

1 **Accounting for 16S rRNA copy number prediction uncertainty and its implications in**
2 **bacterial diversity analyses**

3

4 Yingnan Gao¹, yg5ap@virginia.edu

5 Martin Wu^{1*}, mw4yv@virginia.edu

6 ¹Department of Biology, University of Virginia, 485 McCormick Road, Charlottesville, VA,

7 22904, USA

8 *Corresponding author

9

10 Competing interests

11 The authors declare that they have no competing interests.

12 **Abstract**

13 16S rRNA gene copy number (16S GCN) varies among bacterial species and this variation
14 introduces potential biases to microbial diversity analyses using 16S rRNA read counts. To
15 correct the biases, methods have been developed to predict 16S GCN. A recent study suggests
16 that the prediction uncertainty can be so great that copy number correction is not justified in
17 practice. Here we develop RasperGade16S, a novel method and software to better model and
18 capture the inherent uncertainty in 16S GCN prediction. RasperGade16S implements a
19 maximum likelihood framework of pulsed evolution model and explicitly accounts for
20 intraspecific GCN variation and heterogeneous GCN evolution rates among species. Using cross
21 validation, we show that our method provides robust confidence estimates for the GCN
22 predictions and outperforms other methods in both precision and recall. We have predicted GCN
23 for 592605 OTUs in the SILVA database and tested 113842 bacterial communities that represent
24 an exhaustive and diverse list of engineered and natural environments. We found that the
25 prediction uncertainty is small enough for 99% of the communities that 16S GCN correction
26 should improve their compositional and functional profiles estimated using 16S rRNA reads. On
27 the other hand, we found that GCN variation has limited impacts on beta-diversity analyses such
28 as PCoA, PERMANOVA and random forest test.

29

30 **Introduction**

31 The 16S ribosomal RNA (16S rRNA) gene is the gold standard for bacterial and archaeal
32 diversity study and has been commonly used to estimate the composition of bacterial and
33 archaeal communities through amplicon sequencing. Sequence reads are usually matched to
34 reference databases like SILVA [1] and GreenGenes [2] to determine the presence of taxa and

35 their relative cell abundances. However, the 16S rRNA gene copy number (16S GCN) can vary
36 from 1 to more than 15 [3, 4] and this large copy number variation introduces bias in the relative
37 cell abundance estimated using the gene read counts (thereafter referred to as gene abundance)
38 [5], and consequently it can skew the community profiles, diversity measures and lead to
39 qualitatively incorrect interpretations [5–8]. As a result, it has been argued that 16S GCN
40 variations should be taken into account in 16S rRNA gene-based analyses [5].

41
42 The majority of bacteria species have not been cultured or sequenced and their 16S GCNs are
43 unknown. Studies have shown that 16S GCN exhibits a strong phylogenetic signal [5, 7], and
44 therefore 16S GCN can be inferred from closely related reference bacteria. Based on this
45 principle, software has been developed to predict the 16S GCN [5, 7, 9, 10] in a process often
46 referred to as hidden state prediction [11]. However, a recent study correctly points out that the
47 accuracy of 16S GCN prediction deteriorates as the minimum phylogenetic distance between the
48 query sequence and the reference sequences increases, and the prediction of 16S GCN is still an
49 open question [12].

50
51 The increasing error of 16S GCN prediction with increasing phylogenetic distance roots from the
52 stochastic nature of trait evolution, which leads to inherent uncertainty in the predicted trait
53 values. One way of reducing the inherent uncertainty is to improve taxon sampling in the
54 reference phylogeny to reduce the query's phylogenetic distance to the reference [13]. Another
55 way of addressing the inherent uncertainty is to model the uncertainty directly and have a
56 confidence estimate. By doing so, we will be able to determine how confident we should be
57 about a GCN prediction and make meaningful interpretations. Unfortunately, few 16S GCN

58 prediction tools provide a confidence estimation for the predicted 16S GCN, and uncertainty is
59 mostly ignored when interpreting the results of downstream analyses [5, 7, 10]. For example,
60 PICRUST2 predicts functional profiles of bacterial and archaeal communities from 16S rRNA
61 sequence data. It predicts 16S GCN for each operational taxonomic unit (OTU) in the
62 community and uses the predicted values (point estimates) to estimate “corrected” relative cell
63 abundances and metagenomes, without accounting for the uncertainty of the predictions. As a
64 result, the impact of uncertainty in 16S GCN prediction on bacterial diversity analyses remains
65 unknown and needs to be investigated.

66

67 Several points need to be considered to properly model the prediction uncertainty. First, because
68 the uncertainty roots from the stochastic nature of trait evolution, we need to develop a good
69 model for 16S GCN evolution. Previously the evolution of the 16S GCN trait has been modeled
70 as gradual evolution using the Brownian motion (BM) model [5, 7, 10]. However, alternative
71 models exist and need to be considered [14–16]. For example, pulsed evolution (PE) model
72 postulates that traits evolve by jumps, followed by periods of stasis [14, 17]. It has been shown
73 that pulsed evolution is prevalent in microbial genome trait evolution [18]. 16S GCN of *Bacillus*
74 *subtilis* can jump from 1 to 6 in a matter of days by gene amplification [19]. On the other hand, it
75 is well known that the 16S GCN of some bacterial clades such as the Rickettsiales order, a
76 diverse group of obligate intracellular bacteria, has only one copy of 16S rRNA in their genomes,
77 demonstrating stasis [20, 21]. To develop a proper model for 16S GCN evolution, the tempo and
78 mode of evolution need to be examined.

79

80 Secondly, 16S GCN can vary within the same species [22–25], which introduces uncertainty to
81 GCN prediction that needs to be accounted for. It has been shown that modeling the intraspecific
82 variation is essential for the analysis of comparative trait data and failing to account for this
83 variation can result in model misspecification [14]. Because conspecific strains are usually
84 separated by zero branch length in the phylogeny of the 16S rRNA gene, the intraspecific
85 variation can be modelled as time-independent variation, which can also account for
86 measurement errors [26].

87

88 Thirdly, there is notable rate heterogeneity in 16S GCN evolution. For example, the obligately
89 intracellular bacteria and free-living bacteria with streamlined genomes (e.g., *Rickettsia* and
90 *Pelagibacter*) have elevated molecular evolutionary rates [27, 28] and therefore relatively long
91 branches in the 16S rRNA gene phylogeny [29]. Nevertheless, they have only one copy of 16S
92 rRNA in their genomes and the GCNs rarely change [21]. It is expected that the 16S GCN
93 prediction for this group of bacteria should be accurate despite their large phylogenetic distances
94 to the reference genomes. Such examples suggest that the rate heterogeneity of 16S GCN
95 evolution should be systematically evaluated and modelled properly. However, no previous
96 methods have evaluated and modeled such evolution rate heterogeneity, leading to potential
97 model misspecification in 16S GCN predictions.

98

99 Here, we develop a novel tool *RasperGade16S* that employs a heterogeneous pulsed evolution
100 model for 16S rRNA GCN prediction. Through simulation and cross-validation, we show that
101 *RasperGade16S* outperforms other methods in terms of providing significantly improved
102 confidence estimates. We demonstrate that correcting 16S rRNA GCN improves the relative cell

103 abundance estimates of the bacterial communities and is expected to be beneficial for more than
104 99% of 113842 environmental samples we analyzed. However, our findings suggest that GCN
105 correction may not be necessary for beta-diversity analyses, as it has limited impact on the
106 results.

107

108 **Methods**

109 Compiling 16S GCN data and inferring 16S rRNA reference phylogeny

110 We downloaded annotated RNA gene sequences from 21245 complete bacterial genomes in the
111 NCBI RefSeq database (Release 205) on April 9, 2021. For each genome, we counted the
112 number of annotated 16S rRNA genes. Genomes with questionable 16S GCNs were removed
113 and one representative 16S rRNA sequence from each remaining genome was selected. A 16S
114 rRNA phylogeny (referred to as reference phylogeny hereafter) was inferred from the
115 representative sequences of 6408 genomes. See Supplementary Methods for details.

116

117 Evaluating time-independent variation in 16S GCN

118 To evaluate the extent of 16S GCN time-independent or intraspecific variation, we compared
119 GCN between 5437 pairs of genomes with identical 16S rRNA gene alignments. To formally test
120 whether accounting for time-independent variation is necessary, we modeled time-independent
121 variation as a normal white noise, and fitted the Brownian motion (BM) model to the evolution
122 of 16S GCN in the 6408 reference genomes, with and without time-independent variation. We
123 then calculated the likelihood and chose the best model using the Akaike Information Criterion
124 (AIC).

125

126 Evaluating the rate heterogeneity of 16S GCN evolution

127 We calculated the local average rate of evolution for each genus that contains at least 10
128 genomes in the reference phylogeny and examined the distribution of the average rates among
129 genera. The average rate of a genus is calculated as the variance of phylogenetically independent
130 contrasts (PICs) [30] of GCN within the genus.

131

132 Predicting 16S GCN

133 We developed a heterogeneous pulsed evolution model to model 16S GCN evolution (see
134 Supplementary Methods for details) and a likelihood based R package *RasperGade16S* to predict
135 16S GCN. *RasperGade16S* first assigns the query sequence to either the regularly-evolving or
136 the slowly-evolving group based on where it is inserted in the reference phylogeny. For a query
137 sequence inserted into the slowly-evolving group, its insertion branch length is scaled by the
138 ratio $r_{slow}/r_{regular}$, where r is the rate of evolution in each group. For a query sequence inserted
139 into the regularly-evolving group, a small branch length is added to the insertion branch to
140 represent the estimated time-independent variation. *RasperGade16S* then predicts the GCN of
141 the query using the rescaled reference phylogeny. Because 16S GCN is an integer trait, the
142 continuous prediction from hidden state prediction is rounded and a confidence (probability) that
143 the prediction is equal to the truth is estimated by integrating the predicted uncertainty
144 distribution. We marked the 16S GCN prediction with a confidence smaller than 95% as
145 unreliable, and otherwise as reliable. As a comparison, we also predicted GCN using PICRUST2,
146 which employs multiple hidden state prediction methods in the R package *castor* [31] for 16S
147 GCN predictions. We selected three methods by which confidence can be estimated: the

148 phylogenetically independent contrast (pic) method, the maximum parsimony (mp) method, and
149 the empirical probability (emp) method. Otherwise, we run PICRUST2 using default options and
150 the unscaled reference phylogeny.

151

152 We did not test the tools CopyRighter [7] and PAPRICA [9] in this study because 1) neither
153 provides the option of using a user-supplied reference data, and 2) neither provides uncertainty
154 estimates (i.e., confidence intervals) of its predictions, which is the primary focus of this study.

155

156 Adjust NSTD and NSTI with rate heterogeneity

157 The adjusted nearest-sequenced-taxon-distances (NSTDs) [12] is calculated using the rescaled
158 reference tree. The adjusted nearest-sequenced-taxon-index (NSTI) [10] is calculated as the
159 weighted average of adjusted NSTDs of the community members.

160

161 Validating the quality of predicted 16S GCN and its confidence estimate

162 We used cross-validations to evaluate the quality of 16S GCN prediction and its confidence
163 estimate, and how they vary with NSTD. We randomly selected 2% of the tips in the reference
164 phylogeny as the test set and filtered the remaining reference set by removing tips with a NSTD
165 to any test sequence smaller than a threshold. We then predicted the 16S GCN for each tip in the
166 test set using the filtered reference set. We conducted cross-validation within 9 bins delineated
167 by 10 NSTD thresholds: 0, 0.002, 0.005, 0.010, 0.022, 0.046, 0.100, 0.215, 0.464 and 1.000
168 substitutions/site, and for each bin we repeated the cross-validation 50 times with non-
169 overlapping test sets. We evaluated the quality of the 16S GCN prediction by the coefficient of
170 determination (R^2), the fraction of variance in the true copy numbers explained by the prediction.

171 We evaluated the quality of confidence estimate by precision and recall. Precision is defined as
172 the proportion of accurately predicted 16S GCN in predictions considered as reliable (with $\geq 95\%$
173 confidence), and recall is defined as the proportion of reliable predictions in the accurately
174 predicted 16S GCNs. We averaged the R^2 , precision and recall for the 50 cross-validations in
175 each bin.

176

177 Evaluating the effect of 16S GCN correction on relative cell abundance estimation

178 We simulated bacterial communities with 16S GCN variation (SC1 dataset, see Supplementary
179 Methods). To estimate the confidence interval (CI) of the corrected relative cell abundance of
180 each OTU in a community, we randomly drew 1000 sets of 16S GCNs from their predicted
181 uncertainty distribution. For each set of 16S GCNs, we divided the gene read count of OTUs by
182 their corresponding 16S GCNs to get the corrected cell counts. The median of the corrected cell
183 count for each OTU in the 1000 sets is used as the point estimate of the corrected cell count, and
184 the OTU's relative cell abundance is calculated by normalizing the corrected cell count with the
185 sum of corrected cell counts of all OTUs in the community. The 95% CI for each OTU's relative
186 cell abundance is determined using the 2.5% and 97.5% quantiles of the 1000 sets of corrected
187 relative cell abundances. The support value for the most abundant OTU is calculated as the
188 empirical probability that the OTU has the highest cell abundance in the 1000 sets of corrected
189 cell abundances. We calculated the coverage probability of the CI as the empirical frequency that
190 the relative gene abundance or true relative cell abundance is covered by the estimated CI. We
191 evaluated the effect of 16S GCN correction on relative cell abundance estimation at different
192 NSTD thresholds.

193

194 Evaluating the effect of 16S GCN correction on beta-diversity analyses

195 We used the Bray-Curtis dissimilarity and Aitchison distance for beta-diversity analysis that
196 requires a dissimilarity or distance matrix and evaluated the effect of 16S GCN correction on the
197 simulated bacterial communities (SC2 dataset, see Supplementary Methods). To correct for 16S
198 GCN variation in beta-diversity analyses, we divided the gene abundance of each OTU by its
199 predicted 16S GCN and calculated the corrected relative cell abundance table and the
200 corresponding dissimilarity/distance matrix. We used the corrected cell abundance table to
201 generate the principal coordinates analysis (PCoA) plot and to conduct the permutational
202 multivariate analysis of variance (PERMANOVA) and the random forest test with the R package
203 *vegan* and *randomForest*, respectively.

204

205 Examining the adjusted NSTI of empirical bacterial communities

206 To check the predictability of 16S GCN in empirical data, we examined bacterial communities
207 surveyed by 16S rRNA amplicon sequencing in the MGnify resource platform [32] that were
208 processed with the latest two pipelines (4.1 and 5.0). The MGnify resource platform uses the
209 SILVA database release 132 [1] for OTU-picking in their latest pipelines, and therefore we
210 predicted GCNs for SILVA OTUs (Supplementary Methods). We filtered the surveyed
211 communities from the MGnify platform so that only communities with greater than 80% of their
212 gene reads mapped to the SILVA reference at a similarity of 97% or greater were included. This
213 filtering yielded 113842 bacterial communities representing a broad range of environment types.
214 We calculated the adjusted NSTI for each community and examined the adjusted NSTI
215 distribution in various environmental types.

216

217 **Results**

218 Time-independent variation is present in 16S GCN evolution

219 To evaluate the extent of time-independent or intraspecific variation in 16S GCN, we examined
220 5437 pairs of genomes with identical 16S rRNA gene alignments. The 16S GCN differs in 607
221 (11%) of them, suggesting the presence of significant time-independent variation. For the 6408-
222 genomes in the reference phylogeny, we found that incorporating time-independent variation
223 with the BM model greatly improves the model fit (Table 1), indicating the necessity to take
224 time-independent variation into account in 16S GCN prediction. In addition, we observed that
225 the rate of evolution in the fitted BM model is inflated by 1670 folds when time-independent
226 variation is not included in the model, which will lead to overestimation of uncertainty in BM
227 model-based 16S GCN prediction.

228

229 Pulsed evolution model explains the 16S GCN evolution better than the Brownian motion model

230 When predicting traits using phylogenetic methods, the BM model is commonly assumed to be
231 the model of evolution. We have shown that PE model is a better model for explaining the
232 evolution of bacterial genome size [33], prompting us to test whether pulsed evolution can be
233 applied to explain 16S GCN evolution as well. Using the R package *RasperGade* that
234 implements the maximum likelihood framework of pulsed evolution [14], we fitted the PE model
235 with time-independent variation to the same dataset. Table 1 shows that the PE model provides a
236 significantly better fit than the BM model, indicating that 16S GCN prediction should assume the
237 PE model instead of the BM model. Fitted model parameters are not sensitive to the HMM
238 profiles used for aligning the 16S rRNA sequences (Table S2).

239

240 Substantial rate heterogeneity exists in 16S GCN evolution

241 To systematically examine the rate heterogeneity of 16S GCN evolution in the reference
242 genomes, we first used the variance of PICs as an approximate estimate of the local evolution
243 rate of 16S GCN. We found that the rate of evolution varies greatly among genera (Figure S1),
244 but can be roughly divided into two groups with high and low rates of evolution. Therefore, we
245 developed a heterogeneous pulsed evolution model where all jumps are the same size but the
246 frequency of jumps varies between two groups to accommodate the heterogeneity among
247 different bacterial lineages. Using a likelihood framework and AIC, we classified 3049 and 3358
248 nodes and their descending branches into slowly-evolving and regularly-evolving groups
249 respectively (Figure S2). The frequency of jumps in the regularly-evolving group is 145 folds of
250 the frequency in the slowly-evolving group (Table S3). The heterogeneous PE model provides
251 the best fit among all models tested (Table 1), indicating that a heterogeneous PE model should
252 be assumed in predicting 16S GCN.

253

254 Apart from the rate of pulsed evolution, we also observed heterogeneity in time-independent
255 variation: for the slowly-evolving group, the fitted model parameters indicate no time-
256 independent variation, while for the regularly-evolving group, the magnitude of time-
257 independent variation is approximately 40% of a jump in pulsed evolution (Table S3). The
258 presence of time-independent variation caps the confidence of prediction in the regularly-
259 evolving group at 85%, which can only be achieved when the query has identical 16S rRNA
260 gene alignment to one of the reference genomes.

261

262 RasperGade16S improves confidence estimate for 16S GCN prediction in empirical data

263 Using 16S GCN from the 6408 complete genomes in the reference phylogeny for cross-
264 validation, we compared the performance of various methods in accuracy and confidence
265 estimates. The pic and mp methods produce very large and zero uncertainty respectively (Figure
266 1A), leading to both poor precision and recall rates (Figure 1C and 1D). The emp method
267 performs the worst in terms of accuracy. The PE method produces the best overall precision
268 (Figure 1D), achieving an average precision rate of 0.96, one of the best accuracies (Figure 1B),
269 and the best confidence estimate for 16S GCN prediction (Figure 1C) over the full spectrum of
270 NSTD, and should be preferred when predicting 16S GCN.

271

272 Copy number correction improves relative cell abundance estimation

273 From theoretical calculations, in general, community members with lower relative cell
274 abundances suffer from greater impacts by 16S GCN variation (Figure 2A). If a species has a
275 higher GCN compared to the average GCN of the community, its relative abundance will be
276 overestimated. Otherwise, its presence will be underestimated (Figure 2A). In simulated dataset
277 (SC1), we found that 16S GCN variation has a large detrimental effect on the estimated relative
278 cell abundance (Figure 2B). On average, the relative cell abundance estimated using the gene
279 abundance increased or decreased by 1.8-fold compared to the true relative cell abundance, and
280 the empirical probability of correctly identifying the most abundant OTU based on the gene
281 abundance is only around 13% (Figure 2C). Correcting for 16S GCN improves the estimated
282 relative cell abundance (Figure 2B). As expected, the improvement is greatest when the adjusted
283 NSTI is small (i.e., when there are closely related reference genomes), and it gradually
284 diminishes when the adjusted NSTI increases. At the smallest adjusted NSTI, the average fold

285 change of the estimated relative cell abundance decreases to 1.1-fold after 16S GCN correction
286 and the empirical probability of correctly identifying the most abundant OTU increases to around
287 65% (Figure 2C).

288

289 Because we predict each OTU's 16S GCN with a confidence estimate, we can provide 95%
290 confidence intervals (95% CIs) for their relative cell abundance as well. Ideally, 95% of the true
291 relative cell abundances should be covered by the 95% CIs. Figure 2D shows that the average
292 coverage probability of the true relative cell abundance is about 98% across NSTD cutoffs,
293 indicating that our 95% CIs are slightly over-conservative. Similarly, we can also calculate the
294 coverage probability of our 95% CI to the relative gene abundance. As expected, when the
295 coverage probability to the relative gene abundance increases, the improvement by GCN
296 correction (quantified by the relative reduction in the difference between the estimated and true
297 cell abundances) decreases (Figure 2E), and that when this coverage probability is below 95%,
298 GCN correction always results in strong improvement in relative cell abundance estimates. In
299 empirical studies when the true abundance is unknown, we can use the coverage probability to
300 the relative gene abundance as a conservative statistic to decide if GCN correction for a
301 community will likely improve the relative abundance estimation or not. For the most abundant
302 OTU in the community, we can calculate its support value from the 16S GCN's confidence
303 estimates. We found that the calculated support value matches the empirical probability that the
304 most abundant OTU is correctly identified (Figure 2C).

305

306 To demonstrate the effect of 16S GCN correction in empirical data, we analyzed the data from
307 the first phase of the Human Microbiome Project (HMP1) and the 2000-sample subset of Earth

308 Microbiome Project (EMP). We found that on average the relative cell abundance with and
309 without 16S GCN correction changes around 1.3-fold in HMP1 and 1.6-fold in EMP. Since the
310 true abundance of OTUs is unknown, we use the coverage probability of 95% CIs to the relative
311 gene abundance described above to evaluate the effect of GCN correction. Our results indicate
312 that a majority of HMP1 (over 82%) and EMP (over 90%) samples have a coverage probability
313 below 95% (as shown in Figure 2F). Since our simulations demonstrate that GCN correction
314 improves the accuracy of relative cell abundance estimation in samples with coverage probability
315 less than 95% (as demonstrated in Figure 2E), this suggests that GCN correction will likely
316 improve relative cell abundance estimates in these HMP1 and EMP samples. In terms of the
317 most abundant OTU, we found that the identity of the most abundant OTU changes after copy
318 number correction in around 20% and 31% of the communities in HMP1 and EMP respectively.
319 The support values for the most abundant OTUs are around 0.85 on average in both datasets,
320 indicating high confidence in the identification of the most abundant OTUs.

321

322 Copy number correction provides limited improvements on beta-diversity analyses

323 To examine the effect of 16S GCN variation on beta-diversity analyses, we simulated
324 communities at different turnover rates in two types of environments where 0.25%, 1% or 5% of
325 the OTUs are enriched in one environment compared to the other (the SC2 dataset). We found
326 that when the relative gene abundance is used to calculate the Bray-Curtis dissimilarity or the
327 Aitchison distance, the positions of the samples in the PCoA plot shift from their positions based
328 on the true relative cell abundance (solid lines in Figure 3A and B), although this shift is much
329 smaller if the Aitchison distance is used. Correcting for 16S GCN reduces about 56% and 85%
330 of the shifts in the Bray-Curtis dissimilarity ($P < 0.001$, paired t-test, Figure 3A) and Aitchison

331 distance spaces ($P < 0.001$, paired t-test, Figure 3B) respectively. Despite the shift in the PCoA
332 plot, we found that the clustering of communities does not seem to be affected by the 16S GCN
333 variation.

334

335 We observed a limited effect of 16S GCN variation on PERMANOVA. Depending on the metric
336 used, the signature OTU numbers and turnover rates, the proportion of variance explained (PVE)
337 by the environmental type using the true cell abundances ranges from 5.27% to 17.20% on
338 average. Using gene abundance, the average PVE ranges from 5.27% to 17.22% and the change
339 in PVE is not statistically significant regardless the metric used, the signature OTU numbers, or
340 the turnover rates ($P > 0.002$, paired t-test with Bonferroni correction, $\alpha = 9.26 \times 10^{-4}$, Table S4),
341 indicating that PERMANOVA is not very sensitive to 16S GCN variation.

342

343 It is a common practice to compare the relative cell abundance of OTUs of interest between
344 environments. We found that such comparison is also not sensitive to 16S GCN variation (Table
345 S4), with the fold-change of relative cell abundance estimated using the gene abundance and the
346 truth highly concordant ($R^2 > 0.99$). Random forest identified from 20.0% to 89.0% of the
347 signature OTUs when the true cell abundances were used (Table S4). When the gene abundances
348 were used, this recovery rate changes to from 18.0% to 89.32% (Table S4), and the change is not
349 statistically significant ($P > 0.032$, paired t-test with Bonferroni correction, $\alpha = 1.85 \times 10^{-3}$).

350 Correcting for 16S GCN changes the recovery rate to from 17.8% to 89.2% (Table S4), and the
351 change is not significant either ($P > 0.041$, paired t-test with Bonferroni correction, $\alpha = 1.85 \times 10^{-3}$).

352 Similar results were found when we examined the effect of 16S GCN variation correction on
353 beta-diversity in empirical data (Supplementary Results).

354

355 Vast majorities of bacterial community studies should benefit from copy number correction

356 To examine if analysis of real communities would benefit from 16S GCN correction, we
357 calculated the adjusted NSTI for 113842 communities in the microbiome resource platform
358 MGnify (formerly known as EBI Metagenomics) [32] that passed our quality control. These
359 microbiomes were sampled from various environments and include host-associated microbiomes
360 in animals and plants and free-living microbiomes in soil and aquatic environments (Table S5).
361 The adjusted NSTI varies greatly among samples and the median across all samples is 0.01
362 substitutions/site. In the simulated communities, we observed that GCN correction significantly
363 improves the estimated relative cell abundances ($P < 0.001$, paired t-test) even when the adjusted
364 NSTI reaches 0.3 substitutions/site. We found that more than 99% of the communities from
365 MGnify have an adjusted NSTI less than 0.3 substitutions/site, suggesting that they should
366 benefit from 16S GCN correction when estimating the relative cell abundances. The distribution
367 of adjusted NSTI varies among different environmental types (Figure 4), but the proportion of
368 communities that will likely benefit from 16S GCN correction remains high, ranging from 98%
369 to 100%.

370

371 **Discussion**

372 We address the inherent uncertainty problem in 16S GCN prediction by directly measuring it
373 with confidence estimates. Using simulations and cross-validation, we show that the PE method
374 implemented in *RasperGade16S* outperforms other methods in both the precision and recall rates.
375 This method's strength comes from three features of its modeling of the 16S GCN evolution:
376 implementation of a pulsed evolution model and accounting for the rate heterogeneity and time-

377 independent trait variation. Pulsed evolution model expects no trait changes to occur over a short
378 branch as jumps are not likely to happen on that branch. This leads to a higher confidence to 16S
379 GCN prediction with a short NSTD, and thus improves the recall of the accurate predictions. By
380 incorporating rate heterogeneity, we can make predictions in the slowly-evolving groups with
381 high confidence, even when their NSTDs are large, thereby further improving the overall
382 precision and recall rates. In the reference phylogeny, 48% of branches were estimated to fall
383 within this slowly-evolving group, whose evolution rate is 145 times slower compared to that of
384 the regularly-evolving group. The third source of improvement for *RasperGade16S* comes from
385 accounting for time-independent variation, which can result from measurement error and
386 intraspecific variation. We show that failing to account for time-independent variation results in
387 model misspecification (Table 1) and overestimated rate of evolution for the pic method.

388
389 Having confidence estimates is critical in the presence of inherent uncertainty because they
390 provide direct evaluation of the uncertainty associated with the predictions. Using cross-
391 validation, we show that *RasperGade16S* has high precision (around 0.96), which means for
392 predictions with high confidence ($\geq 95\%$), 96% of the predictions are accurate. Therefore, we can
393 use the confidence score provided by *RasperGade16S* to select high-quality predictions if
394 necessary, or we can draw firm conclusions from the 16S rRNA data when the confidence is
395 high.

396
397 The application of confidence estimation extends beyond the prediction of 16S GCN. Because
398 the uncertainty in the prediction is inherited by statistics derived from the predicted 16S GCN,
399 we can estimate the uncertainty and confidence intervals of important parameters in downstream

400 analyses, such as the relative cell abundance. With confidence intervals, we can draw more
401 meaningful and sound conclusions, such as identifying the most abundant OTU in the
402 community with a support value. Getting confidence estimates of the relative cell abundance is
403 also important for predicting the functional profile of a community based on 16S rRNA
404 sequences. Although PICRUST2 uses an extremely lenient NSTD cut-off to eliminate
405 problematic sequences, it does not provide an accurate confidence measurement of its
406 predictions. As shown in this study, the default maximum parsimony method used by
407 PICRUST2 to predict 16S GCN essentially assumes there is no uncertainty in the predictions,
408 which is unrealistic and leads to poor precision. Incorporation of a more meaningful confidence
409 estimate of 16S GCN prediction in PICRUST2 should make its functional profile prediction
410 more informative.

411
412 Strikingly, 99% of 113842 bacterial communities we examined have an adjusted NSTI less than
413 0.3 substitutions/site, a range where we show that GCN correction improves the accuracy of the
414 relative cell abundance estimation (Figure 2B). Because these communities represent a
415 comprehensive and diverse list of natural and engineered environments, we recommend applying
416 16S GCN correction to practically any microbial community regardless of the environmental
417 type if accurate estimates of relative cell abundance are critical to the study. Our results therefore
418 affirm the conclusion of the previous studies based on analyses of a much smaller number of
419 communities [5, 7].

420
421 Few studies have investigated to what extent the bias introduced by 16S GCN variation will have
422 on the microbiome beta diversity analyses. We show that the effect sizes of 16S rRNA bias on

423 beta-diversity analyses are small. Correcting 16S GCN provides limited improvement on the
424 beta-diversity analyses such as random forest analysis and PERMANOVA test. One possible
425 reason is that for an OTU, the fold change in the relative cell abundance between samples
426 remains more or less the same with or without correcting for the copy number. For example,
427 assuming the estimated relative cell abundances of an OTU in samples A and B are r_a and r_b
428 respectively without copy number correction. When correcting for the copy number, its relative
429 abundance is adjusted with the scaling factor ACN/GCN, where the GCN is the 16S rRNA copy
430 number of the OTU and the ACN is the average copy number of the sample. Assuming the ACN
431 does not vary much between samples, then the scaling factor for the OTU will be roughly the
432 same in samples A and B. So even with copy number correction, the relative abundance change
433 will still be close to r_a/r_b .

434

435 It should be noted that having a confidence associated with the 16S GCN prediction helps to
436 estimate the uncertainty of the prediction, but it does not improve the accuracy of the prediction.
437 Accuracy of the prediction is constrained by the inherent uncertainty, which can only be
438 improved by better sampling the reference genomes. However, as our current sampling is
439 inadequate for accurate 16S GCN prediction of all environmental bacteria, we believe that
440 incorporating confidence estimates is the best practice to control for the uncertainty in the 16S
441 rRNA based bacterial diversity studies, as opposed to not correcting the GCN bias as previously
442 suggested [8, 12].

443

444 **Acknowledgements**

445 None.

446

447 **Competing interests**

448 The authors declare that they have no competing interests.

449

450 **Data Availability Statement**

451 The NCBI accession numbers of the reference genomes, the representative 16S rRNA sequences
452 and alignments, the reference phylogeny, the predicted GCN for OTU99 in the SILVA database,
453 the simulated bacterial community data and scripts to reproduce the figures and tables in this
454 study are available in the Dryad repository

455 (https://datadryad.org/stash/share/OaS9BjM_kIVdJ3WkZRT7KO8fDr8D4k8jy3LsOtlYELM).

456 The R package *RasperGade16S* can be downloaded from <https://github.com/wu-lab->

457 [uva/RasperGade16S](https://github.com/wu-lab-uva/RasperGade16S). The scripts to conduct the analyses in this study are available in the GitHub
458 repository (<https://github.com/wu-lab-uva/16S-rRNA-GCN-Predcition>).

459

460 **References**

- 461 1. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal
462 RNA gene database project: improved data processing and web-based tools. *Nucleic Acids*
463 *Res* 2012; **41**: D590–D596.
- 464 2. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a
465 chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl*
466 *Environ Microbiol* 2006; **72**: 5069–5072.
- 467 3. Klappenbach JA. rrndb: the ribosomal RNA operon copy number database. *Nucleic Acids*
468 *Res* 2001; **29**: 181–184.

- 469 4. Větrovský T, Baldrian P. The variability of the 16S rRNA gene in bacterial genomes and
470 its consequences for bacterial community analyses. *PLoS One* 2013; **8**: e57923.
- 471 5. Kembel SW, Wu M, Eisen JA, Green JL. Incorporating 16S gene copy number
472 information improves estimates of microbial diversity and abundance. *PLoS Comput Biol*
473 2012; **8**: 16–18.
- 474 6. McLaren MR, Willis AD, Callahan BJ. Consistent and correctable bias in metagenomic
475 sequencing experiments. *Elife* 2019; **8**.
- 476 7. Angly FE, Dennis PG, Skarshewski A, Vanwonderghem I, Hugenholtz P, Tyson GW.
477 CopyRighter: a rapid tool for improving the accuracy of microbial community profiles
478 through lineage-specific gene copy number correction. *Microbiome* 2014; **2**: 11.
- 479 8. Starke R, Pylro VS, Morais DK. 16S rRNA gene copy number normalization does not
480 provide more reliable conclusions in metataxonomic surveys. *Microb Ecol* 2021; **81**.
- 481 9. Bowman JS, Ducklow HW. Microbial communities can be described by metabolic
482 structure: A general framework and application to a seasonally variable, depth-stratified
483 microbial community from the coastal west Antarctic peninsula. *PLoS One* 2015; **10**:
484 e0135868.
- 485 10. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes JA, et al.
486 Predictive functional profiling of microbial communities using 16S rRNA marker gene
487 sequences. *Nat Biotechnol* 2013; **31**: 814–821.
- 488 11. Zaneveld JRR, Thurber RL v. Hidden state prediction: a modification of classic ancestral
489 state reconstruction algorithms helps unravel complex symbioses. *Front Microbiol* 2014;
490 **5**: 431.

- 491 12. Louca S, Doebeli M, Parfrey LW. Correcting for 16S rRNA gene copy numbers in
492 microbiome surveys remains an unsolved problem. *Microbiome* 2018; **6**: 41.
- 493 13. Ané C. Analysis of comparative data with hierarchical autocorrelation. *Ann Appl Stat*
494 2008; **2**: 1078–1102.
- 495 14. Landis MJ, Schraiber JG. Pulsed evolution shaped modern vertebrate body sizes.
496 *Proceedings of the National Academy of Sciences* 2017; **114**: 13224–13229.
- 497 15. Joy JB, Liang RH, McCloskey RM, Nguyen T, Poon AFY. Ancestral reconstruction.
498 *PLoS Comput Biol* 2016; **12**: e1004763.
- 499 16. Elliot MG, Mooers AØ. Inferring ancestral states without assuming neutrality or
500 gradualism using a stable model of continuous character evolution. *BMC Evol Biol* 2014;
501 **14**: 226.
- 502 17. Eldredge N, Gould SJ. Punctuated equilibria - an alternative to phyletic gradualism.
503 *Models in Paleobiology*. 1972. pp 82–115.
- 504 18. Gao Y, Wu M. Microbial genomic trait evolution is dominated by frequent and rare pulsed
505 evolution. *Sci Adv* 2022; **8**.
- 506 19. Yano K, Masuda K, Akanuma G, Wada T, Matsumoto T, Shiwa Y, et al. Growth and
507 sporulation defects in *Bacillus subtilis* mutants with a single *rrn* operon can be suppressed
508 by amplification of the *rrn* operon. *Microbiology (N Y)* 2016; **162**: 35–45.
- 509 20. Rastogi R, Wu M, DasGupta I, Fox GE. Visualization of ribosomal RNA operon copy
510 number distribution. *BMC Microbiol* 2009; **9**: 208.
- 511 21. Stoddard SF, Smith BJ, Hein R, Roller BRK, Schmidt TM. rrnDB: Improved tools for
512 interpreting rRNA gene abundance in bacteria and archaea and a new foundation for
513 future development. *Nucleic Acids Res* 2015; **43**.

- 514 22. Sadeghifard N, Guertler V, Beer M, Seviour RJ. The mosaic nature of intergenic 16S-23S
515 rRNA spacer regions suggests rRNA operon copy number variation in *Clostridium*
516 *difficile* strains. *Appl Environ Microbiol* 2006; **72**: 7311–7323.
- 517 23. Lee CM, Sieo CC, Abdullah N, Ho YW. Estimation of 16S rRNA gene copy number in
518 several probiotic *Lactobacillus* strains isolated from the gastrointestinal tract of chicken.
519 *FEMS Microbiol Lett* 2008; **287**: 136–141.
- 520 24. Bodilis J, Nsigure-Meilo S, Besaury L, Quillet L. Variable copy number, intra-genomic
521 heterogeneities and lateral transfers of the 16S rRNA gene in *Pseudomonas*. *PLoS One*
522 2012; **7**: e35647.
- 523 25. Lavrinienko A, Jernfors T, Koskimäki JJ, Pirttilä AM, Watts PC. Does intraspecific
524 variation in rDNA copy number affect analysis of microbial communities? *Trends*
525 *Microbiol* 2021; **29**: 19–27.
- 526 26. Uyeda JC, Hansen TF, Arnold SJ, Pienaar J. The million-year wait for macroevolutionary
527 bursts. *Proceedings of the National Academy of Sciences* 2011; **108**: 15908–15913.
- 528 27. Viklund J, Ettema TJG, Andersson SGE. Independent genome reduction and phylogenetic
529 reclassification of the oceanic SAR11 clade. *Mol Biol Evol* 2012; **29**.
- 530 28. Moran NA. Accelerated evolution and Muller’s ratchet in endosymbiotic bacteria.
531 *Proceedings of the National Academy of Sciences* 1996; **93**: 2873–2878.
- 532 29. Wu D, Hugenholtz P, Mavromatis K, Pukall R, Dalin E, Ivanova NN, et al. A phylogeny-
533 driven genomic encyclopaedia of Bacteria and Archaea. *Nature* 2009; **462**.
- 534 30. Felsenstein J. Phylogenies and the comparative method. *American Naturalist* 1985; **125**:
535 1–15.

- 536 31. Louca S, Doebeli M. Efficient comparative phylogenetics on large trees. *Bioinformatics*
537 2018; **34**: 1053–1055.
- 538 32. Mitchell AL, Almeida A, Beracochea M, Boland M, Burgin J, Cochrane G, et al. MGnify:
539 the microbiome analysis resource in 2020. *Nucleic Acids Res* 2019; **48**: D570–D578.
- 540 33. Gao Y, Wu M. Modeling pulsed evolution and time-independent variation improves the
541 confidence level of ancestral and hidden state predictions. *Syst Biol* 2022.
- 542

543 **Table 1. The AICs of Brownian motion model and pulsed evolution model.**

Model	BM	BM (with time-independent variation)	PE (with time-independent variation)
Homogenous model	34338	18028	-7925
Heterogeneous model	NA	NA	-15395

544

545 **Figure legends**

546 **Figure 1. The performance of prediction on empirical 16S GCN.** Using cross-validation of

547 empirical data, the mean estimated uncertainty and confidence of predictions (A), the mean

548 coefficient of determination R^2 of the predictions (B), and the recall (C) and precision (D) of

549 classification of predictions by their associated confidence estimate, plotted against the mean

550 NSTD. The error bars represent the 95% CI of the mean. The empirical 16S GCN analyzed here

551 are from the 6408 complete genomes in the reference phylogeny.

552

553 **Figure 2. The impact of 16S GCN variation on estimated relative cell abundances.** (A) The

554 impact of GCN variation on estimated relative cell abundance based on theoretical calculations.

555 The color of the lines denotes the ratio of an OTU's of GCN to the average GCN of the

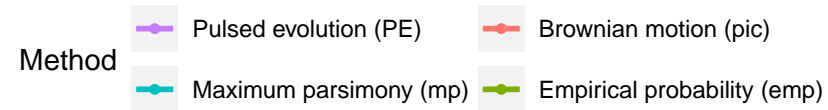
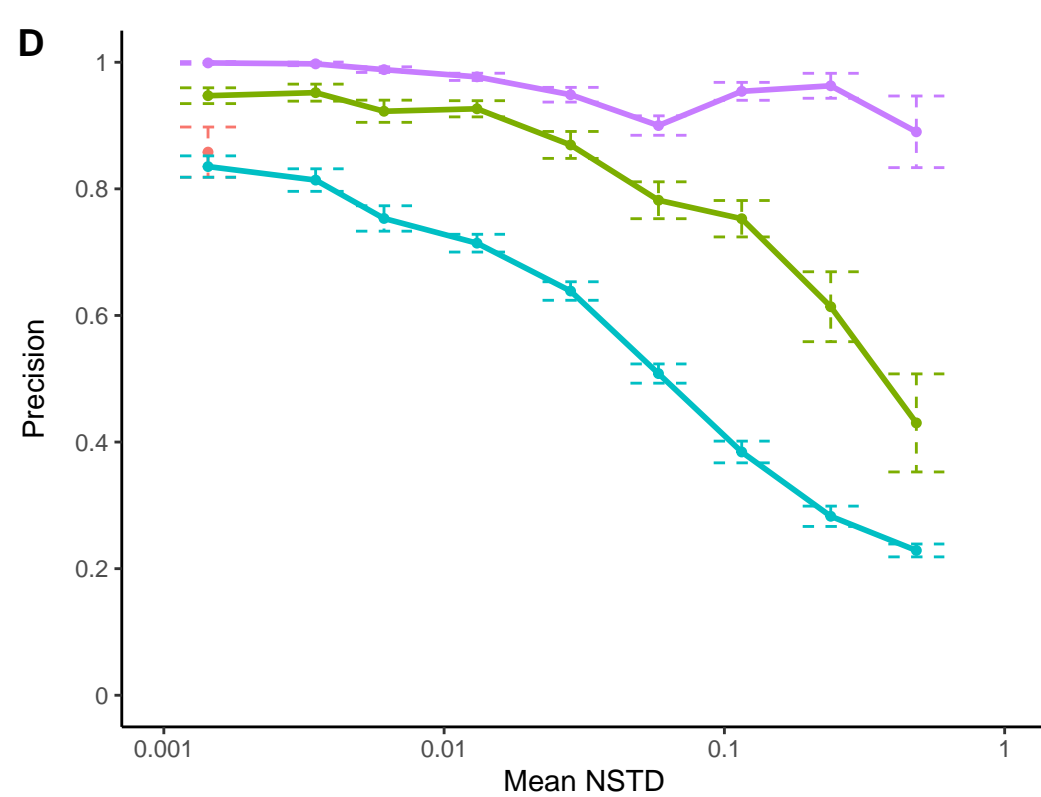
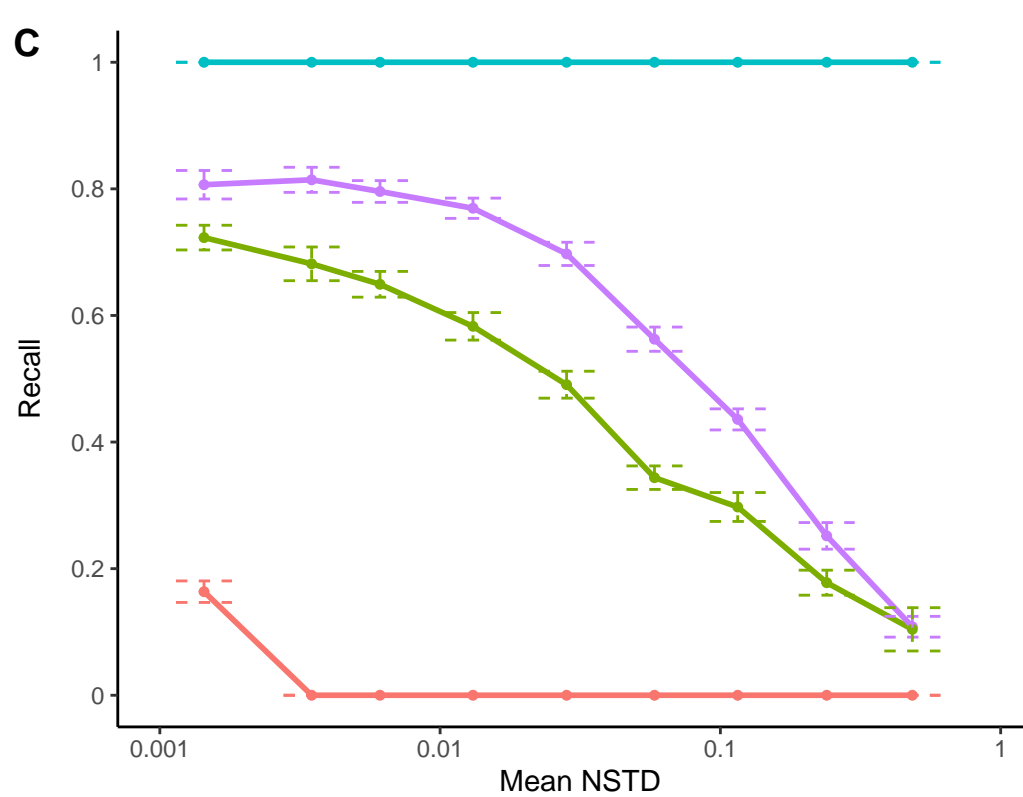
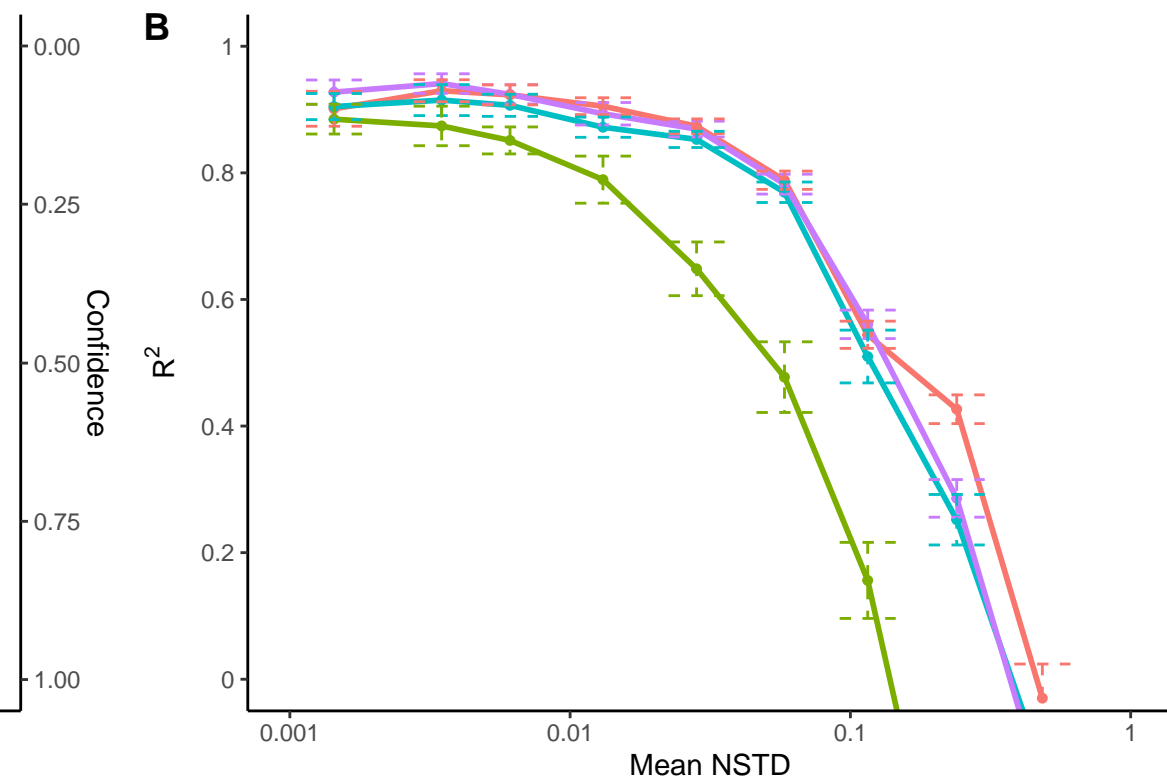
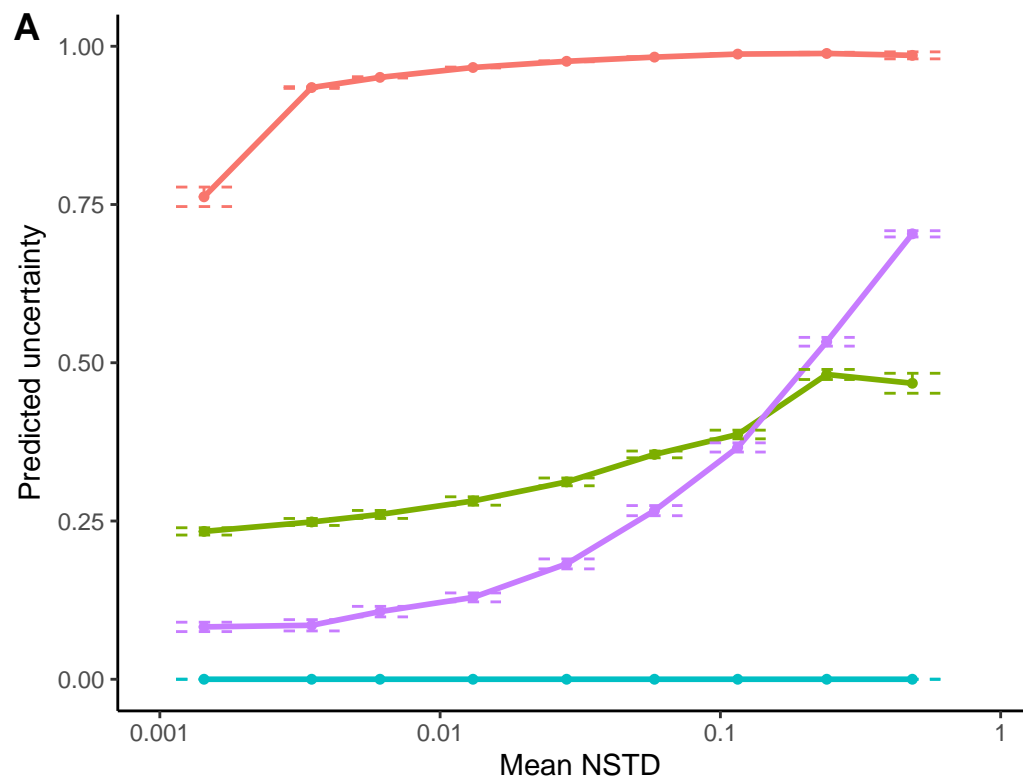
556 community. (B) The average fold-change to the true relative cell abundance. (C) The empirical
557 probability of correctly identifying the most abundant OTU in the community and the support
558 value for the most abundant OTU. (D) The coverage probability of relative cell abundances'
559 estimated 95% CIs to the true relative cell abundance. Accurate confidence estimates (95% CIs)
560 should produce a coverage probability of 95% regardless of the adjusted NSTI (dashed red line).
561 (E) The correlation between the coverage probability to the relative gene abundance and the
562 improvement by GCN correction. The horizontal red dashed line represents no improvement in
563 relative cell abundance estimates; the vertical red dashed line represents 95% coverage
564 probability to the relative gene abundance. The improvement is quantified by the relative
565 reduction in the difference between the estimated and true cell abundances. (F) The empirical
566 cumulative distribution of the coverage probability to the relative gene abundance in 2560
567 samples from the HMP1 dataset and 1856 samples from the EMP dataset. All error bars
568 represent 95% CI of the mean.

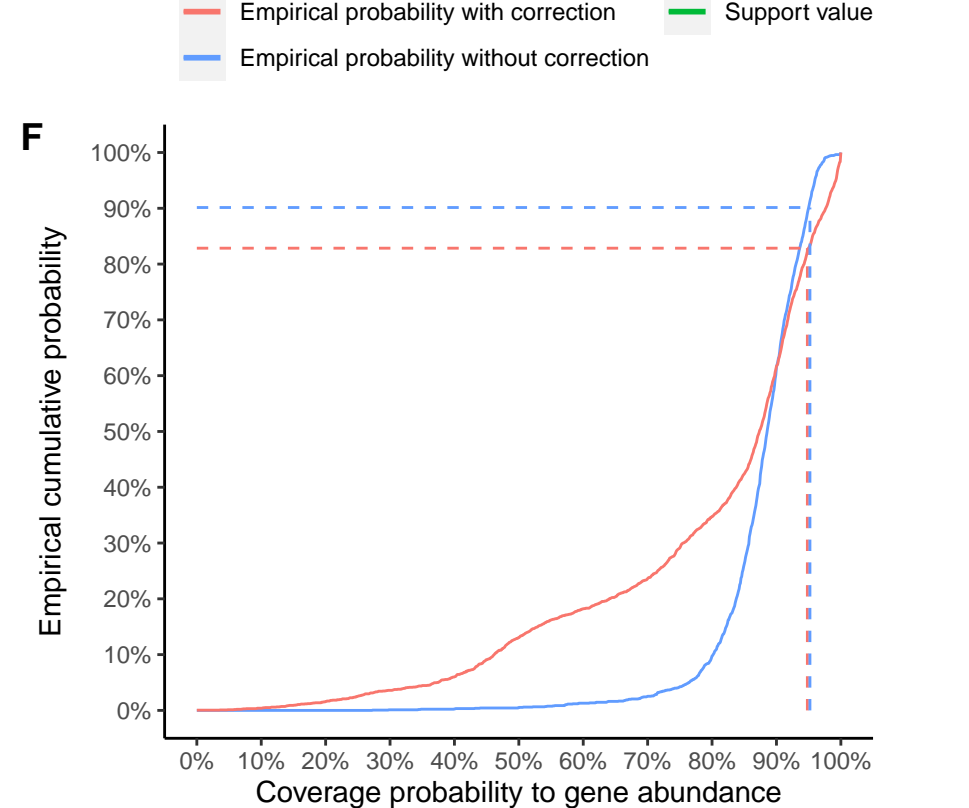
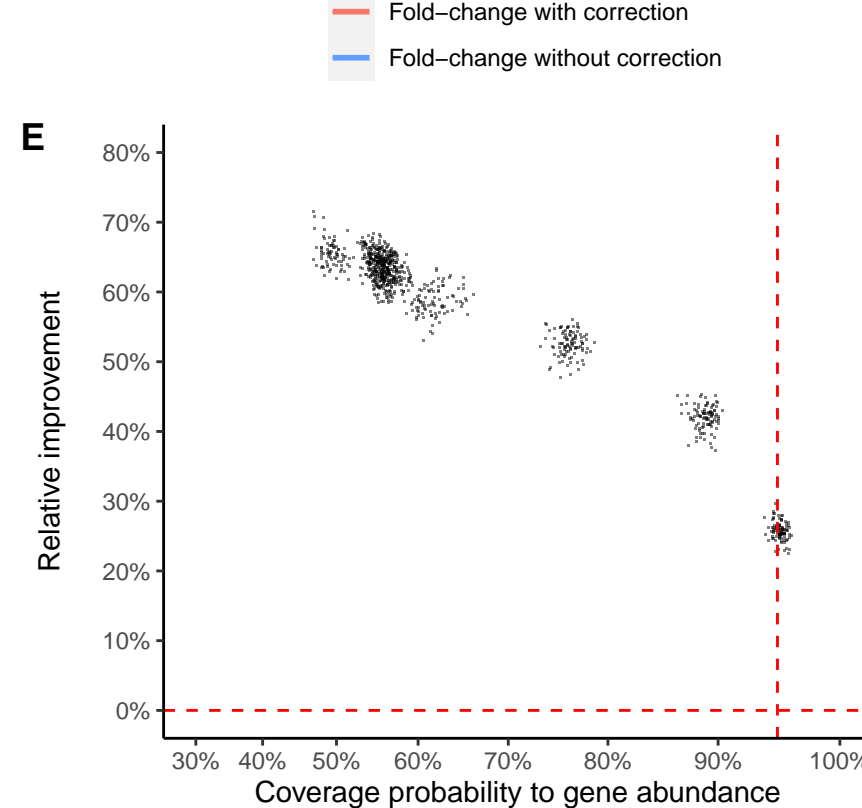
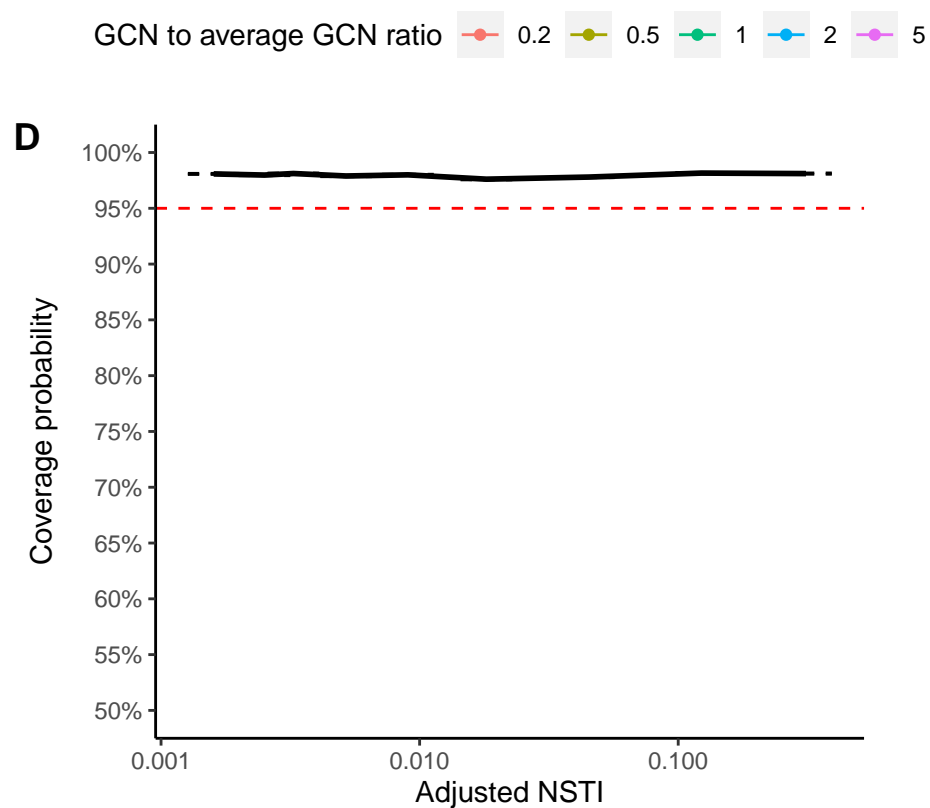
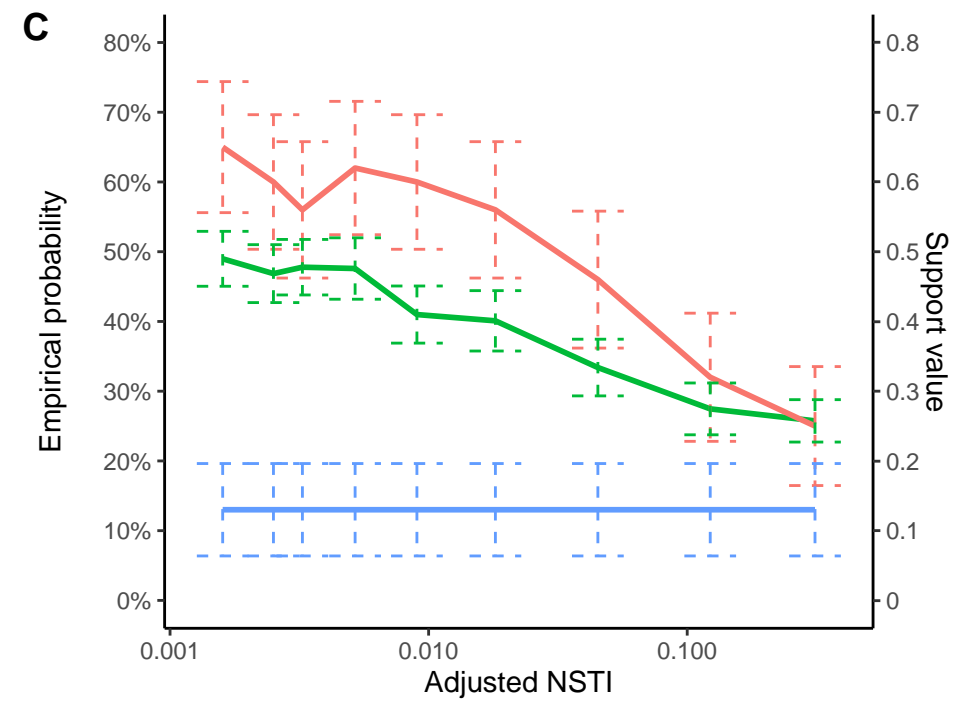
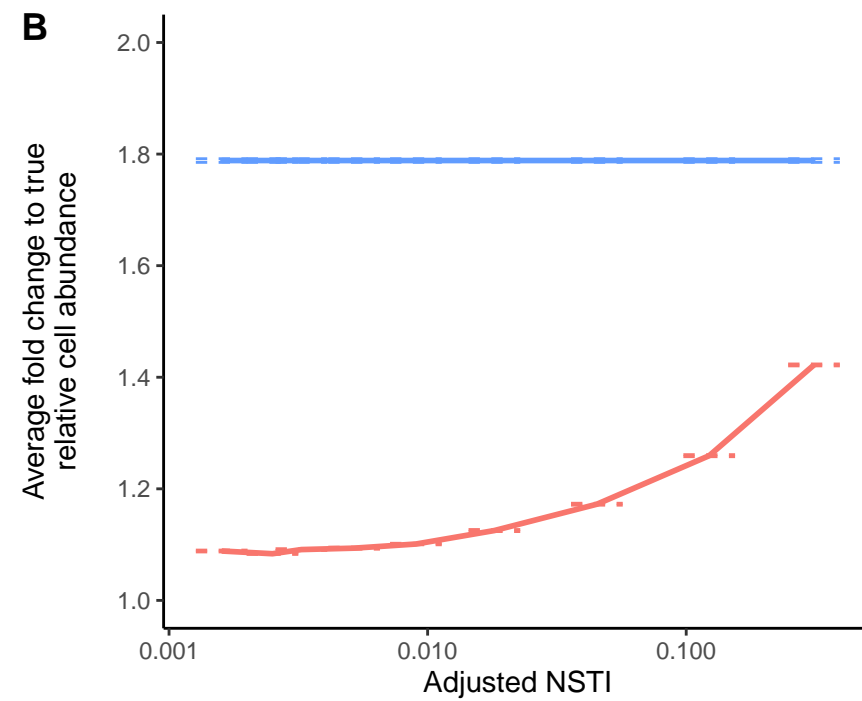
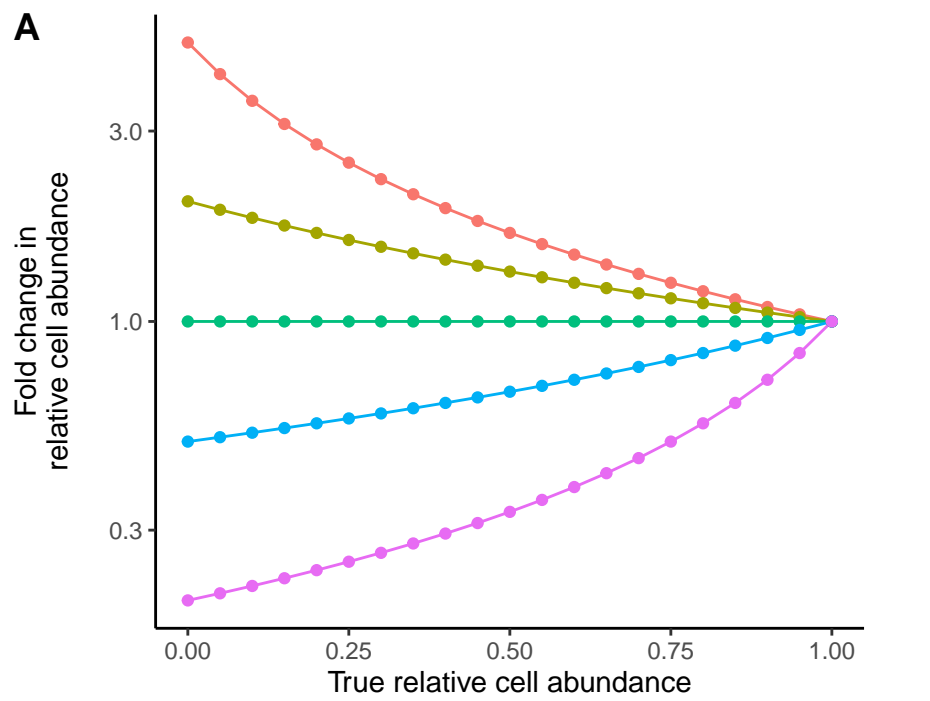
569

570 **Figure 3. The impact of 16S GCN variation on beta-diversity.** Examples of shift in the Bray-
571 Curtis dissimilarity (A) and the Aitchison distance (B) matrices due to 16S GCN variation. The
572 shift for each metric is visualized in a PCoA plot comparing 20 simulated samples from two
573 hypothetical environments with 5 signature OTUs (0.25%) in each environment and a turnover
574 rate of 20%. Solid lines represent the shift of a sample from its true location when using the gene
575 abundance. The results with 1% and 5% enriched signature OTUs are similar to the examples
576 shown in Figure 4.

577

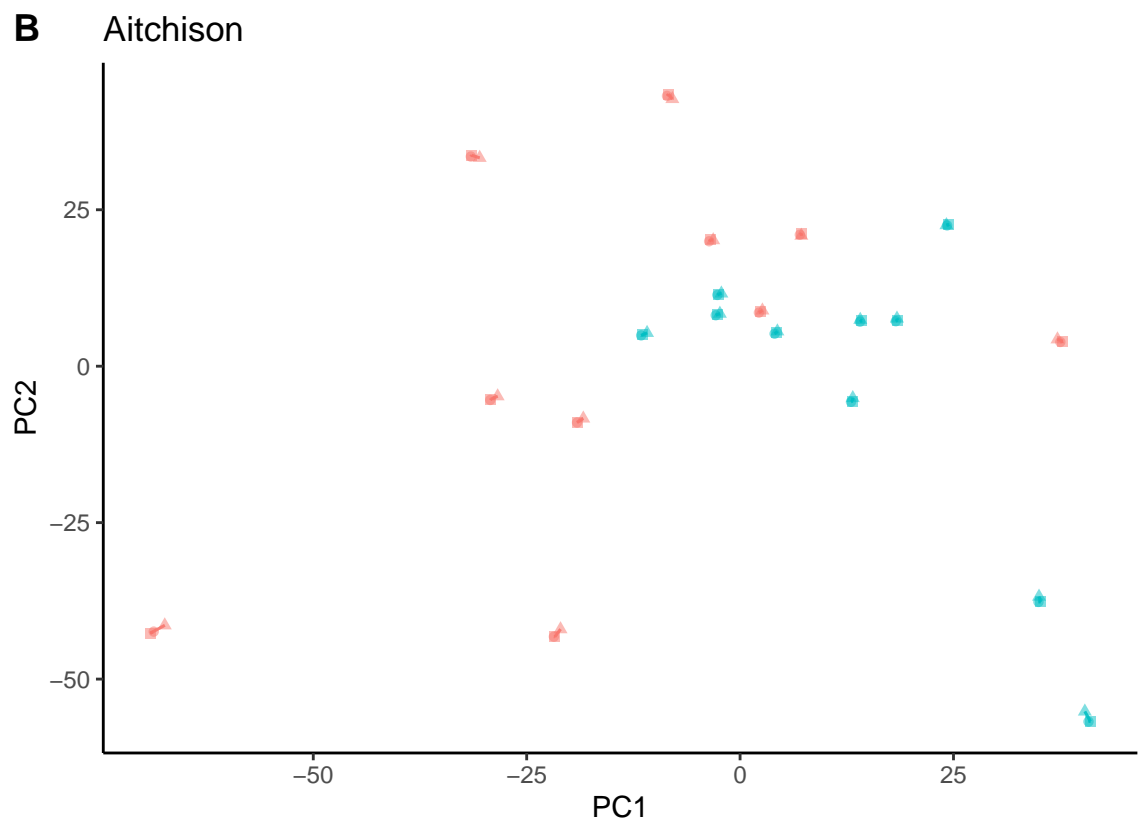
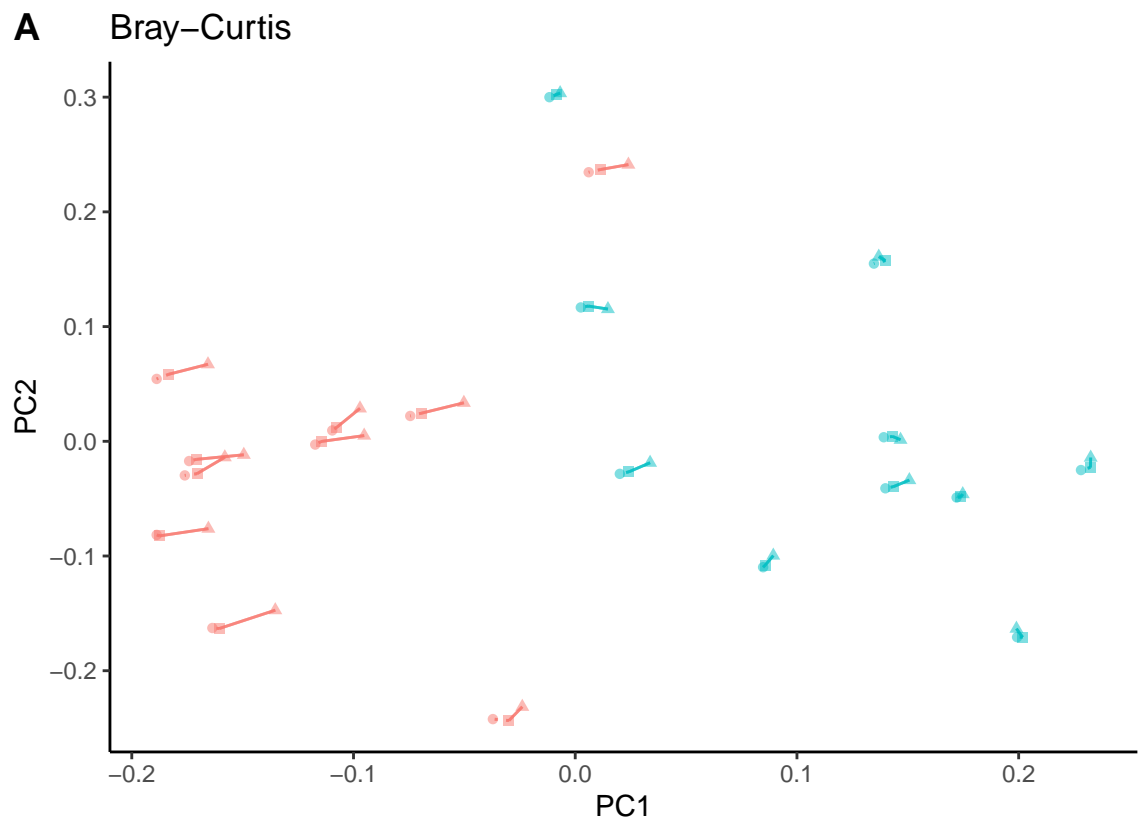
578 **Figure 4. The distribution of adjusted NSTI in empirical data.** The distribution of adjusted
579 NSTI of 113842 communities in the MGnify database representing various environmental types.
580 The red dashed line marks the adjusted NSTI of 0.3 substitutions/site.





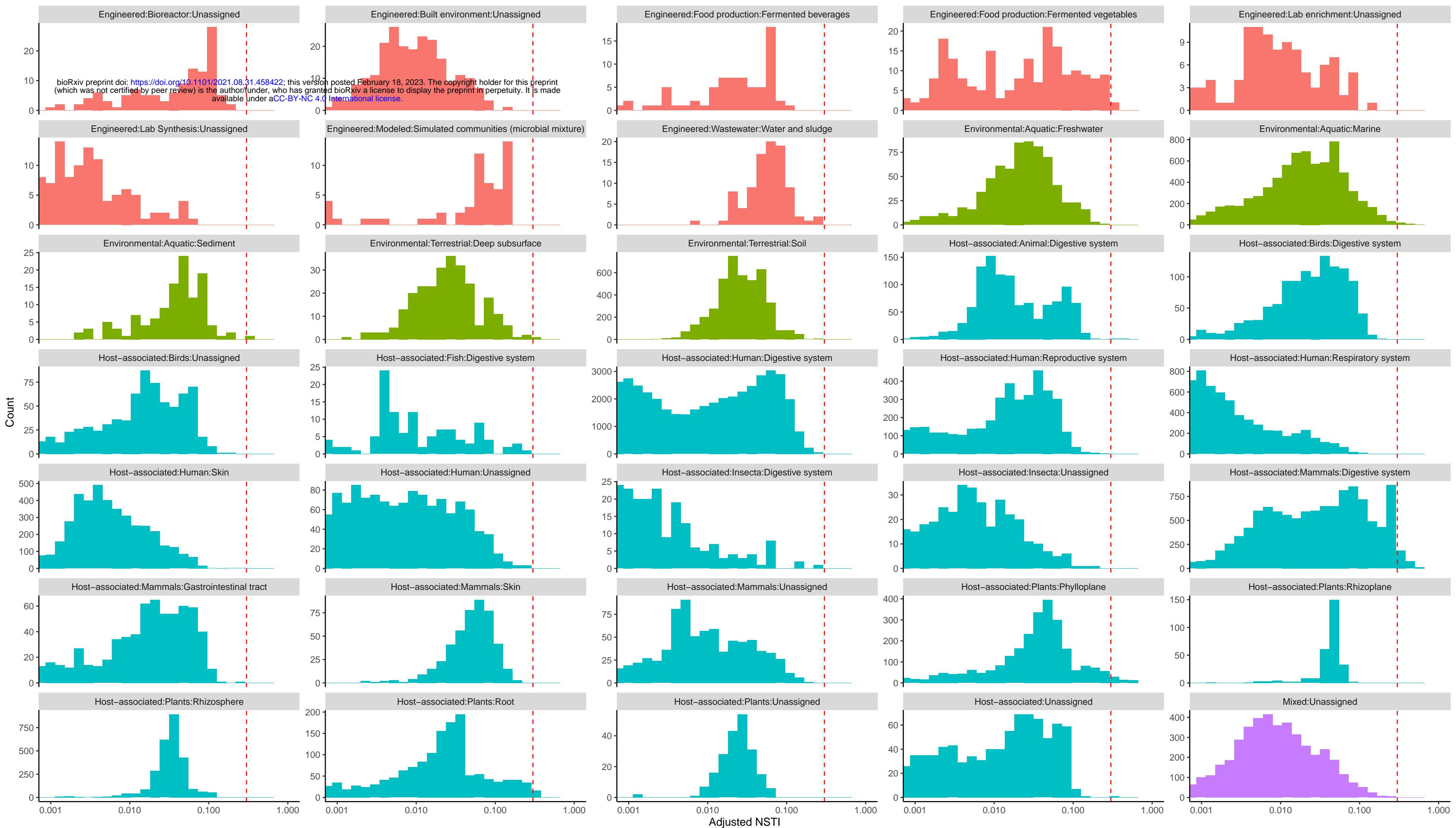
Dataset

- HMP
- EMP



Method ● Corrected abundance ▲ Gene abundance ■ True cell abundance

Environment ● 1 ● 2



Engineered Environmental Host-associated Mixed