

# Looking for a Signal in the Noise: Revisiting Obesity and the Microbiome

Running Title: The Human Microbiome and Obesity

Marc A Sze and Patrick D Schloss<sup>†</sup>

Contributions: Both authors contributed to the planning, design, execution, interpretation, and writing of the analyses.

<sup>†</sup> To whom correspondence should be addressed: [pschloss@umich.edu](mailto:pschloss@umich.edu)

Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI

## 1 **Abstract**

2 Two recent studies have re-analyzed published data and found that when datasets are  
3 analyzed independently there was limited support for the widely accepted hypothesis that  
4 changes in the microbiome are associated with obesity. This hypothesis was reconsidered  
5 by increasing the number of datasets and pooling the results across the individual datasets.  
6 The Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA)  
7 guidelines were applied to identify 10 studies for an updated and more synthetic analysis.  
8 Alpha diversity metrics and the relative risk of obesity based on those metrics were used to  
9 identify a limited number of significant associations with obesity; however, when the results  
10 of the studies were pooled using a random effects model significant associations were  
11 observed between Shannon diversity, number of observed OTUs, and Shannon evenness  
12 and obesity status. They were not observed for the ratio of *Bacteroidetes* and *Firmicutes*  
13 or their individual relative abundances. Although these tests yielded small P-values, the  
14 difference between the Shannon diversity index of non-obese and obese individuals was  
15 2.07%. A power analysis demonstrated that only one of the studies had sufficient power to  
16 detect a 5% difference in diversity. When Random Forest machine learning models were  
17 trained on one dataset and then tested using the other 9 datasets, the median accuracy  
18 varied between 33.01 and 64.77% (median=56.68%). Although there was support for a  
19 relationship between the microbial communities found in human feces and obesity status,  
20 this association was relatively weak and its detection is confounded by large interpersonal  
21 variation and insufficient sample sizes.

## 22 **Importance**

23 As interest in the human microbiome grows there is an increasing number of studies that  
24 can be used to test numerous hypotheses across human populations. The hypothesis

25 that variation in the gut microbiota can explain or be used to predict obesity status has  
26 received considerable attention and is frequently mentioned as an example for the role of  
27 the microbiome in human health. Here we assess this hypothesis using ten independent  
28 studies and find that although there is an association, it is smaller than can be detected  
29 by most microbiome studies. Furthermore, we directly tested the ability to predict obesity  
30 status based on the composition of an individual's microbiome and find that the median  
31 classification accuracy is between 33.01 and 64.77%. This type of analysis can be used to  
32 design future studies and expanded to explore other hypotheses.

## 33 Introduction

34 Obesity is a growing health concern with approximately 20% of the youth (aged 2-19) in  
35 the United States classified as either overweight or obese (1). This number increases  
36 to approximately 35% in adults (aged 20 or older) and these statistics have seen little  
37 change since 2003 (1). Traditionally, the body mass index (BMI) has been used to classify  
38 individuals as non-obese or obese (2). Recently, there has been increased interest in  
39 the role of the microbiome in modulating obesity (3, 4). If the microbiome does affect  
40 obesity status, then manipulating the microbiome could have a significant role in the future  
41 treatment of obesity and in helping to stem the current epidemic.

42 There have been several studies that report observing a link between the composition  
43 of microbiome and obesity in animal models and in humans. The first such study used  
44 genetically obese mice and observed the ratio of the relative abundances of *Bacteroidetes*  
45 to *Firmicutes* (B:F) was lower in obese mice than lean mice (5). Translation of this result  
46 to humans by the same researchers did not observe this effect, but did find that obese  
47 individuals had a lower alpha-diversity than lean individuals (6). They also showed that the  
48 relative abundance of *Bacteroidetes* and *Firmicutes* increased and decreased, respectively,  
49 as obese individuals lost weight while on a fat or carbohydrate restricted diet (7). Two  
50 re-analysis studies by Walters et al. (8) and Finucane et al. (9) interrogated previously  
51 published microbiome and obesity data and concluded that the previously reported  
52 differences in community diversity and B:F among non-obese and obese individuals  
53 could not be generalized. Regardless of the results using human populations, studies  
54 using animal models where the community was manipulated with antibiotics or established  
55 by colonizing germ-free animals with varied communities appear to support the association  
56 since these manipulations yielded differences in animal weight (10–13). The purported  
57 association between the differences in the microbiome and obesity have been widely  
58 repeated with little attention given to the lack of a clear signal in human cohort studies.

59 The recent publication of additional studies that collected BMI data for each subject as  
60 well as other studies that were not included in the earlier re-analysis studies offered the  
61 opportunity to revisit the question relating the structure of the human microbiome to obesity.  
62 One critique of the prior re-analysis studies is that the authors did not aggregate the results  
63 across studies to increase the effective sample size. It is possible that there were small  
64 associations within each study that were not statistically significant because the individual  
65 studies lacked sufficient power. Alternatively, diversity metrics may mask the appropriate  
66 signal and it is necessary to measure the association at the level of microbial populations.  
67 The Walters re-analysis study demonstrated that Random Forest machine learning models  
68 were capable of predicting obesity status within a single cohort, but did not attempt to test  
69 the models on other cohorts. The purpose of this study was to perform a meta-analysis of  
70 the association between differences in the microbiome and obesity status by analyzing  
71 and applying a more systematic and synthetic approach than was used previously.

## 72 **Results**

73 ***Literature Review and Study Inclusion.*** To perform a robust meta-analysis and limit  
74 inclusion bias, we followed the Preferred Reporting Items for Systematic Reviews and  
75 Meta-Analyses (PRISMA) guidelines to identify the studies that we analyzed (14). A  
76 detailed description of our selection process and the exact search terms are provided in  
77 the Supplemental Text and in Figure 1. Briefly, we searched PubMed for original research  
78 studies that involved studying obesity and the human microbiome. The initial search  
79 yielded 187 studies. We identified 10 additional studies that were not designed to explicitly  
80 test for an association between the microbiome and obesity. We then manually curated  
81 the 197 studies to select those studies that included BMI and 16S rRNA gene sequence  
82 data. This yielded 11 eligible studies. An additional study was removed from our analysis  
83 because no individuals in the study had a BMI over 30. Among the final 10 studies, 3 were

84 identified from our PubMed search (10, 15, 16), 5 were originally identified from the 10  
85 studies that did not explicitly investigate obesity but included BMI data (17–21), and two  
86 datasets were used (22, 23) because these publications did not specifically look for any  
87 metabolic or obesity conditions but had control populations and enabled us to help mitigate  
88 against publication biases associated with the bacterial microbiome and obesity. The ten  
89 studies are summarized in Table 1. For comparison, two of these studies were included  
90 in the Finucane re-analysis study (10, 21) and four of these studies were included in the  
91 Walters re-analysis study (10, 15, 20, 21). The 16S rRNA gene sequence data from each  
92 study was re-analyzed using a similar approach based on previously described methods  
93 for reducing the number of chimeric sequences and sequencing errors for 454 and Illumina  
94 MiSeq data (24, 25). The sequences were clustered into operational taxonomic units  
95 (OTUs) using the average neighbor approach (26) and into taxonomic groupings based on  
96 their classification using a naive Bayesian classifier (27).

97 ***Alpha diversity analysis.*** We calculated the Shannon diversity index, observed richness,  
98 and Shannon evenness, the relative abundance of *Bacteroidetes* and *Firmicutes*, and  
99 the ratio of their relative abundance (B:F) for each sample. Once we transformed each  
100 of the six alpha diversity metrics to make them normally distributed, we used a t-test  
101 to identify significant associations between the alpha diversity metric and whether an  
102 individual was obese for each of the ten studies. The B:F and the relative abundance  
103 of *Firmicutes* were not significantly associated with obesity in any study. We identified  
104 seven P-values that were less than 0.05: three studies indicated obese individuals had  
105 a lower richness, two studies indicated a significantly lower diversity, one study indicated  
106 a significantly lower evenness, and one study indicated a significantly higher relative  
107 abundance of *Bacteroidetes* (Figures 2 and S1). These results largely match those of the  
108 Walters and Finucane re-analysis studies. Interestingly, although only two of the ten studies  
109 observed the previously reported association between lower diversity and obesity, the  
110 other studies appeared to have the same trend, albeit the differences were not statistically

111 significant. We used a random effects linear model to combine the studies using the  
112 study as the random effect and found statistical support for decreased richness, evenness,  
113 and diversity among obese individuals (all  $P < 0.011$ ). Although there was a significant  
114 relationship between these metrics and obesity status, the effect size was quite small.  
115 The obese individuals averaged 7.47% lower richness, 0.88% lower evenness, and 2.07%  
116 lower diversity. There were no significant associations when we pooled the phylum-level  
117 metrics across studies. These results indicate that obese individuals do have a statistically  
118 significant lower diversity than non-obese individuals; however, it is questionable whether  
119 the difference is biologically significant.

120 **Relative risk.** Building upon the alpha diversity analysis we calculated the relative risk of  
121 being obese based on an individual's alpha diversity metrics relative to the median metric  
122 for that study. Inspection of funnel plots for each of the metrics suggested that the studies  
123 included in our analysis were not biased (Figure S2). The results using relative risk largely  
124 matched those of using the raw alpha diversity data. Across the ten studies and six metrics,  
125 the only significant relative risk values were the richness, evenness, and diversity values  
126 from the Goodrich study (Figures 3 and S3). Again, although the relative risk values were  
127 not significant for other studies, the values tended to be above one. When we pooled the  
128 data using a random effects model, the relative risk associated with having a richness,  
129 evenness, or diversity below the median for the population was significantly associated  
130 with obesity (all  $P < 0.0044$ ). The relative risks associated with alpha diversity were small.  
131 The relative risk of having a low richness was 1.30 (95% CI: 1.13-1.49), low evenness was  
132 1.20 (95% CI: 1.06-1.37), and low diversity was 1.27 (95% CI: 1.09-1.48). There were no  
133 significant differences in the phylum-level metrics. Again, the relative risk results indicate  
134 that individuals with a lower richness, evenness, or diversity are at statistically significant  
135 increased risk of being obese, it is questionable whether that risk is biologically or clinically  
136 relevant.

137 **Beta diversity analysis.** Following the approach used by the Walters and Finucane  
138 re-analysis studies, for each dataset we calculated a Bray-Curtis distance matrix to  
139 measure the difference in the membership and structure of the individuals from each  
140 study. We then used AMOVA to test for significant differences between the structure  
141 of non-obese and obese individuals (Table 1). The Escobar, Goodrich, and Turnbaugh  
142 datasets indicated a significant difference in community structure (all  $P < 0.05$ ). Because  
143 it was not possible to ascertain the directionality of the difference in community structure  
144 because the samples are arrayed in a non-dimensional space or perform a pooled analysis  
145 using studies that had non-overlapping 16S rRNA gene sequence regions, it is unclear  
146 whether these differences reflect a broader, but perhaps small, shift in community structure  
147 between non-obese and obese individuals.

148 **Development of a microbiome-based classifier of obesity.** The Walters re-analysis  
149 study suggested that it was possible to classify individuals as being non-obese or obese  
150 based on the composition of their microbiota. We repeated this analysis with additional  
151 datasets using OTU and genus-level phylotype data. For each study we developed a  
152 Random Forest machine learning model to classify individuals. Using ten-fold cross  
153 validation, the cross-validated AUC values for the OTU-based models varied between 0.52  
154 and 0.69 indicating a relatively poor ability to classify individuals (Figure 4A). To test models  
155 on other datasets, we trained models using genus-level phylotype data for each dataset.  
156 The cross-validated AUC values for the models applied to the training datasets varied  
157 between 0.51 and 0.65, again indicating a relatively poor ability to classify individuals from  
158 the original dataset (Figure 4B). For each model we identified the probability where the  
159 sum of the sensitivity and specificity was the highest. We then used this probability to  
160 define a threshold for calculating the accuracy of the models when applied to the other  
161 nine datasets (Figure 5). Although there was considerable variation in accuracy values  
162 for each model, the median accuracy for each model varied between 0.33 (Turnbaugh)  
163 and 0.65 (HMP) (median=0.57). We built similar models using taxonomic representation

164 based on phylum, class, order, and family assignments and saw no improvement in the  
165 results (Figure S4). We also attempted to predict individual BMI values as continuous  
166 variables based on the relative abundance of OTUs and genera. The median percent of the  
167 variance explained with the resulting models was 12.9% for the OTU-based models and  
168 8.2% for the genus-based models. When we considered the number of samples, balance  
169 of non-obese and obese individuals, and region within the 16S rRNA gene for each study  
170 it was not possible to identify factors that predictably affected model performance. The  
171 ability to predict obesity status using relative abundance data from the communities was  
172 only marginally better than random. These results suggest that given the large diversity of  
173 microbiome compositions it is difficult to identify a taxonomic signal that can be associated  
174 with obesity.

175 ***Power and Sample Size Estimate Simulations.*** The inability to detect a difference  
176 between non-obese and obese individuals could be due to the lack of a true effect or  
177 because the study had insufficient statistical power to detect a difference because of  
178 insufficient sampling, large interpersonal variation, or unbalanced sampling of non-obese  
179 and obese individuals. To assess these factors, we calculated the power to detect  
180 differences of 1, 5, 10, and 15% in each of the alpha diversity metrics using the sample  
181 sizes used in each of the studies (Figures 6, S5-S10). Although there is no biological  
182 rationale for these effect sizes, they represent a range that includes effect sizes that would  
183 be generally considered to be biologically significant. Only the Goodrich study had power  
184 greater than 0.80 to detect a 5% difference in Shannon diversity and six of the studies had  
185 enough power to detect a 10% difference (Figure 6A). None of the studies had sufficient  
186 power to detect a 15% difference between B:F values (Figure S5). In fact, the maximum  
187 power among any of the studies to detect a 15% difference in B:F values was 0.25. Among  
188 the tests for relative risk, none of the studies had sufficient power to detect a Cohen's d  
189 of 0.10 and only two studies had sufficient power to detect a Cohen's d of 0.15. We next  
190 estimated how many individuals would need to have been sampled to have sufficient power

191 to detect the four effect sizes assuming the observed interpersonal variation from each  
192 study and balanced sampling between the two groups (Figure 6B). To detect a 1, 5, 10, or  
193 15% difference in Shannon diversity, the median required sampling effort per group was  
194 approximately 3,400, 140, 35, or 16 individuals, respectively. To detect a 1, 5, 10, and 15%  
195 difference in B:F values, the median required sampling effort per group was approximately  
196 160,000, 6,300, 1,600, or 700 individuals, respectively. To detect a 1, 5, 10, and 15%  
197 difference in relative risk values using Shannon diversity, the median required sampling  
198 effort per group was approximately 39,000, 1,500, 380, or 170 individuals, respectively.  
199 These estimates indicate that most microbiome studies are underpowered to detect modest  
200 effect sizes using either metric. In the case of obesity, the studies were underpowered to  
201 detect the 0.90 to 6% difference in diversity that was observed across the studies.

## 202 **Discussion**

203 Our meta-analysis helps to provide clarity to the ongoing debate of whether or not there  
204 are specific microbiome-based markers that can be associated with obesity. We performed  
205 an extensive literature review of the existing studies on the microbiome and obesity and  
206 performed a meta-analysis on the studies that remained based on our inclusion and  
207 exclusion criteria. By statistically pooling the data from ten studies, we observed significant,  
208 but small, relationships between richness, evenness, and diversity and obesity status as  
209 well as the relative risk of being obese based on these metrics. We also generated Random  
210 Forest machine learning models trained on each dataset and tested on the remaining  
211 datasets. This analysis demonstrated that the ability to reliably classify individuals as  
212 being obese based solely on the composition of their microbiome was limited. Finally,  
213 we assessed the ability of each study to detect defined differences in alpha diversity and  
214 observed that most studies were underpowered to detect modest effect sizes. Considering  
215 these datasets are among the largest published, it appears that most human microbiome

216 studies are underpowered to detect differences in alpha diversity.

217 Alpha diversity metrics are attractive because they distill a complex dataset to a single  
218 value. For example, Shannon diversity is a measure of the entropy in a community and  
219 integrates richness and evenness information. Two communities with little taxonomic  
220 similarity can have the same diversity. Among ecologists the relevance of these metrics is  
221 questioned because it is difficult to ascribe a mechanistic interpretation to their relationship  
222 with stability or disease. Regardless, the concept of a biologically significant effect size  
223 needs to be developed among microbiome researchers. Alternative metrics could include  
224 the ability to detect a defined difference in the relative abundance of an OTU representing a  
225 defined relative abundance. What makes for a biologically significant difference or relative  
226 abundance is an important point that has yet to be discussed in the microbiome field. The  
227 use of operationally defined effect sizes should be adequate until it is possible to decide  
228 upon an accepted practice.

229 By selecting a range of possible effect sizes, we were able to demonstrate that most studies  
230 are underpowered to detect modest differences in alpha diversity metrics and phylum-level  
231 relative abundances. Several factors interact to limit the power of microbiome studies.  
232 There is wide interpersonal variation in the diversity and structure of the human microbiome.  
233 Some factors such as relationship between subjects could potentially decrease the amount  
234 variation (6) and other factors such as whether one lives in a rural environment could  
235 increase the amount of variation (28). In addition, the common experimental designs limit  
236 their power. As we observed, most of the studies included in our analysis were unbalanced  
237 for the variable that we were interested in. This was also true of those studies that originally  
238 sought to identify associations with obesity. Even with a balanced design, we showed that  
239 it was necessary to obtain approximately 140 and 6,300 samples per group to detect a 5%  
240 difference in Shannon diversity or B:F, respectively. It was interesting that these sample  
241 sizes agreed across studies regardless of their sequencing method, region within the 16S

242 rRNA gene, or subject population (Figure 6). This suggests that regardless of the treatment  
243 or category, these sample sizes represent a good starting point for subject recruitment  
244 when using stool samples. Unfortunately, few studies have been published with this level  
245 of subject recruitment. This is troubling since the positive predictive rate of a significant  
246 finding in an underpowered study is small leading to results that cannot be reproduced  
247 (29). Future microbiome studies should articulate the basis for their experimental design.

248 Two previous re-analysis studies have stated that there was not a consistent association  
249 between alpha diversity and obesity (8, 9); however, neither of these studies made an  
250 attempt to pool the existing data together to try and harness the additional power that  
251 this would give and they did not assess whether the studies were sufficiently powered  
252 to detect a difference. Additionally, our analysis used 16S rRNA gene sequence data  
253 from ten studies whereas the Finucane study used 16S rRNA gene sequence data from  
254 three studies (7, 10, 21) and a metagenomic study (30) and the Walters study used  
255 16S rRNA gene sequence data from five studies (10, 15, 20, 21, 28); two studies were  
256 included in both analyses (10, 21). Our analysis included four of these studies (10, 15,  
257 20, 21) and excluded three of the studies because they were too small (7), only utilized  
258 metagenomic data (30), or used short single read Illumina HiSeq data that has a high  
259 error rate making it intractable for *de novo* OTU clustering (28). The additional seven  
260 datasets were published after the two reviews were performed and include datasets with  
261 more samples than were found in the original studies. Our collection of ten studies allowed  
262 us to largely use the same sequence analysis pipeline for all datasets and relied heavily  
263 on the availability of public data and access to metadata that included variables beyond  
264 the needs of the original study. To execute this analysis, we created an automated data  
265 analysis pipeline, which can be easily updated to add additional studies as they become  
266 available ([https://github.com/SchlossLab/Sze\\_Obesity\\_mBio\\_2016/](https://github.com/SchlossLab/Sze_Obesity_mBio_2016/)). Similarly, it would be  
267 possible to adapt this pipeline to other body sites and treatment or variables (e.g. subject's  
268 sex or age).

269 Similar to our study, the Walters study generated Random Forest machine learning models  
270 to differentiate between non-obese and obese individuals (8). They obtained similar AUC  
271 values to our analysis; however, they did not attempt to test these models on the other  
272 studies in their analysis. When we performed the inter-dataset cross validation the median  
273 accuracy across datasets was only 56.68% indicating that the models did a poor job when  
274 applied to other datasets. This could be due to differences in subject populations and  
275 methods. Furthermore, others have reported improved classification at broader taxonomic  
276 levels (31); we did not find this to be the case across the studies in our analysis (Figure  
277 S4). Considering the median AUC for models trained and tested on the same data with  
278 ten-fold cross validation only varied between 0.51 and 0.65 and that there was not a strong  
279 signal in the alpha diversity data, we suspect that there is insufficient signal to reliably  
280 classify individuals to a BMI category based on their microbiota.

281 Although we failed to find an effect this does not necessarily mean that there is no role  
282 for the microbiome in obesity. There is strong evidence in murine models of obesity that  
283 the microbiome and level of adiposity can be manipulated via genetic manipulation of the  
284 animal and manipulation of the community through antibiotics or colonizing germ free mice  
285 with diverse fecal material from human donors (5, 10–13). These studies appear to conflict  
286 with the observations using human subjects. Recalling the large interpersonal variation in  
287 the structure of the microbiome, it is possible that each individual has their own signatures  
288 of obesity. Alternatively, it could be that the involvement of the microbiome in obesity is not  
289 apparent based on the taxonomic information provided by 16S rRNA gene sequence data.  
290 Rather, the differences could become more apparent at the level of a common set of gene  
291 transcripts or metabolites that can be produced from different structures of the microbiome.

## 292 **Methods**

293 **Sequence Analysis Pipeline.** All sequence data were publicly available and were  
294 downloaded from the NCBI Sequence Read Archive, the European Nucleotide Archive,  
295 or the investigators' personal website ([https://gordonlab.wustl.edu/TurnbaughSE/\\_10/\\_09/](https://gordonlab.wustl.edu/TurnbaughSE/_10/_09/STM/_2009.html)  
296 [STM/\\_2009.html](https://gordonlab.wustl.edu/TurnbaughSE/_10/_09/STM/_2009.html)). In total seven studies used 454 (6, 15, 16, 18, 20–22) and three studies  
297 used Illumina sequencing (17, 19, 23). All of these studies used amplification-based  
298 16S rRNA gene sequencing. Among the studies that sequenced the 16S rRNA gene,  
299 the researchers targeted the V1-V2 (20), V1-V3 (15, 16, 18), V3-V5 (21, 22), V4 [(19);  
300 (23); ], and V3-4 (17) regions. For those studies where multiple regions were sequenced,  
301 we selected the region that corresponded to the largest number of subjects (6, 21). We  
302 processed the 16S rRNA gene sequence data using a standardized mothur pipeline. Briefly,  
303 our pipelines attempted to follow previously recommended approaches for 454 and Illumina  
304 sequencing data (24, 25). All sequences were screened for chimeras using UCHIME and  
305 assigned to operational taxonomic units (OTUs) using the average neighbor algorithm  
306 using a 3% distance threshold (26, 32). All sequence processing was performed using  
307 mothur (v.1.37.0) (33).

308 **Data Analysis.** We split the overall meta-analysis into three general strategies using R  
309 (3.3.0). First, we followed the approach employed by Finucane et al (9) and Walters et al  
310 (8) where each study was re-analyzed separately to identify associations between BMI  
311 and the relative abundance of *Bacteroidetes* and Firmicutes, the ratio of *Bacteroidetes*  
312 and *Firmicutes* relative abundances (B:F), Shannon diversity, observed richness, and  
313 Shannon evenness (34). After each variable was transformed to fit a normal distribution  
314 a two-tailed t-test was performed for comparison of non-obese and obese individuals  
315 (i.e. BMI > 35.0). We performed a pooled analysis on these measured variables using  
316 linear random effect models to correct for study effect to assess differences on the combined  
317 dataset between non-obese and obese groups using the lme4 (v.1.1-12) R package (35).

318 Next, we compared the community structure from non-obese and obese individuals using  
319 analysis of molecular variance (AMOVA) with Bray-Curtis distance matrices (36). This  
320 analysis was performed using the vegan (v.2.3-5) R package. For both analyses, the  
321 datasets were rarefied (N=1000) so that each study had the same number of sequences.  
322 Second, for each study we partitioned the subjects into a low or high group depending  
323 on whether their alpha diversity metrics were below or above the median value for the  
324 study. The relative risk (RR) was then calculated as the ratio of the number of obese  
325 individuals in the low group to the number of obese individuals in the high group. We then  
326 performed a Fisher exact-test to investigate whether the RR was significantly different from  
327 1.0 within each study and across all of the studies using the epiR (0.9-77) and metafor  
328 (1.9-8) packages. Third, we used the AUCRF (1.1) R package to generate Random Forest  
329 models (37). For each study we developed models using either OTUs or genus-level  
330 phylotypes. The quality of each model was assessed by measuring the area under the  
331 curve (AUC) of the Receiver Operating Characteristic (ROC) using ten-fold cross validation.  
332 Because the genus-level phylotype models were developed using a common reference, it  
333 was possible to use one study's model (i.e. the training set) to classify the samples from  
334 the other studies (i.e. the testing sets). The optimum threshold for the training set was set  
335 as the probability threshold that had the highest combined sensitivity and specificity. This  
336 threshold was then used to calculate the accuracy of the model applied to the test studies.  
337 To generate ROC curves and calculate the accuracy of the models we used the pROC (1.8)  
338 R package (38). Finally, we performed power and sample number simulations for different  
339 effect sizes for each study using the pwr (1.1-3) R package and base R functions. We also  
340 calculated the actual sample size needed based on the effect size of each individual study.

341 **Reproducible methods.** A detailed and reproducible description of how the data were  
342 processed and analyzed can be found at [https://github.com/SchlossLab/Sze\\_Obesity\\_](https://github.com/SchlossLab/Sze_Obesity_mBio_2016/)  
343 [mBio\\_2016/](https://github.com/SchlossLab/Sze_Obesity_mBio_2016/).

## 344 **Acknowledgements**

345 The authors would like to thank Nielson Baxter and Shawn Whitefield for their suggestions  
346 on the development of the manuscript. We are grateful to the authors of the studies used  
347 in our meta-analysis who have made their data publicly available or available to us directly.  
348 Without their forethought studies such as this would not be possible. This work was  
349 supported in part by funding from the National Institutes of Health to PDS (U01AI2425501  
350 and P30DK034933).

351 **Figure 1: PRISMA flow diagram of total records searched (39).**

352 **Figure 2: Individual and combined comparison of obese and non-obese groups for**  
353 **Shannon diversity (A) and B:F (B).**

354 **Figure 3: Meta analysis of the relative risk of obesity based on Shannon diversity**  
355 **(A) or B:F (B).**

356 **Figure 4: ROC curves for each study based on classification of non-obese or obese**  
357 **groups using OTUs (A) or genus-level classification (B).**

358 **Figure 5: Overall accuracy of each study to predict non-obese and obese**  
359 **individuals based on that study's Random Forest machine learning model applied**  
360 **to each of the other studies.**

361 **Figure 6: Power (A) and sample size simulations (B) for Shannon diversity for**  
362 **differentiating between non-obese versus obese for effect sizes of 1, 5, 10, and**  
363 **15%. Power calculations use the sampling distribution from the original studies and the**  
364 **sample size estimations assume an equal amount of sampling from each treatment group.**

365 **Figure S1: Individual and Combined comparison of Obese and Non-Obese groups**  
366 **Based on Evenness (A), Richness (B), or the Relative Abundance of *Bacteroidetes***  
367 **(C) and Firmicutes (D).**

368 **Figure S2: Funnel plots depicting the general lack of bias in the selection of**  
369 **datasets included in the analysis.**

370 **Figure S3: Meta Analysis of the Relative Risk of Obesity Based on Evenness (A),**  
371 **Richness (B), or the Relative Abundance of *Bacteroidetes* (C) and Firmicutes (D).**

372 **Figure S4: Overall accuracy of each study to predict non-obese and obese**  
373 **individuals based on that study's Random Forest machine learning model applied**  
374 **to each of the other studies when trained using relative abundance of each phylum,**  
375 **class, order, family, or genus.** The cross-validated AUC values for the training model  
376 are provided for each study and taxonomic level.

377 **Figure S5: Power (A) and sample size simulations (B) for B:F for differentiating**  
378 **between non-obese versus obese for effect sizes of 1, 5, 10, and 15%.** Power  
379 calculations use the sampling distribution from the original studies and the sample size  
380 estimations assume an equal amount of sampling from each treatment group.

381 **Figure S6: Power (A) and sample size simulations (B) for richness for differentiating**  
382 **between non-obese versus obese for effect sizes of 1, 5, 10, and 15%.** Power  
383 calculations use the sampling distribution from the original studies and the sample size  
384 estimations assume an equal amount of sampling from each treatment group.

385 **Figure S7: Power (A) and sample size simulations (B) for evenness for**  
386 **differentiating between non-obese versus obese for effect sizes of 1, 5, 10,**  
387 **and 15%.** Power calculations use the sampling distribution from the original studies and  
388 the sample size estimations assume an equal amount of sampling from each treatment

389 group.

390 **Figure S8: Power (A) and sample size simulations (B) for the relative abundance of**  
391 ***Bacteroidetes* for differentiating between non-obese versus obese for effect sizes**  
392 **of 1, 5, 10, and 15%.** Power calculations use the sampling distribution from the original  
393 studies and the sample size estimations assume an equal amount of sampling from each  
394 treatment group.

395 **Figure S9: Power (A) and sample size simulations (B) for the relative abundance of**  
396 ***Firmicutes* for differentiating between non-obese versus obese for effect sizes of**  
397 **1, 5, 10, and 15%.** Power calculations use the sampling distribution from the original  
398 studies and the sample size estimations assume an equal amount of sampling from each  
399 treatment group.

400 **Figure S10: Power (A) and sample size simulations (B) for relative risk of obesity**  
401 **based on Shannon diversity.** Power calculations use the sampling distribution from the  
402 original studies and the sample size estimations assume an equal amount of sampling  
403 from each treatment group.

## 404 **References**

- 405 1. **Ogden CL, Carroll MD, Kit BK, Flegal KM.** 2014. Prevalence of childhood and adult  
406 obesity in the United States, 2011-2012. *JAMA* **311**:806–814. doi:[http://doi.org/10.1001/](http://doi.org/10.1001/jama.2014.732)  
407 [jama.2014.732](http://doi.org/10.1001/jama.2014.732).
- 408 2. **Lichtash CT, Cui J, Guo X, Chen Y-DI, Hsueh WA, Rotter JI, Goodarzi MO.** 2013.  
409 Body adiposity index versus body mass index and other anthropometric traits as correlates  
410 of cardiometabolic risk factors. *PloS One* **8**:e65954. doi:[http://doi.org/10.1371/journal.](http://doi.org/10.1371/journal.pone.0065954)  
411 [pone.0065954](http://doi.org/10.1371/journal.pone.0065954).
- 412 3. **Brahe LK, Astrup A, Larsen LH.** 2016. Can We Prevent Obesity-Related Metabolic  
413 Diseases by Dietary Modulation of the Gut Microbiota? *Advances in Nutrition* (Bethesda,  
414 Md) **7**:90–101. doi:<http://doi.org/10.3945/an.115.010587>.
- 415 4. **Dror T, Dickstein Y, Dubourg G, Paul M.** 2016. Microbiota manipulation for weight  
416 change. *Microbial Pathogenesis*. doi:<http://doi.org/10.1016/j.micpath.2016.01.002>.
- 417 5. **Ley RE, Bäckhed F, Turnbaugh P, Lozupone CA, Knight RD, Gordon**  
418 **JJ.** 2005. Obesity alters gut microbial ecology. *Proceedings of the National*  
419 *Academy of Sciences of the United States of America* **102**:11070–11075. doi:<http://doi.org/10.1073/pnas.0504978102>.
- 420
- 421 6. **Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, Sogin**  
422 **ML, Jones WJ, Roe BA, Affourtit JP, Egholm M, Henrissat B, Heath AC, Knight R,**  
423 **Gordon JJ.** 2009. A core gut microbiome in obese and lean twins. *Nature* **457**:480–484.  
424 doi:<http://doi.org/10.1038/nature07540>.
- 425 7. **Ley RE, Turnbaugh PJ, Klein S, Gordon JJ.** 2006. Microbial ecology: Human gut  
426 microbes associated with obesity. *Nature* **444**:1022–1023. doi:<http://doi.org/10.1038/>

427 [4441022a](#).

428 **8. Walters WA, Xu Z, Knight R.** 2014. Meta-analyses of human gut microbes associated  
429 with obesity and IBD. FEBS letters **588**:4223–4233. doi:[http://doi.org/10.1016/j.febslet.](http://doi.org/10.1016/j.febslet.2014.09.039)  
430 [2014.09.039](#).

431 **9. Finucane MM, Sharpton TJ, Laurent TJ, Pollard KS.** 2014. A taxonomic signature  
432 of obesity in the microbiome? Getting to the guts of the matter. PloS One **9**:e84689.  
433 doi:<http://doi.org/10.1371/journal.pone.0084689>.

434 **10. Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JI.** 2006.  
435 An obesity-associated gut microbiome with increased capacity for energy harvest. Nature  
436 **444**:1027–31. doi:<http://doi.org/10.1038/nature05414>.

437 **11. Koren O, Goodrich JK, Cullender TC, Spor A, Laitinen K, Bäckhed HK, Gonzalez**  
438 **A, Werner JJ, Angenent LT, Knight R, Bäckhed F, Isolauri E, Salminen S, Ley RE.**  
439 2012. Host remodeling of the gut microbiome and metabolic changes during pregnancy.  
440 Cell **150**:470–480. doi:<http://doi.org/10.1016/j.cell.2012.07.008>.

441 **12. Cox LM, Yamanishi S, Sohn J, Alekseyenko AV, Leung JM, Cho I, Kim SG, Li**  
442 **H, Gao Z, Mahana D, Rodriguez JGZ, Rogers AB, Robine N, Loke P, Blaser MJ.**  
443 2014. Altering the intestinal microbiota during a critical developmental window has lasting  
444 metabolic consequences. Cell **158**:705–721. doi:<http://doi.org/10.1016/j.cell.2014.05.052>.

445 **13. Mahana D, Trent CM, Kurtz ZD, Bokulich NA, Battaglia T, Chung J, Müller CL, Li**  
446 **H, Bonneau RA, Blaser MJ.** 2016. Antibiotic perturbation of the murine gut microbiome  
447 enhances the adiposity, insulin resistance, and liver disease associated with high-fat diet.  
448 Genome Medicine **8**. doi:<http://doi.org/10.1186/s13073-016-0297-9>.

449 **14. Moher D, Liberati A, Tetzlaff J, Altman DG, PRISMA Group.** 2010. Preferred  
450 reporting items for systematic reviews and meta-analyses: The PRISMA statement.

451 International Journal of Surgery (London, England) **8**:336–341. doi:<http://doi.org/10.1016/j.ijso.2010.02.007>.

453 **15. Zupancic ML, Cantarel BL, Liu Z, Drabek EF, Ryan KA, Cirimotich S, Jones**  
454 **C, Knight R, Walters WA, Knights D, Mongodin EF, Horenstein RB, Mitchell BD,**  
455 **Steinle N, Snitker S, Shuldiner AR, Fraser CM.** 2012. Analysis of the gut microbiota  
456 in the old order Amish and its relation to the metabolic syndrome. PLoS One **7**:e43052.  
457 doi:<http://doi.org/10.1371/journal.pone.0043052>.

458 **16. Escobar JS, Klotz B, Valdes BE, Agudelo GM.** 2014. The gut microbiota of  
459 Colombians differs from that of Americans, Europeans and Asians. BMC microbiology  
460 **14**:311. doi:<http://doi.org/10.1186/s12866-014-0311-6>.

461 **17. Zeevi D, Korem T, Zmora N, Israeli D, Rothschild D, Weinberger A, Ben-Yacov**  
462 **O, Lador D, Avnit-Sagi T, Lotan-Pompan M, Suez J, Mahdi JA, Matot E, Malka**  
463 **G, Kosower N, Rein M, Zilberman-Schapira G, Dohnalová L, Pevsner-Fischer M,**  
464 **Bikovsky R, Halpern Z, Elinav E, Segal E.** 2015. Personalized Nutrition by Prediction of  
465 Glycemic Responses. Cell **163**:1079–1094. doi:<http://doi.org/10.1016/j.cell.2015.11.001>.

466 **18. Ross MC, Muzny DM, McCormick JB, Gibbs RA, Fisher-Hoch SP, Petrosino JF.**  
467 2015. 16S gut community of the Cameron County Hispanic Cohort. Microbiome **3**:7.  
468 doi:<http://doi.org/10.1186/s40168-015-0072-y>.

469 **19. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, Beaumont**  
470 **M, Van Treuren W, Knight R, Bell JT, Spector TD, Clark AG, Ley RE.** 2014. Human  
471 genetics shape the gut microbiome. Cell **159**:789–799. doi:<http://doi.org/10.1016/j.cell.2014.09.053>.

473 **20. Wu GD, Chen J, Hoffmann C, Bittinger K, Chen Y-Y, Keilbaugh SA, Bewtra M,**  
474 **Knights D, Walters WA, Knight R, Sinha R, Gilroy E, Gupta K, Baldassano R, Nessel**

475 **L, Li H, Bushman FD, Lewis JD.** 2011. Linking long-term dietary patterns with gut  
476 microbial enterotypes. *Science (New York, NY)* **334**:105–108. doi:[http://doi.org/10.1126/](http://doi.org/10.1126/science.1208344)  
477 [science.1208344](http://doi.org/10.1126/science.1208344).

478 **21. Human Microbiome Project Consortium.** 2012. Structure, function and diversity  
479 of the healthy human microbiome. *Nature* **486**:207–214. doi:[http://doi.org/10.1038/](http://doi.org/10.1038/nature11234)  
480 [nature11234](http://doi.org/10.1038/nature11234).

481 **22. Schubert AM, Rogers MAM, Ring C, Mogle J, Petrosino JP, Young VB,**  
482 **Aronoff DM, Schloss PD.** 2014. Microbiome data distinguish patients with *Clostridium*  
483 *difficile* infection and non-*C. difficile*-associated diarrhea from healthy controls. *mBio*  
484 **5**:e01021–01014. doi:<http://doi.org/10.1128/mBio.01021-14>.

485 **23. Baxter NT, Ruffin MT, Rogers MAM, Schloss PD.** 2016. Microbiota-based model  
486 improves the sensitivity of fecal immunochemical test for detecting colonic lesions. *Genome*  
487 *Medicine* **8**:37. doi:<http://doi.org/10.1186/s13073-016-0290-3>.

488 **24. Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD.** 2013.  
489 Development of a dual-index sequencing strategy and curation pipeline for analyzing  
490 amplicon sequence data on the MiSeq Illumina sequencing platform. *Applied and*  
491 *environmental microbiology* **79**:5112–5120.

492 **25. Schloss PD, Gevers D, Westcott SL.** 2011. Reducing the effects of PCR amplification  
493 and sequencing artifacts on 16S rRNA-based studies. *PLoS ONE* **6**:e27310. doi:<http://doi.org/10.1371/journal.pone.0027310>.

495 **26. Westcott SL, Schloss PD.** 2015. De novo clustering methods outperform  
496 reference-based methods for assigning 16S rRNA gene sequences to operational  
497 taxonomic units. *PeerJ* **3**:e1487. doi:<http://doi.org/10.7717/peerj.1487>.

498 **27. Wang Q, Garrity GM, Tiedje JM, Cole JR.** 2007. Naive bayesian classifier for

- 499 rapid assignment of rRNA sequences into the new bacterial taxonomy. *Applied and*  
500 *Environmental Microbiology* **73**:5261–5267. doi:<http://doi.org/10.1128/aem.00062-07>.
- 501 **28. Yatsunenko T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras**  
502 **M, Magris M, Hidalgo G, Baldassano RN, Anokhin AP, Heath AC, Warner B, Reeder**  
503 **J, Kuczynski J, Caporaso JG, Lozupone CA, Lauber C, Clemente JC, Knights D,**  
504 **Knight R, Gordon JI.** 2012. Human gut microbiome viewed across age and geography.  
505 *Nature* **486**:222–227. doi:<http://doi.org/10.1038/nature11053>.
- 506 **29. Ioannidis JPA.** 2005. Why most published research findings are false. *PLoS Med*  
507 **2**:e124. doi:<http://doi.org/10.1371/journal.pmed.0020124>.
- 508 **30. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons**  
509 **N, Levenez F, Yamada T, Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang**  
510 **H, Zheng H, Xie Y, Tap J, Lepage P, Bertalan M, Batto J-M, Hansen T, Paslier DL,**  
511 **Linneberg A, Nielsen HB, Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu**  
512 **H, Yu C, Li S, Jian M, Zhou Y, Li Y, Zhang X, Li S, Qin N, Yang H, Wang J, Brunak**  
513 **S, Doré J, Guarner F, Kristiansen K, Pedersen O, Parkhill J, Weissenbach J, Antolin**  
514 **M, Artiguenave F, Blottiere H, Borrueil N, Bruls T, Casellas F, Chervaux C, Cultrone**  
515 **A, Delorme C, Denariáz G, Dervyn R, Forte M, Friss C, Guchte M van de, Guedon E,**  
516 **Haimet F, Jamet A, Juste C, Kaci G, Kleerebezem M, Knol J, Kristensen M, Layec**  
517 **S, Roux KL, Leclerc M, Maguin E, Minardi RM, Oozeer R, Rescigno M, Sanchez N,**  
518 **Tims S, Torrejon T, Varela E, Vos W de, Winogradsky Y, Zoetendal E, Bork P, Ehrlich**  
519 **SD, Wang J.** 2010. A human gut microbial gene catalogue established by metagenomic  
520 sequencing. *Nature* **464**:59–65. doi:<http://doi.org/10.1038/nature08821>.
- 521 **31. Sun Y, Cai Y, Mai V, Farmerie W, Yu F, Li J, Goodison S.** 2010. Advanced  
522 computational algorithms for microbial community analysis using massive 16S rRNA  
523 sequence data. *Nucleic Acids Research* **38**:e205–e205. doi:<http://doi.org/10.1093/nar/>

524 [gkq872](#).

525 32. **Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R.** 2011. UCHIME improves  
526 sensitivity and speed of chimera detection. *Bioinformatics* **27**:2194–2200. doi:[http://doi.](http://doi.org/10.1093/bioinformatics/btr381)  
527 [org/10.1093/bioinformatics/btr381](http://doi.org/10.1093/bioinformatics/btr381).

528 33. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB,**  
529 **Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, others.** 2009. Introducing  
530 mothur: open-source, platform-independent, community-supported software for describing  
531 and comparing microbial communities. *Applied and environmental microbiology*  
532 **75**:7537–7541.

533 34. **Magurran AE.** 2003. Measuring biological diversity 264.

534 35. **Bates D, Mächler M, Bolker B, Walker S.** 2015. Fitting linear mixed-effects models  
535 using lme4. *Journal of Statistical Software* **67**:1–48. doi:[http://doi.org/10.18637/jss.v067.](http://doi.org/10.18637/jss.v067.i01)  
536 [i01](http://doi.org/10.18637/jss.v067.i01).

537 36. **Anderson MJ.** 2001. A new method for non-parametric multivariate analysis of  
538 variance. *Austral Ecology* **26**:32–46. doi:[http://doi.org/10.1111/j.1442-9993.2001.01070.](http://doi.org/10.1111/j.1442-9993.2001.01070.pp.x)  
539 [pp.x](http://doi.org/10.1111/j.1442-9993.2001.01070.pp.x).

540 37. **Calle ML, Urrea V, Boulesteix A-L, Malats N.** 2011. AUC-RF: A new strategy for  
541 genomic profiling with random forest. *Human Heredity* **72**:121–132. doi:[http://doi.org/10.](http://doi.org/10.1159/000330778)  
542 [1159/000330778](http://doi.org/10.1159/000330778).

543 38. **Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, Müller M.** 2011.  
544 PROC: An open-source package for r and s+ to analyze and compare rOC curves. *BMC*  
545 *Bioinformatics* **12**:77.

546 39. **Moher D, Liberati A, Tetzlaff J, Altman DG.** 2009. Preferred reporting items for

547 systematic reviews and meta-analyses: The PRISMA statement. PLoS Med **6**:e1000097.

548 doi:<http://doi.org/10.1371/journal.pmed.1000097>.

**Table 1. Summary of obesity, demographic, sequencing, and beta-diversity analysis data for the studies used in the meta-analysis.** NA indicates that those metadata were not available for that study.

Study (Ref.)	Subjects (N)	Obese (%)	Average BMI (Min-Max)	Female (%)	Average Age (Min-Max)	Non-Hispanic White (%)	Sequencing Method	16S rRNA Gene Region	AMOVA (P-value)
Baxter (23)	172	27.3	27.0 (17.5-46.9)	64.5	54.3 (29.0-80.0)	87.8	MiSeq	V4	0.078
Escobar (16)	30	33.3	27.4 (19.5-37.6)	46.7	38.1 (21.0-60.0)	NA	454	V2	0.047
Goodrich (19)	982	19.7	26.3 (16.2-52.4)	98.9	61.0 (23.0-86.0)	NA	MiSeq	V4	<0.001
HMP (21)	287	10.8	24.3 (19.0-34.0)	49.1	26.3 (18.0-40.0)	81.5	454	V3-V5	0.322
Ross (18)	63	60.3	31.6 (22.1-47.9)	76.2	57.0 (33.0-81.0)	0.0	454	V1-V3	0.845
Schubert (22)	104	32.7	28.2 (18.5-62.5)	66.3	52.8 (19.0-88.0)	82.7	454	V3-V5	0.180
Turnbaugh (6)	146	67.8	NA	NA	NA	51.4	454	V2	0.040
Wu (20)	64	7.8	24.3 (14.0-41.3)	53.1	26.3 (2.16-50.0)	NA	454	V1-V2	0.577
Zeevi (17)	731	21.6	26.4 (16.4-47.0)	NA	43.4 (18.0-70.0)	NA	MiSeq	V3-V4	0.135
Zupancic (15)	207	36.2	28.2 (18.2-127.0)	57.0	46.7 (20.0-79.0)	100.0	454	V3-V5	0.206

Records identified through  
database searching  
(n = 187)

Additional records identified  
through other sources  
(n = 10)

Records after duplicates removed  
(n = 197)

Records screened  
(n = 197)

Records excluded  
(n = 184)

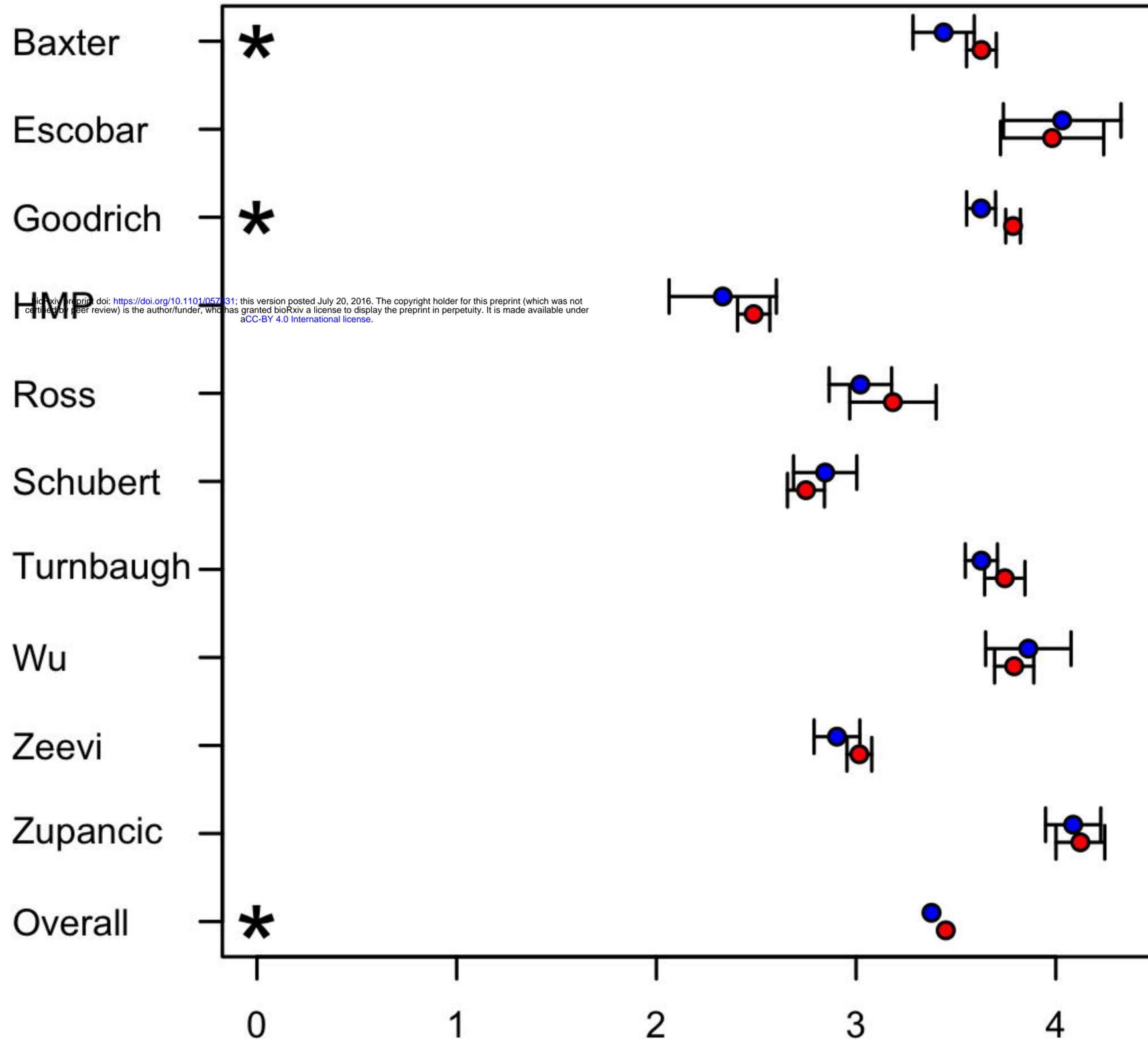
Full-text articles assessed  
for eligibility  
(n = 13)

Full-text articles excluded,  
with reasons  
(n = 2)

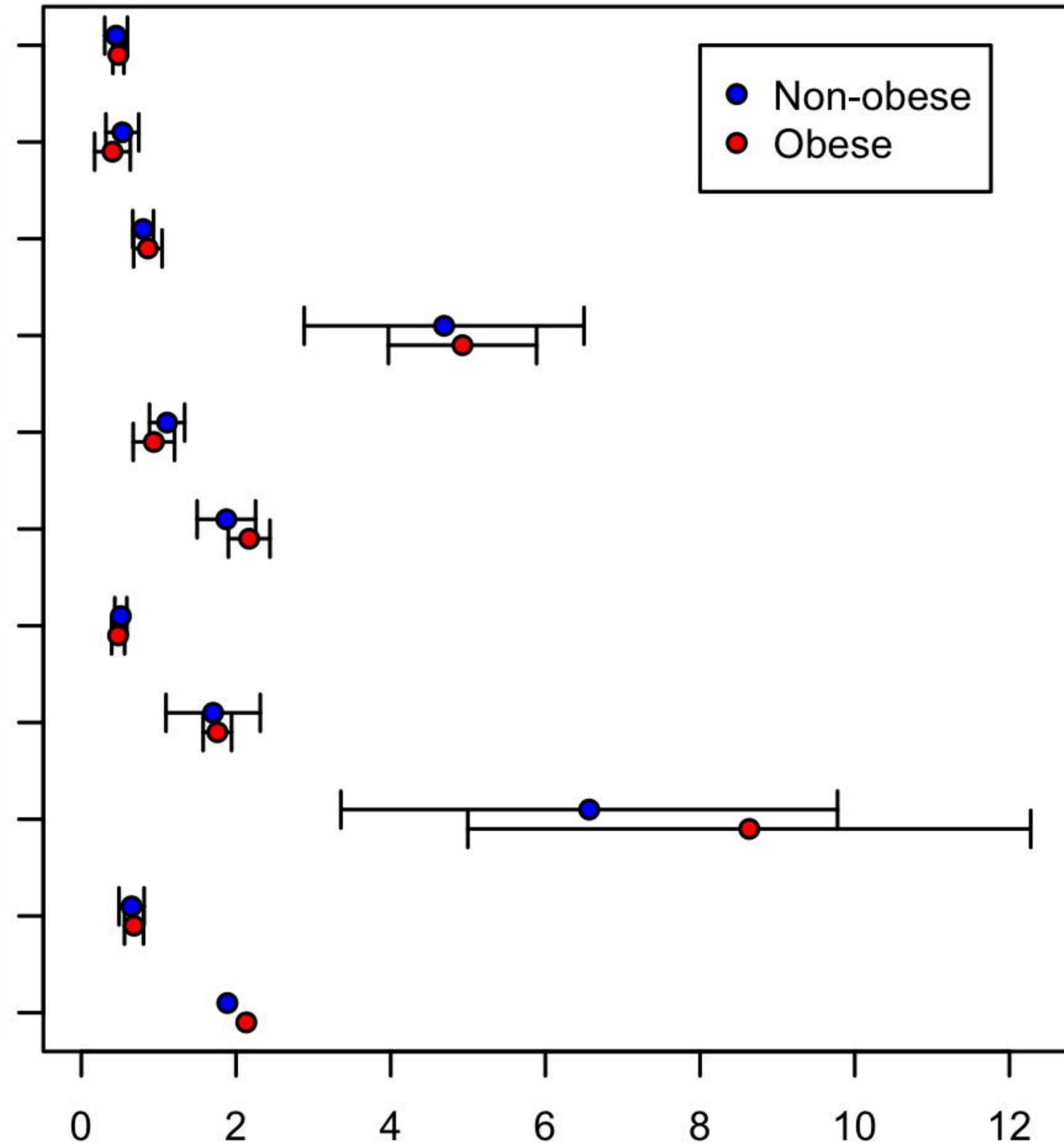
Studies included in  
qualitative synthesis  
(n = 11)

No obese individuals as  
measured by BMI  
(n = 1)

Studies included in  
quantitative synthesis  
(meta-analysis)  
(n = 10)

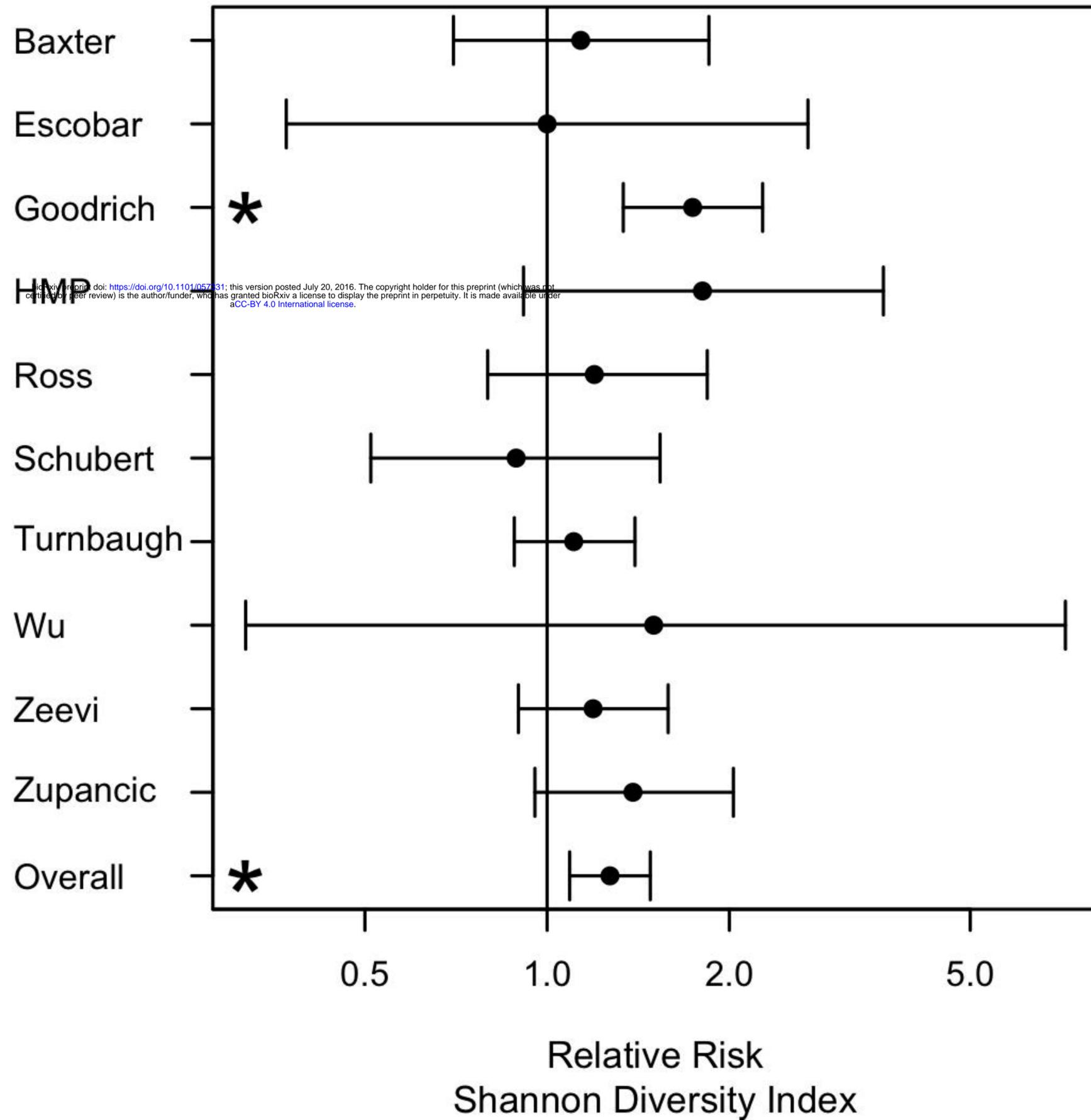
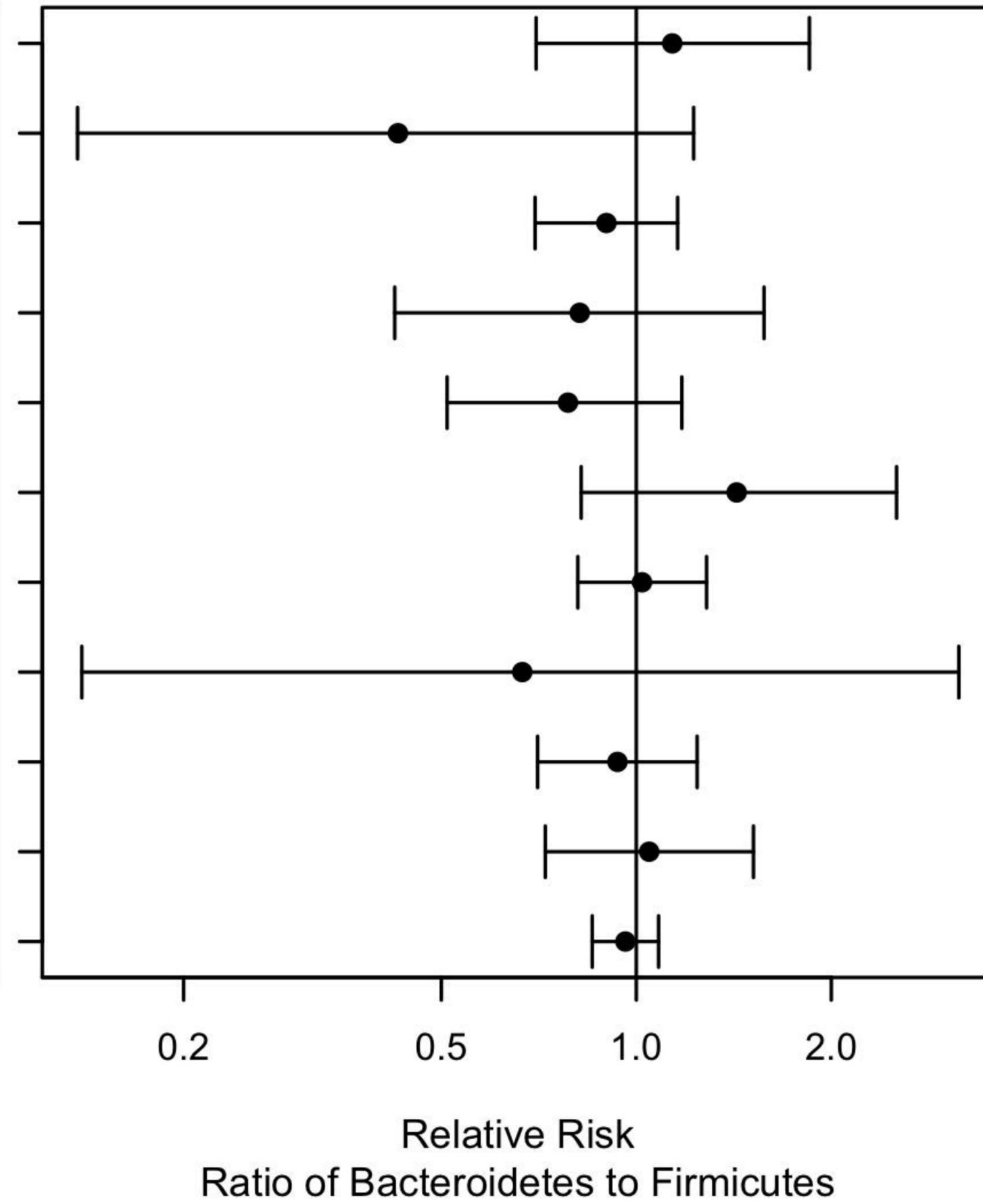
**A**

bioRxiv preprint doi: <https://doi.org/10.1101/057231>; this version posted July 20, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

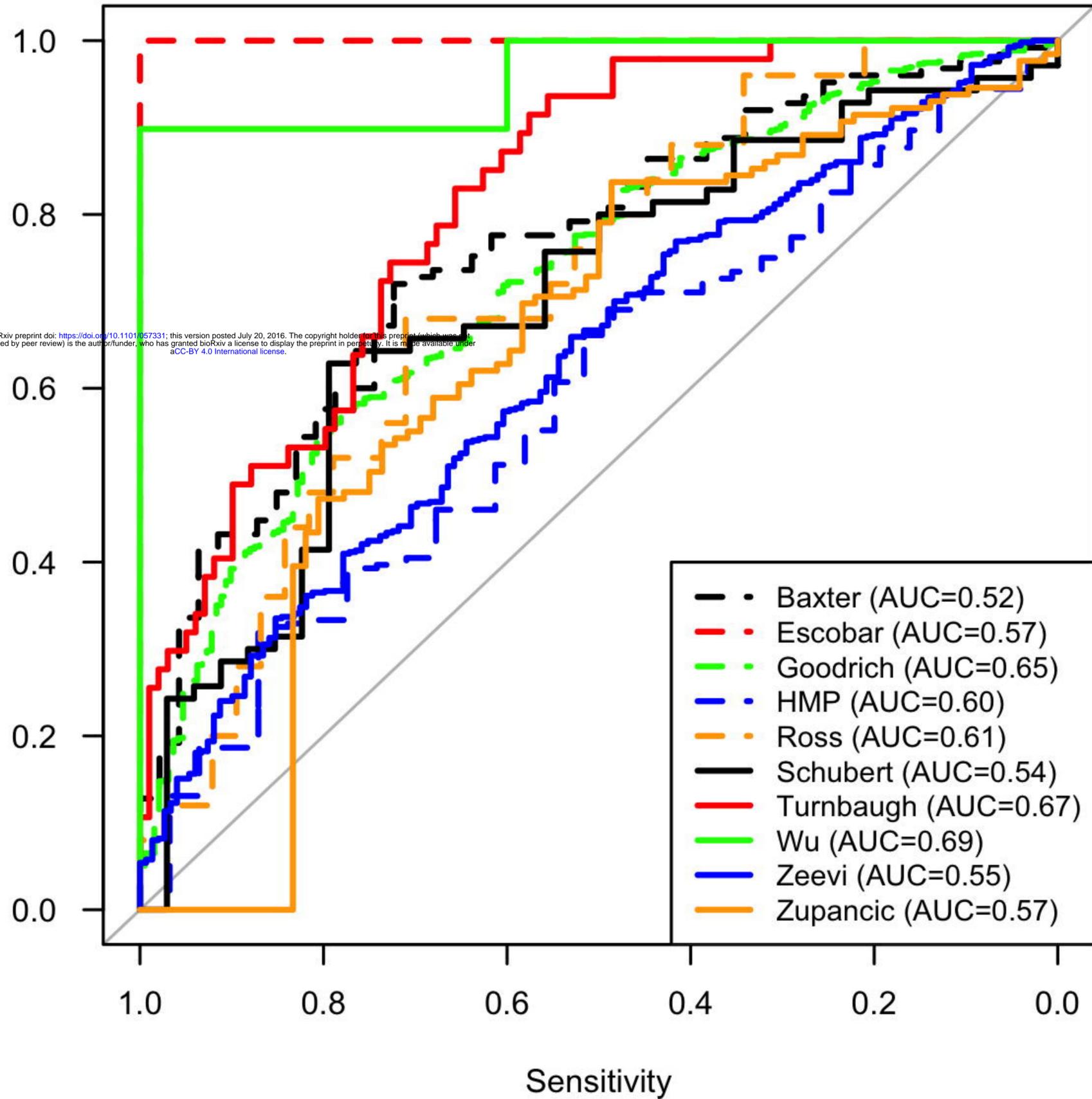
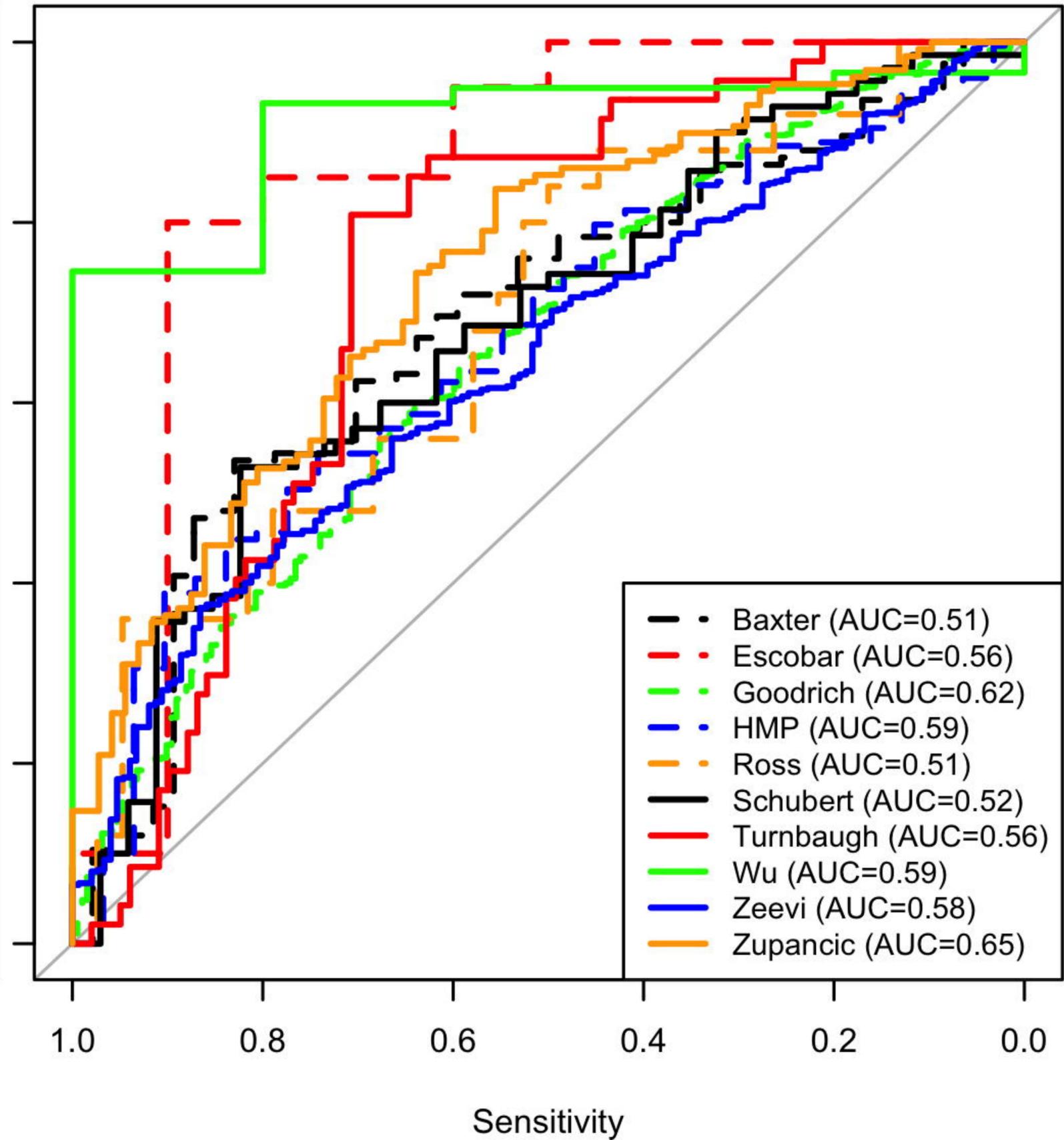
**B**

Shannon Diversity Index

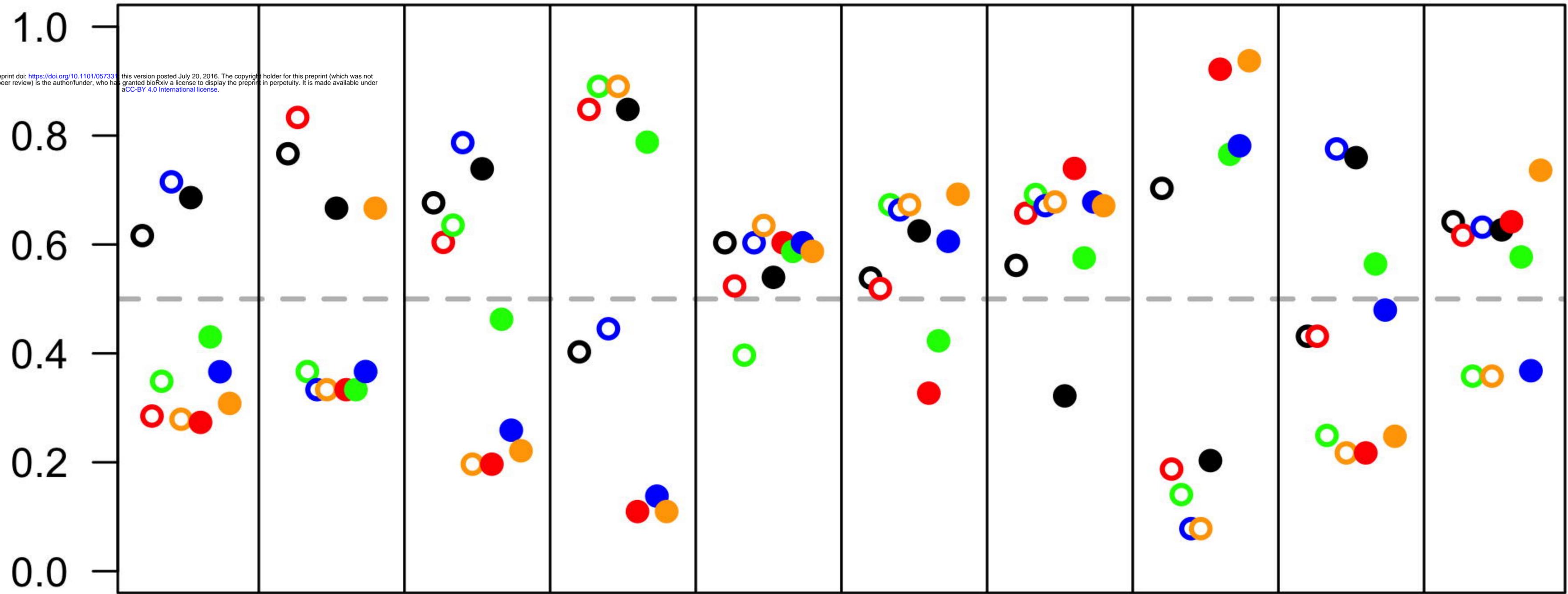
Ratio of Bacteroidetes to Firmicutes

**A****B**

bioRxiv preprint doi: <https://doi.org/10.1101/057331>; this version posted July 20, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

**A****B**

Accuracy



Baxter

Escobar

Goodrich

HMP

Ross

Schubert

Turnbaugh

Wu

Zeevi

Zupancic

