

34 **Author Summary**

35 The evolutionary forces driving genome size in bacteria and archaea have been subject to debate
36 during the last decades. Independent comparative analyses have suggested that unique variables,
37 such as the strength of selection, environmental complexity, and mutation rate, are the main drivers
38 of this trait, which complicates generalizations across the Tree of Life. Here, we applied a
39 phylogeny-based statistical approach to assess how tightly genome size is linked to evolutionary
40 history in bacteria and archaea. Moreover, we also evaluated the predictability of genome size
41 from the strength of purifying selection and ecological strategy on a broad diversity of bacteria
42 and archaea genomes. Our approach indicates that genome size in prokaryotes is strongly
43 dependent on phylogenetic history, and that genome size is the result of the interaction of variables
44 like past events, current selection regimes, and environmental complexity that are clade dependent.

46 **Introduction**

47 Bacterial and Archaeal genomes are densely packed with genes and contain relatively little non-
48 coding DNA, and therefore an increase in genome size is directly translated into more genes [1–
49 3]. In contrast, multicellular eukaryotes generally show genome expansion due to the proliferation
50 of noncoding-DNA as a consequence of high genetic drift [2]. The absence of non-functional
51 elements in prokaryotes is explained through the deletion bias process; newly acquired genes or
52 existing genes are removed through deletions if selection on those genes is ineffective enough due
53 to low selection coefficient [4–6]. Although narrowly constrained when compared with
54 eukaryotes, prokaryotic genome sizes still vary by over an order of magnitude. Assuming an
55 intrinsic deletion bias in prokaryotes, it remains unclear what evolutionary forces determine which
56 genes are maintained and which are lost, and what determines the variability of genome sizes
57 across the broad diversity of bacteria and archaea.

58
59 Multiple individual factors have been hypothesized to be primary drivers of genome size in
60 bacteria and archaea. Early studies suggested that effective population size (N_e) may be the
61 primary force that determines genome size in prokaryotes. For example, genome reduction has
62 been observed in host-dependent microorganisms that have small N_e due to bottlenecks and
63 therefore experience high levels of genetic drift. Under such evolutionary constraints, slightly
64 deleterious deletions are accumulated and cause overall genome reduction [7–11]. Paradoxically,

65 later studies focusing on abundant free-living planktonic lineages in the ocean suggested that
66 genome reduction can also be observed in bacteria with large N_e that experience strong purifying
67 selection [12–15]. Factors other than N_e and the strength of purifying selection have also been
68 postulated to play a role in determining prokaryotic genome size. Recently, one study suggested
69 that environmental stress leads to genome streamlining in soil bacteria [16], and that habitat
70 complexity and ecological strategy therefore may also play a major role in determining genome
71 size. The mutation rate has also been proposed to be a major factor that determines genome size
72 [17,18]. In particular, it was suggested that a high mutation rate would cause the erosion of genes,
73 loss of function, and subsequent reduction in genome size of streamlined and host-dependent
74 microorganisms [17–19]. Given the large number of forces that have been proposed to be primary
75 determinants of genome size, it remains largely unknown whether genome size in prokaryotes is
76 driven by unique variables, their interaction, or variables that have specific influence depending
77 on the lineage.

78
79 In order to explore the evolutionary forces driving genome size in bacteria and archaea, we tested
80 several hypotheses using a collection of genomes encompassing a broad diversity of bacteria and
81 archaea available on the Genome Taxonomy Database (Fig. 1) (Genome Taxonomy Database,
82 GTDB, [20]). We first examined how strongly genome size is linked to prokaryotic phylogeny in
83 order to evaluate whether the recent shared evolutionary history of many lineages may explain
84 why some factors have previously been shown to be correlated with genome size. Because genome
85 size has most commonly been viewed as a result of either effective population size or ecological
86 niche [3,21], we also evaluated the use of dN/dS ratios and *rrn* operon copies as proxies,
87 respectively. Lastly, we then examined the power of predictability of genome size from these
88 variables using a phylogeny-based statistical approach that explicitly accounts for the evolutionary
89 relationships between different taxa. Our work provides important insights into the complex
90 mechanisms that shape genome size in bacteria and archaea.

91
92
93
94
95

96 **Results and discussion**

97 **Genome size distribution across major phyla of bacteria and archaea**

98 In order to explore the distribution of genome size across the Tree of Life of bacteria and archaea,
99 we built a phylogenetic tree using representative genomes of 836 genera and 33 phyla available
100 on the Genome Taxonomy Database [20] (Fig. 1, referred hereafter as GTDB genomes dataset.
101 See methods for details on the criteria for genomes selection). For this phylogeny we used a set of
102 ribosomal proteins and RNA polymerase subunits that we have previously benchmarked [22]. The
103 size of genomes in our analysis and across the phylogeny varied by over one order of magnitude
104 (0.6-14.3 Mbp, Fig. 1). The minimum and maximum corresponded to two bacterial lineages with
105 contrasting lifestyles: the endosymbiont *Buchnera aphidicola* of the phylum Proteobacteria and
106 the free-living Actinobacteria *Nonomuraea sp.* (Fig.1). Regarding the distribution of genome size
107 within each phylum, the greatest variation of genome size was observed for the phyla
108 Actinobacteria and Cyanobacteria (Fig. 2A), whereas the phylum with shortest genomes belonged
109 to symbiotic bacteria of the phylum Patescibacteria.

110
111 Since the strength of selection and ecological strategy have been proposed to be important drivers
112 of genome size in prokaryotes [21,23–25], we calculated the dN/dS and 16S rRNA copy number
113 as approximations to the strength of selection and ecological strategy, respectively. We found that
114 the largest genomes tended to have intermediate dN/dS values, while small genomes were found
115 across a wide range of selection strengths (Fig. 3). These observations are consistent with previous
116 descriptions of genome reduction occurring at high levels of purifying selection (i.e., *Pelagibacter*
117 and *Prochlorococcus*) [15] but also under strong genetic drift (i.e., *Rickettsia*, *Blattabacterium*,
118 and *Buchnera*) (Fig. 3) [21,26], indicating that there is not a strict linear relationship between
119 genome size and selection strength. Similarly, we observed genomes with multiple 16S rRNA
120 copies with variable dN/dS and genome size values (Fig. 3). Although we did not observe linear
121 relationships between genome size and dN/dS or 16S rRNA copy number, we next sought to
122 explore the predictability of this genomic trait from the latter parameters using a quantitative
123 phylogenetic framework.

124

125

126

127 **Genome size is strongly dependent on phylogenetic history and clade-specific factors**

128 Due to the recent shared evolutionary history of many bacteria and archaea, any study involving
129 statistical analyses and species' data potentially violates the assumption of independence of
130 residuals [27,28], and phylogenetic methods are therefore required to analyse evolutionary
131 relationships between traits. We sought to assess whether genome size distributions have a
132 phylogenetic signal (i.e., that genome size is not randomly distributed across the Tree of Life and
133 genome size variation is equivalent to phylogenetic distance). As a first approximation, we
134 estimated Blomberg's K [29] on the genome size of the GTDB genomes dataset (Figs. 1 and 2A).
135 Values of Blomberg's K between 0 and 1 indicate that the genome size of closely related genomes
136 resemble each other, but less than expected under the Brownian Motion model (BM) of trait
137 evolution, whereas K=1 is evidence that genomes size varies according to the Brownian Motion
138 expectation [29]. We observed phylogenetic signal in the data, but less than expected under the
139 Brownian Motion model (BM) (K=0.51, P=0.001), suggesting that although genome size in our
140 data shows phylogenetic signal, variation is not fully explained through phylogenetic distance and
141 relationships [30]. In addition, we tested the fit of different models of evolution for genome size,
142 including Brownian Motion [30], Ornstein-Uhlenbeck [31], Early-Burst [32], a diffusion model,
143 Pagel's model [33], a drift model, and a white-noise model (non-phylogenetic signal) (Table 1).
144 According to a likelihood ratio test (P<0.001 when compared with the next-best likelihood), the
145 Pagel's model showed the best fit (Table 1), supporting our previous finding that genome size in
146 bacteria and archaea shows phylogenetic signal but it is not fully driven by phylogenetic history,
147 which would be expected under the BM model.

148
149 We next used a phylogenetic generalized least squares model (PGLS) under Pagel's approach to
150 control for potential non-independence in the residuals of our regression models derived from the
151 shared ancestry of the genomes analysed [34]. The Pagel's lambda (λ) represents how strongly
152 phylogenetic relationships predict the observed pattern of variation of a trait at the tips of a
153 phylogeny and varies from 0 (no phylogenetic signal) to 1 (phylogenetic signal observed) [33].
154 According to our estimate, $\lambda=0.98$ (95% CI= 0.96-0.99; Table 2), genome size also shows non-
155 phylogenetic independence in the residuals of the regression, confirming the suitability of a PGLS
156 for the purposes of our analyses [35]. These findings indicate that phylogenetic history alone is a
157 strong predictor of genome size, and that genome size in bacteria and archaea does not evolve

158 independently. Similar results were obtained previously in a study analyzing the relationship
159 between $N_e\mu$ and genome size but with a smaller set of prokaryotic and eukaryotic genomes
160 [36,37], suggesting that sample size does not have an effect on the conclusions of our study.
161 Moreover, we estimated kappa (k) and delta (δ), two parameters that describe the mode of
162 evolution of a trait (punctuated vs gradual) and the rate change across the phylogeny (acceleration
163 vs deceleration), respectively [38]. Our estimates ($k=0.48-0.49$ and $\delta=2.44-2.51$, Table 2) are
164 consistent with a gradual and late diversification of genome size in bacteria and archaea, which
165 might indicate species-specific adaptations [33,38].

166

167 **Non-phylogenetic regression overestimates the effect of dN/dS on genome size**

168 Previous studies have suggested that high levels of genetic drift are related to a decrease in genome
169 size in bacteria [21,26]. However, such studies were based on a limited set of genomes available
170 at the time and did not include a broad repertoire of streamlined genomes, which are notable for
171 their small genomes and large effective population sizes [10,39]. In order to investigate whether
172 this trend is maintained when including a broader diversity of taxa, we calculated the dN/dS on
173 the GTDB genomes dataset. Although earlier studies have reported a strong relationship between
174 genome content and dN/dS [21], our non-phylogenetic generalized least squares (GLS) showed a
175 positive and significant but poor predictability of dN/dS when using a broader set of genomes
176 ($P<0.001$, $R^2=0.04$, Table 1, Fig. 4A). Interestingly, when considering phylogeny, PGLS
177 (phylogenetic generalized least squares) showed a non-significant and considerably poorer
178 predictability ($P=0.5$, $R^2=0.0006$, Table 1, Fig. 4A). We also calculated the lambda parameter on
179 dN/dS, and the value found ($\lambda=0.68$; 95% CI= 0.56-0.77) indicates a relatively high phylogenetic
180 signal for this variable, suggesting that phylogenetically related microorganisms tend to experience
181 similar levels of selection. Altogether, these results suggest that correlations between dN/dS and
182 genome size are largely driven by artefacts that arise by not specifically accounting for the recent
183 shared evolutionary history of many lineages [28].

184

185 Although our results indicate that dN/dS is a poor predictor of genome size in bacteria and archaea
186 (Fig. 4A), it is worth mentioning that dN/dS only reflects recent evolutionary constraints due to
187 saturation of substitutions at synonymous sites [40,41]. Therefore, we do not discount the
188 possibility that genome reduction may be driven in part by processes such as population

189 bottlenecks and periods of relaxed selection that happened in the past but are not reflected in dN/dS
190 estimations. This scenario has been suggested for the streamlined autotroph *Prochlorococcus*, in
191 which the genome simplification observed in this clade could be the result of periods of relaxed
192 selection experienced in the past [41].

193

194 **Ecological strategy plays a role on genome size**

195 In addition to testing the effect of strength of selection on genome size, we also assessed the
196 predictability of genome size from 16S rRNA copies as an approximation to ecological strategy.
197 Previous studies have shown that copies of the *rrn* operon can be a predictor of the number of
198 ribosomes that a cell can produce simultaneously, and that this reflects the ecological strategy in
199 microorganisms [25,42]. A large number of *rrn* copies is associated with the ability to adapt
200 quickly to fluctuating environmental conditions (i.e., “boom and bust” strategies) [43], while
201 multiple *rrn* copies would confer a metabolic burden to slow-growing microorganisms living in
202 stable or low-nutrients environments because of ribosome overproduction [25]. Similarly to what
203 we observed for dN/dS, we found a weak, positive, and significant relationship between genome
204 size and 16S rRNA copies when using GLS ($P < 0.001$, $R^2 = 0.01$, Table 2, Fig. 4B). However, our
205 PGLS analysis did not reduce the R^2 estimate when compared with the non-phylogenetic linear
206 model ($P = 0.009$, $R^2 = 0.01$, Table 2, Fig. 4B), probably because the phylogenetic signal of 16S
207 rRNA was relatively low ($\lambda = 0.40$; 95% CI = 0.22-0.57). Although 16S rRNA copies show a poor
208 predictability, our result suggests that larger genomes harbor more copies of the *rrn* open (Fig.
209 4B), consistent with the observation that larger genomes tend to inhabit complex environments in
210 terms of temporal variability and diversity of resources [25]. In addition to fitting our model using
211 dN/dS and 16S rRNA copies individually as predictors, we fitted an additive model with both
212 variables. An ANOVA test showed that a model including both variables does not significantly
213 improve the fit when compared with the model based on 16S rRNA copies as a unique predictor
214 variable ($P = 0.48$).

215

216 **A hypothesis for the evolutionary processes that shape genome size in bacteria and archaea**

217 According to our PGLS analysis (Table 1-2, Fig. 4), evolutionary history is a dominant variable
218 determining genome size in bacteria and archaea, meaning that genomes with recent shared
219 evolutionary history tend to maintain similar sizes since their divergence from their common

220 ancestor. Nevertheless the pattern of variation in genome size differs from what would be expected
221 under the Brownian Motion model, and microorganisms of the same clade can show variable
222 genome size (Fig. 1-3). Based on our results we propose that genome size in prokaryotes is the
223 result of a complex interplay of multiple variables, including evolutionary history, past events such
224 as population bottlenecks, and environmental complexity (substrates available, variability in
225 environmental factors, biotic pressure, etc.). The strong dependence of genome size on
226 phylogenetic history suggests that when a group diverges, the resulting clades deviate from the
227 genome size of the ancestor as a result of the colonization of new habitats, niche-specific
228 adaptations, and/or population processes like bottlenecks or long population stability. Although
229 several factors have been proposed to be singular drivers of genome size in prokaryotes, such as
230 effective population size [44], ecological strategy [23], and mutation rate [17–19], our findings
231 strongly suggest that genome size is a complex trait determined by lineage-specific factors that
232 vary from group to group.

233
234 The phylogenetic signal detected in genome size does not discount that current and past processes
235 like bottlenecks have a relevant role in the genome reduction of some bacteria and archaea. This
236 is particularly expected in endosymbionts like *Buchnera* and *Blattabacterium*, which are thought
237 to derive from a large-genome ancestor [8], and are frequently going through bottlenecks and
238 periods of diversity loss [7,8,45]. Such exacerbated loss of diversity is enhanced by the nearly
239 absent homologous recombination found in vertically transmitted endosymbionts [46]. These
240 observations are consistent with the relatively high dN/dS value and small genome size that we
241 observed for *Buchnera* and *Blattabacterium* (Fig. 3). In contrast, some abundant marine clades
242 inhabiting the open ocean such as *Prochlorococcus* and *Pelagibacter* have undergone long periods
243 of adaptation and specialization to their stable environments [47,48]. The open ocean is
244 characterized by chronically-oligotrophic nutrient conditions that are stable throughout the year
245 [49], and genes that are under relaxed selection are therefore pseudogenized and lost [10]. The
246 latter is supported by the unusual growth requirements and low number of transcriptional
247 regulators found in *Pelagibacter*, which is expected to limit its response to changing environmental
248 conditions [50,51]. Consistent with these observations, we observed low dN/dS values, small
249 genome size, and fewer 16S rRNA for these streamlined bacteria (Fig. 3). The small genomes

250 observed in both endosymbionts and free-living planktonic lineages are therefore likely the result
251 of distinct evolutionary processes, as previously proposed [15].

252

253 In contrast to the genome simplification observed in host-dependent and streamlined prokaryotes,
254 genome expansion is expected in free-living lineages that inhabit complex environments like soils
255 or sediments, where microenvironments with strikingly different abiotic conditions can be found
256 millimeters apart [52]. Although temporal diversity declines and sweeps for specific gene variants
257 are likely to occur in soil prokaryotes due to rapidly changing environmental conditions [52,53],
258 larger genomes may be positively selected in these environmental realms due to variable abiotic
259 and biotic constraints. Indeed, a study exploring the genes enriched in larger genomes of soil
260 prokaryotes found a larger proportion of genes involved in regulation and secondary metabolism,
261 and were depleted in genes related with translation, replication, cell division, and nucleotides
262 metabolism when compared with smaller genomes [23]. These environmental and genomic
263 findings are consistent with the large genome sizes, high dN/dS, and multiple 16S rRNA copies
264 estimated in our study for soil microorganisms of the genera *Streptacidiphilus*, *Actinomyces*,
265 *Conexibacter*, *Actinoplanes*, and *Myxococcus* (Fig. 3), the latter showing complex fruiting body
266 development [54]. It is interesting to note that the largest genomes analyzed in our study (>6 Mpb)
267 tend to experience intermediate levels of purifying selection (dN/dS), suggesting that either
268 extremely high or low purifying selection are not conducive to genomic expansion events.

269

270 **Conclusions**

271 Despite the increase of genomes available on publicly available databases, the evolutionary
272 processes and factors driving genome size and content in bacteria and archaea are continuously
273 debated. Several studies have proposed ecological strategies, the strength of purifying selection,
274 and mutation rate as prominent forces that determine prokaryotic genome size, but our study shows
275 that these factors likely vary in importance depending on the lineage. Moreover, our statistical
276 approach showed that evolutionary history plays a large role in structuring genome size
277 distributions across the Tree of Life, and that genome size is not a phylogeny independent trait.
278 The significant but poor relationship between genome size and 16S rRNA copies suggest that
279 besides phylogenetic history, ecological strategy plays a role in shaping genome size in bacteria
280 and archaea, although this single trait is insufficient to completely represent ecological strategies.

281 Future studies will be necessary to evaluate this in detail on a lineage-by-lineage basis. The strong
282 phylogenetic signal observed in genome size data indicates that analyses involving this trait cannot
283 consider species as phylogenetically independent, therefore phylogenetic relatedness should be
284 taken into account when studying the evolutionary forces driving genome size in order to avoid
285 biased association between traits and simplified models.

286

287 **Material and methods**

288 **Genomes compilation, dN/dS estimation, and *rrn* genes identification**

289 In order to assess the predictability of genome size (response variable) from dN/dS and 16S rRNA
290 copies (predictor variables), all the bacteria and archaea representative genomes available on the
291 Genome Taxonomy Database (GTDB) (Release 05-RS95; 17th July 2020) [55] were filtered based
292 on completeness ($\geq 95\%$) and contamination ($\leq 5\%$) and then classified at the Class levels. In
293 order to include the phylum *Patescibacteria* in our analysis (also known as Candidate Phyla
294 Radiation or CPR), we used completeness $\geq 80\%$ and contamination $\leq 5\%$ for this taxa. Classes
295 having more than 500 genomes were randomly downsampled to 500 genomes. The resulting
296 genomes were clustered based on their taxonomic identity at the genus level. Genera with fewer
297 than two genomes after filtering and clustering were discarded from further analyses. To estimate
298 the strength and direction of selection on the genomes analysed, we calculated the ratio of
299 synonymous and nonsynonymous substitutions (dN/dS) within each genus cluster using two sets
300 of conserved marker genes, checkm_bact and checkm_arch for bacteria and archaea, respectively
301 [56]. Genomes used to calculate the dN/dS for each genus cluster are reported in Supplemental
302 File 1. The open reading frames (ORFs) retrieved from the GTDB were compared to the HMMs
303 of the checkm_bact (120 marker genes) and checkm_arch marker (122 marker genes) sets using
304 the hmmsearch tool available in HMMER v. 3.2.1 with the reported model-specific cutoffs [57].
305 We aligned the amino acid sequences for each marker gene and each cluster individually using
306 ClustalOmega [58] and then converted amino acid alignments into codon alignments using
307 PAL2NAL with the parameter --nogap [59]. We used the resulting codon alignments to estimate
308 the ratio of synonymous and nonsynonymous substitutions for each pair of genomes using the
309 maximum likelihood approximation (codeML) available on PAML 4.9h [60]. In order to avoid
310 bias associated with divergence, dN/dS estimates with $dS > 1.5$ were removed due to potential
311 saturation. We also discarded pairwise comparisons with $dS < 0.1$ because these might represent

312 dN/dS values calculated from genomes of the same population. Moreover, dN/dS values >10 were
313 considered artifactual [39]. Genomes with fewer than 25 marker genes remaining after filtering
314 were discarded. After dN/dS estimation, we randomly selected one representative genome for each
315 genus for further analyses (GTDB genomes dataset). We predicted ribosomal RNA genes in our
316 selected genomes using Barrnap (barrnap 0.9: rapid ribosomal RNA prediction;
317 <https://github.com/tseemann/barrnap>), with the default parameters. Genome size, 16S rRNA
318 copies, and dNdS values for the GTDB representative genomes dataset are reported in
319 Supplemental File S2.

320

321 **Statistical analyses**

322 Due to the tendency of related species to resemble each other because of their shared phylogenetic
323 ancestry, we assessed the suitability of a phylogeny-based method for our regression analyses by
324 first estimating Blomberg's K on genome size data [29] using the phylosignal function on R [61].
325 This parameter represents the phylogenetic signal in a continuous trait, and goes from 0 (no
326 phylogenetic signal) to ∞ (phylogenetic signal) with the null hypothesis (K=1) meaning that the
327 trait analysed evolves under Brownian Motion (BM, variation of the trait is proportional to the
328 distance between species [30,62]. In addition, we also tested the fit of different trait evolution
329 models, including including Brownian Motion [30], Ornstein-Uhlenbeck [31], Early-Burst [32], a
330 diffusion model, Pagel's model [33], a drift model, and a white-noise model (non-phylogenetic).
331 We also performed a Generalized Least Square analysis to explore the predictability of genome
332 size using dN/dS and 16S rRNA copies as predictor variables using the "glm" function available
333 on R. Since we detected phylogenetic signal in genome size data, we additionally accounted for
334 potential phylogenetic nonindependence in the residuals using the PGLS method with the function
335 pglsl on the R package Caper [63] and the Pagel's model [33]. We calculated the lambda (λ)
336 parameter (which showed phylogenetic signal in the residuals), delta (δ) and kappa (κ) (pattern of
337 evolution of trait) through maximum likelihood. The best fitting model according to AIC and
338 likelihood was checked visually using diagnostic plots (residuals vs. fitted values, and QQ plots to
339 check normality) (Fig. S3).

340

341

342

343 **Phylogenetic reconstruction**

344 To perform a Phylogenetic Generalized Least Square analysis (PGLS), we reconstructed a
345 phylogenetic tree using the GTDB genomes dataset described above. We used the MarkerFinder
346 pipeline reported previously [22], consisting in the identification of 27 ribosomal proteins and
347 three RNA polymerase genes [64] using HMMER3 and the resulting individual sequences aligned
348 with ClustalOmega and concatenated. In addition, the concatenated alignment was trimmed with
349 trimAl [65] using the option -gt 0.1. The Ribosomal-RNAP alignment was then used to build the
350 phylogenetic tree with IQ-TREE 1.6.12 [66] with the substitutions model LG+R10 and the options
351 -wbt, -bb 1000, and --runs 10 [67–69]. The resulting phylogeny was manually inspected on iTOL
352 [70] (Fig. 1).

353

354 **Acknowledgments**

355 We acknowledge the use of the Virginia Tech Advanced Research Computing Center for
356 bioinformatic analyses performed in this study. This investigation was supported by grants from
357 the Institute for Critical Technology and Applied Science and the National Science Foundation
358 (IIBR-1918271), and a Simons Early Career Award in Marine Microbial Ecology and Evolution
359 to F.O.A. We kindly thank members of the Aylward Lab for their insightful comments on an earlier
360 version of this manuscript and Prof. Josef Uyeda for advice on phylogeny-based statistical
361 methods.

362

363 **Figure legends**

364 **Figure 1.** Genome size distribution across the Tree of Life of bacteria and archaea. Phylogenetic
365 tree was built using a concatenated alignment of ribosomal and RNA polymerase sequences
366 through a maximum likelihood approach and the substitution model LG+R10. Abbreviations:
367 TACK = Thaumarchaeota, Aigarchaeota, Crenarchaeota and Korarchaeota; TDS = Thermotogota,
368 Deinococcota, and Synergistota; AMND = Acidobacteriota, Methyloirabilota, Nitrospirota,
369 Deferribacterota. Raw data for genome size can be found in Supplemental File S2.

370

371 **Figure 2.** Distribution of genome size within bacteria and archaea taxonomic groups. Genome size
372 grouping based on phylum. First, third quantile, and median are shown for each phylum
373 distribution. Abbreviations: TDS = Thermotogota, Deinococcota, and Synergistota. Raw data for
374 genome size can be found in Supplemental File S2.

375

376 **Figure 3.** Relationship between genome size and dN/dS. dN/dS values were log transformed. Dots
377 size is equivalent to the number of 16S rRNA gene copies. Raw data can be found in Supplemental
378 File S2.

379
380 **Figure 4.** Relationship between genome size and genomic traits for bacteria and archaea. A)
381 Regression line of the relationship between genome size and dN/dS ratio before (dashed line) and
382 after (solid line) taking phylogenetic relationships into account. B) Regression line of the
383 relationship between genome size and 16S rRNA copies before (dashed line) and after (solid line)
384 taking phylogenetic relationships into account. Parameters of the regression equation for both
385 relationships can be found in Table 2. Raw data can be found in Supplemental File S2.

386
387 **Table 1.** Summary of model fitting for genome size data. We highlighted the model that showed
388 the highest likelihood and the lowest corrected AIC.

389
390 **Table 2.** Statistics of the regression models relating genome size and dN/dS and 16S rRNA as
391 predictor variables using Generalized Least Square and Phylogenetic Least Square analyses. We
392 highlighted the models that were statistically significant ($\alpha = 0.05$).

393
394 **Supplementary File 1.** Genomes used to calculate pairwise dN/dS within each genus cluster.

395
396 **Supplementary File 2.** dN/dS, genome size, and *rrn* operon copies for each genus representative.

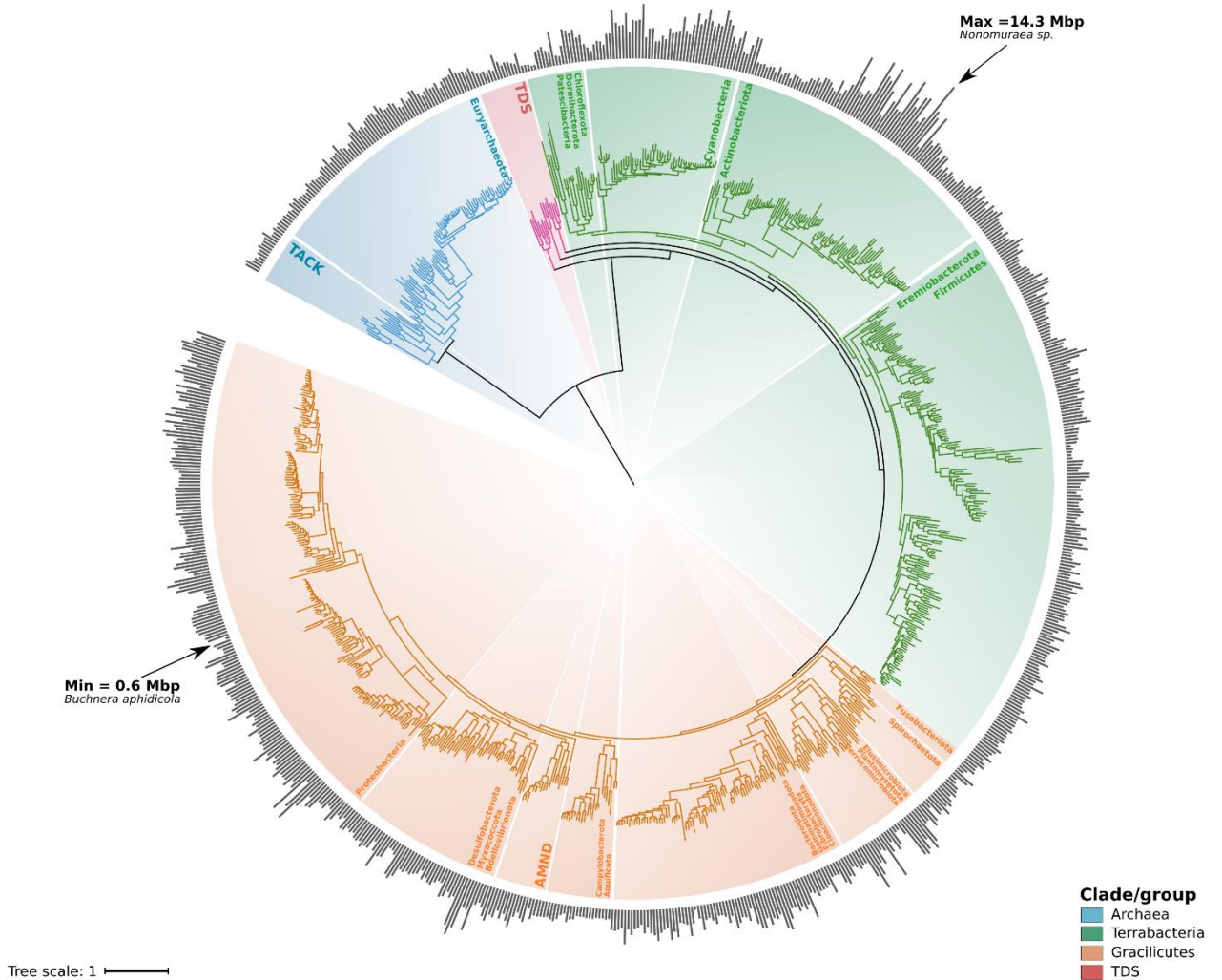
397
398 **Supplementary Figure 3.** Diagnostic plots for the PGLS model genome size ~ 16S rRNA copies.

399
400
401
402
403
404
405
406
407
408
409
410
411
412
413

414 **Figures and Tables**

415 **Figure 1**

416



417

418

419

420

421

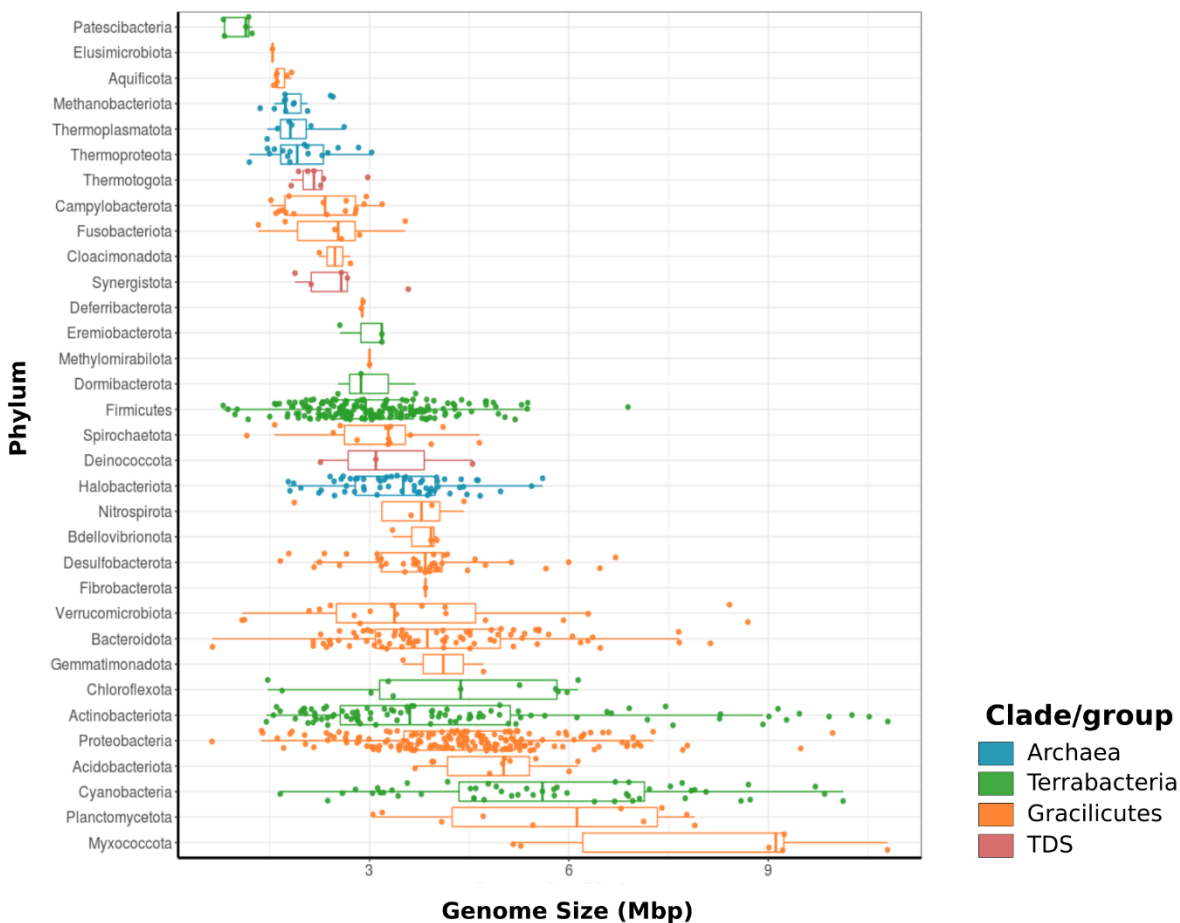
422

423

424

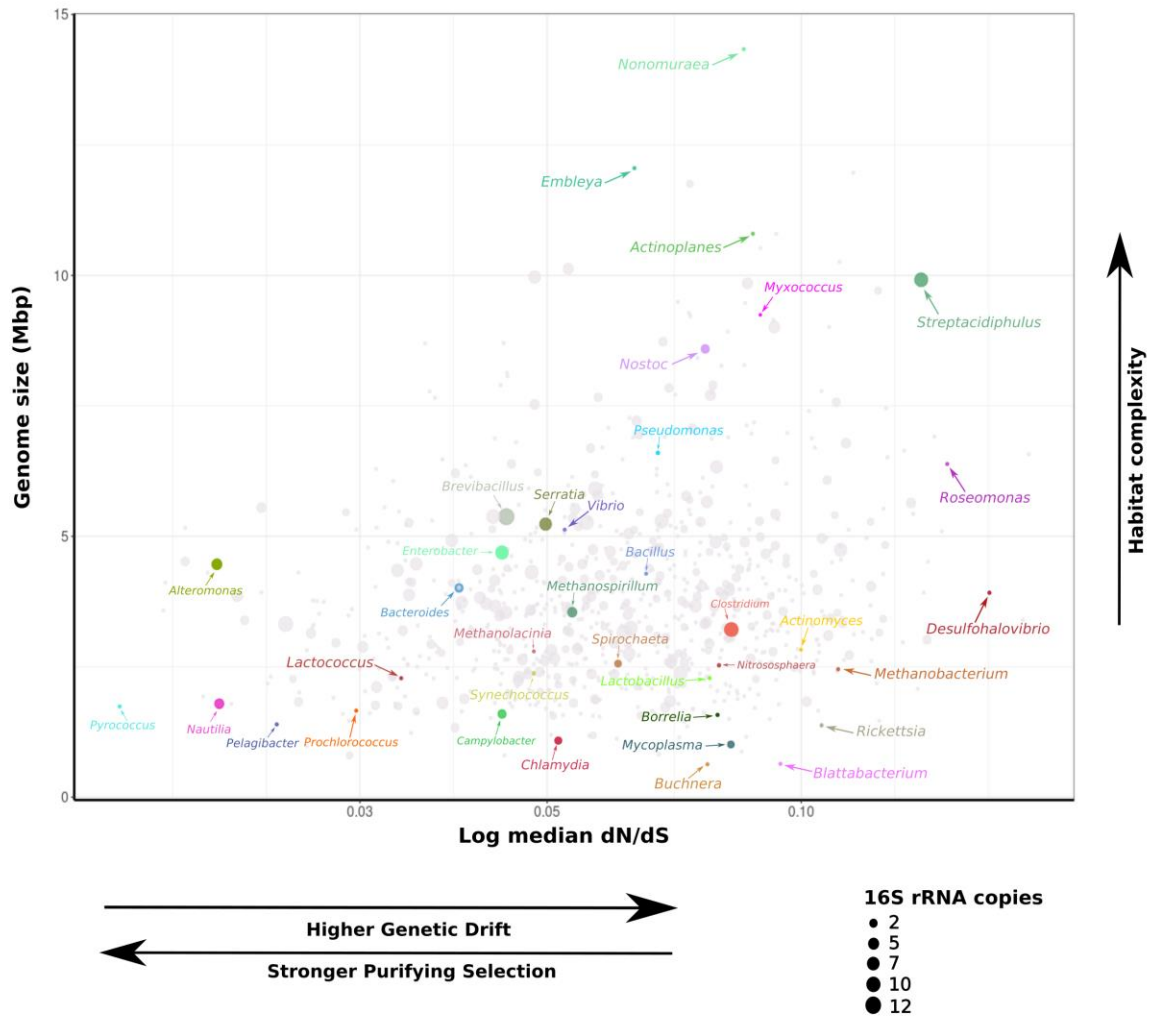
425

426 **Figure 2**



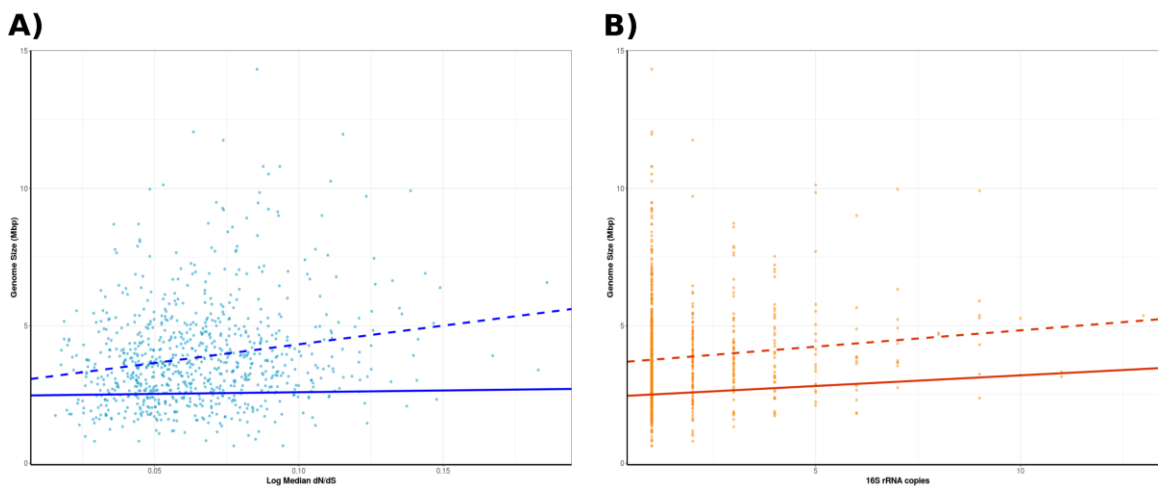
427
428
429
430
431
432
433
434
435
436
437
438
439
440

441 **Figure 3**



442

443 **Figure 4**



444

445 **Table 1**

Model	Loglik	AICc
Brownian motion	-1463.3	2930.7
Ornstein-Uhlenbeck	-1420.7	2847.4
Early-Burst	-1463.3	2932.7
Pagel's lambda*	-1415.6	2837.3
Trend diffusion	-1447.7	2901.4
Drift	-1463.3	2932.7
White-noise	-1695.6	3395.2

446 *Significantly higher likelihood when compared with the rest of the models tested according to the chisq test

447

448

449 **Table 2**

450

Model	Predictor variable	Kappa (95% CI)	Lambda (95% CI)	Delta (95% CI)	Slope	Intercept	P-val	AIC	R2*
Generalized Least Square									
Genome Size ~ Median dN/dS	dN/dS	-	-	-	13.57	2.97	<0.001	3366.2	0.04
Genome Size ~ 16S rRNA copies	16S rRNA copies	-	-	-	0.12	3.65	0.002	3387.7	0.01
Genome Size ~ Median dN/dS + 16S rRNA copies	dN/dS + 16S rRNA copies	-	-	-	14.11/0.13	2.7	<0.001	3355.9	0.05
Phylogenetic Generalized Least Square									

Genome Size ~ Median dN/dS	dN/dS	0.48 (0.39-0.58)	0.98 (0.96-0.99)	2.44 (2.01-2.85)	1.26	2.46	0.5	2748.8	0.0006
Genome Size ~ 16S rRNA copies	16S rRNA copies	0.49 (0.34-0.59)	0.98 (0.96-0.99)	2.49 (2.06-2.9)	0.08	2.42	0.003	2740.268	0.01
Genome Size ~ Median dN/dS + 16S rRNA copies**	dN/dS + 16S rRNA copies	0.49 (0.40-0.59)	0.98 (0.96-0.99)	2.51 (2.08-2.93)	1.29/0.08	2.35	0.009	2741.79	0.01

451 *R2 calculated based on residuals, likelihood, and predicted data (multiple R-squared)

452 **Anova did not show significant differences between models Genome Size ~ 16S rRNA copies and Genome Size ~
453 Median dN/dS + 16S rRNA copies (P=0.48)

454

455

456

457 **References**

458 1. Mira A, Ochman H, Moran NA. Deletional bias and the evolution of bacterial genomes.
459 Trends Genet. 2001;17: 589–596.

460 2. Lynch M. Streamlining and simplification of microbial genome architecture. Annu Rev
461 Microbiol. 2006;60: 327–349.

462 3. Koonin EV. Evolution of genome architecture. Int J Biochem Cell Biol. 2009;41: 298–306.

463 4. Lawrence JG, Hendrix RW, Casjens S. Where are the pseudogenes in bacterial genomes?
464 Trends Microbiol. 2001;9: 535–540.

465 5. Wolf YI, Aravind L, Grishin NV, Koonin EV. Evolution of aminoacyl-tRNA synthetases--
466 analysis of unique domain architectures and phylogenetic trees reveals a complex history of
467 horizontal gene transfer events. Genome Res. 1999;9: 689–710.

468 6. Bobay L-M, Ochman H. The Evolution of Bacterial Genome Architecture. Front Genet.
469 2017;8: 72.

470 7. van Ham RCHJ, Kamerbeek J, Palacios C, Rausell C, Abascal F, Bastolla U, et al.
471 Reductive genome evolution in *Buchnera aphidicola*. Proc Natl Acad Sci U S A. 2003;100:
472 581–586.

- 473 8. Moran NA, Mira A. The process of genome shrinkage in the obligate symbiont *Buchnera*
474 *aphidicola*. *Genome Biol.* 2001;2: RESEARCH0054.
- 475 9. Chong RA, Park H, Moran NA. Genome Evolution of the Obligate Endosymbiont *Buchnera*
476 *aphidicola*. *Mol Biol Evol.* 2019;36: 1481–1489.
- 477 10. Batut B, Knibbe C, Marais G, Daubin V. Reductive genome evolution at both ends of the
478 bacterial population size spectrum. *Nat Rev Microbiol.* 2014;12: 841–850.
- 479 11. Woolfit M, Bromham L. Increased rates of sequence evolution in endosymbiotic bacteria
480 and fungi with small effective population sizes. *Mol Biol Evol.* 2003;20: 1545–1555.
- 481 12. Biller SJ, Berube PM, Lindell D, Chisholm SW. *Prochlorococcus*: the structure and function
482 of collective diversity. *Nat Rev Microbiol.* 2015;13: 13–27.
- 483 13. Kashtan N, Roggensack SE, Rodrigue S, Thompson JW, Biller SJ, Coe A, et al. Single-cell
484 genomics reveals hundreds of coexisting subpopulations in wild *Prochlorococcus*. *Science.*
485 2014;344: 416–420.
- 486 14. Grote J, Thrash JC, Huggett MJ, Landry ZC, Carini P, Giovannoni SJ, et al. Streamlining
487 and core genome conservation among highly divergent members of the SAR11 clade.
488 *MBio.* 2012;3. doi:10.1128/mBio.00252-12
- 489 15. Giovannoni SJ, Cameron Thrash J, Temperton B. Implications of streamlining theory for
490 microbial ecology. *ISME J.* 2014;8: 1553–1565.
- 491 16. Simonsen AK. Environmental stress leads to genome streamlining in a widely distributed
492 species of soil bacteria. *The ISME Journal.* 2021. doi:10.1038/s41396-021-01082-x
- 493 17. Bourguignon T, Kinjo Y, Villa-Martín P, Coleman NV, Tang Q, Arab DA, et al. Increased
494 Mutation Rate Is Linked to Genome Reduction in Prokaryotes. *Curr Biol.* 2020;30: 3848–
495 3855.e4.
- 496 18. Marais GAB, Calteau A, Tenaillon O. Mutation rate and genome reduction in
497 endosymbiotic and free-living bacteria. *Genetica.* 2008;134: 205–210.
- 498 19. Marais GAB, Batut B, Daubin V. Genome Evolution: Mutation Is the Main Driver of
499 Genome Size in Prokaryotes. *Curr Biol.* 2020;30: R1083–R1085.
- 500 20. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. A complete
501 domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol.* 2020;38: 1079–
502 1086.
- 503 21. Sela I, Wolf YI, Koonin EV. Theory of prokaryotic genome evolution. *Proc Natl Acad Sci*
504 *U S A.* 2016;113: 11399–11407.
- 505 22. Martinez-Gutierrez CA, Aylward FO. Phylogenetic Signal, Congruence, and Uncertainty
506 across Bacteria and Archaea. *Mol Biol Evol.* 2021. doi:10.1093/molbev/msab254

- 507 23. Konstantinidis KT, Tiedje JM. Trends between gene content and genome size in prokaryotic
508 species with larger genomes. *Proc Natl Acad Sci U S A*. 2004;101: 3160–3165.
- 509 24. Guieysse B, Wuertz S. Metabolically versatile large-genome prokaryotes. *Curr Opin*
510 *Biotechnol*. 2012;23: 467–473.
- 511 25. Klappenbach JA, Dunbar JM, Schmidt TM. rRNA operon copy number reflects ecological
512 strategies of bacteria. *Appl Environ Microbiol*. 2000;66: 1328–1333.
- 513 26. Kuo C-H, Moran NA, Ochman H. The consequences of genetic drift for bacterial genome
514 complexity. *Genome Res*. 2009;19: 1450–1454.
- 515 27. Freckleton, Freckleton, Harvey, Pagel. Phylogenetic Analysis and Comparative Data: A
516 Test and Review of Evidence. *The American Naturalist*. 2002. p. 712. doi:10.2307/3078855
- 517 28. Felsenstein J. Phylogenies and the Comparative Method. *The American Naturalist*. 1985.
518 pp. 1–15. doi:10.1086/284325
- 519 29. Blomberg SP, Garland T Jr, Ives AR. Testing for phylogenetic signal in comparative data:
520 behavioral traits are more labile. *Evolution*. 2003;57: 717–745.
- 521 30. Felsenstein J. Maximum-likelihood estimation of evolutionary trees from continuous
522 characters. *Am J Hum Genet*. 1973;25: 471–492.
- 523 31. Butler MA, King AA. Phylogenetic Comparative Analysis: A Modeling Approach for
524 Adaptive Evolution. *Am Nat*. 2004;164: 683–695.
- 525 32. Harmon LJ, Losos JB, Jonathan Davies T, Gillespie RG, Gittleman JL, Bryan Jennings W,
526 et al. Early bursts of body size and shape evolution are rare in comparative data. *Evolution*.
527 2010;64: 2385–2396.
- 528 33. Pagel M. Inferring the historical patterns of biological evolution. *Nature*. 1999. pp. 877–
529 884. doi:10.1038/44766
- 530 34. Mundry R. Statistical Issues and Assumptions of Phylogenetic Generalized Least Squares.
531 *Modern Phylogenetic Comparative Methods and Their Application in Evolutionary*
532 *Biology*. 2014. pp. 131–153. doi:10.1007/978-3-662-43550-2_6
- 533 35. Revell LJ. Phylogenetic signal and linear regression on species data. *Methods in Ecology*
534 *and Evolution*. 2010. pp. 319–329. doi:10.1111/j.2041-210x.2010.00044.x
- 535 36. Whitney KD, Garland T Jr. Did genetic drift drive increases in genome complexity? *PLoS*
536 *Genet*. 2010;6. doi:10.1371/journal.pgen.1001080
- 537 37. Whitney KD, Boussau B, Baack EJ, Garland T Jr. Drift and genome complexity revisited.
538 *PLoS Genet*. 2011;7: e1002092.
- 539 38. Hernández CE, Rodríguez-Serrano E, Avaria-Llautureo J, Inostroza-Michael O, Morales-

- 540 Pallero B, Boric-Bargetto D, et al. Using phylogenetic information and the comparative
541 method to evaluate hypotheses in macroecology. *Methods in Ecology and Evolution*. 2013.
542 pp. 401–415. doi:10.1111/2041-210x.12033
- 543 39. Martinez-Gutierrez CA, Aylward FO. Strong Purifying Selection Is Associated with
544 Genome Streamlining in Epipelagic Marinimicrobia. *Genome Biol Evol*. 2019;11: 2887–
545 2894.
- 546 40. Rocha EPC, Smith JM, Hurst LD, Holden MTG, Cooper JE, Smith NH, et al. Comparisons
547 of dN/dS are time dependent for closely related bacterial genomes. *J Theor Biol*. 2006;239:
548 226–235.
- 549 41. Luo H, Huang Y, Stepanauskas R, Tang J. Excess of non-conservative amino acid changes
550 in marine bacterioplankton lineages with reduced genomes. *Nat Microbiol*. 2017;2: 17091.
- 551 42. Niederdorfer R, Besemer K, Battin TJ, Peter H. Ecological strategies and metabolic trade-
552 offs of complex environmental biofilms. *NPJ Biofilms Microbiomes*. 2017;3: 21.
- 553 43. Condon C, Liveris D, Squires C, Schwartz I, Squires CL. rRNA operon multiplicity in
554 *Escherichia coli* and the physiological implications of *rrn* inactivation. *J Bacteriol*.
555 1995;177: 4152–4156.
- 556 44. Lynch M, Conery JS. The Origins of Genome Complexity. *Science*. 2003. pp. 1401–1404.
557 doi:10.1126/science.1089370
- 558 45. Tamas I, Klasson L, Canbäck B, Näslund AK, Eriksson A-S, Wernegreen JJ, et al. 50
559 million years of genomic stasis in endosymbiotic bacteria. *Science*. 2002;296: 2376–2379.
- 560 46. McCutcheon JP, Moran NA. Extreme genome reduction in symbiotic bacteria. *Nat Rev*
561 *Microbiol*. 2011;10: 13–26.
- 562 47. López-Pérez M, Haro-Moreno JM, Coutinho FH, Martinez-Garcia M, Rodriguez-Valera F.
563 The Evolutionary Success of the Marine Bacterium SAR11 Analyzed through a
564 Metagenomic Perspective. *mSystems*. 2020;5. doi:10.1128/mSystems.00605-20
- 565 48. Giovannoni SJ. SAR11 Bacteria: The Most Abundant Plankton in the Oceans. *Ann Rev*
566 *Mar Sci*. 2017;9: 231–255.
- 567 49. Partensky F, Garczarek L. *Prochlorococcus*: advantages and limits of minimalism. *Ann Rev*
568 *Mar Sci*. 2010;2: 305–331.
- 569 50. Cottrell MT, Kirchman DL. Transcriptional Control in Marine Copiotrophic and
570 Oligotrophic Bacteria with Streamlined Genomes. *Appl Environ Microbiol*. 2016;82: 6010–
571 6018.
- 572 51. Carini P, Steindler L, Beszteri S, Giovannoni SJ. Nutrient requirements for growth of the
573 extreme oligotroph “*Candidatus Pelagibacter ubique*” HTCC1062 on a defined medium.
574 *The ISME Journal*. 2013. pp. 592–602. doi:10.1038/ismej.2012.122

- 575 52. Fierer N. Embracing the unknown: disentangling the complexities of the soil microbiome.
576 Nat Rev Microbiol. 2017;15: 579–590.
- 577 53. Takeuchi N, Cordero OX, Koonin EV, Kaneko K. Gene-specific selective sweeps in
578 bacteria and archaea caused by negative frequency-dependent selection. BMC Biol.
579 2015;13: 20.
- 580 54. Goldman B, Bhat S, Shimkets LJ. Genome evolution and the emergence of fruiting body
581 development in *Myxococcus xanthus*. PLoS One. 2007;2: e1329.
- 582 55. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify
583 genomes with the Genome Taxonomy Database. Bioinformatics. 2019.
584 doi:10.1093/bioinformatics/btz848
- 585 56. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the
586 quality of microbial genomes recovered from isolates, single cells, and metagenomes.
587 Genome Res. 2015;25: 1043–1055.
- 588 57. Eddy SR. Accelerated Profile HMM Searches. PLoS Comput Biol. 2011;7: e1002195.
- 589 58. Sievers F, Higgins DG. Clustal Omega for making accurate alignments of many protein
590 sequences. Protein Science. 2018. pp. 135–145. doi:10.1002/pro.3290
- 591 59. Suyama M, Torrents D, Bork P. PAL2NAL: robust conversion of protein sequence
592 alignments into the corresponding codon alignments. Nucleic Acids Research. 2006. pp.
593 W609–W612. doi:10.1093/nar/gkl315
- 594 60. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol. 2007;24:
595 1586–1591.
- 596 61. Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD, et al.
597 Picante: R tools for integrating phylogenies and ecology. Bioinformatics. 2010;26: 1463–
598 1464.
- 599 62. O’Meara BC, Ané C, Sanderson MJ, Wainwright PC. Testing for different rates of
600 continuous trait evolution using likelihood. Evolution. 2006;60: 922–933.
- 601 63. Website. Available: Orme D, Freckleton R, Thomas G, Petzold T, Fritz S, Isaac N, Pears
602 W. 2012. Caper: Comparative Analyses of Phylogenetics and Evolution in R. Version 0.5.
603 [WWW document] URL <http://cran.r-project.org/web/packages/caper/caper.pdf>.
- 604 64. Sunagawa S, Mende DR, Zeller G, Izquierdo-Carrasco F, Berger SA, Kultima JR, et al.
605 Metagenomic species profiling using universal phylogenetic marker genes. Nat Methods.
606 2013;10: 1196–1199.
- 607 65. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated
608 alignment trimming in large-scale phylogenetic analyses. Bioinformatics. 2009;25: 1972–
609 1973.

- 610 66. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
611 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.*
612 2015;32: 268–274.
- 613 67. Minh BQ, Nguyen MAT, von Haeseler A. Ultrafast approximation for phylogenetic
614 bootstrap. *Mol Biol Evol.* 2013;30: 1188–1195.
- 615 68. Le SQ, Gascuel O. An improved general amino acid replacement matrix. *Mol Biol Evol.*
616 2008;25: 1307–1320.
- 617 69. Soubrier J, Steel M, Lee MSY, Der Sarkissian C, Guindon S, Ho SYW, et al. The influence
618 of rate heterogeneity among sites on the time dependence of molecular rates. *Mol Biol Evol.*
619 2012;29: 3345–3358.
- 620 70. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new
621 developments. *Nucleic Acids Res.* 2019;47: W256–W259.
- 622