

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

Lack of power and small effect size confounds the ability to differentiate non-obese and obese individuals using gut microbiome data

Running Title: The Human Microbiome and Obesity

Marc A Sze and Patrick D Schloss[†]

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

[†] To whom correspondence should be addressed: pschloss@umich.edu

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

Abstract

Two recent studies have re-analyzed published data and found that when datasets are analyzed independently there was limited support the widely accepted hypothesis that changes in the microbiome are associated with obesity. This hypothesis was reconsidered by increasing the number of data sets and pooling the results across the individual datasets. The Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines were applied to identify 10 studies for an updated and more synthetic analysis. Alpha diversity metrics and the relative risk of obesity based on those metrics were used to identify a limited number of significant associations with obesity; however, when the results of the studies were pooled using a random effects model significant associations were observed between Shannon diversity, number of observed OTUs, and Shannon evenness and obesity status. They were not observed for the ratio of *Bacteroidetes* and *Firmicutes* or their individual relative abundances. Although these tests yielded small P-values, the difference between the Shannon diversity index of non-obese and obese individuals was 2.07%. A power analysis demonstrated that one of the studies had sufficient power to detect a 5% difference in diversity. When models trained on one dataset were then tested using the other 9 datasets, the median accuracy varied between 33.01 and 64.77% (median=56.67%). Although there is statistical support for a relationship between the microbial communities found in human feces and obesity status, this association is relatively weak and its detection is confounded by large interpersonal variation and insufficient sample sizes.

Importance

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

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

Introduction

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

There have been several studies that report observing a link between the composition of microbiome and obesity in animal models and in humans. The first such study used genetically obese mice and observed the ratio of the relative abundances of Bacteroidetes to Firmicutes (B:F) was lower in obese mice than lean mice (5). Translation of this result to humans by the same researchers did not observe this effect, but did find that obese individuals had a lower diversity than lean individuals (6). They also showed that the relative abundance of Bacteroidetes and Firmicutes increased and decreased, respectively, as obese individuals lost weight while on a fat or carbohydrate restricted diet (7). Two re-analysis studies interrogated previously published microbiome and obesity data and concluded that the previously reported differences in community diversity and B:F among non-obese and obese individuals could not be generalized (8, 9). Regardless of the results using human populations, mechanistic studies using animal models that were manipulated with antibiotics or colonization with varied communities appears to support the association since direct manipulation of the communities yielded variation in animal weight (10–13). The purported association between the differences in the microbiome and obesity have been widely repeated with little attention given to the lack of a clear signal in human cohort

studies.

The recent publication of additional studies that collected BMI data for each subject as well as other studies that were not included in the earlier re-analysis studies offered the opportunity to revisit the question relating the structure of the human microbiome to obesity (14–22). One critique of the prior re-analysis studies is that the authors did not aggregate the results across studies to increase the effective sample size. It is possible that there were small associations within each study that were not statistically significant because the individual studies lacked sufficient power. Alternatively, diversity metrics may mask the appropriate signal and it is necessary to measure the association at the level of microbial populations. Walters et al. (8) demonstrated that Random Forest machine learning models were capable of predicting obesity status within a single cohort, but did not attempt to test the models on other cohorts. The purpose of this study was to perform a meta-analysis of the association between differences in the microbiome and obesity status by analyzing and applying a more systematic and synthetic approach than was used previously.

Methods

Literature Review and Study Inclusion. We followed the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines to identify studies to include in our meta-analysis (23). A detailed description of our selection process and the exact search terms are provided in the Supplement and in figure 1. Briefly, we searched PubMed for original research studies that involved studying obesity and the human microbiome. The initial search yielded 187 studies. We identified number_string[n] additional studies that were not designed to explicitly test for an association between the microbiome and obesity. We then manually curated the 196 studies to select those studies that included BMI and sequence data. This yielded 10 eligible studies. An additional study was removed

from our analysis because no individuals in the study had a BMI over 30. Among the final 10 studies, 3 were from identified from our PubMed search (6, 15, 20), 5 were originally identified from the 9 studies that did not explicitly investigate obesity but included BMI data (14, 18, 19, 24, 25), and two datasets were used (21, 22) because these publications did not specifically look for any metabolic or obesity conditions but had control populations and enabled us to help mitigate against publication biases associated with the bacterial microbiome and obesity. The number_string[n] studies are summarized in Tables 1 and 2.

Sequence Analysis Pipeline. All sequence data were publicly available and were downloaded from the NCBI Sequence Read Archive, the European Nucleotide Archive, or the investigators' personal website (https://gordonlab.wustl.edu/TurnbaughSE/_10/_09/STM/_2009.html). In total seven studies used 454 (6, 14, 15, 19, 20, 22, 25) and three studies used Illumina sequencing (18, 21, 24). All of these studies used amplification-based 16S rRNA gene sequencing. Among the studies that sequenced the 16S rRNA gene, the researchers targeted the V1-V2 (19), V1-V3 (14, 15, 20), V3-V5 (22, 25), V4 [(18); (21);], and V3-4 (24) regions. For those studies where multiple regions were sequenced, we selected the region that corresponded to the largest number of subjects (6, 25). We processed the 16S rRNA gene sequence data using a standardized mothur pipeline. Briefly, our pipelines attempted to follow previously recommended approaches for 454 and Illumina sequencing data (26, 27). All sequences were screened for chimeras using UCHIME and assigned to operational taxonomic units (OTUs) using the average neighbor algorithm using a 3% distance threshold (28, 29). All sequence processing was performed using mothur (v.1.37.0) (30).

Data Analysis. We split the overall meta-analysis into three general strategies using R (3.3.0). First, we followed the approach employed by Finucane et al (9) and Walters et al (8) where each study was re-analyzed separately to identify associations between BMI and the relative abundance of Bacteroidetes and Firmicutes, the ratio of Bacteroidetes and

Firmicutes relative abundances (B:F), Shannon diversity, observed richness, and Shannon evenness. After each variable was transformed to fit a normal distribution a two-tailed t-test was performed for comparison of non-obese and obese individuals (i.e. BMI > 35.0). We performed a pooled analysis on these measured variables using linear random effect models to correct for study effect to assess differences on the combined dataset between non-obese and obese groups using the lme4 (v.1.1-12) R package. Next, we compared the community structure from non-obese and obese individuals using PERMANOVA analysis of Bray-Curtis distance matrices. This analysis was performed using the vegan (v.2.3-5) R package. For both analyses, the datasets were rarefied (N=1000) so that each study within a study had the same number of sequences. Second, for each study we partitioned the subjects into a low or high group depending on whether their alpha diversity metrics were below or above the median value for the study. The relative risk (RR) was then calculated as the ratio of the number of obese individuals in the low group to the number of obese individuals in the high group. We then performed a Fisher exact-test to investigate whether the RR was significantly different from 1.0 within each study and across all of the studies using the epiR (0.9-77) and metafor (1.9-8) packages. Third, we used the AUCRF (1.1) R package to generate Random Forest models. For each study we developed models using either OTUs or genus-level phylotypes. The quality of each model was assessed by measuring the area under the curve (AUC) of the Receiver Operating Characteristic (ROC) using ten-fold cross validation. Because the genus-level phylotype models were developed using a common reference, it was possible to use one study's model (i.e. the training set) to classify the samples from the other studies (i.e. the testing sets). The optimum threshold for the training set was set as the probability threshold that had the highest combined sensitivity and specificity. This threshold was then used to calculate the accuracy of the model applied to the test studies. To generate Receiver Operator Characteristic (ROC) curves and calculate the accuracy of the models we used the pROC (1.8) R package. Finally, we performed power and sample number simulations for different effect sizes for

each study using the pwr (1.1-3) R package and base R functions. We also calculated the actual sample size needed based on the effect size of each individual study.

Reproducible methods. A detailed and reproducible description of how the data were processed and analyzed can be found at https://github.com/SchlossLab/Size_Obesity_mBio_2016/.

Results

Alpha diversity analysis. We calculated the Shannon diversity index, richness, and Shannon evenness, the relative abundance of *Bacteroidetes* and *Firmicutes*, and the ratio of their relative abundance (B:F) for each sample. Once we transformed each of the six alpha diversity metrics to make them normally distributed, we used a t-test to identify significant associations between the alpha diversity metric and whether an individual was obese for each of the ten studies. The B:F and the relative abundance of *Firmicutes* were not significantly associated with obesity in any study. We identified 7 P-values less than 0.05: three studies indicated obese individuals had a lower richness, two studies indicated a significantly lower diversity, one study indicated a significantly lower evenness, and one study indicated a significantly higher relative abundance of *Bacteroidetes* (Figures 2 and S1). These results largely match those of the Walters and Finucane re-analysis studies. Interestingly, although only two of the 10 studies observed the previously reported association between lower diversity and obesity, the other studies appeared to have the same trend, albeit the differences were not statistically significant. We used a random effects linear model to combine the studies using the study as the random effect and found statistical support for decreased richness, evenness, and diversity among obese individuals (all $P < 0.011$). Although there was a significant relationship between these metrics and diversity and obesity status, the effect size was quite small. The obese individuals averaged

160 7.47% lower richness, 0.88% lower evenness, and 2.07% lower diversity. There were no
161 significant associations when we pooled the phylum-level metrics across studies. These
162 results indicate that obese individuals do have a statistically significant lower diversity than
163 non-obese individuals; however, it is questionable whether the difference is biologically
164 significant.

165 **Relative risk.** Building upon the alpha diversity analysis we calculated the relative risk of
166 being obese based on whether an individual's alpha diversity metrics were below or above
167 the median metric for that study. The results using relative risk largely matched those of
168 using the untransformed alpha diversity data. Across the number_string[n] studies and six
169 metrics, the only significant relative risk values were the richness, evenness, and diversity
170 values from the Goodrich study (Figures 3 and S2). Again, although the relative risk values
171 were not significant for other studies, the values tended to be above one. When we pooled
172 the data using a random effects model, the relative risk associated with having a richness,
173 evenness, or diversity below the median for the population was significantly associated
174 with obesity (all $P < 0.0044$). The relative risks associated with alpha diversity were small.
175 The relative risk of having a low richness was 1.30 (1.13-1.49), low evenness was 1.20
176 (1.06-1.37), and low diversity was 1.27 (1.09-1.48). There were no significant difference
177 in the phylum-level metrics. Again, the relative risk results indicate that individuals with a
178 lower richness, evenness, or diversity are at statistically significant increased risk of being
179 obese, it is questionable whether that risk is biologically or clinically relevant.

180 **Beta diversity analysis.** Following the approach used by the Walters and Finucane
181 re-analysis studies, for each dataset we calculated a Bray-Curtis distance matrix to
182 measure the difference in the membership and structure of the individuals from each
183 study. We then used PERMANOVA to test for a significant differences between the
184 structure of non-obese and obese individuals. The Escobar, Goodrich, and Turnbaugh
185 datasets indicated a significant difference in community structure (all $P < 0.05$). Because

186 it was not possible to ascertain the directionality of the difference in community structure
187 nor perform a pooled analysis using studies that had non-overlapping 16S rRNA gene
188 sequence regions it is unclear whether these differences reflect a broader, but perhaps
189 small, shift in community structure between non-obese and obese individuals.

190 ***Development of a microbiome-based classifier of obesity.*** The Walters re-analysis
191 study suggested that it was possible to classify individuals as being obese or non-obese
192 based on the composition of their microbiota. We repeated this analysis with additional
193 datasets using OTU and genus-level phylotype data. For each study we developed a
194 Random Forest machine learning model to classify individuals. Using ten-fold cross
195 validation, the observed AUC values varied between 0.52 and 0.69 indicating a relatively
196 poor ability to classify individuals (Figure 4A). So that we could test models on other
197 datasets, we trained models using genus-level phylotype data for each dataset. The the
198 observed AUC values for the models applied to the training datasets varied between
199 0.51 and 0.65, again indicating a relatively poor ability to classify individuals from the
200 original dataset (Figure 4B). For each model we identified the probability where the sum of
201 the sensitivity and specificity was the highest. We then used this to identify a threshold
202 to calculate the accuracy of the models when applied to the other number_string[n-1]
203 datasets (Figure 5). Although there considerable variation in accuracy values for each
204 model, the median accuracy for each model varied between 0.33 (Turnbaugh) and 0.65
205 (HMP) (median=0.57). When we considered the number of samples, balance of non-obese
206 and obese individuals, and region within the 16S rRNA gene it was not possible to identify
207 factors that predictably affected model performance. The ability to predict obesity status
208 using the relative abundance of OTUs and genera in the communities is only marginally
209 better than random. These results suggest that given the large diversity of microbiome
210 compositions it is difficult to identify a taxonomic signal that can be associated with obesity.

211 ***Power and Sample Size Estimate Simulations.*** The inability to detect a difference

between non-obese and obese individuals could be due to the lack of a true effect or because the study had insufficient statistical power to detect a difference because of insufficient sampling, large interpersonal variation, and unbalanced sampling of non-obese and obese individuals. To assess this, we calculated the power to detect differences of 1, 5, 10, and 15% in each of the alpha diversity metrics using the sample sizes used in each of the studies (Figures 6, S3-S8). Although there is no biological rationale for these effect sizes, they represent a range that is plausible. Only the Goodrich study had power greater than 0.80 to detect a 5% difference in Shannon diversity and six of the studies had enough power to detect a 10% difference (Figure 6). None of the studies had sufficient power to detect a 15% difference between B:F values (Figure S5). In fact, the maximum power among any of the studies to detect a 15% difference in B:F values was 0.25. Among the tests for relative risk, none of the studies had sufficient power to detect a Cohen's d of 0.10 and only two studies had sufficient power to detect a Cohen's d of 0.15. We next estimated how many individuals would need to have been sampled to have sufficient power to detect the four effect sizes assuming the observed interpersonal variation from each study and balanced sampling between the two groups. To detect a 1, 5, 10, and 15% difference in Shannon index, the median sampling effort per group was approximately 3,400, 140, 35, and 16 individuals, respectively. To detect a 1, 5, 10, and 15% difference in B:F values, the median sampling effort per group was approximately 160,000, 6,300, 1,600, and 700 individuals, respectively. To detect a 1, 5, 10, and 15% difference in relative risk values using Shannon diversity, the median sampling effort per group was approximately 39,000, 1,500, 380, and 170 individuals, respectively. These estimates indicate that most microbiome studies are underpowered to detect modest effect sizes using either metric. In the case of obesity, the studies were underpowered to detect the `signif(range_effect_size, 1)[1]` to 6% difference in diversity that was observed across the studies.

Discussion

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

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

upon an accepted practice.

By selecting a range of possible effect sizes, we were able to demonstrate that most studies are underpowered to detect modest differences in alpha diversity metrics and phylum-level relative abundances. Several factors interact to limit the power of microbiome studies. There is wide interpersonal variation in the diversity and structure of the human microbiome. In addition, the common experimental designs limit their power. As we observed, most of the studies included in our analysis were unbalanced for the variable that we were interested in. This was also true of those studies that originally sought to identify associations with obesity. Even with a balanced design, we showed that it was necessary to obtain approximately 140 and 6,300 sequences per sample to detect a 5% difference in Shannon diversity or B:F, respectively. It was interesting that these sample sizes agreed across studies regardless of their sequencing method, region within the 16S rRNA gene, or subject population (Figure 6). This suggests that regardless of the treatment or category, these sample sizes represent a good starting point for subject recruitment when using stool samples. Unfortunately, few studies have been published with this level of subject recruitment. This is troubling since the positive predictive rate of a significant finding in an underpowered study is small leading to results that cannot be reproduced (31). Future microbiome studies should articulate the basis for their experimental design.

Two previous reviews (8, 9) have stated that there was not a consistent association between alpha diversity and obesity; however, neither of these studies made an attempt to pool the existing data together to try and harness the additional power that this would give and they did not assess whether the studies were sufficiently powered to detect a difference. Our analysis also used 16S rRNA gene sequence data from ten studies whereas the Finucane study used 16S rRNA gene sequence data from 3 studies (7, 10, 25) and a metagenomic study (32) and the Walters study used 16S rRNA gene sequence data from 5 studies (10, 15, 19, 25, 33); two studies were included in both analyses (10, 25). Our analysis included

4 of these studies (10, 15, 19, 25) and excluded 3 of the studies because they were too small (7), only utilized metagenomic data (32), or used short single read Illumina HiSeq data that has a high error rate making it untractable for *de novo* OTU clustering (33). The additional seven datasets were published after the two reviews were performed and include datasets with more samples than were found in the original studies. Our collection of ten studies allowed us to largely use the same sequence analysis pipeline for all datasets and relied heavily on the availability of public data and access to metadata that included variables beyond the needs of the original study. To execute this analysis, we created an automated data analysis pipeline, which can be easily updated to add additional studies as they become available. Similarly, it would be possible to adapt this pipeline to other body sites and treatment or variables (e.g. subject's sex or age).

Similar to our study, the Walters analysis (8), the authors generated Random Forest machine learning models to differentiate between non-obese and obese individuals. They obtained similar AUC values to our analysis; however, they did not attempt to test these models on the other studies in their analysis. When we performed the inter-dataset cross validation the median accuracy across datasets was only 56.67% indicating that the models did a poor job when applied to other datasets. This could be due to differences in subject populations and methods. Considering the median AUC for models trained and tested on the same data with ten-fold cross validation only varied between 0.51 and 0.65 and there was not a strong signal in the alpha diversity data, we suspect that there is insufficient signal to reliably classify individuals.

Although we failed to find an effect it is not realistic to necessarily state that there is no microbiome impact on obesity. There is strong evidence in murine models of obesity that the microbiome and level of adiposity can be manipulated via genetic manipulation of the animal and manipulation of the community through antibiotics or colonizing germ free mice with diverse fecal material from human donors (5, 10–13). These studies appear to conflict

with the observations using human subjects. Recalling the large interpersonal variation in the structure of the microbiome, it is possible that each individual has their own signatures of obesity. Alternatively, it could be that the involvement of the microbiome in obesity is at the level of a common set of metabolites that can be produced from different structures of the microbiome.

Acknowledgements

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

Table 1. Summary Demographics of Individuals used in the Meta-analysis.

```

##
## > capwords <- function(s, strict = FALSE) {
## +   s <- as.character(s)
## +   cap <- function(s) paste(toupper(substring(s, 1, 1)), {
## +     s <- su .... [TRUNCATED]
##
## > make_study_label <- function(dataset) {
## +   study <- ifelse(dataset == "hmp", "HMP", capwords(dataset))
## + }
##
## > format_p <- function(p_value) {
## +   char_p <- format(round(p_value, 3), nsmall = 3)
## +   if (p_value < 0.001) {
## +     char_p <- "<0.001"
## +     .... [TRUNCATED]
##
## > get_study_summary <- function(study, beta = beta_summary) {
## +   metadata_file <- paste0("data/", study, "/", study, ".metadata")
## +   metadata <- .... [TRUNCATED]
##
## > beta_summary <- read.table("data/process/beta_tests.summary",
## +   header = T, row.names = 1)
##
## > datasets <- sort(rownames(beta_summary))
##

```



```

354 ## > study_summary <- t(sapply(datasets, get_study_summary))
355 ##
356 ## > table1 <- kable(study_summary, row.names = FALSE,
357 ## +      col.names = c("Study", "Subjects (N)", "Obese (%)", "Average BMI (Min-Max)",
358 ## +      "Fe ..." ... [TRUNCATED])

```

Study	Subjects (N)	Obese (%)	Average BMI (Min-Max)	Female (%)	Average Age (Min-Max)
Baxter	172	27.3	27.0 (17.5-46.9)	64.5	54.3 (29.0-80.0)
Escobar	30	33.3	27.4 (19.5-37.6)	46.7	38.1 (21.0-60.0)
Goodrich	982	19.7	26.3 (16.2-52.4)	98.9	61.0 (23.0-86.0)
Hmp	287	10.8	24.3 (19.0-34.0)	49.1	26.3 (18.0-40.0)
Ross	63	60.3	31.6 (22.1-47.9)	76.2	57.0 (33.0-81.0)
Schubert	104	32.7	28.2 (18.5-62.5)	66.3	52.8 (19.0-88.0)
Turnbaugh	146	67.8	NA	NA	NA
Wu	64	7.8	24.3 (14.0-41.3)	53.1	26.3 (2.16-50.0)
Zeevi	731	NA	26.4 (16.4-47.0)	NA	43.4 (18.0-70.0)
Zupancic	207	36.2	28.2 (18.2-127.0)	57.0	46.7 (20.0-79.0)

Figure 1: PRISMA flow diagram of total records searched (34).

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

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

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

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

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

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

Figure S2: Meta Analysis of the Relative Risk of Obesity Based on Evenness (A), Richness (B), or the Relative Abundance of Bacteroidetes (C) and Firmicutes (D).

Figure S3: Power (A) and sample size simulations (B) for B:F for differentiating between non-obese versus obese for effect sizes of 1, 5, 10, and 15%. Power calculations use the sampling distribution from the original studies and the sample size estimations assume an equal amount of sampling from each treatment group.

Figure S4: Power (A) and sample size simulations (B) for richness for differentiating between non-obese versus obese for effect sizes of 1, 5, 10, and 15%. Power calculations use the sampling distribution from the original studies and the sample size estimations assume an equal amount of sampling from each treatment group.

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

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

Figure S7: Power (A) and sample size simulations (B) for the relative abundance of

Firmicutes for differentiating between non-obese versus obese for effect sizes of 1, 5, 10, and 15%. Power calculations use the sampling distribution from the original studies and the sample size estimations assume an equal amount of sampling from each treatment group.

Figure S8: Power (A) and sample size simulations (B) for relative risk of obesity based on Shannon diversity. Power calculations use the sampling distribution from the original studies and the sample size estimations assume an equal amount of sampling from each treatment group.

References

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

428 4441022a.

429 8. **Walters WA, Xu Z, Knight R.** 2014. Meta-analyses of human gut microbes associated
430 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)
431 [2014.09.039](http://doi.org/10.1016/j.febslet.2014.09.039).

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

435 10. **Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JI.** 2006.
436 An obesity-associated gut microbiome with increased capacity for energy harvest. *Nature*
437 **444**:1027–31. doi:<http://doi.org/10.1038/nature05414>.

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

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

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

450 14. **Ross MC, Muzny DM, McCormick JB, Gibbs RA, Fisher-Hoch SP, Petrosino JF.**
451 2015. 16S gut community of the Cameron County Hispanic Cohort. *Microbiome* **3**:7.

doi:<http://doi.org/10.1186/s40168-015-0072-y>.

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

16. Nam Y-D, Jung M-J, Roh SW, Kim M-S, Bae J-W. 2011. Comparative analysis of Korean human gut microbiota by barcoded pyrosequencing. *PloS One* 6:e22109. doi:<http://doi.org/10.1371/journal.pone.0022109>.

17. Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, Fernandes GR, Tap J, Bruls T, Batto J-M, Bertalan M, Borruel N, Casellas F, Fernandez L, Gautier L, Hansen T, Hattori M, Hayashi T, Kleerebezem M, Kurokawa K, Leclerc M, Levenez F, Manichanh C, Nielsen HB, Nielsen T, Pons N, Poulain J, Qin J, Sicheritz-Ponten T, Tims S, Torrents D, Ugarte E, Zoetendal EG, Wang J, Guarner F, Pedersen O, Vos WM de, Brunak S, Doré J, MetaHIT Consortium, Antolín M, Artiguenave F, Blottiere HM, Almeida M, Brechot C, Cara C, Chervaux C, Cultrone A, Delorme C, Denariáz G, Dervyn R, Foerstner KU, Friss C, Guchte M van de, Guedon E, Haimet F, Huber W, Hylckama-Vlieg J van, Jamet A, Juste C, Kaci G, Knol J, Lakhdari O, Layec S, Le Roux K, Maguin E, Mérieux A, Melo Minardi R, M'rini C, Muller J, Oozeer R, Parkhill J, Renault P, Rescigno M, Sanchez N, Sunagawa S, Torrejon A, Turner K, Vandemeulebrouck G, Varela E, Winogradsky Y, Zeller G, Weissenbach J, Ehrlich SD, Bork P. 2011. Enterotypes of the human gut microbiome. *Nature* 473:174–180. doi:<http://doi.org/10.1038/nature09944>.

18. Goodrich JK, Waters JL, Poole AC, Sutter JL, Koren O, Blekhman R, Beaumont M, Van Treuren W, Knight R, Bell JT, Spector TD, Clark AG, Ley RE. 2014. Human

genetics shape the gut microbiome. *Cell* **159**:789–799. doi:<http://doi.org/10.1016/j.cell.2014.09.053>.

19. **Wu GD, Chen J, Hoffmann C, Bittinger K, Chen Y-Y, Keilbaugh SA, Bewtra M, Knights D, Walters WA, Knight R, Sinha R, Gilroy E, Gupta K, Baldassano R, Nessel L, Li H, Bushman FD, Lewis JD.** 2011. Linking long-term dietary patterns with gut microbial enterotypes. *Science (New York, NY)* **334**:105–108. doi:<http://doi.org/10.1126/science.1208344>.

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

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

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

23. **Moher D, Liberati A, Tetzlaff J, Altman DG, PRISMA Group.** 2010. Preferred reporting items for systematic reviews and meta-analyses: The PRISMA statement. *International Journal of Surgery (London, England)* **8**:336–341. doi:<http://doi.org/10.1016/j.ijssu.2010.02.007>.

24. **Zeevi D, Korem T, Zmora N, Israeli D, Rothschild D, Weinberger A, Ben-Yacov O, Lador D, Avnit-Sagi T, Lotan-Pompan M, Suez J, Mahdi JA, Matot E, Malka G, Kosower N, Rein M, Zilberman-Schapira G, Dohnalová L, Pevsner-Fischer M,**

Bikovsky R, Halpern Z, Elinav E, Segal E. 2015. Personalized Nutrition by Prediction of Glycemic Responses. *Cell* **163**:1079–1094. doi:<http://doi.org/10.1016/j.cell.2015.11.001>.

25. Human Microbiome Project Consortium. 2012. Structure, function and diversity of the healthy human microbiome. *Nature* **486**:207–214. doi:<http://doi.org/10.1038/nature11234>.

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

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

28. Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R. 2011. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics* **27**:2194–2200. doi:<http://doi.org/10.1093/bioinformatics/btr381>.

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

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

31. Ioannidis JPA. 2005. Why most published research findings are false. *PLoS Med*

2:e124. doi:<http://doi.org/10.1371/journal.pmed.0020124>.

32. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez F, Yamada T, Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H, Xie Y, Tap J, Lepage P, Bertalan M, Batto J-M, Hansen T, Paslier DL, Linneberg A, Nielsen HB, Pelletier E, Renault P, Sicheritz-Ponten T, Turner K, Zhu H, Yu C, Li S, Jian M, Zhou Y, Li Y, Zhang X, Li S, Qin N, Yang H, Wang J, Brunak S, Doré J, Guarner F, Kristiansen K, Pedersen O, Parkhill J, Weissenbach J, Antolin M, Artiguenave F, Blottiere H, Borruel N, Bruls T, Casellas F, Chervaux C, Cultrone A, Delorme C, Denariáz G, Dervyn R, Forte M, Friss C, Guchte M van de, Guedon E, Haimet F, Jamet A, Juste C, Kaci G, Kleerebezem M, Knol J, Kristensen M, Layec S, Roux KL, Leclerc M, Maguin E, Minardi RM, Oozeer R, Rescigno M, Sanchez N, Tims S, Torrejon T, Varela E, Vos W de, Winogradsky Y, Zoetendal E, Bork P, Ehrlich SD, Wang J. 2010. A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* **464**:59–65. doi:<http://doi.org/10.1038/nature08821>.

33. Yatsunenko T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M, Magris M, Hidalgo G, Baldassano RN, Anokhin AP, Heath AC, Warner B, Reeder J, Kuczynski J, Caporaso JG, Lozupone CA, Lauber C, Clemente JC, Knights D, Knight R, Gordon JI. 2012. Human gut microbiome viewed across age and geography. *Nature* **486**:222–227. doi:<http://doi.org/10.1038/nature11053>.

34. Moher D, Liberati A, Tetzlaff J, Altman DG. 2009. Preferred reporting items for systematic reviews and meta-analyses: The PRISMA statement. *PLoS Med* **6**:e1000097. doi:<http://doi.org/10.1371/journal.pmed.1000097>.