

1 **Title:** Correcting for batch effects in case-control microbiome studies.

2

3 **Authors:** Sean M. Gibbons^{1,2,3}, Claire Duvallet^{1,2}, and Eric J. Alm^{1,2,3}

4

5 **Affiliations:** ¹ Department of Biological Engineering, Massachusetts Institute of
6 Technology, Cambridge, MA, USA; ² Center for Microbiome Informatics and
7 Therapeutics, Cambridge, MA, USA; ³ The Broad Institute of MIT and Harvard,
8 Cambridge, MA, USA

9

10 **Abstract**

11

12 High-throughput data generation platforms, like mass-spectrometry, microarrays,
13 and second-generation sequencing are susceptible to batch effects due to run-to-
14 run variation in reagents, equipment, protocols, or personnel. Currently, batch
15 correction methods are not commonly applied to microbiome sequencing
16 datasets. In this paper, we compare multiple batch-correction methods applied to
17 microbiome case-control studies. We introduce a model-free normalization
18 procedure where features (i.e. bacterial taxa) in case samples are converted to
19 percentiles of the equivalent features in control samples within a study prior to
20 pooling data across studies. We look at how this percentile-normalization method
21 compares to ComBat, a widely used batch-correction model developed for RNA
22 microarray data, and traditional meta-analysis methods for combining
23 independent p-values. Overall, we show that percentile-normalization is a simple,
24 model-free approach for removing batch effects and improving sensitivity in case-
25 control meta-analyses.

26

27 **Author Summary**

28

29 Batch effects present a significant obstacle to comparing results across
30 independent studies. Traditional meta-analysis techniques for combining p-
31 values from independent studies, like Fisher's method, are effective, but
32 statistically conservative. If batch-effects can be corrected, then statistical tests
33 can be performed on data pooled across studies, increasing sensitivity to detect
34 differences between treatment groups. Here, we show how a simple, model-free
35 approach corrects for batch effects in case-control datasets.

36

37 **Introduction**

38 Data generated by high throughput methods like mass-spectrometry, second-
39 generation sequencing, or microarrays are sensitive to experimental and
40 computational processing [1]. This sensitivity gives rise to 'batch effects' between
41 independent runs of an experiment. Even when different research groups adhere

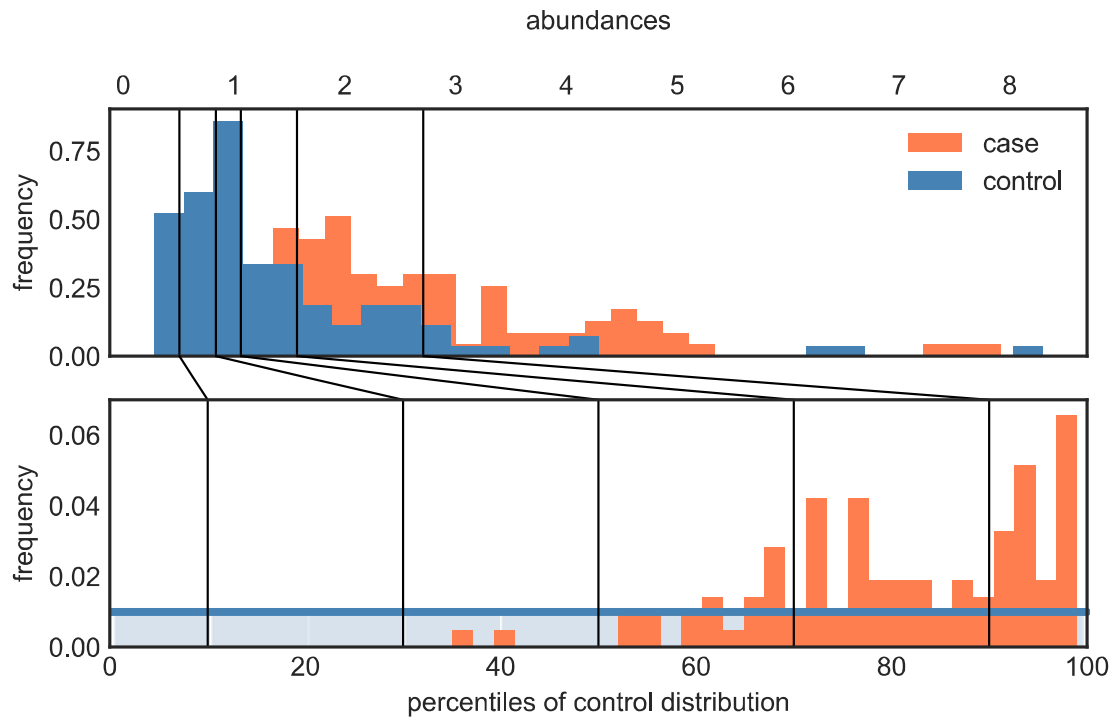
42 to the same methodologies, these effects can arise due to slight differences in
43 hardware, reagents, or personnel [2]. Thus, it is inadvisable to directly compare
44 non-normalized data across studies.

45 Several tools for reducing batch effects in RNA expression microarray
46 data have been developed. For example, surrogate variable analysis (SVA)
47 estimates a set of inferred variables (eigenvectors) that explain variance
48 associated with putative batch effects [3]. These inferred variables are then
49 incorporated into a linear model to correct downstream significance tests. SVA is
50 part of a family of batch-correction methods that use different varieties of factor
51 analysis or singular value decomposition [3-5]. The most relied upon method to
52 date [6], called ComBat, uses a Bayesian approach to estimate location and
53 scale parameters for each feature within a batch [7]. These methods are most
54 effective when batch effects are not conflated with the true biological effects [1].
55 Furthermore, these methods work best when batch effects are not diffuse and
56 can be projected onto a low-dimensional manifold.

57 Unfortunately, batch effects are often diffuse and conflated with biological
58 signals [8-10], limiting the usefulness of these methods. This is especially true for
59 low-biomass samples in microbiome sequencing studies, like samples taken from
60 the built environment [11], where the biological signal is relatively weak and
61 batch effects can be quite large [12]. One way to get around this issue is to
62 calculate statistics within a given batch, and then compare significant features
63 across batches using classic meta-analysis techniques for combining p-values,
64 like Fisher's and Stouffer's methods [13, 14]. These meta-analysis techniques

65 are robust to batch effects across independent studies. However, in cases where
66 pooling raw data across studies might increase statistical power to detect subtle
67 differences or in cases where batches are not statistically independent of one
68 another (e.g. multiple sequencing runs within the same study), these methods fall
69 short.

70 Here, we describe a simple data-normalization procedure for controlling
71 batch effects in case-control microbiome studies. Case-control studies include a
72 built-in population of control samples (e.g. healthy subjects) that can be used to
73 normalize the case samples (e.g. diseased subjects). For every feature (e.g.
74 bacterial taxon) with sufficient representation in the data, the case abundance
75 distributions can be converted to percentiles of the equivalent control abundance
76 distributions (Fig. 1). Study-specific batch effects present in the case samples will
77 also be present in the control samples, and by converting the case data into
78 percentiles of the control distribution these effects are effectively removed. Upon
79 conversion to percentiles of the within-study controls, percentile-normalized
80 samples from multiple studies with similar case-control definitions can be
81 appropriately pooled for statistical testing. We show that this approach controls
82 batch effects in microbiome case-control studies and we compare this method to
83 pooling non-normalized relative abundance data, pooling ComBat-corrected
84 data, and to Fisher's and Stouffer's methods for combining p-values from
85 unpooled analyses.



86

87 **Figure 1.** Theoretical feature abundance distributions for the control samples (blue) and case
88 samples (orange) are shown in the upper panel. Converting the control distribution into
89 percentiles of itself naturally gives rise to a uniform distribution (blue horizontal line in lower
90 panel), while converting the case distribution into percentiles of the control distribution produces a
91 non-uniform distribution when these two distributions differ (lower panel). Black lines show where
92 control distribution percentiles lie on the original and percentile-normalized histograms (10th, 30th,
93 50th, 70th, and 90th percentiles). The control distribution was produced by randomly sampling 100
94 times from a lognormal distribution with parameters $\mu = 0.1$ and $\sigma = 0.7$. The case distribution
95 was produced in a similar fashion, with distribution parameters $\mu = 0.8$ and $\sigma = 0.5$.

96

97 **Methods**

98 *Datasets*

99 We used a collection of case-control datasets obtained from the MicrobiomeHD
100 database [15] to validate our batch-normalization method. We focused our

101 analyses on studies spanning four diseases: colorectal cancer (CRC), Crohn's
102 Disease (CD), Ulcerative Colitis (UC), and *Clostridium difficile* induced diarrhea
103 (CDI). For a subset of three CRC studies [16-18], we were able to obtain
104 sequence data from the same region of the 16S gene so that these data could be
105 processed together. The remaining MicrobiomeHD case-control datasets were
106 processed separately using the same pipeline, and then Operational Taxonomic
107 Units (OTUs) were summarized at the genus level for comparison across studies.

108

109 *Sequence Data Processing*

110 To perform OTU-level analyses across the CRC studies, we downloaded the raw
111 data from all of the MicrobiomeHD datasets that sequenced the V4 region of the
112 16S gene. We quality filtered and length trimmed each V4 dataset as described
113 in [15] and concatenated these raw, trimmed FASTQ files into one file. We
114 removed any unique sequences that did not appear more than 20 times and
115 clustered the remaining reads with USEARCH [19] at 97% similarity. We
116 assigned these OTUs taxonomic identifiers using the RDP classifier [20] with a
117 cutoff of 0.5.

118 For genus-level analyses, data and metadata were acquired from the
119 MicrobiomeHD database (<https://doi.org/10.5281/zenodo.569601>). Raw data
120 were downloaded from the original studies and processed through our in-house
121 16S-processing pipeline
122 (https://github.com/thomasgurry/amplicon_sequencing_pipeline) as described in
123 [15]. Each study's OTU table was converted to relative abundance by dividing

124 each sample by the total number of reads and collapsed to genus level by
125 summing all OTUs with the same genus, throwing out any OTUs which did not
126 have a genus label.

127 To plot data in ordination space, Bray-Curtis distances were calculated
128 from relative abundance data using Scikit-learn
129 (`sklearn.metrics.pairwise.pairwise_distances`; `metric='braycurtis'`) [21]. Non-metric
130 multidimensional scaling (NMDS) coordinates were calculated for two axes
131 based on Bray-Curtis distances using Scikit-learn (`sklearn.manifold.MDS`;
132 `n_components=2`, `metric=False`, `max_iter=500`, `eps=1e-12`,
133 `dissimilarity='precomputed'`).

134

135 *Percentile Normalization*

136 Empirical relative abundance distributions were converted to percentiles using
137 the SciPy v 0.19.0 [22] `stats.percentileofscore` method (`kind='mean'`). Within each
138 study, control distributions for each individual OTU or genus were converted into
139 percentiles of themselves and case distributions were converted into percentiles
140 of their corresponding control distribution (Fig. 1). We restricted our analysis to
141 OTUs that occurred in at least one third of control *or* one third of case samples in
142 order to avoid statistical artifacts due to sampling effects. We have written a
143 python script that performs percentile-normalization given an OTU table, a list of
144 case sample IDs, and a list of control sample IDs as inputs
145 (https://github.com/seangibbons/percentile_normalization)

146

147 *ComBat*

148 For each disease, we applied ComBat [6] to the case-control data sets analyzed
149 in this study. Relative abundances (OTUs in the CRC analysis or OTUs
150 collapsed to the genus level in the genus-level analysis) were log-transformed
151 prior to running ComBat, adding a pseudocount of 1.0 to replace zeros in the
152 OTU count matrix. ComBat-corrected data were then transformed back from log-
153 space (i.e. exponential transformation) prior to downstream analyses.

154

155 *Statistical Analysis*

156 We used the Wilcoxon rank-sum test, as implemented in SciPy v0.19.0
157 (`sicipy.stats.ranksums`) [22], to determine significant differences between
158 independent groups of samples. Wilcoxon tests were calculated either within or
159 across studies. In order to calculate statistics across datasets, case and control
160 samples from multiple studies of the same disease were combined together into
161 the same OTU table. Hereafter, combining datasets is referred to as 'pooling.' P-
162 values were multiple-test corrected using the Benjamini-Hochberg False
163 Discovery Rate (FDR) procedure, as implemented in StatsModels v 0.8.0
164 (`statsmodels.sandbox.stats.multicomp.multipletests`) [23]. Differences in overall
165 community structure were assessed using the Permutational Multivariate
166 Analysis of Variance (PERMANOVA) test in R's `vegan` package [24] as
167 implemented in `scikit-bio` (`skbio.stats.distance.permanova`). Fisher's and
168 Stouffer's methods for combining p-values were performed using SciPy v0.19.0
169 (`scipy.stats.combine_pvalues`; `method='fisher'` or `method='stouffer'`). For

170 Stouffer's method, weights for each study were defined as the square root of the
171 number of cases plus the number of controls.

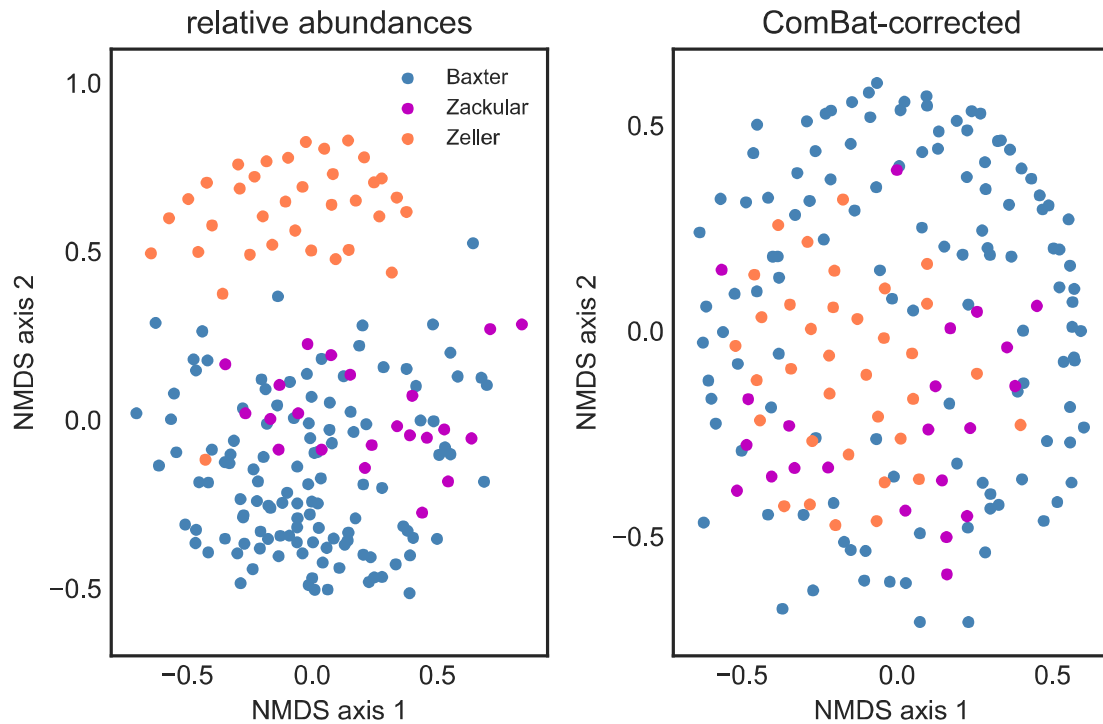
172

173 **Results**

174 *Batch effects at OTU-level resolution*

175 To minimize possible biases across data sets, we identified three colorectal
176 cancer (CRC) studies that sequenced the same region of the 16S gene (V4). We
177 reprocessed the raw sequence data from each study in the same quality filtering
178 and OTU picking pipeline to obtain bioinformatically-standardized results. OTUs
179 that occurred in at least a third of case or a third of control samples (i.e. either
180 within individual studies or across studies) were retained for all downstream
181 statistical analyses. Despite standardizing the computational processing of these
182 data, we saw significant batch effects in healthy patients across studies
183 (PERMANOVA $p < 0.001$; Fig. 2). The similarity between samples from the
184 Baxter and Zackular studies is due to the fact that they were sourced from the
185 same patient cohort, making this comparison a good pseudo-negative control for
186 batch effects [16, 18]. There was an apparent reduction in the batch effect after
187 applying ComBat, although differences between batches remained statistically
188 significant (PERMANOVA $p < 0.001$; Fig. 2) [6]. Due to the non-independence
189 between the Baxter and Zackular patient cohorts, we removed the smaller of the
190 two studies (Zackular) from all downstream analyses. Out of a total of 5,585
191 OTUs found in healthy controls, 725 OTUs differed significantly in relative

192 abundance between the Baxter et al. (2016) and Zeller et al. (2014) controls
193 (FDR $q \leq 0.05$).



194

195 **Figure 2.** Non-metric multimensional scaling (NMDS) plot showing the distribution of healthy
196 controls from three colorectal cancer studies in ordination space (Bray-Curtis distances of relative
197 abundance data). Despite standardized bioinformatic processing, healthy patients differed
198 significantly in their gut microbiomes across studies (PERMANOVA $p < 0.001$). Studies were still
199 significantly different even after applying ComBat, an established batch-correction method
200 (PERMANOVA $p < 0.001$).

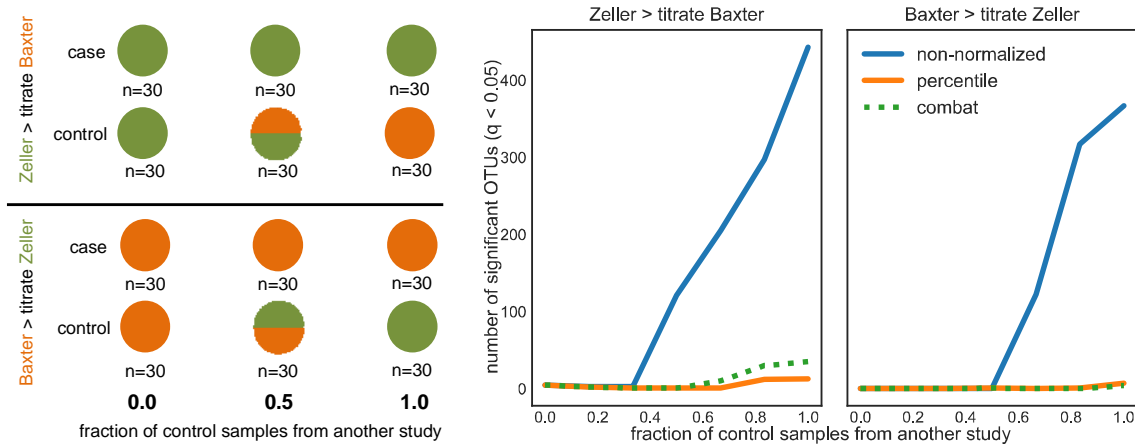
201

202 As expected for the Wilcoxon rank-based statistical test, within-study
203 results at the OTU level were identical before and after percentile-normalization.
204 In addition, these within-study results were also identical with the results from
205 ComBat-corrected data. In the Baxter study, there were 172 healthy (control)
206 samples and 120 CRC (case) samples, with 3 OTUs (from *Parvimonas*,

207 *Porphyromonas*, and *Peptostreptococcus* genera) showing significant differences
208 in abundance between cases and controls for all analyses (FDR $q \leq 0.05$). For
209 Zeller, there were 71 control and 40 case samples, with 4 OTUs (from
210 *Fusobacterium*, *Clostridium XIVa*, *Peptostreptococcus*, and *Dialister* genera) that
211 differed significantly across cases and controls for all analyses (FDR $q \leq 0.05$).

212 We ran an *in silico* titration experiment to simulate pooling of control
213 samples from different datasets before calculating significant differences. Healthy
214 samples from one study were mixed with healthy samples from another study at
215 different proportions prior to calculating significant differences in OTU
216 frequencies between cases and controls (Fig. 3). Case and control groups were
217 subsampled to 30 samples each. Control samples were substituted by samples
218 from another study along a fractional gradient (0-100% control samples from
219 another study; see conceptual outline in Fig. 3). For the relative abundance data
220 (non-normalized), the number of significant OTUs greatly increased due to batch
221 effects as more control samples were substituted in from another study.
222 However, the ComBat-corrected and percentile-normalized results were almost
223 totally unaffected by the proportion of control samples coming from another
224 study, indicating that batch effects were no longer driving spurious associations
225 in the normalized data.

226



227

228 **Figure 3.** *In silico* titration experiment, where the control group from one study is gradually
229 substituted with randomly chosen control samples from another study (non-normalized,
230 percentile-normalized, and ComBat-corrected), keeping the total number of case and control
231 samples fixed at n=30 (see conceptual illustration on the left). Mixing non-normalized data from
232 control samples from another study often gave rise to spurious significant results due to technical
233 differences across studies (blue lines). However, when the data were percentile-normalized or
234 ComBat-corrected, we did not see a large increase in significant OTUs as control samples from
235 different studies were mixed in (solid orange and green dashed lines).

236

237 In the absence of batch effects, pooling data across datasets of the
238 same disease should increase sensitivity to detect significant associations. We
239 pooled relative abundances, percentile-normalized abundances, and ComBat-
240 corrected abundances, respectively, across the Baxter and Zeller studies to look
241 for OTUs that differed significantly across cases and controls. These pooled
242 results were then compared to classic methods for combining p-values from each
243 dataset's individual results. For the relative abundance data, we found six OTUs
244 (from *Porphyromonas*, *Fusobacterium*, *Clostridium XIVa*, *Peptostreptococcus*,
245 *Dialister*, and *Parvimonas* genera) that differed significantly across cases and

246 controls (FDR $q \leq 0.05$). After pooling the percentile-normalized data, we found
247 seven OTUs that were significantly enriched in cancer patients relative to
248 controls -- two OTUs from the *Clostridium XIVa* genus, one from *Parvimonas*,
249 one from *Peptostreptococcus*, one from *Porphyromonas*, one from *Dialister*, and
250 one from *Fusobacterium* (FDR $q \leq 0.05$). The pooled ComBat-corrected results
251 included the same significant hits identified in the percentile-normalization
252 results. Fisher's method identified just two significant OTUs from the
253 *Peptostreptococcus* and *Parvimonas* genera, which were also found in the
254 pooled results. Stouffer's method identified one significant OTU from the
255 *Peptostreptococcus* genus, which was also identified in the pooled results.
256 Overall, the pooled methods maximize statistical power to detect significant
257 OTUs over traditional meta-analysis methods. For example, OTUs from
258 *Fusobacterium*, *Porphyromonas*, *Clostridium XIVa* and *Dialister* genera were
259 identified as significantly enriched in CRC patients by the normalization methods
260 but not by Fisher's or Stouffer's methods. Previous meta-analyses of CRC
261 microbiome studies have shown these genera to be consistently associated with
262 CRC, which supports our findings [15, 25].

263

264 *Batch effects at genus-level resolution across multiple diseases*

265 In order to assess the performance of different meta-analysis techniques across
266 a larger set of studies and diseases, we summarized OTU abundances at the
267 genus level for four diseases - *Clostridium difficile* induced diarrhea (CDI),
268 Crohn's disease (CD), ulcerative colitis (UC), and CRC - across 11 case-control
269 studies. There were a total of 306 unique genera detected across studies. There

270 were two CDI case-control studies: Schubert et al. (2014) had 154 control and 93
271 case samples [26]; Vincent et al. (2013) had 25 control and 25 case samples
272 [27]. There were four inflammatory bowel disease (IBD) studies that included CD
273 patients and three that also included UC patients: Papa et al. (2012) had 24 non-
274 IBD control sample, 23 CD samples, and 43 UC samples [28]; Morgan et al.
275 (2012) had 18 control, 61 CD and 47 UC samples [29]; Willing et al. (2010) had
276 35 control, 16 UC and 29 CD samples [30]; Gevers et al. (2014) had 16 non-IBD
277 control and 146 CD samples, with no UC samples [31]. There were four
278 independent CRC studies, including the Baxter and Zeller studies listed in the
279 OTU-level analysis (see above for sample sizes). The remaining two CRC
280 studies added to the genus-level analysis are Wang et al. (2012), which had 54
281 control and 44 case samples [32], and Chen et al. (2012), which had 22 controls
282 and 21 cases [33].

283 The number of genera that differed significantly across cases and controls
284 changed depending on how the data were normalized, pooled, and analyzed
285 (Table 1). Wilcoxon rank-sum tests yielded identical within-study results for non-
286 normalized and percentile-normalized data. However, unlike the OTU-level
287 analysis, within-study ComBat-corrected results showed fewer significant genera
288 than the non-normalized results for unpooled, within-study tests (Table 1). Thus,
289 in correcting for batch effects, ComBat appears to smooth out some biological
290 signal. While pooling non-normalized data across studies is technically
291 inappropriate, it frequently resulted in significant hits that were consistent with
292 percentile-normalized results, suggesting that the biological signal was often

293 stronger than the batch effect. Except in the case of UC, pooling percentile-
 294 normalized data consistently yielded more significant hits than pooling non-
 295 normalized data (see ‘across’ column in Table 1). ComBat-correction generally
 296 resulted in many fewer significant genera after pooling, especially for CD and
 297 UC. Half of the IBD studies included non-IBD patients with inflammatory
 298 symptoms as controls rather than clinically healthy patients. These biologically
 299 relevant differences in inflammatory symptoms between control cohorts were
 300 conflated with batches and were likely smoothed out by ComBat. In all cases,
 301 Fisher’s and Stouffer’s methods identified fewer significant results than pooling
 302 the percentile-normalized data. These results illustrate that pooling data is more
 303 sensitive than classic meta-analysis techniques [34] and that percentile-
 304 normalization further increases the statistical power to detect differences while
 305 controlling for batch effects.

306

disease	normalization	studies	within	across	shared	Fisher	Stouffer
CDI	none	2	31	33	30	30	29
	ComBat	2	29	33	29	29	27
	percentile	2	31	33	30	30	29
CD	none	4	3	16	1	0	0
	ComBat	4	0	0	0	0	0
	percentile	4	3	18	3	0	0
UC	none	3	8	10	6	3	0
	ComBat	3	3	3	1	1	0
	percentile	3	8	9	6	0	0
CRC	none	4	7	4	2	3	0
	ComBat	4	2	4	1	0	0
	percentile	4	7	7	3	3	0

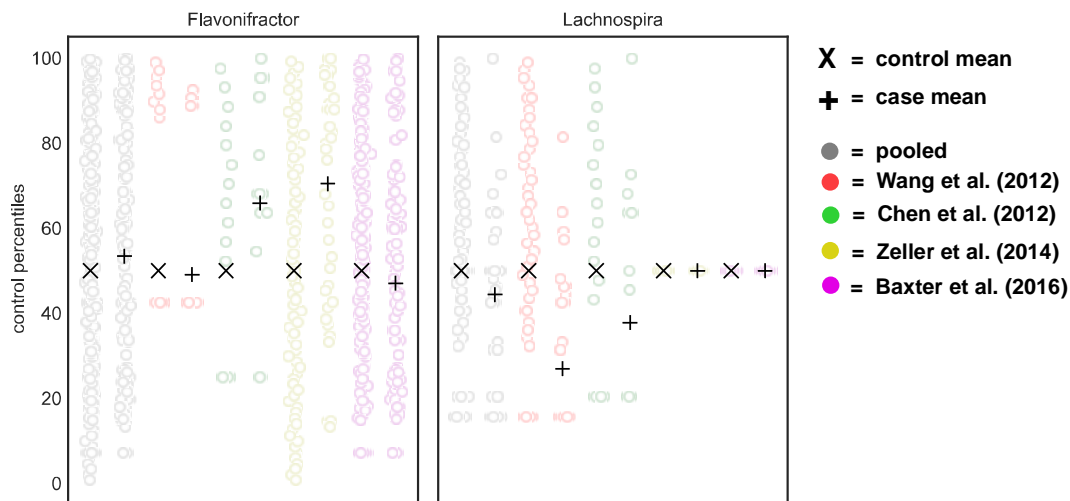
307

308 **Table 1.** Numbers of taxa that differ significantly between cases and controls for four diseases.
309 The normalization column indicates how the data were treated prior to running significance tests
310 (non-normalized, ComBat -corrected, or percentile-normalization). In the 'disease' column, 'CD' =
311 Crohn's Disease, UC = Ulcerative Colitis, CRC = Colorectal Cancer, and CDI = *Clostridium*
312 *difficile* induced diarrhea. The significance threshold used was $q \leq 0.05$ (FDR). The 'within'
313 column shows how many significant taxa were identified when running the statistics for each
314 study independently, while the 'pooled' column shows the number of significant taxa identified
315 when running the statistics on the combined datasets. The 'shared' column shows how many taxa
316 overlap between the 'within' and 'pooled' columns. The 'Fisher' and 'Stouffer' columns show the
317 number of significant taxa identified using Fisher's and Stouffer's methods for combining p-values
318 from the independent within-study tests.

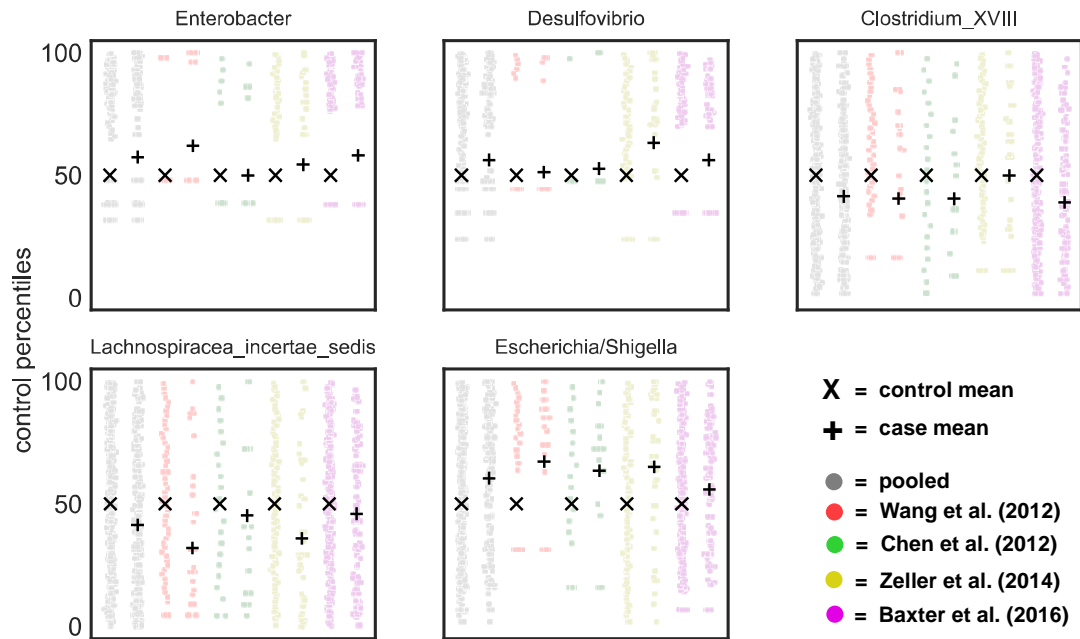
319

320 To better assess how percentile normalization impacted the pooled
321 results, we looked at genera that were significant within a single-study but not
322 across studies after pooling and also at OTUs that were significant across pooled
323 studies but not within a given study. We ran this analysis on the CRC data,
324 where we had the largest number of independent studies with consistently
325 defined healthy control cohorts ($n = 4$). There were two genera that were
326 significant within a subset of studies, but not across all studies after pooling (Fig.
327 4). *Lachnospira* was absent in three out of the five CRC studies and was
328 enriched in controls in the two studies where it was detected. *Flavonifractor* was
329 more abundant in cases for two studies, but this signal was not consistent across
330 all studies. Thus, these taxa were either too rare or sensitive to different
331 experimental and/or processing techniques to be reliable biomarkers. There were
332 five genera that showed significant differences after pooling studies together but

333 were not significant in any individual study. *Escherichia/Shigella*, *Enterobacter*,
334 and *Desulfovibrio* genera were slightly enriched in CRC patients in most studies,
335 but did not show a statistically significant enrichment in any individual study (Fig.
336 5). Conversely, *Clostridium XVIII* and *Lachnospiraceae incertae sedis* genera
337 were enriched in controls across most studies. These OTUs show small, yet
338 consistent effect sizes across independent studies that can only be detected after
339 pooling (Fig. 5).



340
341 **Figure 4.** The *Flavonifractor* and *Lachnospira* genera showed significant differences between
342 cases and controls within a study (FDR $q \leq 0.05$), but not after pooling across CRC studies.
343



344

345 **Figure 5.** Five genera did not show a significant difference between cases and controls within an
346 individual study, but were significantly different when pooling across CRC studies (FDR $q \leq$
347 0.05). Scatter plots show distributions of percentile-normalized data for case and control samples
348 across studies.

349

350 Discussion

351 Batch effects are unavoidable when working with high-throughput data
352 generation platforms. The RNA microarray community has been proactive in the
353 development of tools for dealing with these effects [1, 6]. However, these tools
354 are not as effective when batch effects are confounded with biological signals, or
355 when these effects cannot be projected onto a small number of dimensions,
356 which is often the case in microbiome case-control studies [35-37]. Fortunately,
357 case-control studies can be internally normalized by their own control samples.
358 Any study-specific batch effects in the case samples will be present in the control

359 samples, and by converting the case data into percentiles of the control
360 distribution these effects are removed.

361 Relative abundance data -- but not the percentile-normalized or ComBat-
362 corrected data -- quickly gave spurious results when cases from one study were
363 tested against controls from another (Fig. 3). For studies with small numbers of
364 control and/or case samples, it is tempting to pool with other datasets to increase
365 statistical power. In the past, pooling of non-normalized data from different
366 studies has been done [31, 35, 38], but as demonstrated above, this is
367 inadvisable. In these scenarios, datasets can first be percentile-normalized and
368 then appropriately combined without introducing batch-related artifacts.

369 We found substantial overlap in the relative abundance and percentile-
370 normalized results when calculating significance across studies. This overlap is
371 expected when there is a strong biological signal that overrides batch effects
372 [39]. Despite the similarity between pooled relative abundance and percentile-
373 normalized results, there were several cases where the percentile-normalized
374 results identified significant differences between cases and controls that were
375 missed in the non-normalized results and there was one instance (UC) where
376 one fewer significant difference was found in the percentile-normalized results
377 (Table 1). Percentile-normalization also identified more significant hits than
378 ComBat-corrected data in the genus-level pooled analyses, especially for UC
379 and CD (IBD). The reduced number of significant hits from ComBat-corrected
380 data for IBD was likely due to heterogeneous control cohorts across these

381 studies (i.e. healthy patients vs. non-IBD patients), which likely smoothed-out
382 inflammation-associated signals.

383 We compared normalization and pooling methods (i.e. percentile-
384 normalization and ComBat) to Fisher's and Stouffer's methods for combining p-
385 values. Stouffer's method is similar to Fisher's, but includes weights for each p-
386 value based on the number of samples in a study. For all diseases, the pooling
387 methods identified a larger number of significant hits than Stouffer's and Fisher's
388 methods, indicating that pooling provides more sensitivity to detect differences
389 between cases and controls. The bacterial taxa identified as significant by the
390 percentile-normalization method were largely consistent with prior results [15].

391 In conclusion, we present a robust, model-free procedure for transforming
392 each feature in a case-control dataset into percentiles of its control distribution
393 (Fig. 1). These percentile-normalized features can be pooled across independent
394 studies for non-parametric, univariate statistical testing, circumventing the batch
395 effect problem. We find that this procedure allows us to identify differences
396 between cases and controls that are often missed by more conservative meta-
397 analysis techniques. Methods developed for batch-correction in microarray data,
398 like ComBat, can reduce batch effects in microbiome studies (Fig. 2-3), but may
399 also obscure real patterns if batch effects are not totally independent of biological
400 signals. We suggest that ComBat and other similar methods are useful for
401 studies without case and control groups. However, when studies have internal
402 controls, percentile-normalization should be the preferred batch correction
403 approach.

404

405 **Acknowledgements**

406 This work was supported by the Center for Microbiome Informatics and
407 Therapeutics. We thank the members of the Alm lab for helpful feedback.

408

409 **References**

- 410 1. Leek JT, Scharpf RB, Bravo HC, Simcha D, Langmead B, Johnson WE, et
411 al. Tackling the widespread and critical impact of batch effects in high-throughput
412 data. *Nat Rev Genet.* 2010;11(10):733-9.
- 413 2. Schloss PD, Gevers D, Westcott SL. Reducing the Effects of PCR
414 Amplification and Sequencing Artifacts on 16S rRNA-Based Studies. *PLoS One.*
415 2011;6(12):e27310. doi: 10.1371/journal.pone.0027310.
- 416 3. Leek JT, Storey JD. Capturing heterogeneity in gene expression studies
417 by surrogate variable analysis. *PLoS Genet.* 2007;3(9):e161.
- 418 4. Alter O, Brown PO, Botstein D. Singular value decomposition for genome-
419 wide expression data processing and modeling. *Proc Natl Acad Sci USA.*
420 2000;97(18):10101-6.
- 421 5. Benito M, Parker J, Du Q, Wu J, Xiang D, Perou CM, et al. Adjustment of
422 systematic microarray data biases. *Bioinformatics.* 2004;20(1):105-14.
- 423 6. Chen C, Grennan K, Badner J, Zhang D, Gershon E, Jin L, et al.
424 Removing Batch Effects in Analysis of Expression Microarray Data: An
425 Evaluation of Six Batch Adjustment Methods. *PLoS One.* 2011;6(2):e17238. doi:
426 10.1371/journal.pone.0017238.
- 427 7. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray
428 expression data using empirical Bayes methods. *Biostatistics.* 2007;8(1):118-27.
- 429 8. Weiss S, Amir A, Hyde ER, Metcalf JL, Song SJ, Knight R. Tracking down
430 the sources of experimental contamination in microbiome studies. *Genome Biol.*
431 2014;15(12):564. doi: 10.1186/s13059-014-0564-2.
- 432 9. Shen H, Rogelj S, Kieft TL. Sensitive, real-time PCR detects low-levels of
433 contamination by *Legionella pneumophila* in commercial reagents. *Mol Cell*
434 *Probes.* 2006;20. doi: 10.1016/j.mcp.2005.09.007.
- 435 10. Nguyen NH, Smith D, Peay K, Kennedy P. Parsing ecological signal from
436 noise in next generation amplicon sequencing. *New Phytol.* 2015;205(4):1389-
437 93. doi: 10.1111/nph.12923.
- 438 11. Gibbons SM. The Built Environment Is a Microbial Wasteland. *mSystems.*
439 2016;1(2):e00033-16.
- 440 12. Chase J, Fouquier J, Zare M, Sonderegger DL, Knight R, Kelley ST, et al.
441 Geography and Location Are the Primary Drivers of Office Microbiome
442 Composition. *mSystems.* 2016;1(2). doi: 10.1128/mSystems.00022-16.

- 443 13. Fisher RA. Statistical methods for research workers: Genesis Publishing
444 Pvt Ltd; 1925.
- 445 14. Stouffer SA. Adjustment during army life: Princeton University Press;
446 1949.
- 447 15. Duvallet C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta Analysis Of
448 Microbiome Studies Identifies Shared And Disease-Specific Patterns. bioRxiv.
449 2017:134031.
- 450 16. Baxter NT, Ruffin MT, Rogers MA, Schloss PD. Microbiota-based model
451 improves the sensitivity of fecal immunochemical test for detecting colonic
452 lesions. *Genome Med.* 2016;8(1):37.
- 453 17. Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, et al.
454 Potential of fecal microbiota for early stage detection of colorectal cancer. *Molec*
455 *Sys Biol.* 2014;10(11):766.
- 456 18. Zackular JP, Rogers MA, Ruffin MT, Schloss PD. The human gut
457 microbiome as a screening tool for colorectal cancer. *Cancer Prev Res.*
458 2014;7(11):1112-21.
- 459 19. Edgar RC. Search and clustering orders of magnitude faster than BLAST.
460 *Bioinformatics.* 2010;26(19):2460-1.
- 461 20. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naïve Bayesian Classifier for
462 Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Appl*
463 *Environ Microbiol.* 2007;73(16):5261-7. doi: 10.1128/aem.00062-07.
- 464 21. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et
465 al. Scikit-learn: Machine learning in Python. *J Mach Learn Res.*
466 2011;12(Oct):2825-30.
- 467 22. Jones E, Oliphant T, Peterson P. SciPy: Open source scientific tools for
468 Python, 2001–. URL [http://www](http://www.scipy.org) scipy org. 2007;73:86.
- 469 23. Seabold S, Perktold J, editors. Statsmodels: Econometric and statistical
470 modeling with python. Proceedings of the 9th Python in Science Conference;
471 2010.
- 472 24. Oksanen J. Multivariate analysis of ecological communities in R: vegan
473 tutorial. R package version. 2011;1(7).
- 474 25. Tjalsma H, Boleij A, Marchesi JR, Dutilh BE. A bacterial driver–passenger
475 model for colorectal cancer: beyond the usual suspects. *Nat Rev Microbiol.*
476 2012;10(8):575-82.
- 477 26. Schubert AM, Rogers MA, Ring C, Mogle J, Petrosino JP, Young VB, et
478 al. Microbiome data distinguish patients with *Clostridium difficile* infection and
479 non-*C. difficile*-associated diarrhea from healthy controls. *mBio.*
480 2014;5(3):e01021-14.
- 481 27. Vincent C, Stephens DA, Loo VG, Edens TJ, Behr MA, Dewar K, et al.
482 Reductions in intestinal Clostridiales precede the development of nosocomial
483 *Clostridium difficile* infection. *Microbiome.* 2013;1(1):18.
- 484 28. Papa E, Docktor M, Smillie C, Weber S, Preheim SP, Gevers D, et al.
485 Non-invasive mapping of the gastrointestinal microbiota identifies children with
486 inflammatory bowel disease. *PLoS One.* 2012;7(6):e39242.

- 487 29. Morgan XC, Tickle TL, Sokol H, Gevers D, Devaney KL, Ward DV, et al.
488 Dysfunction of the intestinal microbiome in inflammatory bowel disease and
489 treatment. *Genome Biol.* 2012;13(9):R79.
- 490 30. Willing BP, Dicksved J, Halfvarson J, Andersson AF, Lucio M, Zheng Z, et
491 al. A pyrosequencing study in twins shows that gastrointestinal microbial profiles
492 vary with inflammatory bowel disease phenotypes. *Gastroenterology.*
493 2010;139(6):1844-54. e1.
- 494 31. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W,
495 Ren B, et al. The treatment-naive microbiome in new-onset Crohn's disease. *Cell*
496 *Host Microbe.* 2014;15(3):382-92.
- 497 32. Wang T, Cai G, Qiu Y, Fei N, Zhang M, Pang X, et al. Structural
498 segregation of gut microbiota between colorectal cancer patients and healthy
499 volunteers. *ISME J.* 2012;6(2):320-9.
- 500 33. Chen W, Liu F, Ling Z, Tong X, Xiang C. Human intestinal lumen and
501 mucosa-associated microbiota in patients with colorectal cancer. *PLoS One.*
502 2012;7(6):e39743.
- 503 34. Glass GV. Primary, secondary, and meta-analysis of research. *Educ Res.*
504 1976;5(10):3-8.
- 505 35. Ross MC, Muzny DM, McCormick JB, Gibbs RA, Fisher-Hoch SP,
506 Petrosino JF. 16S gut community of the Cameron County Hispanic Cohort.
507 *Microbiome.* 2015;3(1):7.
- 508 36. Pasolli E, Truong DT, Malik F, Waldron L, Segata N. Machine learning
509 meta-analysis of large metagenomic datasets: tools and biological insights. *PLoS*
510 *Comput Biol.* 2016;12(7):e1004977.
- 511 37. Escobar JS, Klotz B, Valdes BE, Agudelo GM. The gut microbiota of
512 Colombians differs from that of Americans, Europeans and Asians. *BMC*
513 *Microbiol.* 2014;14(1):311.
- 514 38. Dubourg G, Lagier J-C, Hüe S, Surenaud M, Bachar D, Robert C, et al.
515 Gut microbiota associated with HIV infection is significantly enriched in bacteria
516 tolerant to oxygen. *BMJ Open Gastroenter.* 2016;3(1).
- 517 39. Forslund K, Hildebrand F, Nielsen T, Falony G, Le Chatelier E, Sunagawa
518 S, et al. Disentangling type 2 diabetes and metformin treatment signatures in the
519 human gut microbiota. *Nature.* 2015;528(7581):262-6.
- 520