-USA and China-UK comparisons to investigate the regional difference of BMI in the
233 China group compared to the USA and UK groups, respectively.

234 **American Gut Project**

235 The AGP project is a crowd-sourcing citizen science cohort to describe the comprehensive
236 characterization of human gut microbiota and to identify factors being linked to human microbiota. The
237 AGP includes 16S rRNA V4 gene sequences from more than 8,000 fecal samples using standard
238 pipelines, and host metadata. Detailed descriptions can be found in Liu et al. and Hu et al. [1, 33]. Our
239 primary investigation is on the disparity of BMI between Asian or Pacific Islander (API) and non-
240 Hispanic Caucasian adults. We selected a subset of the AGP data based on the following inclusion
241 criteria: 1) USA resident; 2) Asian or Pacific Islander or Caucasian race; 3) no missing values in gender,
242 age, and BMI; 4) age ≥ 18 ; 5) $80 \leq \text{BMI}$; 6) $210\text{cm} \geq \text{height} \geq 80\text{cm}$; 7) $200\text{kg} \geq \text{weight} \geq 35\text{kg}$; 8)
243 fecal sample with more than 1,250 sample reads; 9) not duplicate sample; and 10) no self-reported history
244 of inflammatory bowel disease, diabetes, or antibiotic use in the past year. The subjects are filtered out
245 when the reported BMIs are not consistent with the calculated BMI based on the reported heights and
246 weights, i.e. $(|\text{BMI}_{\text{reported}} - \text{BMI}_{\text{calculated}}| / \text{BMI}_{\text{calculated}} > 5\%)$. A dataset with 130 API and 2,263
247 Caucasian adults then is used in this paper (Figure S1a).

248 **Statistical Analysis**

249 Data pre-processing and PSM were conducted in three BMI studies. Specifically, for the China-USA and
250 China-UK comparisons, we performed PSM with the parameters described in Section S2 to control for age

251 and gender. For the API-Caucasian comparison, as the AGP includes more than 400 covariates that were
252 collected through self-reported surveys, we first implemented several pre-processing steps to prepare the
253 self-reported covariates for the subsequent analysis, including cleaning up the inconsistent definition of
254 variables, and collapsing the sparse categorical variables into fewer and less sparse categories. Details are
255 provided in Section S3. Forty-four covariates were retained for PSM. We performed univariate linear
256 regressions to identify the potential confounding variables for the relationship among race, microbiome,
257 and BMI. Twenty-three covariates ($p\text{-value} \leq 0.05$; Figure S1b) were identified as confounders that need
258 to be controlled further based on PSM.

259 With the matched data, alpha (Observed, Shannon, and Simpson indices) and beta diversities (Bray–Curtis
260 dissimilarity and Jensen–Shannon divergence) were used to estimate microbial community-level diversity.
261 T tests were used for group comparisons of BMI and alpha diversity. Permutational multivariate analysis
262 of variance (PERMANOVA) [34] was used to assess group difference of beta diversity. We performed the
263 proposed SparseMCMC_HD framework at the species rank (Section S4) to quantify RDM, MDM, and
264 ODM, and examine whether the microbiome could explain the health disparity between two non-
265 manipulable exposure groups. The proposed SparseMCMC_HD pipeline was implemented through an
266 interactive web app (https://chanw0.shinyapps.io/sparsemcmc_hd/) for easy exploration.

267 **Results**

268 **Results for curatedMetagenomicDataV3.4.2**

269 **Matched datasets.** With the healthy adults included in the China-USA and China-UK comparisons, we
270 identified 328 matched Chinese-USA subject pairs, and 559 matched Chinese-UK subject pairs,
271 separately. Figures S2 and S3 show that both matched datasets have comparable propensity scores. The
272 SMDs decrease dramatically on the matched subjects (SMD=0.036 and 0.033), from using all subjects
273 (SMD=0.302 and 0.470) in both China-USA and China-UK datasets. This indicates that PSM has
274 effectively evened the distribution of confounders between two exposure groups in our studies and

275 practically eliminated or controlled the influence of the confounders. In the well-matched datasets, the
276 China group still has significantly lower average BMIs compared to the matched USA (mean [standard
277 deviation]: 22.64 [3.77] vs. 25.77 [4.56]) and the matched UK (22.98 [4.48] vs. 25.77 [4.79]) groups
278 (Figure 2a and 2d).

279 **Community level results.** The Chinese group has distinctive microbial community diversities, compared
280 to the matched USA or UK group. For alpha diversity, samples from China have lower Shannon and
281 Simpson diversities and a higher observed diversity than the matched USA or UK samples (Figure 2b and
282 2e). For beta diversity, Bray-Curtis dissimilarity and Jensen-Shannon divergence both indicate that the
283 Chinese group is significantly different in community structure from the matched USA or UK groups
284 (PERMANOVA [34] all p-values $< 1.0 \times 10^{-4}$. Figure 2c and 2f).

285 **Taxon-level analysis.** After implementing the filtering criteria described in Section S4, 25 species
286 remained in both matched datasets (China vs. USA and China vs. UK). The testing results for OMD and
287 CMD show that the overall and component-wise MDMs through microbiome are significant in both data
288 sets for regional differences in BMI (all p-values < 0.001 based on 1,000 permutations). Figure 3a shows
289 that the ODM of BMI are 3.17 and 2.79, respectively, for the matched Chinese and USA subjects, and the
290 matched Chinese and UK subjects; the corresponding MDMs due to microbiome are 0.38 and 0.36. These
291 results suggest that 11.99% and 12.90% of the disparity in BMI between the Chinese and matched USA
292 and UK groups, respectively, would be eliminated if the between-group microbiome profiles were
293 equalized.

294 Significant CMD testing results show that there is at least one species playing a mediating role in the
295 disparity of BMI between Chinese and USA subjects, and Chinese and UK subjects. Figure 3b reports 15
296 species and 21 species further identified by SparseMCMC_HD, with the point and 95% confidence
297 interval (CI) estimates for their mediation effects on the regional differences of BMI between China and
298 USA, and between China and UK, respectively. Among the twelve overlapping species identified in both
299 matched datasets (Figure 3b and 3c), five species—*Anaerostipes hadrus*, *Bacteroides plebeius*,

300 *Bacteroides thetaiotaomicron*, *Bacteroides uniformis*, and *Escherichia coli*—play consistent positive
301 mediating roles in regional disparity in BMI for Chinese compared to USA subjects, and for Chinese
302 compared to UK subjects. The relative evaluation of these five species in terms of their relative
303 abundances (Figure 4a) and their associations with BMI (Figure 4b) are quite similar between two
304 independent studies: China-USA comparison and China-UK comparison, which validates their mediating
305 roles in the regional disparity on BMI. Confirming with the published studies, *B. plebeius*, *B.*
306 *thetaitotaomicron*, and *B. uniformis* belong to the same genus *Bacteroides*, and all play important roles in
307 human metabolism and have been linked with diet-induced obesity, by improving whole-body glucose
308 disposal, promoting lipid digestion and absorption, and degrading host-derived carbohydrates [35-38]. *B.*
309 *thetaitotaomicron* also possesses glycine lipid biosynthesis pathway (Figure S4). *A. hadrus*, and *E. coli*
310 also have been reported by multiple studies that they contribute to or are associated with the BMI or
311 obesity [39-41]. On the other hand, 12 species play mediating roles in BMI but with the opposite
312 directions between China-USA comparison and China-UK comparison, that reflects the distinguishing
313 characteristics between USA and UK (Figure S5). This is not surprising considering the microbial profile
314 is inherently dynamic and racially or geographically specific. Moreover, there are three and nine unique
315 species identified in the China-USA and China-UK comparisons respectively (Figures S6 and S7). Most
316 of these study-specific species have been reported being associated with BMI, obesity or metabolic
317 disorders [41-50]. Notably, *Anaerostipes hadrus*, *Fusicatenibacter saccharivorans*, *Lachnospira*
318 *pectinoschiza*, and *Roseburia inulinivorans* belong to family *Lachnospiraceae* (Figure 5d), which is
319 related to metabolic syndrome and obesity and whose controversial role has been discussed across
320 different studies [51].

321 **Results for AGP**

322 **Matched dataset.** After performing PSM, as described in Section S2, 98 Caucasians and 98 APIs are
323 matched. Figures S8 and S9 show that the matched Caucasians and APIs have very similar propensity
324 scores (SMD=0.005 for the matched subjects vs. SMD=1.033 for the raw subjects), indicating that the

325 confounding effects are well controlled. With this well-matched dataset, Figure 5a shows that the
326 Caucasian group has a significantly higher BMI (23.96 [3.92]), compared to the API group (22.38 [3.59]),
327 as observed in the other studies [52, 53].

328 **Community level results.** Caucasians and APIs have distinct microbial profiles in terms of community
329 diversity. For alpha diversity, Caucasians have higher microbial richness and evenness as measured by
330 Observed, Shannon, and Simpson diversities (p-value = 3.1×10^{-5} , 1.5×10^{-4} , and 3.9×10^{-3} ,
331 respectively. Figure S10a). For Beta diversity, Bray-Curtis dissimilarity and Jensen-Shannon divergence
332 both show that Caucasian samples have different community structures compared to API samples
333 (PERMANOVA p-value=0.0036 and 0.0012, respectively. Figure S10b).

334 **Taxon-level analysis.** The above community level results indicate that the microbiome may play a
335 mediating role in the racial diversity of BMI. To investigate this assumption, we perform the proposed
336 SparseMCMC_HD on this matched dataset. With the filtering criteria described in Section S4, 28 species
337 are included in the following taxon-level analysis.

338 We found that the ODM of BMI between Caucasians and APIs is 1.63 (Figure 5b). Microbiome plays a
339 significant role in mediating the racial disparity of BMI indicated by the test results of both OMD (p-
340 value=0.038) and CMD (p-value=0.048). The microbial manipulable disparity measure MDM is 0.12.
341 This suggests that the difference of microbiome profiles contributes to 7.4% of ODM, which would be
342 eliminated if the microbiome profiles between the Caucasians and APIs were identical.

343 We further identified 12 species playing mediating roles in the racial disparity of BMI between the
344 Caucasians and APIs (Figure 5c). Eight species (*[Ruminococcus] gnavus*, *Faecalibacterium prausnitzii*,
345 *Bacteroides uniformis*, *[Eubacterium] bifforme*, *Bacteroides fragilis*, *Prevotella copri*, *Bacteroides ovatus*,
346 *Haemophilus parainfluenzae*) mediate positively on the racial disparity of BMI, meanwhile, four species
347 (*Bifidobacterium adolescentis*, *Bacteroides plebeius*, *Parabacteroides distasonis*, *Staphylococcus aureus*)
348 play negative mediating roles. Remarkably, there are six common species *B. ovatus*, *B. plebeius*, *B.*

349 *uniformis*, *B. adolescentis*, *F. prausnitzii*, *P. distasonis*, and *P. copri* identified by both comparisons:
350 China-USA and China-UK illustrated in the previous subsection (Figure 5d). Literature reveals that all
351 identified species are associated with the BMI or obesity [41-49].

352 Collectively, the findings in the matched China vs. USA, China vs. UK, and API vs. Caucasian datasets
353 show that the microbiome is an important mediator in the regional or racial disparity of BMI and they
354 substantially shed light on how to reduce the disparity of BMI. The identified microbial agents can be
355 used as the potential therapeutic target for the treatment based on microbiota modulation in the future.

356 **Discussion**

357 The emerging evidence highlights the potential of microbiome in understanding health disparity. In this
358 paper, we proposed a mediation analytical framework, SparseMCMM_HD, to investigate the
359 microbiome's role in health disparity. Considering a health disparity framework with three components:
360 non-manipulable exposure (e.g. race or region), microbiome as mediator, and outcome, the proposed
361 SparseMCMM_HD deciphers the overall health disparity of the non-manipulable exposure on the
362 outcome into two components: MDM that would be eliminated by equalizing microbiome profiles and
363 RDM that would remain and could not be explained through microbiome. Remarkably, MDM paves a
364 viable path towards reduction of health disparity with microbial modulation. Similar to SparseMCMM,
365 SparseMCMM_HD can be used to identify the signature causal microbes and examine whether the
366 overall or component-wise MDM is significantly non-zero.

367 It is vital to control confounding effects beforehand in the real data analysis to satisfy the identifiability
368 assumptions of the proposed SparseMCMM_HD. In three BMI applications, we first employed PSM to
369 remove the confounding effects by selecting matched subsets in which the distributions of confounders
370 were notably comparable between two exposure groups, and then performed the proposed
371 SparseMCMM_HD framework. The utilization of SparseMCMM_HD in two datasets, the
372 curatedMetagenomicData 3.4.2 package and the AGP dataset, depicts an explicit causal path among

373 region or race, microbiome, and BMI. These findings confirm not only that microbiome is differentially
374 distributed across races or regions, but also that the differential microbiome profile contributes to the
375 disparities in BMI across races or regions. The identified microbial signatures potentially aid in
376 developing personalized medication or nutrition to reduce obesity disparity.

377 It is not surprising that the proportion of disparities in BMI explained by the microbiome profiles is not
378 large (~10%) in all three applications, due to the heritable and polygenic nature of BMI [54, 55]. Further
379 investigations to integrate the microbiome profile and genetic factors are necessary to better understand
380 disparity in BMI. However, we here emphasize that the proposed SparseMCMM_HD is a rigorous and
381 validated causal mediation framework and has preeminent potential to identify the microbiome's roles in
382 much broader health disparity studies.

383 Recently, several other microbial mediation methods have been proposed, such as CMM [56], MedTest
384 [57], Zhang, et al. [58], LDM-med [59], and MarZIC [60], in a typical three-factor (manipulable
385 exposure, microbiome as mediator, and outcome) study design. Considering distinct model assumptions
386 and characteristics, a few recent benchmark studies [12, 56-60] show that there is no method performing
387 consistently and accurately better than others in all circumstances. However, since the assumptions for
388 model identification in health disparity are weaker than those for the causal mediation effects in the
389 manipulable exposure-mediator-outcome framework [23], it is expected that the idea of how the proposed
390 SparseMCMM_HD framework rigorously defines, quantifies, and tests health disparity measures as an
391 extension of SparseMCMM [12] can provide insight into extending these available mediation models to
392 investigate the microbiome's role in health disparity. Then, a useful path forward will be to mutually
393 employ these multiple and complimentary methods to better characterize the microbiome's role in health
394 disparity by capitalizing their distinct assumptions and strengths.

395 Our study has several limitations. First, similar to discussions in SparseMCMM [12], SparseMCMM_HD
396 takes microbiome data at a fixed time point into the proposed frame and is limited to accommodate the
397 dynamic nature of microbiome. Second, the proposed SparseMCMM_HD currently deals with disparity

398 in a continuous outcome. Given the fact that multiple binary or categorical outcomes are
399 disproportionately prevalent across races or regions [61-63], it will be worthwhile to extend the current
400 framework to handle categorical outcomes. Third, microbiome studies typically characterize both
401 taxonomic and functional profiles of microbial communities. Functional profile is generally thought to be
402 more closely linked with human health and disease. Identifying the role of microbiome in terms of gene
403 function in health disparity is of high practical value [64].

404 **Conclusions**

405 This paper elucidates the role of microbiome in health disparity by providing a causal mediation analytic
406 framework for investigating the relationship among race or region, microbiome, and outcome under the
407 counterfactual framework. The proposed SparseMCMM_HD framework is a useful tool to investigate the
408 underlying biological mechanism of health disparity and disentangles the substantial contributions of
409 microbiome to health disparity. The applications of SparseMCMM_HD in the disparity of BMI across
410 races and regions uncover the microbial mediating roles in reducing the disparities of BMI and improving
411 health equality.

412

413 **List of abbreviations**

414 AGP: American gut project; API: Asian or Pacific Islander; BMI: body mass index; MDM: manipulable
415 disparity measure; ODM: overall disparity measure; PERMANOVA: permutational multivariate analysis
416 of variance; PSM: propensity score matching; RDM: residual disparity measure; SparseMCMM: sparse
417 microbial causal mediation model; SparseMCMM_HD: SparseMCMM for health disparity; SMD:
418 standardized mean difference; UK: United Kingdom; USA: United States.

419

420 **Declarations**

421 **Ethics approval and consent to participate**

422 All utilized microbiome datasets are publicly available. No ethics approval or consent to participate was
423 required for this study.

424 **Consent for publication**

425 Not applicable: All utilized microbiome datasets are publicly available. No consent for publication was
426 required for this study.

427 **Availability of data and materials**

428 All relevant datasets are publicly available. The data used in investigations of the regional difference of
429 BMI in the China group compared to the United States (USA) and United Kingdom (UK) groups can be
430 downloaded from the curatedMetagenomicData 3.4.2 package [28]. The data used in investigations of the
431 racial difference in BMI between Caucasians and Asian or Pacific Islanders are from the American Gut
432 Project. Their raw data and metadata are publicly available on the FTP website
433 (<ftp://ftp.microbio.me/AmericanGut/>). Version 07/29/2016 is used in our analyses.

434 SparseMCMM R package is available at <https://github.com/chanw0/SparseMCMM>. The interactive web
435 app for the proposed SparseMCMM_HD framework is available at
436 <https://chanw0.shinyapps.io/sparsemcm HD/>.

437 **Competing interests**

438 The authors declare that they have no competing interests.

439 **Funding**

440 The study was supported in part by grant number U54MD000538 from the National Institutes of Health
441 (NIH) National Institute on Minority Health and Health Disparities, and grant number P20CA252728

442 from the National Cancer Institute. The contents of this publication are solely the responsibility of the
443 authors and do not necessarily represent the official views of the NIH.

444 **Authors' contributions**

445 C.W. developed the microbial causal mediation analytic framework, performed data analyses, and wrote
446 the manuscript. J.A., T.T., R.B.H., and S.S.Y. contributed to the biological insights and interpretation, and
447 to manuscript writing. H.L. contributed to the methodological ideas for the proposed framework,
448 simulations, real data analyses, and manuscript writing. All authors read and approved the final
449 manuscript.

450 **Acknowledgements**

451 Not applicable.

452

453 **References**

- 454 1. Hu J, Koh H, He L, Liu M, Blaser MJ, Li H: **A two-stage microbial association mapping framework**
455 **with advanced FDR control.** *Microbiome* 2018, **6**(1):1-16.
- 456 2. Gilbert JA, Quinn RA, Debelius J, Xu ZZ, Morton J, Garg N, Jansson JK, Dorrestein PC, Knight R:
457 **Microbiome-wide association studies link dynamic microbial consortia to disease.** *Nature*
458 2016, **535**(7610):94-103.
- 459 3. Koh H, Livanos AE, Blaser MJ, Li H: **A highly adaptive microbiome-based association test for**
460 **survival traits.** *BMC genomics* 2018, **19**(1):1-13.
- 461 4. Koh H, Blaser MJ, Li H: **A powerful microbiome-based association test and a microbial taxa**
462 **discovery framework for comprehensive association mapping.** *Microbiome* 2017, **5**(1):1-15.
- 463 5. Ahn J, Sinha R, Pei Z, Dominianni C, Wu J, Shi J, Goedert JJ, Hayes RB, Yang L: **Human gut**
464 **microbiome and risk for colorectal cancer.** *Journal of the National Cancer Institute* 2013,
465 **105**(24):1907-1911.
- 466 6. Kostic AD, Xavier RJ, Gevers D: **The microbiome in inflammatory bowel disease: current status**
467 **and the future ahead.** *Gastroenterology* 2014, **146**(6):1489-1499.
- 468 7. Hoffmann AR, Proctor L, Surette M, Suchodolski J: **The microbiome: the trillions of**
469 **microorganisms that maintain health and cause disease in humans and companion animals.**
470 *Veterinary Pathology* 2016, **53**(1):10-21.
- 471 8. Kelly TN, Bazzano LA, Ajami NJ, He H, Zhao J, Petrosino JF, Correa A, He J: **Gut microbiome**
472 **associates with lifetime cardiovascular disease risk profile among bogalusa heart study**
473 **participants.** *Circulation research* 2016, **119**(8):956-964.
- 474 9. Gilbert JA, Blaser MJ, Caporaso JG, Jansson JK, Lynch SV, Knight R: **Current understanding of the**
475 **human microbiome.** *Nature medicine* 2018, **24**(4):392-400.

- 476 10. Fattorusso A, Di Genova L, Dell'Isola GB, Mencaroni E, Esposito S: **Autism spectrum disorders**
477 **and the gut microbiota.** *Nutrients* 2019, **11**(3):521.
- 478 11. Integrative H, Proctor LM, Creasy HH, Fettweis JM, Lloyd-Price J, Mahurkar A, Zhou W, Buck GA,
479 Snyder MP, Strauss III JF: **The integrative human microbiome project.** *Nature* 2019,
480 **569**(7758):641-648.
- 481 12. Wang C, Hu J, Blaser MJ, Li H: **Estimating and testing the microbial causal mediation effect with**
482 **high-dimensional and compositional microbiome data.** *Bioinformatics (Oxford, England)* 2020,
483 **36**(2):347-355.
- 484 13. Gupta VK, Paul S, Dutta C: **Geography, ethnicity or subsistence-specific variations in human**
485 **microbiome composition and diversity.** *Frontiers in microbiology* 2017, **8**:1162.
- 486 14. Dehingia M, Adak A, Khan MR: **Ethnicity-influenced microbiota: a future healthcare**
487 **perspective.** *Trends in microbiology* 2019, **27**(3):191-193.
- 488 15. Findley K, Williams DR, Grice EA, Bonham VL: **Health disparities and the microbiome.** *Trends in*
489 *microbiology* 2016, **24**(11):847-850.
- 490 16. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, Costea PI, Godneva A,
491 Kalka IN, Bar N: **Environment dominates over host genetics in shaping human gut microbiota.**
492 *Nature* 2018, **555**(7695):210-215.
- 493 17. Schulz MD, Atay Ç, Heringer J, Romrig FK, Schwitalla S, Aydin B, Ziegler PK, Varga J, Reindl W,
494 Pommerenke C: **High-fat-diet-mediated dysbiosis promotes intestinal carcinogenesis**
495 **independently of obesity.** *Nature* 2014, **514**(7523):508-512.
- 496 18. Zhang X, Zhao Y, Zhang M, Pang X, Xu J, Kang C, Li M, Zhang C, Zhang Z, Zhang Y: **Structural**
497 **changes of gut microbiota during berberine-mediated prevention of obesity and insulin**
498 **resistance in high-fat diet-fed rats.** 2012.
- 499 19. Cox LM, Blaser MJ: **Antibiotics in early life and obesity.** *Nature Reviews Endocrinology* 2015,
500 **11**(3):182-190.
- 501 20. Taur Y, Pamer EG: **Microbiome mediation of infections in the cancer setting.** *Genome medicine*
502 2016, **8**(1):1-7.
- 503 21. Lv B-M, Quan Y, Zhang H-Y: **Causal inference in microbiome medicine: Principles and**
504 **applications.** *Trends in microbiology* 2021, **29**(8):736-746.
- 505 22. Ananthakrishnan AN, Bernstein CN, Iliopoulos D, Macpherson A, Neurath MF, Ali RAR, Vavricka
506 SR, Fiocchi C: **Environmental triggers in IBD: a review of progress and evidence.** *Nature Reviews*
507 *Gastroenterology & Hepatology* 2018, **15**(1):39-49.
- 508 23. VanderWeele TJ, Robinson WR: **On causal interpretation of race in regressions adjusting for**
509 **confounding and mediating variables.** *Epidemiology (Cambridge, Mass)* 2014, **25**(4):473.
- 510 24. Haththotuwa RN, Wijeyaratne CN, Senarath U: **Worldwide epidemic of obesity.** In: *Obesity and*
511 *obstetrics.* Elsevier; 2020: 3-8.
- 512 25. Gaulke CA, Sharpton TJ: **The influence of ethnicity and geography on human gut microbiome**
513 **composition.** *Nature medicine* 2018, **24**(10):1495-1496.
- 514 26. Deschasaux M, Bouter KE, Prodan A, Levin E, Groen AK, Herrema H, Tremaroli V, Bakker GJ,
515 Attaye I, Pinto-Sietsma S-J: **Depicting the composition of gut microbiota in a population with**
516 **varied ethnic origins but shared geography.** *Nature medicine* 2018, **24**(10):1526-1531.
- 517 27. He Y, Wu W, Zheng H-M, Li P, McDonald D, Sheng H-F, Chen M-X, Chen Z-H, Ji G-Y, Zheng Z-D-X:
518 **Regional variation limits applications of healthy gut microbiome reference ranges and disease**
519 **models.** *Nature medicine* 2018, **24**(10):1532-1535.
- 520 28. Pasolli E, Schiffer L, Manghi P, Renson A, Obenchain V, Truong DT, Beghini F, Malik F, Ramos M,
521 Dowd JB: **Accessible, curated metagenomic data through ExperimentHub.** *Nature methods*
522 2017, **14**(11):1023-1024.

- 523 29. Hijazi RH, Jernigan RW: **Modelling compositional data using Dirichlet regression models.** *Journal of Applied Probability & Statistics* 2009, **4**(1):77-91.
- 524
- 525 30. Naimi AI, Schnitzer ME, Moodie EE, Bodnar LM: **Mediation analysis for health disparities**
- 526 **research.** *American journal of epidemiology* 2016, **184**(4):315-324.
- 527 31. Rosenbaum PR, Rubin DB: **The central role of the propensity score in observational studies for**
- 528 **causal effects.** *Biometrika* 1983, **70**(1):41-55.
- 529 32. Austin PC: **An introduction to propensity score methods for reducing the effects of**
- 530 **confounding in observational studies.** *Multivariate behavioral research* 2011, **46**(3):399-424.
- 531 33. Liu M, Koh H, Kurtz ZD, Battaglia T, PeBenito A, Li H, Nazzari L, Blaser MJ: **Oxalobacter**
- 532 **formigenes-associated host features and microbial community structures examined using the**
- 533 **American Gut Project.** *Microbiome* 2017, **5**(1):1-17.
- 534 34. Anderson MJ: **Permutational multivariate analysis of variance (PERMANOVA).** *Wiley statsref:*
- 535 *statistics reference online* 2014:1-15.
- 536 35. López-Almela I, Romaní-Pérez M, Bullich-Vilarrubias C, Benítez-Páez A, Gómez Del Pulgar EM,
- 537 Francés R, Liebisch G, Sanz Y: **Bacteroides uniformis combined with fiber amplifies metabolic**
- 538 **and immune benefits in obese mice.** *Gut Microbes* 2021, **13**(1):1-20.
- 539 36. Cho S-H, Cho Y-J, Park J-H: **The human symbiont Bacteroides thetaiotaomicron promotes diet-**
- 540 **induced obesity by regulating host lipid metabolism.** *Journal of Microbiology* 2022, **60**(1):118-
- 541 127.
- 542 37. Hehemann J-H, Kelly AG, Pudlo NA, Martens EC, Boraston AB: **Bacteria of the human gut**
- 543 **microbiome catabolize red seaweed glycans with carbohydrate-active enzyme updates from**
- 544 **extrinsic microbes.** *Proceedings of the National Academy of Sciences* 2012, **109**(48):19786-
- 545 19791.
- 546 38. Thomas F, Hehemann J-H, Rebuffet E, Czejek M, Michel G: **Environmental and gut**
- 547 **bacteroidetes: the food connection.** *Frontiers in microbiology* 2011, **2**:93.
- 548 39. Holmes ZC, Silverman JD, Dressman HK, Wei Z, Dallow EP, Armstrong SC, Seed PC, Rawls JF,
- 549 David LA: **Short-chain fatty acid production by gut microbiota from children with obesity**
- 550 **differs according to prebiotic choice and bacterial community composition.** *MBio* 2020,
- 551 **11**(4):e00914-00920.
- 552 40. Million á, Angelakis E, Maraninchi M, Henry M, Giorgi R, Valero R, Vialettes B, Raoult D:
- 553 **Correlation between body mass index and gut concentrations of Lactobacillus reuteri,**
- 554 **Bifidobacterium animalis, Methanobrevibacter smithii and Escherichia coli.** *International*
- 555 *journal of obesity* 2013, **37**(11):1460-1466.
- 556 41. Ignacio A, Fernandes M, Rodrigues V, Groppo F, Cardoso A, Avila-Campos M, Nakano V:
- 557 **Correlation between body mass index and faecal microbiota from children.** *Clinical*
- 558 *microbiology and infection* 2016, **22**(3):258. e251-258. e258.
- 559 42. Journey EK, Ortega-Santos CP, Bruening M, Whisner CM: **Changes in weight status and the**
- 560 **intestinal microbiota among college freshman, aged 18 years.** *Journal of Adolescent Health*
- 561 2020, **66**(2):166-171.
- 562 43. Palmas V, Pisanu S, Madau V, Casula E, Deledda A, Cusano R, Uva P, Vascellari S, Loviselli A,
- 563 Manzin A: **Gut microbiota markers associated with obesity and overweight in Italian adults.**
- 564 *Scientific reports* 2021, **11**(1):1-14.
- 565 44. Chatel J-M, Maioli TU, Borrás-Nogues E, Barbosa SC, Martins VD, Torres L, Azevedo VADC:
- 566 **Possible benefits of Faecalibacterium prausnitzii for obesity-associated gut disorders.** *Frontiers*
- 567 *in Pharmacology* 2021:2869.
- 568 45. Duan M, Wang Y, Zhang Q, Zou R, Guo M, Zheng H: **Characteristics of gut microbiota in people**
- 569 **with obesity.** *Plos one* 2021, **16**(8):e0255446.

- 570 46. Li Y, Yang Y, Wang J, Cai P, Li M, Tang X, Tan Y, Wang Y, Zhang F, Wen X: **Bacteroides ovatus-**
571 **mediated CD27– MAIT cell activation is associated with obesity-related T2D progression.**
572 *Cellular & Molecular Immunology* 2022:1-14.
- 573 47. Assmann TS, Cuevas-Sierra A, Riezu-Boj JI, Milagro FI, Martínez JA: **Comprehensive analysis**
574 **reveals novel interactions between circulating MicroRNAs and gut microbiota composition in**
575 **human obesity.** *International journal of molecular sciences* 2020, **21**(24):9509.
- 576 48. Befus M, Lowy FD, Miko BA, Mukherjee DV, Herzig CT, Larson EL: **Obesity as a determinant of**
577 **Staphylococcus aureus colonization among inmates in maximum-security prisons in New York**
578 **State.** *American journal of epidemiology* 2015, **182**(6):494-502.
- 579 49. Yan H, Qin Q, Chen J, Yan S, Li T, Gao X, Yang Y, Li A, Ding S: **Gut microbiome alterations in**
580 **patients with visceral obesity based on quantitative computed tomography.** *Frontiers in*
581 *Cellular and Infection Microbiology* 2022:1451.
- 582 50. Yang M, Bose S, Lim S, Seo J, Shin J, Lee D, Chung W-H, Song E-J, Nam Y-D, Kim H: **Beneficial**
583 **effects of newly isolated Akkermansia muciniphila strains from the human gut on obesity and**
584 **metabolic dysregulation.** *Microorganisms* 2020, **8**(9):1413.
- 585 51. Vacca M, Celano G, Calabrese FM, Portincasa P, Gobetti M, De Angelis M: **The controversial**
586 **role of human gut lachnospiraceae.** *Microorganisms* 2020, **8**(4):573.
- 587 52. McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, Aksenov AA, Behsz B,
588 Brennan C, Chen Y: **American gut: an open platform for citizen science microbiome research.**
589 *Msystems* 2018, **3**(3):e00031-00018.
- 590 53. Obana KK, Davis J: **Racial disparities in the prevalence of arthritis among native Hawaiians and**
591 **Pacific Islanders, Whites, and Asians.** *Hawai'i Journal of Medicine & Public Health* 2016,
592 **75**(6):155.
- 593 54. Bouchard C: **Genetics of obesity: what we have learned over decades of research.** *Obesity*
594 2021, **29**(5):802-820.
- 595 55. Loos RJ, Yeo GS: **The genetics of obesity: from discovery to biology.** *Nature Reviews Genetics*
596 2022, **23**(2):120-133.
- 597 56. Sohn MB, Li H: **Compositional mediation analysis for microbiome studies.** *The Annals of Applied*
598 *Statistics* 2019, **13**(1):661-681.
- 599 57. Zhang J, Wei Z, Chen J: **A distance-based approach for testing the mediation effect of the**
600 **human microbiome.** *Bioinformatics* 2018, **34**(11):1875-1883.
- 601 58. Zhang H, Chen J, Feng Y, Wang C, Li H, Liu L: **Mediation effect selection in high-dimensional and**
602 **compositional microbiome data.** *Statistics in medicine* 2021, **40**(4):885-896.
- 603 59. Yue Y, Hu Y: **Testing Mediation Effects in High-Dimensional Microbiome Data with False**
604 **Discovery Rate Control.** 2021.
- 605 60. Wu Q, O'malley J, Datta S, Gharaibeh RZ, Jobin C, Karagas MR, Coker MO, Hoen AG, Christensen
606 BC, Madan JC: **MarZIC: A Marginal Mediation Model for Zero-Inflated Compositional**
607 **Mediators with Applications to Microbiome Data.** *Genes* 2022, **13**(6):1049.
- 608 61. Royston KJ, Adedokun B, Olopade OI: **Race, the microbiome and colorectal cancer.** *World*
609 *Journal of Gastrointestinal Oncology* 2019, **11**(10):773.
- 610 62. Siddharth S, Sharma D: **Racial disparity and triple-negative breast cancer in African-American**
611 **women: a multifaceted affair between obesity, biology, and socioeconomic determinants.**
612 *Cancers* 2018, **10**(12):514.
- 613 63. Johnson JR, Kittles RA: **Genetic ancestry and racial differences in prostate tumours.** *Nature*
614 *Reviews Urology* 2022, **19**(3):133-134.
- 615 64. Tian L, Wang X-W, Wu A-K, Fan Y, Friedman J, Dahlin A, Waldor MK, Weinstock GM, Weiss ST,
616 Liu Y-Y: **Deciphering functional redundancy in the human microbiome.** *Nature communications*
617 2020, **11**(1):1-11.

618 65. Foster ZS, Sharpton TJ, Grünwald NJ: **Metacoder: An R package for visualization and**
619 **manipulation of community taxonomic diversity data.** *PLoS computational biology* 2017,
620 **13(2):e1005404.**

621

622

623 **Figure 1.** Microbiome (M) may play a mediating role in the health disparity of the continuous outcome
624 (Y) between two categories of a non-manipulable exposure group (e.g. race or region) (R). We aim to
625 investigate how much disparity of the outcome Y can be reduced by manipulating microbiome profiles.

626

627 **Figure 2.** Association analyses in two matched datasets from the curatedMetagenomicData package [28].
628 a Violin plots of BMI in matched Chinese vs. USA subjects. b Violin plots of alpha diversities (Observed,
629 Shannon, and Simpson indices) in matched Chinese vs. USA samples. c PCoA plots using Bray–Curtis
630 dissimilarity and Jensen–Shannon divergence in matched Chinese and USA samples. d Violin plots of
631 BMI in matched Chinese vs. UK subjects. e Violin plots of alpha diversities (Observed, Shannon, and
632 Simpson indices) in matched Chinese and UK samples. f PCoA plots using Bray–Curtis dissimilarity and
633 Jensen–Shannon divergence in matched Chinese vs. UK samples.

634

635 **Figure 3.** Health disparity analyses in two matched datasets from the curatedMetagenomicData package
636 [28]. a Manipulable disparity measure (MDM) and residual disparity measure (RDM) of BMI in the
637 China-USA comparison and China-UK comparison, respectively. b Component-wise point and 95% CI
638 estimates of MDM_j for the identified species that have mediation effects on the differences of BMI
639 between matched Chinese vs. USA subjects and between matched Chinese vs. UK subjects, respectively.
640 95% CI estimates of MDM_j were calculated by bootstrapping procedure, and the number of bootstrapping
641 is 50. c Venn diagram to show the relationship of the species playing mediation effects in the disparity of
642 BMI among China-USA, China-UK, and API-Caucasian comparisons. API: Asian or Pacific Islander.

643

644 **Figure 4.** Five species who play positive mediation roles in the disparity of BMI in both China-USA and
645 China-UK comparisons. a Violin plots illustrating the relative abundances of these five identified species
646 in the matched Chinese and USA samples, and the matched Chinese and UK samples, respectively. b
647 Scatterplots of BMI and the relative abundances of these five identified species in the matched Chinese
648 and USA subjects, and the matched Chinese and UK subjects, respectively.

649

650 **Figure 5.** Health disparity analyses in the matched APIs and Caucasians from the AGP dataset. a Violin
651 plots of BMI in the matched APIs and Caucasians from the AGP dataset. b MDM and RDM of BMI in
652 the API- Caucasian comparison. c Component-wise point and 95% CI estimates of MDM_j for the
653 identified species that have mediation effects on the differences of BMI between matched APIs and
654 Caucasians from the AGP dataset. 95% CI estimates of MDM_j were calculated by bootstrapping
655 procedure, and the number of bootstrapping is 50. d The taxonomic relationship of the species playing
656 mediation effects in the disparity of BMI among China-USA, China-UK, and API-Caucasian
657 comparisons. The tree figure was generated by Metacoder [65]. From the outer to the center, taxonomic
658 ranks are species, genus, family, order, class, phylum, and kingdom (Bacteria), respectively. For each

659 species, color represents the number of comparisons that identify it among China-USA, China-UK, and
660 API-Caucasian comparisons. APIs: Asian or Pacific Islanders.

661

662 **Additional material**

663

664 **Additional file 1: Section S1** Derivations for MDM and RDM expressions.

665 **Section S2** Propensity score matching (PSM).

666 **Section S3** Metadata curation in the AGP.

667 **Section S4** Taxon-level alignment.

668

669 **Additional file 2: Figure S1.** Flowcharts for data pre-processing in the AGP dataset. a Pre-processing for
670 all covariates. b The sample breakdown for the disparity analysis.

671 **Figure S2.** Plots of standardized mean differences before and after propensity score matching for the
672 datasets from the curatedMetagenomicData package [28]. a Comparison between Chinese and USA
673 subjects. b Comparison between Chinese and UK subjects.

674 **Figure S3.** Histogram plots of propensity score before and after propensity score matching for the
675 datasets from the curatedMetagenomicData package [28]. a Comparison between Chinese and USA
676 subjects. b Comparison between Chinese and UK subjects.

677 **Figure S4.** Glycine lipid biosynthesis pathway generated based on MetaCyc database
678 (<https://metacyc.org/>). The gene from *B.thetaiotaomicro* is located in an operon together with a second
679 gene, *glsA*, which encodes the second enzyme of the pathway, an *O*-acyltransferase that forms the
680 diacylated compound.

681 **Figure S5.** The species with opposite mediation directions in the disparity of BMI between China-USA
682 and China-UK comparisons. a Violin plots illustrating the relative abundances of these identified species
683 in the matched Chinese and USA samples, and the matched Chinese and UK samples, respectively. b
684 Scatterplots of BMI and the relative abundances of these identified species in the matched Chinese and
685 USA samples, and the matched Chinese and UK samples, respectively.

686 **Figure S6.** The species playing mediation roles in the disparity of BMI in the comparison between
687 Chinese and USA subjects only. a Violin plots illustrating the relative abundances of these identified
688 species in the matched Chinese and USA samples. b Scatterplots of BMI and the relative abundances of
689 these identified species in the matched Chinese and USA samples.

690 **Figure S7.** The species playing mediating roles in the disparity of BMI in the comparison between
691 Chinese and UK subjects only. a Violin plots illustrating the relative abundances of these identified
692 species in the matched Chinese and UK samples. b Scatterplots of BMI and the relative abundances of
693 these identified species in the matched Chinese and UK samples.

694 **Figure S8.** Plots of standardized mean differences before and after propensity score matching for the
695 comparison between the API and Caucasian samples from the AGP dataset. API: Asian or Pacific
696 Islander.

697 **Figure S9.** Histogram plots of propensity score before and after propensity score matching for the
698 comparison between the API and Caucasian samples from the AGP dataset. API: Asian or Pacific
699 Islander.

700 **Figure S10.** Association analyses in the AGP dataset. a Violin plots of alpha diversities including
701 Observed, Shannon, and Simpson indices in the matched API and Caucasian samples. b PCoA plots using
702 Bray–Curtis dissimilarity and Jensen–Shannon divergence in the matched API and Caucasian samples.
703 API: Asian or Pacific Islander.

Figures

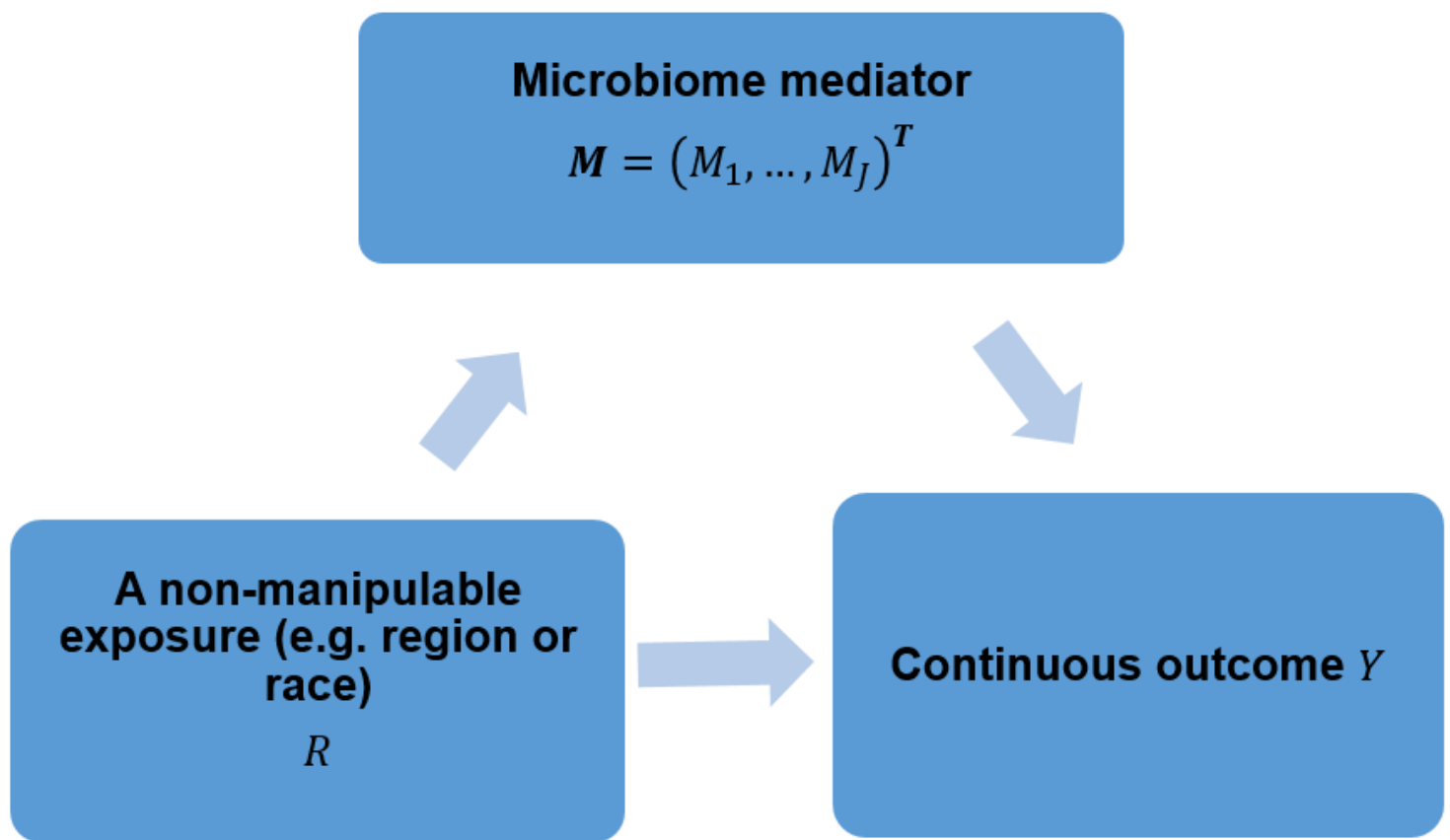


Figure 1

Microbiome (M) may play a mediating role in the health disparity of the continuous outcome (Y) between two categories of a non-manipulable exposure group (e.g. race or region) (R). We aim to investigate how much disparity of the outcome Y can be reduced by manipulating microbiome profiles.

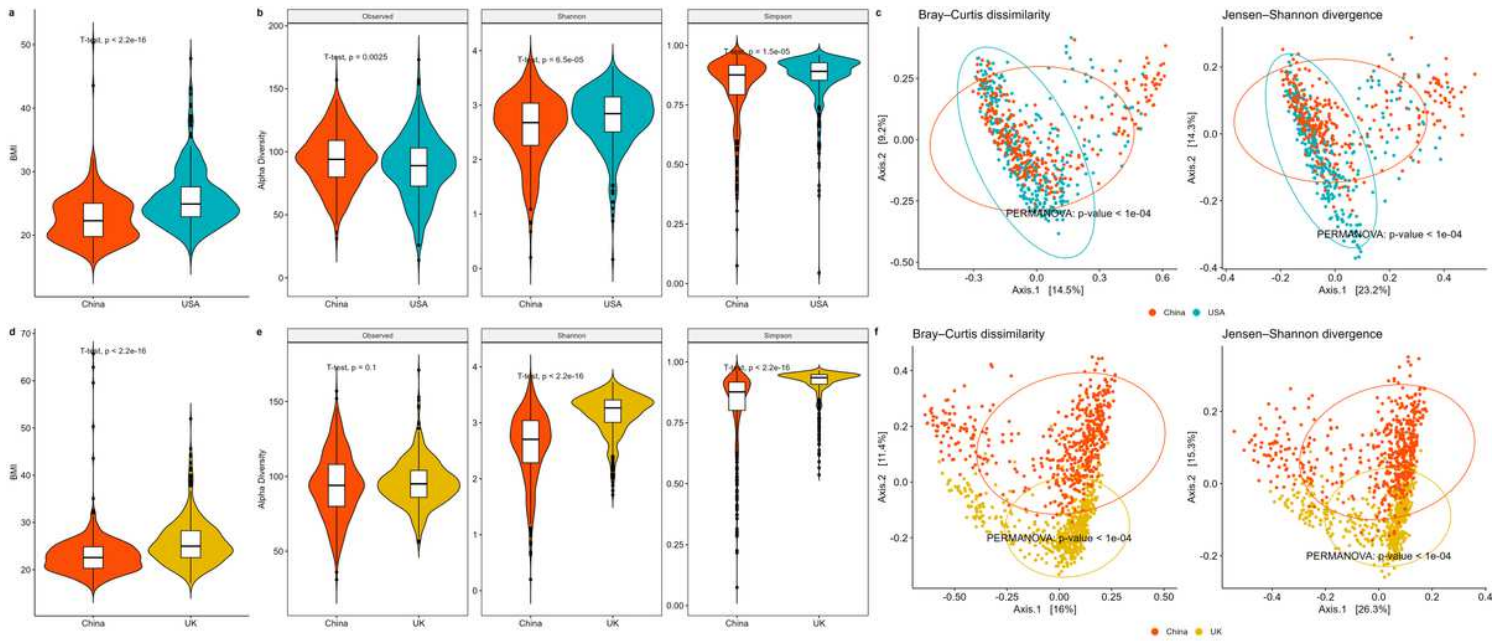


Figure 2

Association analyses in two matched datasets from the curatedMetagenomicData package [28]. a Violin plots of BMI in matched Chinese vs. USA subjects. b Violin plots of alpha diversities (Observed, Shannon, and Simpson indices) in matched Chinese vs. USA samples. c PCoA plots using Bray-Curtis dissimilarity and Jensen-Shannon divergence in matched Chinese and USA samples. d Violin plots of BMI in matched Chinese vs. UK subjects. e Violin plots of alpha diversities (Observed, Shannon, and Simpson indices) in matched Chinese and UK samples. f PCoA plots using Bray-Curtis dissimilarity and Jensen-Shannon divergence in matched Chinese vs. UK samples.

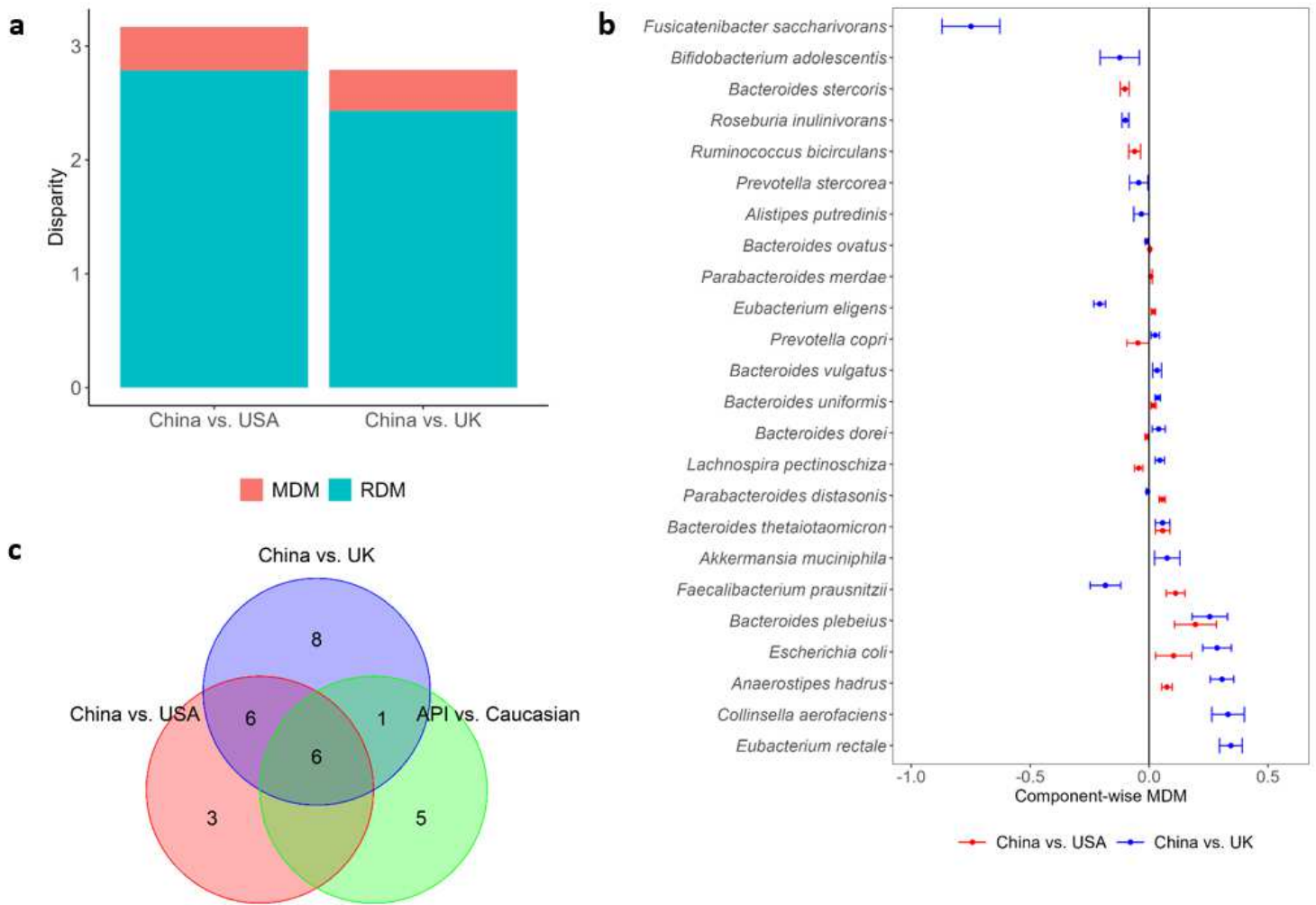


Figure 3

Health disparity analyses in two matched datasets from the curatedMetagenomicData package [28]. a Manipulable disparity measure (MDM) and residual disparity measure (RDM) of BMI in the China-USA comparison and China-UK comparison, respectively. b Component-wise point and 95% CI estimates of MDM_j for the identified species that have mediation effects on the differences of BMI between matched Chinese vs. USA subjects and between matched Chinese vs. UK subjects, respectively. 95% CI estimates of MDM_j were calculated by bootstrapping procedure, and the number of bootstrapping is 50. c Venn diagram to show the relationship of the species playing mediation effects in the disparity of BMI among China-USA, China-UK, and API-Caucasian comparisons. API: Asian or Pacific Islander.

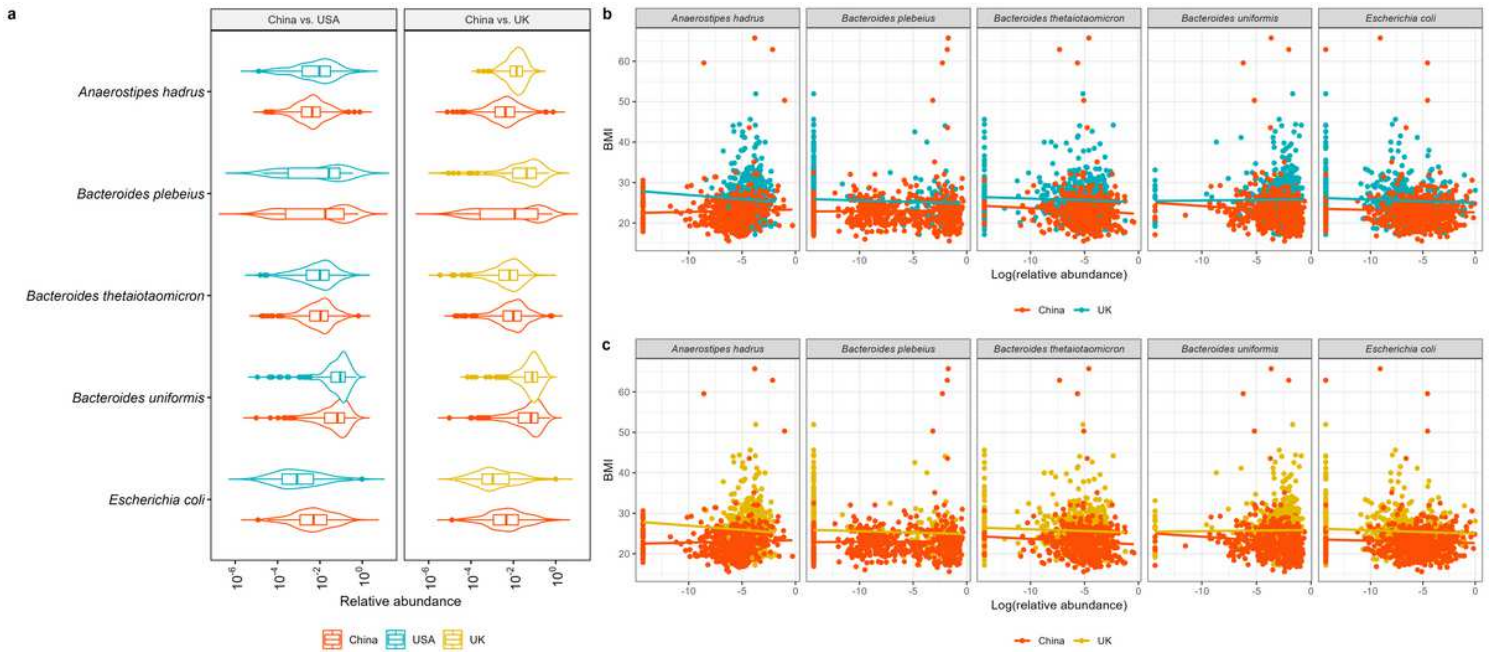


Figure 4

Five species who play positive mediation roles in the disparity of BMI in both China-USA and China-UK comparisons. **a** Violin plots illustrating the relative abundances of these five identified species in the matched Chinese and USA samples, and the matched Chinese and UK samples, respectively. **b** Scatterplots of BMI and the relative abundances of these five identified species in the matched Chinese and USA subjects, and the matched Chinese and UK subjects, respectively.

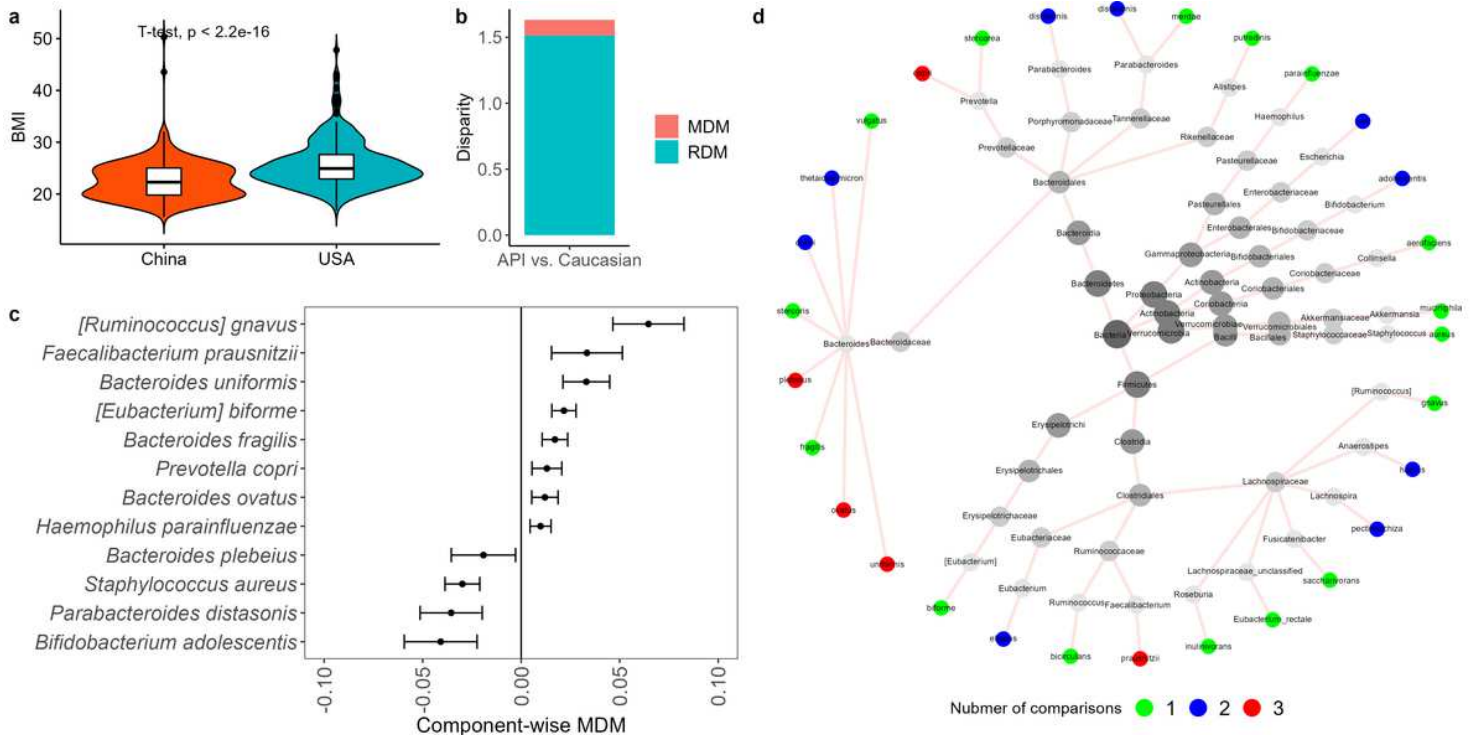


Figure 5

Health disparity analyses in the matched APIs and Caucasians from the AGP dataset. a Violin plots of BMI in the matched APIs and Caucasians from the AGP dataset. b MDM and RDM of BMI in the API-Caucasian comparison. c Component-wise point and 95% CI estimates of MDM_j for the identified species that have mediation effects on the differences of BMI between matched APIs and Caucasians from the AGP dataset. 95% CI estimates of MDM_j were calculated by bootstrapping procedure, and the number of bootstrapping is 50. d The taxonomic relationship of the species playing mediation effects in the disparity of BMI among China-USA, China-UK, and API-Caucasian comparisons. The tree figure was generated by Metacoder [65]. From the outer to the center, taxonomic ranks are species, genus, family, order, class, phylum, and kingdom (Bacteria), respectively. For each species, color represents the number of comparisons that identify it among China-USA, China-UK, and API-Caucasian comparisons. APIs: Asian or Pacific Islanders.

Supplementary Files

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

- [Additionalfile1.docx](#)
- [Additionalfile2.docx](#)