

# 1 Pseudogenes as a neutral reference for detecting selection in 2 prokaryotic pangenomes

3  
4 Gavin M. Douglas<sup>a,b,\*</sup>, W. Ford Doolittle<sup>c</sup>, B. Jesse Shapiro<sup>a,b,d,\*</sup>

5  
6 <sup>a</sup>*Department of Microbiology and Immunology, McGill University, Montréal, QC, Canada*

7 <sup>b</sup>*McGill Genome Centre, McGill University, Montréal, QC, Canada*

8 <sup>c</sup>*Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax, NS, Canada*

9 <sup>d</sup>*McGill Centre for Microbiome Research, McGill University, Montréal, QC, Canada*

10  
11 \*Emails for correspondence: [gavin.douglas@mcgill.ca](mailto:gavin.douglas@mcgill.ca) and [jesse.shapiro@mcgill.ca](mailto:jesse.shapiro@mcgill.ca)

12  
13 Keywords: Pseudogenes, pangenome, evolution, horizontal gene transfer, lateral gene transfer, mobile genes,  
14 mobilome, adaptation

## 16 **Abstract**

17 A long-standing question is to what degree genetic drift vs. selection drives the divergence in  
18 rare accessory gene content between closely related bacteria. Rare genes, including singletons,  
19 make up a large proportion of pangenomes (the set of all genes in a set of genomes), but it  
20 remains unclear how many such genes are adaptive, deleterious, or neutral to their host genome.  
21 Estimates of species' effective population sizes ( $N_e$ ) are positively associated with pangenome  
22 size and fluidity, which has independently been interpreted as evidence for both neutral and  
23 adaptive pangenome models. We hypothesised that these models could be distinguished if  
24 measures of pangenome diversity were normalized by pseudogene diversity as a proxy for  
25 neutral genic diversity. To this end, we defined the ratio of singleton intact genes to singleton  
26 pseudogenes ( $s_i/s_p$ ) within a pangenome, which shows a signal across prokaryotic species  
27 consistent with the relative adaptive value of many rare accessory genes. We also identified  
28 differences in functional annotations between intact genes and pseudogenes. For instance,  
29 transposons are highly enriched among pseudogenes, while most other functional categories are  
30 more often intact. Our work demonstrates that including pseudogenes as a neutral reference leads  
31 to improved inferences of the evolutionary forces driving pangenome variation.

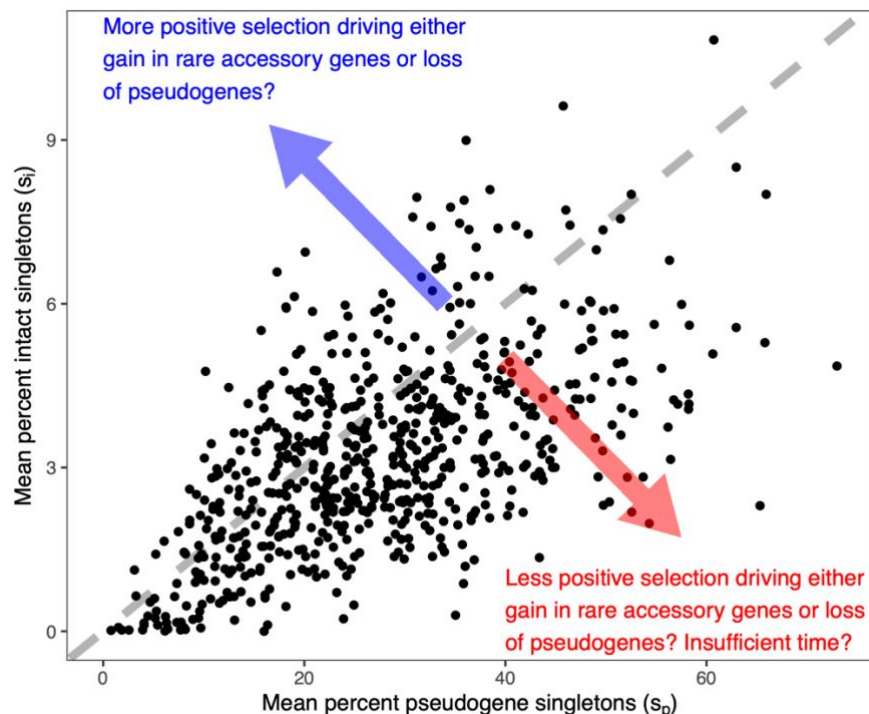
## 32 **Main**

33 A long-standing question is to what degree genetic drift vs. selection drives the divergence in  
34 rare accessory gene content between closely related bacteria<sup>1-4</sup>. Rare genes, including singletons,  
35 make up a large proportion of pangenomes (the set of all genes in a set of genomes), but it  
36 remains unclear how many such genes are adaptive, deleterious, or neutral to their host genome.  
37 Estimates of species' effective population sizes ( $N_e$ ) are positively associated with pangenome  
38 size and fluidity<sup>5-7</sup>, which has independently been interpreted as evidence for both neutral<sup>6</sup> and  
39 adaptive<sup>5,7</sup> pangenome models. We hypothesised that these models could be distinguished if  
40 measures of pangenome diversity were normalized by pseudogene diversity as a proxy for  
41 neutral genic diversity. To this end, we defined the ratio of singleton intact genes to singleton  
42 pseudogenes ( $s_i/s_p$ ) within a pangenome, which shows a signal across prokaryotic species  
43 consistent with the relative adaptive value of many rare accessory genes. We also identified  
44 differences in functional annotations between intact genes and pseudogenes. For instance,  
45 transposons are highly enriched among pseudogenes, while most other functional categories are  
46 more often intact. Our work demonstrates that including pseudogenes as a neutral reference leads  
47 to improved inferences of the evolutionary forces driving pangenome variation. (Please note that  
48 the first paragraph was duplicated as the above abstract to allow bioRxiv to correctly parse it: the  
49 original version of this preprint did not have a separate abstract).

50 These evolutionary forces have been investigated through several approaches, such as  
51 analysing gene frequency distributions,<sup>8,9</sup> gene co-occurrence<sup>10</sup>, and patterns of nucleotide  
52 variation within transferred genes<sup>11</sup>. This work has primarily provided insight into the higher  
53 frequency accessory genes, rather than rare genes that make up the largest fraction of  
54 pangenomes<sup>12</sup>. These rare genes (often singletons sequenced in just one genome) are frequently  
55 mobile genes with high turnover rates and dubious adaptive value to their bacterial hosts.  
56 Nonetheless, rare genes have been hypothesised to provide adaptative benefits in rare ecological  
57 niches<sup>13-15</sup>, although this hypothesis remains largely untested. Here, we propose a new metric of  
58 selection on rare accessory genes, which we apply to a dataset of >600 prokaryotic species. We  
59 then analyse a subset of well-sampled bacterial species to identify functional categories that are  
60 enriched in intact genes compared to pseudogenes. Our results provide strong evidence that an  
61 identifiable subset of rare accessory genes likely provide adaptive value to their hosts.

62 Pseudogenization – gene degeneration through the introduction of mutations, such as  
63 premature stop codons, insertions, and deletions – can occur when genetic drift overcomes  
64 purifying selection to retain a gene<sup>16</sup>, or through positive selection to eliminate a deleterious  
65 gene<sup>17</sup>. We reasoned that rare accessory gene families that tend to remain intact are under  
66 stronger positive selection than those that tend to be pseudogenized. We expressed this by  
67 calculating the mean percentages of intact singleton genes ( $s_i$ ) and pseudogenes ( $s_p$ ) within a  
68 species' pangenome. We analysed 668 named prokaryotic species represented by at least nine  
69 genomes in the Genome Taxonomy Database<sup>18</sup> and found that the mean values of  $s_i$  and  $s_p$  were  
70 correlated (Spearman's  $\rho=0.57$ ;  $P < 0.001$ ), with deviations suggesting species-specific  
71 differences in selection on rare accessory genes (**Figure 1**). For example, *Escherichia coli* has a  
72 high  $s_i/s_p$  ratio, consistent with selection to retain rare accessory genes, while the obligate  
73 intracellular bacteria *Chlamydia trachomatis* and *Rickettsia prowazekii* exhibit the lowest ratios,  
74 which could indicate less selective constraint on their rare genes.

75 The analysed species span substantial prokaryotic diversity (**Extended Data Table 1**) but  
76 were biased towards Gammaproteobacteria (286 species) and Bacilli (161 species). We identified  
77 pseudogenes with Pseudofinder<sup>19</sup>, which identifies several classes of potential pseudogenes. We  
78 focused on intergenic pseudogenes, which represent significant matches to database sequences  
79 outside of gene calls, as this class is more likely to represent degenerating gene sequences  
80 compared to the other candidate pseudogene classes (see Online Methods). We filtered out  
81 pseudogenes based on several criteria, including restricting analysed pseudogenes to those  $\geq$   
82 100 bp and  $\leq$  5000 bp in length. Based on all criteria, a mean of 11.90% (standard deviation  
83 [SD]: 5.78%) of pseudogenes were excluded per species. After this filtering, intergenic  
84 pseudogenes represented a mean of 4.52% (SD: 2.96%) of called elements per genome (range:  
85 0.30-19.81%). These elements comprised an even smaller portion of overall genome size (mean:  
86 1.42%, SD: 0.99%) compared to intact genes (mean: 87.34%, SD: 2.77%) because pseudogenes  
87 are generally shorter than intact genes.  
88



89  
90 **Figure 1:** Mean percentage of intact genes and pseudogenes that are singletons (i.e. genome-specific)  
91 per species. Each point represents one of 668 prokaryotic species ( $\geq$  nine genomes each). The mean  
92 percent singletons (for both intact genes and pseudogenes) per species was based on repeated  
93 subsampling to nine genomes (for up to 100 replicates). Possible (but non-exhaustive) drivers of higher or  
94 lower  $s_i/s_p$  ratios are indicated alongside coloured arrows.  
95

96 Species' pangenome size and complexity have been characterised based on different  
97 metrics, including the mean number of genes per genome<sup>5</sup> and genomic fluidity<sup>6,20</sup>. We  
98 computed these metrics for all species based on both intact genes and pseudogenes. As we were  
99 especially interested in rare elements, we computed the mean numbers and percentages of  
100 singleton genes and pseudogenes per species (i.e. those present in a single genome per species),  
101 based on repeated subsampling to nine genomes. Larger genomes tend to encode more  
102 singletons, both in mean number and percentage (**Extended Data Fig. 1a,b**). In addition, the  
103 percentage of intact singletons ( $s_i$ ) is highly correlated with genomic fluidity, but the traditional  
104 fluidity metric is sensitive to intermediate frequency accessory genes (**Extended Data Fig.**  
105 **1c,d**). We therefore focused on the percentage of intact ( $s_i$ ) and pseudogene ( $s_p$ ) singletons for  
106 most analyses. All metrics ranged substantially across species for both intact genes (fluidity:  
107 0.003-0.246; mean number: 836.4-8692.7; mean number of singletons: 0.00-581.29;  $s_i$ : 0.00-  
108 10.83%) and pseudogenes (fluidity: 0.014-0.513; mean number: 8.1-922.5; mean number of  
109 singletons: 0.78-325.17;  $s_p$ : 0.78-72.97%). These results highlight that, as expected<sup>21</sup>,  
110 pseudogenes are frequently genome-specific.

111 We next recapitulated the previously observed association between genome-wide non-  
112 synonymous to synonymous substitution rates (dN/dS) and pangenome diversity<sup>5,7</sup>, and then  
113 explored whether dN/dS is also associated with  $s_i/s_p$ . We computed dN/dS across the core  
114 genome of each species, based on the mean values of all pairwise strain comparisons. This  
115 metric is often taken as a proxy for selection efficacy: lower dN/dS values indicate increased  
116 efficacy of purifying selection (which is associated with higher  $N_e$ ) against non-synonymous  
117 changes, which tend to be deleterious. However, within-species dN/dS values are highly  
118 dependent on strain divergence times, with recent divergences enriched in higher dN/dS due to  
119 insufficient time for purifying selection to purge deleterious non-synonymous mutations<sup>22,23</sup>.  
120 Using within-species dS as a proxy for the molecular clock, we also observed a time-dependence  
121 of dN/dS in our data (**Extended Data Fig. 2a**).

122 Due to this relationship between within-species dN/dS and dS, we included dS as a  
123 covariate when computing partial Spearman correlations between measures of pangenome  
124 diversity and dN/dS. Based on this approach, the mean number of genes, genomic fluidity, and  $s_i$   
125 were all significantly negatively correlated with dN/dS across species (**Figure 2 a-c**; Partial  
126 Spearman correlations,  $P < 0.05$ ). This observation agrees with past work<sup>5,7</sup>, which has been  
127 taken as evidence for an adaptive pangenome model. However,  $N_e$ , which determines selection  
128 efficacy and thus core genome dN/dS (assuming equal selection pressure), is also associated with  
129 higher standing levels of neutral variation, due to less variation being lost through genetic drift in  
130 larger populations. Accordingly, a positive association between pangenome diversity and  $N_e$  can  
131 be explained by both an adaptive or neutral model.

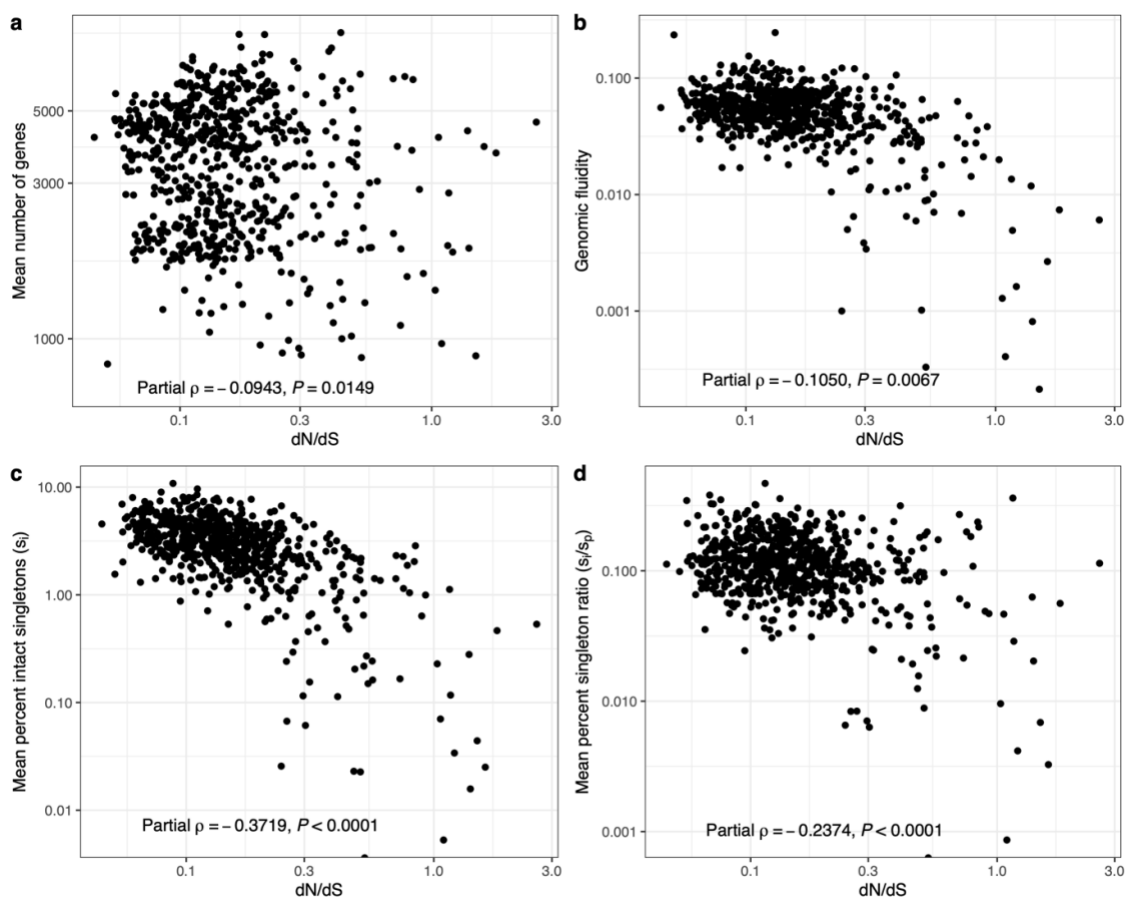
132 To disentangle these models, we explored whether our new metric,  $s_i/s_p$ , is differently  
133 associated with dS and dN/dS compared to the unnormalized measures of pangenome diversity.  
134 Based on a partial Spearman correlation, we found  $s_i/s_p$  to be significantly associated with dN/dS  
135 (partial Spearman's  $\rho=0.237$ ;  $P < 0.001$ ; **Figure 2d**), although less so than  $s_i$  alone (partial  
136 Spearman's  $\rho=0.372$ ;  $P < 0.001$ ). This result highlights that  $s_i$  remains associated with dN/dS  
137 even after normalization by  $s_p$ . If pseudogene diversity is assumed to be a proxy for neutral genic  
138 diversity, this finding suggests that intact singleton gene prevalence is particularly associated  
139 with selection efficacy, and not simply with standing neutral variation. In other words, there is a  
140 role for natural selection in maintaining even very rare intact genes within pangenomes.

141           Although it is difficult to prove that most rare pseudogenes are evolving neutrally, it is  
142 possible to test for signals expected if there is positive selection for pseudogene loss. If this were  
143 true, pseudogene content would be expected to be lower in species with higher efficacy of  
144 selection. Contrary to this prediction, the mean percent of species' genomes covered by  
145 pseudogenes was not significantly associated with dN/dS (partial Spearman's  $\rho=0.005$ ;  $P =$   
146  $0.8972$ ; **Extended Data Fig. 2b**), which is inconsistent with a model of widespread slightly  
147 deleterious pseudogenes that are purged only in species with sufficiently high  $N_e$ .

148           A limitation of our partial correlation analyses is that they did not control for systematic  
149 differences across taxonomic groups. In addition, they provide no insight into the relative  
150 explanatory power of dN/dS vs. dS for explaining pangenome diversity. To address these points,  
151 we conducted a complementary linear modelling analysis, where a separate model was generated  
152 with each of the four pangenome diversity measures as the response, and dS, dN/dS, and  
153 taxonomic class as predictors. Continuous variables were converted to standard units so that  
154 coefficients could be compared across models. All models were highly significant ( $P<0.001$ ;  
155 **Figure 3**) and ranged in adjusted  $R^2$  values from 0.197 to 0.420 for the  $s_i/s_p$  and  $s_i$  models,  
156 respectively.

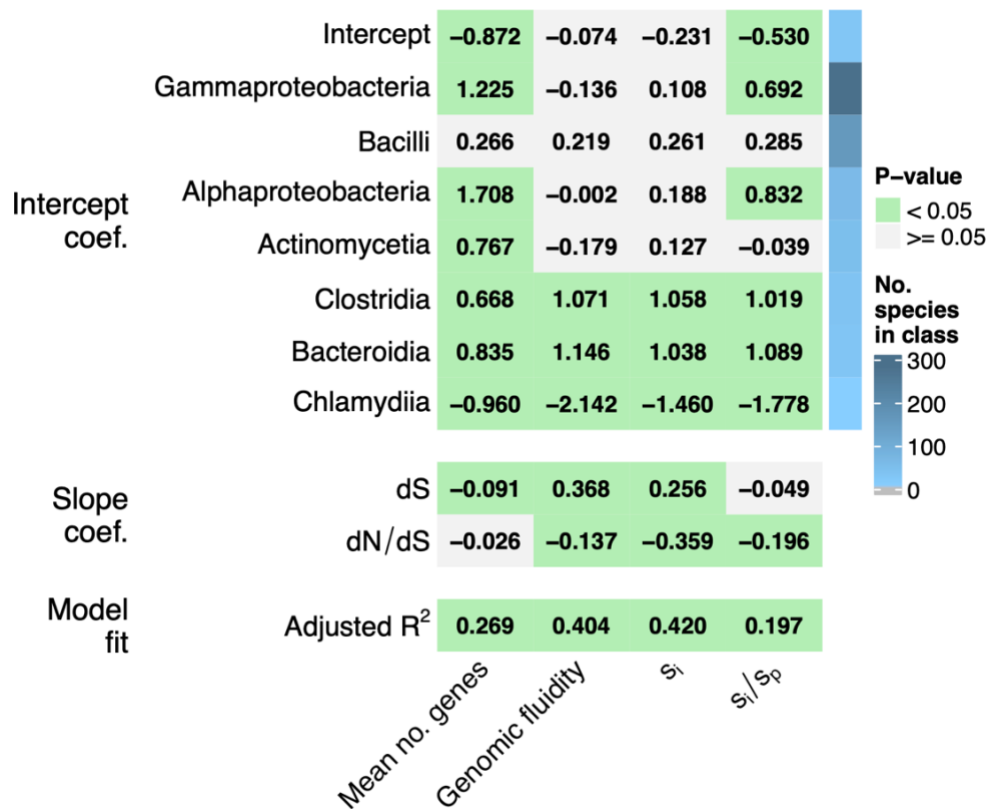
157           All but one class (Bacilli) were significant predictors in at least one model, and  
158 Clostridia, Bacteroidia, and Chlamydia were significant predictors across all four models.  
159 Similarly, dS was a significant predictor of all pangenome diversity metrics except for  $s_i/s_p$ . In  
160 contrast, dN/dS was a significant predictor for all pangenome diversity metrics except for the  
161 mean number of genes, which could indicate that gene number is an overly simplistic measure of  
162 pangenome diversity. Most pertinently, these results highlight that dN/dS, a proxy for selection  
163 efficacy, remains a significant predictor of  $s_i/s_p$ . In addition, dS, a measure that incorporates both  
164 divergence time and the species-wide level of standing neutral variation, is a predictor of  $s_i$ , but  
165 not  $s_i/s_p$ , which would be unexpected were singleton intact genes and pseudogenes both evolving  
166 neutrally. Instead, these results are consistent with  $s_i/s_p$  behaving somewhat analogously to dN/dS  
167 as a measure of the efficacy of selection. As a higher fraction of rare genes (relative to  
168 pseudogenes) are retained when selection is more effective, this is consistent with many  
169 singleton genes conferring adaptive benefits, and/or some singleton pseudogenes being slightly  
170 deleterious. As the latter effect is undetectable in our data (**Extended Data Fig. 2b**), we favour  
171 the hypothesis that rare intact genes tend to provide benefits to their host genomes.

172



173

174 **Figure 2: Associations between pangenome diversity metrics and estimated efficacy of selection**  
175 **(dN/dS).** Each panel presents the association between the ratio of non-synonymous to synonymous  
176 substitution rates (dN/dS; across each species' core genome) and one of the following measures: (a) the  
177 mean number of genes per genome, (b) genomic fluidity, (c) the mean percent of intact singletons, and  
178 the percentage of singleton intact genes normalized by the percentage of singleton pseudogenes per  
179 species. Each point is one of 668 prokaryotic species, plotted on log<sub>10</sub> scales. The partial Spearman  
180 correlation coefficients (which control for dS) and  $P$ -values are indicated in the bottom left corners. In both  
181 panels c and d, one species (*Rickettsia prowazekii*) contained no singleton intact genes and is indicated  
182 by the point intersecting the x-axis in both panels.



183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197

**Figure 3: The  $s_i/s_p$  metric varies across taxa and is correlated with the efficacy of selection.**

Summaries of four pangenome diversity linear models are shown. One model was fit for each pangenome diversity metric: the mean number of genes, genomic fluidity, the percentage of singleton intact genes ( $s_i$ ), and the ratio of the percentages of singleton intact genes vs. pseudogenes ( $s_i/s_p$ ). All continuous response and predictor variables were standardized (i.e. converted to z-scores) prior to building models. Most continuous variables were also transformed to normal distributions prior to this standardization (see Online Methods). Coefficients are displayed for each model, split by those that affect the intercept vs. the slope. The adjusted R<sup>2</sup> is also indicated for each model, and the cell colouring indicates whether each value is statistically significant ( $P < 0.05$ ). The number of genomes per taxonomic class is indicated by the blue bar. The category used to infer the overall intercept was based on a combination of all classes with  $\leq 5$  species present. These models were built based on 667 species, after excluding one species with no singleton intact genes.

198

199

200

201

202

203

204

205

Having established  $s_i/s_p$  as a measure of selection on rare accessory genes, we asked how selection varies across different functional categories of rare genes. To answer this question, we used a dataset of 10 species with a relatively high number of genomes, including highly sampled human pathogens and bacteria with other lifestyles: *Agrobacterium tumefaciens* (223 genomes), *Enterococcus faecalis* (1,298 genomes), *Escherichia coli* (2,955 genomes), *Lactococcus lactis* (135 genomes), *Pseudomonas aeruginosa* (4,115 genomes), *Sinorhizobium meliloti* (166 genomes), *Staphylococcus epidermidis* (447 genomes), *Streptococcus pneumoniae* (6,845 genomes), *Wolbachia pipientis* (716 genomes), and *Xanthomonas oryzae* (326 genomes). We



206 called intact genes and intergenic pseudogenes across these genomes as described above, but  
207 performed joint clustering of intact genes and pseudogenes, to ensure that differences in how  
208 sequence clusters are defined do not influence the results. These 10 species substantially varied  
209 in genome content and characteristics (**Extended Data Table 2**); for example, *Wolbachia*  
210 *pipientis* genomes encoded a mean of 897.0 intact genes (SD: 25.1) and 55.4 pseudogenes (SD:  
211 20.8), while *Sinorhizobium meliloti* genomes encoded a mean of 6032.8 intact genes (SD: 205.7)  
212 and 489.7 pseudogenes (SD: 53.4).

213 We annotated each sequence cluster using eggNOG-mapper<sup>24</sup> to identify Clusters of  
214 Orthologous Genes (COG) annotations<sup>25</sup>. This tool annotates protein sequences, which is  
215 problematic for most pseudogenes as the protein-coding information is generally lost. Instead,  
216 we annotated all proteins (i.e. those from a larger database used to define pseudogenes  
217 originally) that matched each pseudogene sequence. We identified a mean of 57.94% (SD:  
218 7.06%) of intact gene clusters and 49.46% (SD: 7.09%) of pseudogene clusters as COG-  
219 annotated. The ratio of the percent COG-annotated intact genes vs. pseudogenes was  
220 significantly higher than one in 6/10 of species and lower than one in 2/10 species (Fisher's exact  
221 tests,  $P < 0.05$ ). We separated all clusters into three pangenome partitions, based on their  
222 frequency within a species: cloud ( $\leq 15\%$ ), shell ( $> 15\%$  and  $< 95\%$ ), and soft-core ( $\geq 95\%$ ). We  
223 also further partitioned cloud clusters into ultra-rare, including clusters found in only one or two  
224 genomes (singletons and doubletons), and other-rare, referring to higher-frequency cloud  
225 clusters. As expected, most pseudogene clusters were within the cloud partitions: mean of  
226 95.46% (SD: 3.78%) vs. a mean of 84.01% (SD: 8.34%) for intact genes (**Extended Data**  
227 **Figure 3a**). Some pseudogene clusters were in the soft-core partition (mean: 0.54%, SD: 0.66%),  
228 which primarily lacked COG annotations (**Extended Data Figure 3b**). For the subsequent  
229 analyses we proceeded with COG-annotated clusters only (**Extended Data Figure 4**).

230 We applied generalized linear mixed models, for each pangenome partition separately  
231 (excluding soft-core elements), to investigate which factors best explain whether an element is  
232 intact or a pseudogene. These models included 213,912, 3,650,010, and 12,234,597 elements for  
233 the ultra-rare, other-rare, and shell partitions, respectively. The fixed effects included each  
234 element's COG category and whether the element was redundant with an intact gene with the  
235 same COG ID in the same genome. We included the 'redundancy' effect because adaptive genes  
236 might neutrally degenerate if they are complemented by an intact copy of the same gene family

237 in the genome. The interaction between COG category and functional redundancy was also  
238 included as a fixed effect. Last, we also included species names, the interaction between COG  
239 category and species, and the interaction between functional redundancy and species random  
240 effects. All variables added significant information to these models, but there were some slight  
241 differences in their relative contributions. For instance, species identity and element functional  
242 redundancy were particularly informative in the ultra-rare model compared to the more frequent  
243 categories of genes (**Extended Data Figure 5**), and certain species displayed different  
244 associations with pseudogenization by pangenome partition (**Extended Data Figure 6**).

245 We identified significant coefficients in the ultra-rare model (**Figure 4**), which provided  
246 insight into what factors were most associated with pseudogene status ( $P < 0.05$ ). These  
247 coefficients represent decreased log-odds (logit) probabilities of an element being a pseudogene.  
248 Five COG categories were positively associated with pseudogenization: ‘energy production and  
249 conversion’ (C), ‘nucleotide transport and metabolism’ (F), ‘translation, ribosomal structure and  
250 biogenesis’ (J), ‘function unknown’ (S), and – most strongly – ‘mobilome: prophages,  
251 transposons’ (X). ‘Cell cycle control, cell division, chromosome partitioning’ (D), was the sole  
252 COG category specifically associated with decreased pseudogenization. Non-redundant elements  
253 were highly associated with decreased pseudogenization, over most COG categories. This  
254 indicates that even very rare accessory genes are often under selection to maintain a functional  
255 copy in the genome. Non-redundant elements were also depleted for pseudogenes in the other-  
256 rare and shell models, but different COG categories were associated with pseudogenization  
257 overall (**Extended Data Figure 7**). The exception was an enrichment of pseudogenes in  
258 mobilome-associated elements in the other-rare partition.

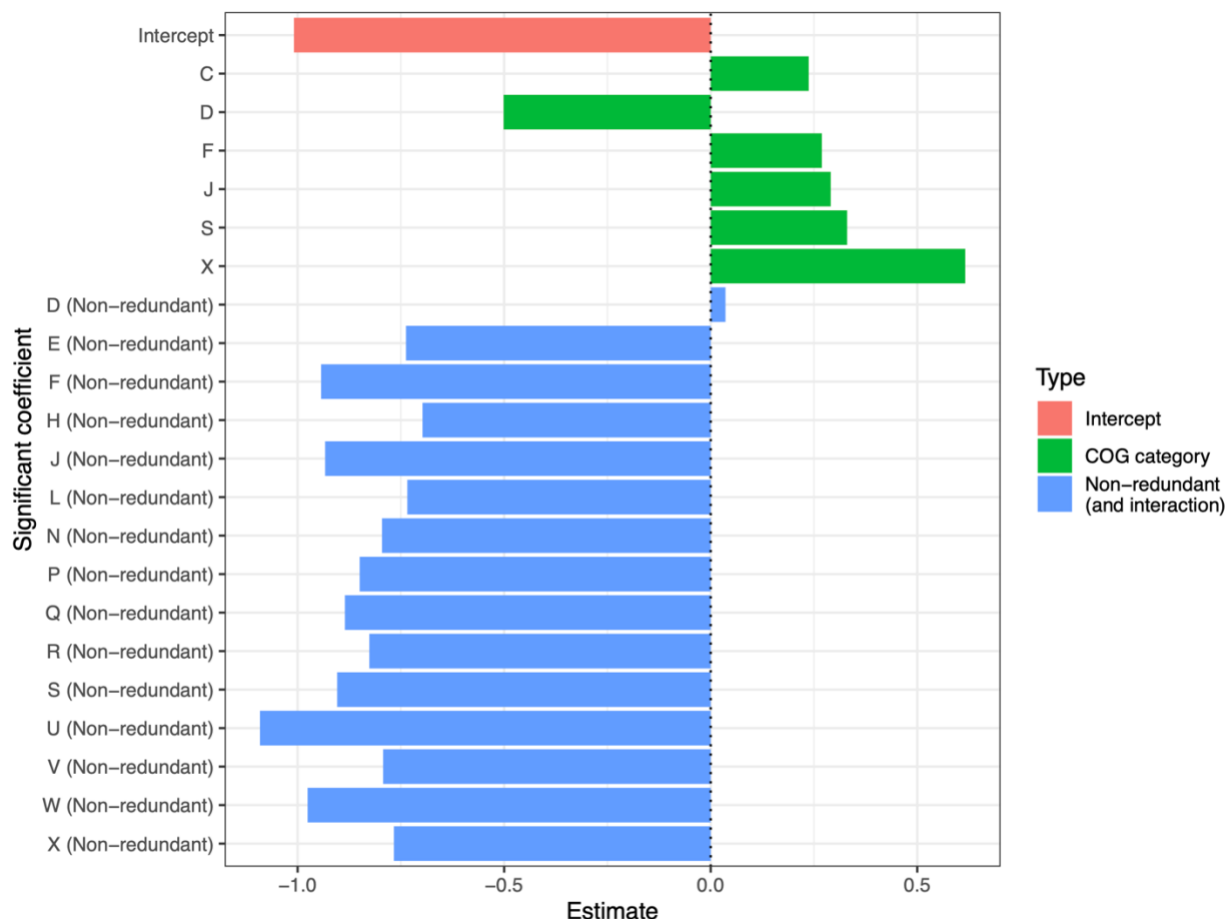
259 In the study of pangenome evolution, a key question is what proportion of rare genes are  
260 under selection or subject to genetic drift. This question is challenging to answer precisely; yet  
261 our models yield estimates of the percentage of genes found in functional groupings depleted for  
262 pseudogenes, providing a lower bound for the percentage of adaptive genes. For instance, genes  
263 in COG category D and non-redundant genes in COG category E are two such pseudogene-  
264 depleted groupings. Based on these definitions, a mean of 19.41% (SD: 5.27%), 20.32% (SD:  
265 6.84%), and 26.02% (SD: 7.05%) of intact genes are found in pseudogene-depleted groupings  
266 across the ultra-rare, other-rare, and shell partitions, respectively. The increasing percentage of  
267 genes classified as pseudogene-depleted as gene frequency increases from ultra-rare to shell is

268 consistent with more frequent genes being more likely adaptive to their host. Nevertheless, an  
269 appreciable percentage (>19%) of ultra-rare genes are likely adaptive according to this estimate.  
270 Note that although element COG non-redundancy was highly negatively associated with  
271 pseudogenization, only 24.39% of elements were non-redundant, which accounts for why only a  
272 minority of intact genes were categorized into pseudogene-depleted groupings. Conversely,  
273 18.68% (SD: 5.62%), 13.29% (SD: 7.69%), and 3.65% (SD: 0.74%) of intact genes are found in  
274 groupings enriched for pseudogenes across these three partitions. The decreasing percentages as  
275 gene frequency increases is consistent with rarer genes being more likely deleterious to their  
276 host. Therefore, although rare accessory genes may on average be adaptive to their host  
277 genomes, a substantial fraction may also be deleterious. Most intact genes do not fall cleanly into  
278 either the pseudogene-enriched or -depleted category, meaning that these estimates represent  
279 rough lower bounds of how many genes are likely adaptive or deleterious.

280 Several COG categories were significant in our models, but these are broad groupings  
281 that can be difficult to biologically interpret. We investigated which individual COG IDs within  
282 significant COG categories were driving the overall signals in the ultra-rare model (see Online  
283 Methods). The clearest signal was of transposase-associated COGs being highly enriched among  
284 pseudogenes (mean of significant odds ratios: 5.10, SD: 6.86), which contrasted with other  
285 mobilome-associated COGs (**Extended Data Fig. 8**). We also identified several COGs highly  
286 associated with pseudogenization in specific species. For instance, anaerobic selenocysteine-  
287 containing dehydrogenases (COG0243, category C), were highly enriched for pseudogenes  
288 across multiple species, particularly in *Agrobacterium tumefaciens* (odds ratio: 103.6,  $P <$   
289 0.001). In addition, several COGs in category D involved in cell division and chromosome  
290 segregation were significantly depleted for pseudogenes, including BcsQ (COG1192), a ParA-  
291 like ATPase, which was significantly depleted for pseudogenes in six species (false discovery  
292 rate < 0.05).

293

294



295  
296 **Figure 4:** Summary of significant coefficients ( $P < 0.05$ ) in generalized linear mixed model with singleton  
297 and doubleton (ultra-rare) element state (intact or pseudogene) as the response. The predictors were  
298 each element's annotated COG category (indicated by single-letter codes), whether the element is  
299 redundant with an intact gene of the same COG ID (i.e. gene family, not COG category) in the same  
300 genome, and the interaction between these variables. The non-redundant coefficients represent the sum  
301 of the overall non-redundant coefficient and the interaction of non-redundancy and each COG category.  
302 Estimates correspond to logit (log-odds) values: estimates  $> 0$  indicate an increased probability of an  
303 element being classified as a pseudogene. Significant COG categories (excluding those significant when  
304 non-redundant) include: energy production and conversion (C), cell cycle control, cell division,  
305 chromosome partitioning (D), nucleotide transport and metabolism (F), translation, ribosomal structure  
306 and biogenesis (J), function unknown (S), and mobilome: prophages, transposons (X).

307  
308 The ability to distinguish neutral and adaptive models of pangenome evolution has been  
309 hindered by a lack of tools to test for selection acting on gene content. This contrasts with an  
310 established toolkit of tests for selection at the nucleotide or protein level, including dN/dS and its  
311 extensions. Here we propose pseudogene diversity as a reference for distinguishing neutral and  
312 adaptive forces acting on pangenomes – particularly rare genes. We showed that the association  
313 between pangenome diversity and synonymous-site variation disappears after correcting for  
314 pseudogene diversity with the  $s_i/s_p$  metric, while the association with dN/dS is maintained. This

315 indicates that a higher proportion of intact singleton genes (relative to singleton pseudogenes) are  
316 present when selection is more effective. This would be unexpected if all rare intact genes were  
317 evolving neutrally, and so is strong evidence against a fully neutral model of prokaryotic  
318 pangenome diversity. Instead, it is consistent with a model where rare intact genes confer slightly  
319 adaptive functions, which are more likely to be preserved by selection given higher selection  
320 efficacy<sup>7</sup> (such as in *E. coli*), but that may degenerate neutrally and become pseudogenes in  
321 species with lower  $N_e$  (such as obligate intracellular bacteria). It would also be consistent with a  
322 model where there are widespread slightly deleterious rare pseudogenes, which can be purged  
323 only in species with high  $N_e$ , but we did not detect a significant association between dN/dS and  
324 pseudogene content, making this less likely.

325 A common explanation for widespread selection on rare accessory genes is adaptation to  
326 highly specialized niches<sup>13-15</sup>. While genes recently acquired through horizontal gene transfer are  
327 often hypothesised to be niche-specific adaptations<sup>26</sup>, it is challenging to make high-confidence  
328 inferences without knowing the background of all recently transferred genes that were not  
329 retained – and are thus unobservable by definition. By focusing on pseudogenes, which are  
330 observable but likely to evolve mostly by drift, we can establish a (nearly) neutral background  
331 against which to discern potentially niche-specific adaptations.

332 We relied on the assumption that any selection pressures acting upon pseudogenes overall  
333 are of much lower magnitude compared to intact genes. In other words, we assumed that, overall,  
334 the pseudogenization instances we identified do not reflect adaptive gene loss<sup>27</sup> (which is  
335 unlikely to substantially increase with selection efficacy, as described above), nor do they  
336 represent adaptive regulatory informative transferred between bacteria through HGT<sup>28</sup>. This  
337 second possibility would be inconsistent with the positive association we observed between  $s_i/s_p$   
338 and selection efficacy. Instead, our results are consistent with rare pseudogenes evolving under a  
339 regime closer to neutrality relative to rare intact genes.

340 Our enrichment test results highlight that a significant proportion of rare accessory genes  
341 are under selection. Notably, 19% of ultra-rare intact genes are in COG categories significantly  
342 depleted for pseudogenes. We hypothesise that many such genes are under purifying selection,  
343 while relaxed purifying selection could account for the observed enrichment of transposons  
344 among pseudogenes. The clear enrichment of selenocysteine-containing dehydrogenases could  
345 similarly reflect relaxed, or sporadic, purifying selection on these elements, which is interesting

346 as selenium, selenocysteine's defining component, is sporadically used across the prokaryotic  
347 tree<sup>29</sup>.

348 Gene-level selection could also account for certain observations. For instance, the DNA  
349 partitioning protein highly enriched in intact ultra-rare genes, COG1192, is a known plasmid-  
350 encoded element predicted to be involved with plasmid partitioning<sup>30</sup>. It is possible that there is  
351 an ascertainment bias in identifying such genes as intact, because were they pseudogenized or  
352 lost the entire plasmid might not be transferred to daughter cells. Similar biases could also  
353 account for why prophage and plasmid-associated elements in the mobilome more generally are  
354 depleted among pseudogenes, although these elements are also more likely to be adaptive to the  
355 host genome<sup>31,32</sup>.

356 Another caveat is that pseudogene diversity can be influenced by many factors, including  
357 life history. For instance, obligate intracellular bacteria are characterized by widespread  
358 degeneration of their genome, followed by streamlining<sup>33</sup>. Depending on a species' stage in this  
359 evolutionary process, its genome could be enriched or depleted for pseudogenes relative to other  
360 bacteria. This likely accounts for certain  $s_i/s_p$  outliers we observed, such as the obligate  
361 intracellular bacteria *Rickettsia prowazekii*, which had the lowest  $s_i/s_p$  ratio. Accordingly, our  
362 framework could be improved by incorporating per-species parameters of pseudogene gain and  
363 loss dynamics.

364 Despite these caveats, our work highlights the value of using pseudogene diversity as a  
365 neutral null<sup>34</sup> for evaluating the evolutionary forces acting upon intact accessory genes.  
366 Establishing true neutrality in microbial genomes is challenging<sup>35</sup>, but the clear association we  
367 identified between  $dN/dS$  and  $s_i/s_p$  suggests that pseudogene diversity can provide insight into  
368 how rare accessory genes evolve. Using this approach, we show that a purely neutral pangenome  
369 model can be rejected and identify which types of rare genes, based on their functional  
370 annotation and what species encodes them, are more likely to be retained by selection.

371

### 372 **Code and data availability**

373 The code used for the analyses in this manuscript is located at  
374 [https://github.com/gavinmdouglas/pangenome\\_pseudogene\\_null](https://github.com/gavinmdouglas/pangenome_pseudogene_null) and the key datafiles are  
375 available on Zenodo (DOI: [10.5281/zenodo.7942837](https://doi.org/10.5281/zenodo.7942837)). All analysed genomes are publicly  
376 available as part of NCBI RefSeq/GenBank.

## 377 **Acknowledgements**

378 We would like to thank Louis-Marie Bobay for reading a draft of this manuscript and providing  
379 feedback. GMD is supported by a Natural Sciences and Engineering Research Council of Canada  
380 (NSERC) Postdoctoral Fellowship. WFD is funded by the Gordon and Betty Moore Foundation.  
381 BJS is supported by an NSERC Discovery Grant.

382

## 383 **Ethics declarations**

384 The authors declare that they have no competing interests related to the content of this article.

385

## 386 **Online Methods**

### 387 *Dataset processing – broad pangenome analysis*

388 We downloaded all genomes used in this study from the Genome Taxonomy Database<sup>18</sup> release  
389 202. We identified all species in this database with at least ten high quality genomes, based on  
390 these criteria: (1) marked as passing the minimum information about a metagenome-assembled  
391 genome<sup>36</sup> check; (2) CheckM<sup>37</sup> completeness > 98% and contamination < 1%; (3) fewer than  
392 1000 contigs; (4) contig N50 > 5000; (6) fewer than 100,000 ambiguous bases. We also  
393 restricted our analyses to genomes in RefSeq (rather than those in GenBank only), except for  
394 *Wolbachia pipientis* genomes, which were numerous but primarily limited to GenBank. For  
395 species with more than twenty genomes, we randomly sampled down to twenty genomes. We  
396 identified 670 species that fit these criteria and downloaded the corresponding genomes. Certain  
397 genomes had been relabelled or removed from NCBI since the release of Genome Taxonomy  
398 Database release 202, which resulted in a minimum of nine genomes per species (we eliminated  
399 two species with fewer than nine genomes). We annotated all genomes with Prokka<sup>38</sup> version  
400 1.14.5 with the `-kingdom`, `--compliant`, and `-rfam` options. We also specified the `—metagenome`  
401 flag for all genomes with 50 or more contigs. We ran Panaroo<sup>39</sup> version 1.3.0 on all output GFFs,  
402 with the `-remove-invalid-genes` and `--clean-mode strict` options. We then ran Pseudofinder<sup>19</sup> on  
403 the Prokka-output GenBank files to identify all putative pseudogenes, using protein sequences  
404 from the UniRef90 database<sup>40</sup> (UniProt KB release 2022\_01) as a reference database. We  
405 restricted the output to intergenic pseudogenes specifically, as the other pseudogene types  
406 identified by Pseudofinder correspond to divergent intact coding sequences (in length or  
407 modularity), which are difficult to interpret as truly degenerating sequences, and could simply

408 represent functionally divergent proteins. We performed three filtering steps on the output  
409 intergenic pseudogenes. Specifically, we excluded all (1) pseudogene calls within 500 bp of  
410 contig ends, (2) pseudogenes of called length  $< 100$  bp or  $> 5000$  bp, and (3) pseudogenes that  
411 substantially differed from the mean size of all matching database hits (mean database size –  
412 pseudogene size was inclusively required to be between -500 bp and 2000 bp). Pseudogenes  
413 were clustered with cd-hit<sup>41</sup> version 4.8.1 with an identity cut-off of 95% over at least 90% of  
414 both compared sequences. The mean numbers of genes and singletons per species were identified  
415 by repeated subsampling to nine strains per species and then comparing Panaroo gene sets. This  
416 procedure was repeated for up to 100 replicates (or until the maximum number of strain  
417 combinations was reached) and the mean number of genes and singletons per genome was  
418 computed across all replicates. This same procedure was repeated for computing the pseudogene  
419 statistics, and the mean percentage of singletons per species was calculated by dividing the mean  
420 number of singletons by the mean number of genes per species (and multiplying by 100). To be  
421 clear, this computation means that the  $s_i/s_p$  metric corresponds to a comparison of the percentage  
422 of singleton intact and pseudogene calls overall per species, rather than of calls within each  
423 individual genome. Where possible, these commands were parallelized with GNU Parallel<sup>42</sup>  
424 version 20161222.

425

### 426 *Metric computation*

427 We performed codon-aware multiple-sequence alignment of all ubiquitous and single-copy genes  
428 sequences per-species with muscle<sup>43</sup> version 3.8.1551, based on the HyPhy<sup>44</sup> version 2.5.36  
429 codon-aware workflow (<https://github.com/veg/hyphy-analyses/tree/master/codon-msa>). We then  
430 concatenated the core gene alignments per species with a Python script  
431 (cat\_core\_genome\_msa.py) and computed pairwise dN/dS and dS for each combination of strain  
432 pairs per species with an additional script (mean\_pairwise\_dnds.py). Both scripts, and the bash  
433 commands for running the codon-aware alignments, are available in v1.1.0 of this repository:  
434 [https://github.com/gavinmdouglas/handy\\_pop\\_gen](https://github.com/gavinmdouglas/handy_pop_gen). The latter script identifies potential non-  
435 synonymous and synonymous mutation sites between each sequence pair using the NG86  
436 approach<sup>45</sup>. We computed the mean values across all pairwise strain comparisons, resulting in a  
437 single measure of dN/dS and dS per species.

438



### 439 ***Linear models***

440 We built linear models using the `lm` function in R to predict pangenome diversity, based on (per  
441 species) either the mean number of genes, the genomic fluidity,  $s_i$ , or  $s_i/s_p$ . The predictors  
442 included  $dS$ ,  $dN/dS$ , and taxonomic class. Classes with  $\leq 5$  member species were collapsed into  
443 the “Other” category, which was set as the intercept for the models. One species, *Rickettsia*  
444 *prowazekii*, was excluded from this analysis due to values of zero for  $s_i$  and  $s_i/s_p$ . We transformed  
445 all continuous variables to be normally distributed, except for the mean number of genes, which  
446 was already normally distributed. We performed a square-root transformation of the genomic  
447 fluidity,  $s_i$ ,  $s_i/s_p$ , and  $dS$  values. The  $dN/dS$  values were especially right skewed and required a  
448 negative inverse transformation ( $-1 * 1/(x)$ , where  $x$  is each  $dN/dS$  value) to be normalized. We  
449 then converted each continuous variable to standardized units, by mean-centring and dividing by  
450 the standard deviation. This step means that the model outputs refer to units of standard deviation  
451 per variable, which makes it possible to compare the magnitude of coefficients across models  
452 with different response variables.

453

### 454 ***Dataset processing – In-depth pangenome analysis***

455 We conducted a subsequent analysis on 10 bacterial species with a relatively high number of  
456 genomes (ranging from 135-6,916). We selected these species from our original set as those with  
457  $> 100$  genomes that were not phylogenetically redundant. For these data, we clustered both intact  
458 genes and pseudogenes with `cd-hit`, using the same settings as above. This clustering was  
459 performed on all genes and pseudogenes across all ten species. We functionally annotated each  
460 resulting cluster with COG IDs and categories<sup>25</sup> using `eggNOG-mapper`<sup>24</sup> version 2.1.6 (based on  
461 `eggNOG` orthology data<sup>46</sup> version 5.0.2) with `DIAMOND`<sup>47</sup> version 2.0.14 and these parameter  
462 options: `--score 60`, `--pident 40`, `--query_cover 20`, `--subject_cover 20`, `--tax_scope auto`, and `--`  
463 `target_orthologs all`. This was performed for individual elements separately (i.e. the original  
464 sequences rather than the cluster representatives), and for database sequence matches to  
465 pseudogene hits. We used majority rule of all member sequences per cluster to assign individual  
466 COG IDs and categories, and the same approach for assigning functions to individual  
467 pseudogene sequences based on database sequence annotations. We manually assigned COG  
468 categories based on a mapping of COG IDs from the COG 2020 database release. This was

469 performed as the raw output COG categories were based on an earlier version of the database  
470 that did not include mobilome (category X) annotations.

471

### 472 *Generalized linear mixed models*

473 Generalized linear mixed models were fit in R using the `glmmTMB`<sup>48</sup> package v1.1.5, one for the  
474 ultra-rare, other-rare, and shell pangenome partitions, respectively. Only COG-annotated  
475 elements were included in these models, excluding those annotated by the (rare) A, B, Y, and Z  
476 COG categories only. We used the binomial family and `nlminb` optimization algorithm with  
477 1000 set for both `iter.max` and `eval.max`. The full R-style formula for each model was:

478

479  $\text{pseudogene} \sim \text{COG-category} + \text{non-redundant-status} + \text{COG-category}:\text{non-redundant-status} + (1$   
480  $| \text{species}) + (1 | \text{COG-category}:\text{species}) + (1 | \text{non-redundant-status}:\text{species})$

481

482 In this formula, random effects are specified as those in parentheses including “1|” and  
483 interaction terms are indicated with “:”. The response was a Boolean variable indicating whether  
484 each element is a pseudogene. The COG-category variable is categorical indicating the one-letter  
485 COG category code that each element belongs to. In cases where elements were members of  
486 multiple categories, duplicate rows were created for each category. The Transcription category  
487 (K) was selected as the first level, to be used for the intercept, as it was the most consistently  
488 abundant COG category across all three partitions (third in the other-rare and shell, and fourth in  
489 ultra-rare). The non-redundant-status variable was a Boolean variable indicating whether each  
490 element was not redundant with another intact element of the same COG ID (gene family, not  
491 category) in the same genome. This negative formulation of redundancy (i.e. whether an element  
492 is not redundant, rather than whether it is redundant) was chosen as most elements were  
493 redundant, and so we decided to set the default level in each model (False) to be more  
494 representative. The species variable corresponded to the name of the species encoding each  
495 element.

496

497 We also fit simpler models with subsets of these variables and computed Akaike Information  
498 Criterion (AIC) values for each model, that allowed us to compare across models and investigate  
499 whether more complex models provide significantly more information. We visualized the AICs

500 per model based on normalized scores that transformed the minimum model AIC per partition to  
501 be 0 and the maximum model AIC per partition to be 1.

502

503 Finally, for each significant COG category in the ultra-rare generalized linear model (excluding  
504 those interacting with non-redundancy), we systematically tested whether individual COG IDs  
505 were enriched for pseudogenes based on Fisher's exact tests comparing the number of  
506 pseudogene and intact genes within each COG ID (and with the same redundancy status and in  
507 the same species) compared to the background of all other elements with the same redundancy  
508 status in the same species.

509

### 510 ***General analyses***

511 No tests for statistical power were conducted to determine the sample sizes required for this  
512 study, but we used genomes from all available species in the Genome Taxonomy Database of  
513 sufficient quality. All analyses were conducted in R v4.2.2. Figures were generated with  
514 ggplot2<sup>49</sup> v3.4.0, with the exception of the heatmaps, which were created with the  
515 ComplexHeatmap<sup>50</sup> package v2.14.0.

516

### 517 **References**

- 518 1. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus*  
519 *agalactiae*: Implications for the microbial 'pan-genome'. *Proc. Natl. Acad. Sci. U. S. A.* **102**,  
520 3950–13955 (2005).
- 521 2. Vos, M., Hesselman, M. C., te Beek, T. A., van Passel, M. W. J. & Eyre-Walker, A. Rates of  
522 Lateral Gene Transfer in Prokaryotes: High but Why? *Trends Microbiol.* **23**, 598–605  
523 (2015).
- 524 3. Innamorati, K. A., Earl, J. P., Aggarwal, S. D., Ehrlich, G. D. & Hiller, N. L. The Bacterial  
525 Guide to Designing a Diversified Gene Portfolio. in *The Pangenome: Diversity, Dynamics*  
526 *and Evolution of Genomes* (eds. Tettelin, H. & Medini, D.) 51–87 (Springer, 2020).  
527 doi:10.1007/978-3-030-38281-0\_3.

- 528 4. Novick, A. & Doolittle, W. F. Horizontal persistence and the complexity hypothesis. *Biol.*  
529 *Philos.* **35**, 2 (2020).
- 530 5. Sela, I., Wolf, Y. I. & Koonin, E. V. Theory of prokaryotic genome evolution. *Proc. Natl.*  
531 *Acad. Sci. U. S. A.* **113**, 11399–11407 (2016).
- 532 6. Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent on effective  
533 population size. *ISME J.* **11**, 1719–1721 (2017).
- 534 7. Bobay, L. M. & Ochman, H. Factors driving effective population size and pan-genome  
535 evolution in bacteria. *BMC Evol. Biol.* **18**, 153 (2018).
- 536 8. Haegeman, B. & Weitz, J. S. A neutral theory of genome evolution and the frequency  
537 distribution of genes. *BMC Genomics* **13**, 196 (2012).
- 538 9. Lobkovsky, A. E., Wolf, Y. I. & Koonin, E. V. Gene frequency distributions reject a neutral  
539 model of genome evolution. *Genome Biol. Evol.* **5**, 233–242 (2013).
- 540 10. Fiona J Whelan, Rebecca J Hall, & James O McInerney. Evidence for Selection in the  
541 Abundant Accessory Gene Content of a Prokaryote Pangenome. *Mol. Biol. Evol.* **38**, 3697–  
542 3708 (2021).
- 543 11. N’Guessan, A., Brito, I. L., Serohijos, A. W. R. & Shapiro, J. Mobile Gene Sequence  
544 Evolution within Individual Human Gut Microbiomes Is Better Explained by Gene-Specific  
545 Than Host-Specific Selective Pressures. *Genome Biol. Evol.* **13**, (2021).
- 546 12. Wolf, Y. I., Makarova, K. S., Lobkovsky, A. E. & Koonin, E. V. Two fundamentally  
547 different classes of microbial genes. *Nat. Microbiol.* **2**, 1–6 (2016).
- 548 13. Boucher, Y. *et al.* Local Mobile Gene Pools Rapidly Cross Species Boundaries To Create  
549 Endemicity within Global *Vibrio cholerae* Populations. *mBio* **2**, e00335-10 (2011).

- 550 14. Niehus, R., Mitri, S., Fletcher, A. G. & Foster, K. R. Migration and horizontal gene transfer  
551 divide microbial genomes into multiple niches. *Nat. Commun.* **6**, 8924 (2015).
- 552 15. Smillie, C. S. *et al.* Ecology drives a global network of gene exchange connecting the human  
553 microbiome. *Nature* **480**, 241–244 (2011).
- 554 16. Danneels, B., Pinto-Carbó, M. & Carlier, A. Patterns of nucleotide deletion and insertion  
555 inferred from bacterial pseudogenes. *Genome Biol. Evol.* **10**, 1792–1802 (2018).
- 556 17. Kuo, C. H. & Ochman, H. The extinction dynamics of bacterial pseudogenes. *PLoS Genet.* **6**,  
557 e1001050 (2010).
- 558 18. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny  
559 substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
- 560 19. Syberg-Olsen, M. J., Garber, A. I., Keeling, P. J., McCutcheon, J. P. & Husnik, F.  
561 Pseudofinder: Detection of Pseudogenes in Prokaryotic Genomes. *Mol. Biol. Evol.* **39**,  
562 msac153 (2022).
- 563 20. Kislyuk, A. O., Haegeman, B., Bergman, N. H. & Weitz, J. S. Genomic fluidity: An  
564 integrative view of gene diversity within microbial populations. *BMC Genomics* **12**, (2011).
- 565 21. Lerat, E. & Ochman, H.  $\Psi$ - $\Phi$ : Exploring the outer limits of bacterial pseudogenes. *Genome*  
566 *Res.* **14**, 2273–2278 (2004).
- 567 22. Rocha, E. P. C. *et al.* Comparisons of dN/dS are time dependent for closely related bacterial  
568 genomes. *J. Theor. Biol.* **239**, 226–235 (2006).
- 569 23. Kryazhimskiy, S. & Plotkin, J. B. The Population Genetics of dN/dS. *PLoS Genet.* **4**,  
570 e1000304 (2008).

- 571 24. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J. eggNOG-  
572 mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the  
573 Metagenomic Scale. *Mol. Biol. Evol.* **38**, 5825–5829 (2021).
- 574 25. Galperin, M. Y. *et al.* COG database update: focus on microbial diversity, model organisms,  
575 and widespread pathogens. *Nucleic Acids Res.* **49**, D274–D281 (2021).
- 576 26. McInerney, J. O., McNally, A. & O’Connell, M. J. Why prokaryotes have pangenomes. *Nat.*  
577 *Microbiol.* **2**, 170402 (2017).
- 578 27. Hottes, A. K. *et al.* Bacterial Adaptation through Loss of Function. *PLoS Genet.* **9**, e1003617  
579 (2013).
- 580 28. Oren, Y. *et al.* Transfer of noncoding DNA drives regulatory rewiring in bacteria. *Proc. Natl.*  
581 *Acad. Sci. U. S. A.* **111**, 16112–16117 (2014).
- 582 29. Peng, T., Lin, J., Xu, Y.-Z. & Zhang, Y. Comparative genomics reveals new evolutionary  
583 and ecological patterns of selenium utilization in bacteria. *ISME J.* **10**, 2048–2059 (2016).
- 584 30. A Schlüter *et al.* Erythromycin Resistance-Confering Plasmid pRSB105, Isolated from a  
585 Sewage Treatment Plant, Harbors a New Macrolide Resistance Determinant, an Integron-  
586 Containing Tn402-Like Element, and a Large Region of Unknown Function. *Appl. Environ.*  
587 *Microbiol.* **73**, (2007).
- 588 31. Bobay, L. M., Rocha, E. P. C. & Touchon, M. The adaptation of temperate bacteriophages to  
589 their host genomes. *Mol. Biol. Evol.* **30**, 737–751 (2013).
- 590 32. McKerral, J. C. *et al.* The Promise and Pitfalls of Prophages. *bioRxiv* 2023.04.20.537752  
591 (2023) doi:10.1101/2023.04.20.537752.
- 592 33. Giovannoni, S. J., Cameron Thrash, J. & Temperton, B. Implications of streamlining theory  
593 for microbial ecology. *ISME J.* **8**, 1553–1565 (2014).

- 594 34. Koonin, E. V. Splendor and misery of adaptation, or the importance of neutral null for  
595 understanding evolution. *BMC Biol.* **14**, 114 (2016).
- 596 35. Rocha, E. P. C. Neutral Theory, Microbial Practice: Challenges in Bacterial Population  
597 Genetics. *Mol. Biol. Evol.* **35**, 1338–1347 (2018).
- 598 36. Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and a  
599 metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat. Biotechnol.* **35**,  
600 725–731 (2017).
- 601 37. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
602 assessing the quality of microbial genomes recovered from isolates, single cells, and  
603 metagenomes. *Genome Res.* **25**, 1043–55 (2015).
- 604 38. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069  
605 (2014).
- 606 39. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo pipeline.  
607 *Genome Biol.* **21**, 180 (2020).
- 608 40. The UniProt Consortium. The Universal Protein Resource. *Nucleic Acids Res.* **36**, D190–  
609 D195 (2008).
- 610 41. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein  
611 or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 612 42. Tange, O. GNU Parallel: the command-line power tool. *Linux USENIX Mag.* **36**, 42–47  
613 (2011).
- 614 43. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high  
615 throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).

- 616 44. Kosakovsky Pond, S. L. *et al.* HyPhy 2.5—A Customizable Platform for Evolutionary  
617 Hypothesis Testing Using Phylogenies. *Mol. Biol. Evol.* **37**, 295–299 (2020).
- 618 45. Nei, M. & Gojobori, T. Simple methods for estimating the numbers of synonymous and  
619 nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.* **3**, 418–426 (1986).
- 620 46. Huerta-Cepas, J. *et al.* eggNOG 5.0: A hierarchical, functionally and phylogenetically  
621 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.*  
622 **47**, D309–D314 (2019).
- 623 47. Buchfink, B., Reuter, K. & Drost, H.-G. Sensitive protein alignments at tree-of-life scale  
624 using DIAMOND. *Nat. Methods* **18**, 366–368 (2021).
- 625 48. Brooks, M., E. *et al.* glmmTMB Balances Speed and Flexibility Among Packages for Zero-  
626 inflated Generalized Linear Mixed Modeling. *R J.* **9**, 378 (2017).
- 627 49. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York,  
628 2016).
- 629 50. Zuguang Gu. Complex heatmap visualization. *iMeta* **1**, e43 (2022).

630