

1 **Supplementary Methods**

2 16S rRNA GCN quality control

3 For genomes with multiple copies of 16S rRNA gene, we aligned the 16S rRNA sequences using
4 MAFFT [1] (with parameters: --maxiterate 1000 --globalpair) and picked the 16S rRNA gene
5 sequence that has the highest average similarity (calculated as the proportion of identical bases in
6 the alignment) to other 16S rRNA gene sequences in the genome as the representative sequence.
7 To remove potential errors introduced by mis-assembled genomes [2], we removed genomes
8 whose 16S rRNA GCN differs from their 5S rRNA GCN by greater than 2 copies, genomes
9 whose 16S rRNA sequence contains ambiguous bases, or genomes on the list of withheld
10 genomes in the curated ribosomal RNA operon copy number database rrnDB [3]. The 17
11 genomes in the rrnDB withheld list are rejected from rrnDB because their 16S rRNA genes are
12 missing, the 16S rRNA GCNs are too high, or the genomes have inconsistent meta data
13 (<https://rrndb.umms.med.umich.edu/withheld/>).

14

15 Reconstruction of the 16S rRNA phylogeny

16 We aligned the remaining representative 16S rRNA gene sequences using HMMER version 3.2
17 [4] (hmmalign with parameters: --trim --dna --mapali) with the hidden Markov model (HMM)
18 built from the GreenGenes 13.8 16S rRNA gene alignment (hmmbuild with default parameters),
19 and trimmed the alignment with a mask from the GreenGenes database [5]. The HMM, profile
20 alignment and the alignment mask are included in the R package RasperGade16S. After
21 collapsing identical 16S rRNA alignments, 6408 representative sequences remained. They serve
22 as the reference sequences and their taxonomies of are summarized in Table S1. We built a
23 reference tree from the trimmed alignment using RAxML version 8.2 [6] with options -f d -m

24 GTRGAMMA. We used the Deinococcus-Thermus group to root this reference phylogeny. To
25 examine the effect of sequence alignment on model fitting, we also used the 16S rRNA HMM
26 profile from the software Barnap [7] to align the 16S rRNA genes (hmmalign with default
27 parameters). We trimmed the alignment using a consensus posterior probability threshold of 0.95
28 (esl-alimask with parameters: -p --ppcons 0.95) and made a 16S rRNA phylogeny as described
29 above.

31 Modeling 16S rRNA GCN evolution with homogeneous and heterogeneous pulsed evolution 32 models

33 Using the R package *RasperGade* [8], we fitted one PE model to the entire reference phylogeny
34 and calculated the likelihood of this homogeneous PE model. An analysis of the variance of the
35 PICs associated with each genus indicated that there is a slowly-evolving group and a regularly-
36 evolving group, with the average rate of the slowly-evolving group estimated to be at least 100-
37 fold lower than that of the regularly-evolving group (Fig. 1). To model the rate heterogeneity, we
38 created two PE models: $PE_{regular}$ for the regularly-evolving group and PE_{slow} for the slowly-
39 evolving group. We then use a two-step iterative binning procedure to estimate the parameters of
40 $PE_{regular}$ and PE_{slow} (i.e., jump size and frequency). The $PE_{regular}$ model was initiated to take the
41 parameter values of the homogeneous PE model. PE_{slow} was initiated to have a jump size equal
42 to that of $PE_{regular}$ but a jump frequency 100-fold lower. In our first round of binning, from the
43 root to the tip of the reference phylogeny, we classified each node into the regularly- or slowly-
44 evolving group by testing which model ($PE_{regular}$ or PE_{slow}) provided a better fit. We merged
45 neighboring nodes belonging to the same group into one neighborhood and flipped neighborhood
46 assignment if the flip resulted in an improved overall AIC value. After the first round of binning,

47 we updated $PE_{regular}$ and PE_{slow} by fitting $PE_{regular}$ to nodes that were classified as regularly-
48 evolving and PE_{slow} to slowly-evolving nodes. We used the updated models to perform a second
49 round of binning to assign each node in the phylogeny to a group. Finally, we calculated r , the
50 rate of evolution in each group, as the process variance per unit branch length defined in a
51 previous study [9]. We then rescaled the reference tree by multiplying the branches in the
52 slowly-evolving group by the ratio $r_{slow}/r_{regular}$. To accommodate time-independent variation in
53 the tip trait values, we calculated a branch length over which the process variance of the fitted
54 pulsed evolution model is equal to the model's time-independent variation, and added this branch
55 length to each tip branch. We compared the homogeneous and heterogeneous PE models by AIC.
56

57 Simulating bacterial communities with 16S rRNA GCN variation

58 To evaluate the effect of 16S rRNA GCN correction on bacterial diversity analyses, we
59 simulated two sets of bacterial communities using the reference genomes: one set for relative cell
60 abundance analyses (SC1) and the other set for beta-diversity analyses (SC2).

61
62 For SC1, we simulated a total of 100 communities. For each simulated community, we randomly
63 selected 2000 OTUs from the 6408 reference genomes, treating each reference genome as one
64 OTU, and assigned each OTU a cell abundance randomly drawn from a log-series species
65 abundance distribution with the expected number of individuals in the community set to 40000
66 and Fisher's α set to 400.

67
68 In SC2, we simulated communities in two environmental types to evaluate the effect of 16S
69 rRNA GCN correction on beta diversity analyses. We simulated 10 communities per

70 environmental type and 2000 OTUs per community. The 16S rRNA GCN of each OTU was
71 assigned randomly from the reference genomes' GCN. We controlled the community turnover
72 rate by controlling the number of unique OTUs in each community. For example, at a turnover
73 rate of 10%, a community would have 200 unique OTUs and 1800 core OTUs that are shared
74 among all communities across all environmental types. We varied the turnover rate from 10% to
75 90% at 10% intervals. To control for the effect size of environmental type, we assigned 5
76 (0.25%), 20 (1%) or 100 (5%) signature OTUs from the core OTUs to each environmental type.
77 These signature OTUs were twice more likely to be placed in top ranks of the log-series
78 distribution (i.e., to be more abundant) than the non-signature OTUs in their corresponding
79 environmental type. We simulated 50 batches of communities for each combination of 9 turnover
80 rates and 3 signature OTU numbers, resulting in 27000 (10 communities/type × 2 types × 50 × 9
81 × 3) simulated communities in SC2.

82

83 Evaluating the effect of GCN correction in HMP1 and EMP dataset

84 To check the effect of 16S rRNA GCN correction in empirical data, we analyzed the 16S rRNA
85 V1-V3 amplicon sequence data of the first phase of Human Microbiome Project (HMP1) [10]
86 and the sequence data processed by Deblur [11] in the first release of the Earth Microbiome
87 Project (EMP) [12]. The 16S rRNA GCN for each OTU in the HMP1 and EMP datasets was
88 predicted using *RasperGade16S*. We picked 2560 samples in the HMP1 dataset with complete
89 metadata and used the 2000-sample subset of EMP, and determined the adjusted NSTI and
90 relative cell abundance in each community as described above. For beta-diversity, we randomly
91 picked 100 representative samples from each of the 5 body sites in the HMP1 dataset and
92 analyzed their beta-diversity as described above. For the EMP dataset, we analyzed the beta-

93 diversity within each level-2 EMP ontology (EMPO) category (around 400 to 600 samples per
94 category).

95

96 Predicting 16S rRNA GCN for SILVA OTUs

97 We downloaded 592605 full-length representative bacterial 16S rRNA sequences of non-
98 redundant OTUs at 99% similarity (OTU99) in the SILVA release 132 [13]. We aligned and
99 trimmed the sequences using the method described above. We then inserted the OTUs into the
100 reference phylogeny using the evolutionary placement algorithm (EPA-ng) [14] with the model
101 parameters estimated by RAxML when building the reference phylogeny. We limited the
102 maximum number of placements per SILVA representative sequence to 1. We predicted the 16S
103 rRNA GCN for each SILVA OTU99 as described above using the heterogeneous pulsed
104 evolution model and calculated adjusted NSTDs.

105

106 **Supplementary Results**

107 Copy number correction provides limited improvements on beta-diversity analyses in empirical 108 data

109 We analyzed the beta-diversity using the HMP1 and EMP datasets. Because we observed that the
110 effect of GCN correction is independent of the metric used in beta-diversity analyses, we only
111 used Bray-Curtis dissimilarity in HMP1 and EMP datasets. We found that correction of 16S
112 rRNA GCN does not seem to affect the clustering of communities by body sites in the HMP1
113 PCoA plot. Pairwise PERMANOVA shows that the mean PVE by the body site in HMP1 is
114 14.9% before 16S rRNA GCN correction and decreases marginally to 14.6% after correction,
115 and the PVEs using the gene abundance and the corrected cell abundance are also highly

116 concordant ($R^2 > 0.98$). In EMP, within each level-2 environment (EMPO2), the average PVE by
117 level-3 environment (EMPO3) remains at 7.7% before and after 16S GCN correction and the
118 PVEs using the gene abundance and the corrected cell abundance are highly concordant
119 ($R^2 > 0.99$) as well. On the other hand, pairwise random forest tests yield similar results before
120 and after 16S rRNA GCN correction, with around 9 out of the top 10 features identified by the
121 random forest test remaining unchanged before and after correction in HMP1 and around 8 out
122 of the top 10 unchanged in EMP. In terms of the fold-change of relative cell abundances between
123 body sites, we found that copy number correction has little impact as the estimated fold-change
124 before and after correction are highly similar ($R^2 > 0.95$) in both datasets.

125

126 Predicting 16S rRNA GCNs for SILVA OTUs

127 Using *RasperGade16S*, we predicted the 16S rRNA GCN for 592605 bacterial OTUs (99%
128 identity) in the release 132 of the SILVA database. Overall, the median adjusted NSTD for all
129 bacterial OTUs is 0.070 substitutions/site, and 34.7% of the predictions have a high confidence
130 of 95% or greater, and 74.9% of the predictions have a moderate confidence of 50% or greater
131 (Table S2). This shows that for most OTUs in the SILVA database, the phylogenetic distance to
132 a reference 16S rRNA is small enough that we can have reasonable confidence in the predictions.
133 In comparison, randomly guessing has a null confidence of around 6.7% (1 out of 15 possible
134 GCNs). Among major phyla with more than 10000 OTUs, the proportion of highly confident
135 predictions varies greatly (Table S2), with Cyanobacteria having the lowest proportion of 19.1%
136 and Acidobacteria having the highest proportion of 50.4%. Similarly, the proportion of
137 moderately confident predictions varies from 58.3% to 89.5% among these phyla. Interestingly,

138 the proportions of highly confident predictions closely match the proportions of slowly-evolving
139 OTUs in each phylum (Table S2), suggesting a causal relationship between them.

140

141 **Figure S1. The distribution of rate groups along the 16S rRNA reference phylogeny.** The
142 distribution of rate groups is denoted by colors. Red color represents slowly-evolving group and
143 black color represents regularly-evolving group. The branch lengths displayed in the figure are
144 not scaled by the GCN evolution rate.

145

146 **Figure S2. The impact of 16S GCN variation on NMDS analysis.** Simulated samples from
147 two hypothetical environments are plotted using the Bray-Curtis dissimilarity (top row),
148 weighted UniFrac distance (middle row), and the Aitchison distance (bottom row) matrices, and
149 the true cell abundance (left column), gene abundance (middle column) and corrected abundance
150 (right column). In each plot, there are 20 simulated samples from two hypothetical environments
151 with 20 signature OTUs (1%) in each environment and a turnover rate of 20%.

152

153 **Table S1. The taxonomic composition of genomes in the reference phylogeny.** The count of
154 reference genomes within a clade is listed at phylum, class, order, and family level.

155

156 **Table S2. Summary of SILVA 16S rRNA GCN predictions.** Highly confident predictions are
157 defined as predictions with a confidence of 95% or greater. Moderately confident predictions are
158 defined as predictions with a confidence of 50% or greater.

159

160 **Table S3. The effect of HMM profiles on model fitting.** The AIC and parameters of fitted
161 Brownian motion (BM) and pulsed evolution (PE) models when the alignment of 16S rRNA
162 genes uses different HMM profiles are listed.

163

164 **Table S4. Fitted parameter of homogeneous and heterogeneous pulsed evolution models.**

165 The jump frequency, jump size and the magnitude of time-independent variation of the fitted
166 homogeneous and heterogeneous pulsed evolution models are listed. The unit of jump frequency
167 is jump per unit branch length.

168

169 **Table S5. The effect of 16S rRNA GCN correction on beta-diversity analyses.** The

170 performance statistics of random forest tests, PERMANOVA, and abundance comparison before
171 and after 16S rRNA GCN correction are listed at different signature OTU number and turnover
172 rate.

173

174 **Table S6. Count summary of environmental types in MGnify dataset.** The count number of
175 samples within each biome (environmental type) is listed at the first and second biome level.

176

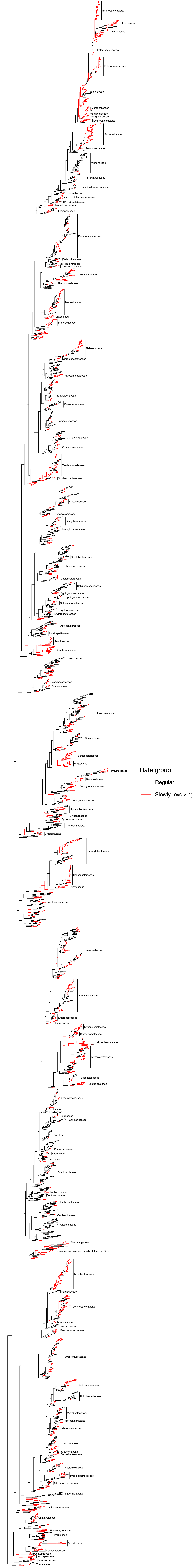
177 **References**

- 178 1. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
179 Improvements in performance and usability. *Mol Biol Evol* 2013; **30**: 772-80.
- 180 2. Perisin M, Vetter M, Gilbert JA, Bergelson J. 16Stimator: statistical estimation of
181 ribosomal gene copy numbers from draft genome assemblies. *ISME J* 2016; **10**: 1020–
182 1024.

- 183 3. Klappenbach JA. rrndb: the ribosomal RNA operon copy number database. *Nucleic Acids*
184 *Res* 2001; **29**: 181–184.
- 185 4. Eddy SR. Accelerated profile HMM searches. *PLoS Comput Biol* 2011; **7**: e1002195.
- 186 5. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a
187 chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl*
188 *Environ Microbiol* 2006; **72**: 5069–5072.
- 189 6. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
190 large phylogenies. *Bioinformatics* 2014; **30**: 1312–1313.
- 191 7. Torsten Seemann. Barrnap. <https://github.com/tseemann/barrnap>. Accessed 12 Mar 2022.
- 192 8. Gao Y, Wu M. Modeling pulsed evolution and time-independent variation improves the
193 confidence level of ancestral and hidden state predictions. *Syst Biol* 2022; **71**:1225-1232.
- 194 9. Landis MJ, Schraiber JG. Pulsed evolution shaped modern vertebrate body sizes.
195 *Proceedings of the National Academy of Sciences* 2017; **114**: 13224–13229.
- 196 10. The Human Microbiome Project Consortium. A framework for human microbiome
197 research. *Nature* 2012; **486**: 215–221.
- 198 11. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, et al.
199 Deblur rapidly resolves single-nucleotide community sequence patterns. *mSystems* 2017;
200 **2**: e00191-16.
- 201 12. Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, et al. A communal
202 catalogue reveals Earth’s multiscale microbial diversity. *Nature* 2017; **551**: 457–463.
- 203 13. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal
204 RNA gene database project: improved data processing and web-based tools. *Nucleic Acids*
205 *Res* 2012; **41**: D590–D596.

206 14. Barbera P, Kozlov AM, Czech L, Morel B, Darriba D, Flouri T, et al. EPA-ng: Massively
207 parallel evolutionary placement of genetic sequences. *Syst Biol* 2019; **68**: 365–369.

208



Rate group
 — Regular
 — Slowly-evolving

