This version of the article has been accepted for publication, after peer review (when applicable) but is not the Version of Record and does not reflect postacceptance improvements, or any corrections. The Version of Record is available online at: https://doi.org/10.1038/s41559-023-02268-6. Use of this Accepted Version is subject to the publisher's Accepted Manuscript terms of use https://www.springernature.com/gp/open-research/policies/accepted-manuscript-terms.

1	Title
2	Pseudogenes act as a neutral reference for detecting selection in prokaryotic pangenomes
3	
4	Author list
5	Gavin M. Douglas <sup>1,2,*</sup> and B. Jesse Shapiro <sup>1,2,3,*</sup>
6	
7	Affiliations
8	<sup>1</sup> Department of Microbiology and Immunology, McGill University, Montréal, QC, Canada
9	<sup>2</sup> McGill Genome Centre, McGill University, Montréal, QC, Canada
10	<sup>3</sup> McGill Centre for Microbiome Research, McGill University, Montréal, QC, Canada
11	
12	*Emails for correspondence: gavin.douglas@mcgill.ca and jesse.shapiro@mcgill.ca
13	
14	Keywords: Pseudogenes, pangenome, horizontal gene transfer, lateral gene transfer, mobile
15	genes, mobilome, adaptation
16	
17	Licensing note
18	This version of the article has been accepted for publication, after peer review (when applicable)
19	but is not the Version of Record and does not reflect post-acceptance improvements, or any
20	corrections. The Version of Record is available online at: <u>https://doi.org/10.1038/s41559-023-</u>
21	02268-6. Use of this Accepted Version is subject to the publisher's Accepted Manuscript terms
22	of use https://www.springernature.com/gp/open-research/policies/accepted-manuscript-terms.
23	
24	Abstract
25	A long-standing question is to what degree genetic drift and selection drive the divergence in
26	rare accessory gene content between closely related bacteria. Rare genes, including singletons,
27	make up a large proportion of pangenomes (all genes in a set of genomes), but it remains unclear
28	how many such genes are adaptive, deleterious, or neutral to their host genome. Estimates of
29	species' effective population sizes (Ne) are positively associated with pangenome size and

- <sup>30</sup> fluidity, which has independently been interpreted as evidence for both neutral and adaptive
- pangenome models. We hypothesised that pseudogenes, used as a neutral reference, could be

used to distinguish these models. We find that most functional categories are depleted for rare 32 pseudogenes when a genome encodes only a single intact copy of a gene family. In contrast, 33 transposons are enriched in pseudogenes, suggesting they are mostly neutral or deleterious to the 34 host genome. Thus, even if individual rare accessory genes vary in their effects on host fitness, 35 we can confidently reject a model of entirely neutral or deleterious rare genes. We also define the 36 ratio of singleton intact genes to singleton pseudogenes (s<sub>i</sub>/s<sub>p</sub>) within a pangenome, compare this 37 measure across 668 prokaryotic species, and detect a signal consistent with the adaptive value of 38 39 many rare accessory genes. Taken together, our work demonstrates that comparing to pseudogenes can improve inferences of the evolutionary forces driving pangenome variation. 40

41

## 42 Main text

## 43 Introduction

Bacterial strains within the same species often encode substantially different genes. This has
been established through genome analyses where the entire set of ubiquitous ('core') and
variably present ('accessory') genes across strains are taken to encompass a single 'pangenome'.
Based on such analyses, the percentage of accessory genes within pangenomes varies from 4080% across prokaryotic species<sup>1</sup>. Many individual accessory genes have been shown to be
adaptive, but it remains controversial whether genetic drift or natural selection are responsible
for driving overall pangenome variation across species.

This has been investigated by comparing species pangenome diversity to measures of 51 effective population size (Ne). Ne represents the population size under idealized conditions and is 52 the key parameter determining the efficacy of selection vs. genetic drift. Several methods can be 53 used to estimate, or be used as proxies for, Ne. One common proxy is the ratio of non-54 synonymous to synonymous substitution rates (dN/dS) in core genes. Assuming stronger 55 purifying selection against non-synonymous mutations relative to synonymous mutations in core 56 genes, a lower dN/dS ratio indicates higher selection efficacy and thus a higher Ne. Lower dN/dS 57 ratios are associated with higher pangenome diversity across prokaryotic species<sup>2,3</sup>, which has 58 been interpreted as evidence for higher selection efficacy to retain slightly adaptive accessory 59 genes, resulting in larger pangenomes<sup>3,4</sup>. If there was strong selection acting on most accessory 60 genes, then they would be retained regardless of species  $N_e$  (except perhaps for those with very 61 low Ne). 62

One issue with this adaptive explanation is that neutral genetic variation is also higher in 63 species with higher N<sub>e</sub> due to weaker genetic drift<sup>5</sup> (e.g., fewer population bottlenecks). Indeed, 64 nucleotide diversity at neutral sites is the standard method for calculating Ne, assuming equal 65 mutation rates across species. Pangenome diversity is also positively associated with Ne based on 66 this approach<sup>6</sup>. Accordingly, neutral and adaptive explanations for pangenome diversity cannot 67 be distinguished based on associations with Ne alone, especially for rare accessory genes that 68 could represent segregating neutral variation. Others have remarked that this remains an issue as 69 70 it is unclear how to partition genes into categories expected to experience distinct selective pressures<sup>7</sup>. 71

We hypothesised that pseudogenes – genes degenerating through the introduction of 72 mutations such as premature stop codons, insertions, and deletions - could be used as an 73 approximately neutral reference for detecting selection on rare intact accessory genes. 74 Pseudogenes can arise when genetic drift overcomes purifying selection to retain a gene<sup>8</sup>, or 75 through positive selection to eliminate a deleterious gene<sup>9</sup>. We reasoned that accessory gene 76 families that tend to remain intact are likely under stronger purifying selection than those that 77 tend to be pseudogenized. This insight is particularly relevant for rare accessory genes, which 78 make up the largest fraction of pangenomes<sup>10</sup>, and for which the evolutionary forces are most 79 controversial. We investigated this idea by first comparing the functional categories of rare intact 80 genes and pseudogenes across ten well-sampled bacterial species. We then conducted a broader 81 assessment of intact genes and pseudogenes in the pangenomes of nearly 700 different 82 prokaryotic species. 83

84

## 85 **Results**

## 86 Functions of rare elements across well-sampled species

If rare accessory genes were effectively neutral to host fitness, then we would expect no
difference in the functional annotations between intact genes and pseudogenes within a species.
In contrast, differences in functional annotations between these rare element types could suggest
functions that tend to be beneficial in a genome.

We conducted a comparison of the functional annotations of rare intact genes and pseudogenes in a dataset of 10 bacteria species with a relatively high number of sequenced genomes (135-6,845 genomes per species), including highly sampled human pathogens and

bacteria with other lifestyles: Agrobacterium tumefaciens, Enterococcus faecalis, Escherichia 94 coli, Lactococcus lactis, Pseudomonas aeruginosa, Sinorhizobium meliloti, Staphylococcus 95 epidermidis, Streptococcus pneumoniae, Wolbachia pipientis, and Xanthomonas oryzae. We 96 performed joint clustering of intact genes and pseudogenes, to ensure that differences in how 97 sequence clusters are defined did not influence the results. These 10 species varied widely in 98 genome content and characteristics (Extended Data Table 1); for example, Wolbachia pipientis 99 genomes encoded a mean of 897.0 intact genes (SD: 25.1) and 55.4 pseudogenes (SD: 20.8), 100 101 while Sinorhizobium meliloti encoded a mean of 6032.8 intact genes (SD: 205.7) and 489.7 pseudogenes (SD: 53.4). 102

We separated gene/pseudogene clusters into three pangenome partitions, based on their 103 frequency within a species: cloud ( $\leq 15\%$ ), shell ( $\geq 15\%$  and  $\leq 95\%$ ), and soft-core ( $\geq 95\%$ ). We 104 also further partitioned cloud clusters into ultra-rare, including clusters found in only one or two 105 genomes (singletons and doubletons), and other-rare, referring to higher-frequency cloud 106 clusters. Most pseudogene clusters were within the cloud partitions: mean of 95.46% (SD: 107 3.78%) vs. a mean of 84.01% (SD: 8.34%) for intact genes (Extended Data Figure 1a). This 108 would be expected under the common assumption that pseudogenes are under relaxed purifying 109 selection, such that they rapidly accumulate mutations which cause them to be split into separate 110 sequence clusters. Some pseudogene clusters were in the soft-core partition (mean: 0.54%, SD: 111 0.66%), which often could not be annotated with Clusters of Orthologous Genes<sup>11</sup> (COG) 112 identifiers (Extended Data Figure 1b). For subsequent analyses we proceeded with COG-113 annotated clusters only (Figure 1). 114

We applied generalized linear mixed models, for each pangenome partition separately 115 (excluding soft-core elements), to investigate which factors best explain whether a genetic 116 element is an intact gene or a pseudogene. These models included 213,912, 3,650,010, and 117 12,234,597 elements for the ultra-rare, other-rare, and shell partitions, respectively. The fixed 118 effects included each element's COG functional category and whether the element was 119 redundant with an intact gene with the same COG ID in the same genome. We included this 120 'redundancy' effect because adaptive genes might neutrally degenerate only if they are 121 complemented by an intact copy of the same gene family in the genome. The interaction between 122 COG category and functional redundancy was also included as a fixed effect. Last, we also 123 included species names, the interaction between COG category and species, and the interaction 124

between functional redundancy and species random effects. All variables added significant

information to the models, but there were some slight differences in their relative contributions.

<sup>127</sup> For instance, species identity and functional redundancy were particularly informative in the

<sup>128</sup> ultra-rare model compared to the more frequent categories of genes (Extended Data Figure 2),

and certain species displayed different associations with pseudogenization by pangenome

130 partition (**Extended Data Figure 3**).

We identified significant coefficients in the ultra-rare model (Figure 2), which provided 131 132 insight into what factors were most associated with pseudogene status (P < 0.05). These coefficients represent decreased log-odds (logit) probabilities of an element being a pseudogene. 133 Five COG categories were positively associated with pseudogenization: 'energy production and 134 conversion' (C), 'nucleotide transport and metabolism' (F), 'translation, ribosomal structure and 135 biogenesis' (J), 'function unknown' (S), and – most strongly – 'mobilome: prophages, 136 transposons' (X). In contrast, 'Cell cycle control, cell division, chromosome partitioning' (D), 137 was the sole COG category specifically associated with decreased pseudogenization. However, 138 non-redundant elements were highly associated with decreased pseudogenization, over most 139 COG categories. Non-redundant elements were also depleted for pseudogenes in the other-rare 140 and shell models, but different COG categories were associated with pseudogenization overall 141 (Extended Data Figure 4). The exception was an enrichment of pseudogenes in mobilome-142 associated elements in the other-rare partition. 143

In the study of pangenome evolution, a key question is what proportion of rare genes are 144 under selection or subject to genetic drift. This question is challenging to answer precisely; yet 145 our models yield estimates of the percentage of genes found in functional categories depleted for 146 pseudogenes, providing a lower bound for the percentage of adaptive genes. For instance, genes 147 in COG category D and non-redundant genes in COG category E are two such pseudogene-148 depleted groupings. Based on these definitions, a mean of 19.41% (SD: 5.27%), 20.32% (SD: 149 6.84%), and 26.02% (SD: 7.05%) of intact genes are found in pseudogene-depleted groupings 150 across the ultra-rare, other-rare, and shell partitions, respectively. The increasing percentage of 151 genes classified as pseudogene-depleted as gene frequency increases from ultra-rare to shell is 152 consistent with more frequent genes being more likely adaptive to their host. Nevertheless, an 153 appreciable percentage (>19%) of ultra-rare genes are likely adaptive according to this estimate. 154 Although non-redundancy was strongly and negatively associated with pseudogenization, only 155

24.39% of elements were non-redundant, which explains why only a minority of intact genes 156 were categorized into pseudogene-depleted groupings. Conversely, 18.68% (SD: 5.62%), 157 13.29% (SD: 7.69%), and 3.65% (SD: 0.74%) of intact genes are found in groupings enriched for 158 pseudogenes across these three partitions. The decreasing percentages as gene frequency 159 increases is consistent with rarer genes being more likely deleterious to their host. Therefore, 160 although rare accessory genes may on average be adaptive to their host genomes, a substantial 161 fraction may also be deleterious. Most intact genes do not fall cleanly into either the pseudogene-162 163 enriched or -depleted category, meaning that these estimates represent rough lower bounds of how many genes are likely adaptive or deleterious. 164

Several COG categories were significantly enriched or depleted in pseudogenes, but these 165 are broad functional groupings that can be difficult to biologically interpret. We investigated 166 which individual COG IDs within significant COG categories were driving the overall signals in 167 the ultra-rare model (see Online Methods). The clearest signal was of transposase-associated 168 COGs being highly enriched among pseudogenes (mean of significant odds ratios: 5.10, SD: 169 6.86), which contrasted with other mobilome-associated COGs (Extended Data Fig. 5). We also 170 identified several COGs highly associated with pseudogenization in specific species. For 171 instance, anaerobic selenocysteine-containing dehydrogenases (COG0243, category C), were 172 highly enriched for pseudogenes across multiple species, particularly in Agrobacterium 173 *tumefaciens* (odds ratio: 103.6, P < 0.001). In addition, several COGs in category D involved in 174 cell division and chromosome segregation were significantly depleted for pseudogenes, 175 including BcsQ (COG1192), a ParA-like ATPase, which was significantly depleted for 176 pseudogenes in six species (false discovery rate < 0.05). 177

This in-depth analysis of 10 species highlighted several functional categories enriched within rare pseudogenes, particularly for mobilome-related genes. Conversely, we identified a clear depletion of pseudogenes among non-redundant elements, which strongly suggests that even very rare accessory genes are often under selection to maintain a working copy in the genome. Taken together, these comparisons serve as a proof-of-concept that comparing extremely rare intact genes and pseudogenes can be useful for disentangling the action of evolutionary forces.

185

### 186 Comparisons of pangenome diversity and N<sub>e</sub>

We next investigated whether the inclusion of pseudogenes can help resolve the prior conflicting interpretations of the association between pangenome diversity and  $N_e$ . To this end, we analysed 668 named prokaryotic species represented by at least nine genomes in the Genome Taxonomy Database<sup>12</sup>.

We first summarized pangenome diversity across these species. Species' pangenome size 191 and complexity have been characterised previously based on different metrics, including the 192 mean number of genes per genome<sup>2</sup> and genomic fluidity<sup>6,13</sup>. We computed these metrics for all 193 194 species based on both intact genes and pseudogenes. In addition, as we were especially interested in rare elements, we computed the percentages of singleton genes and pseudogenes per species 195 (i.e. those present in a single genome per species), based on repeated subsampling to nine 196 genomes. Larger genomes tend to encode more singletons, both in mean number and percentage 197 (Extended Data Fig. 6a,b). In addition, the percentage of intact singletons is highly correlated 198 with genomic fluidity, but the traditional fluidity metric is sensitive to intermediate frequency 199 accessory genes (Extended Data Fig. 6c,d), which can be driven by inconsistent species 200 definitions or population structure within species. We therefore focused on the percentage of 201 intact (s<sub>i</sub>) and pseudogene (s<sub>p</sub>) singletons for most analyses. All metrics ranged substantially 202 across species for both intact genes (fluidity: 0.00-0.246; mean number: 836.4-8692.7; si: 0.00-203 10.83%) and pseudogenes (fluidity: 0.014-0.513; mean number: 8.1-922.5; sp: 0.78-72.97%). 204

Values of  $s_i$  and  $s_p$  were positively correlated (Spearman's  $\rho$ =0.57; P < 0.001), with deviations suggesting species-specific differences in selection on rare accessory genes (**Figure 3a**). For example, *Escherichia coli* has a relatively higher  $s_i$  value, consistent with selection to retain rare accessory genes, while the obligate intracellular bacteria *Chlamydia trachomatis* and *Rickettsia prowazekii* have lower values, suggesting less selective constraint on their rare genes. To summarize the  $s_i$  and  $s_p$  values per species, we focused on the  $s_i/s_p$  ratio as a metric encompassing both intact gene and pseudogene pangenome diversity.

For each species, we computed two proxies for  $N_e$ : dN/dS and dS. These metrics were negatively correlated (**Figure 3b**), which could be due to the expected impact of  $N_e$  on each metric, and of course their direct dependence since dS is the denominator of dN/dS. However, the dependence structure may be more complex because recently diverged strains are biased towards higher dN/dS ratios due to insufficient time for purifying selection to purge deleterious non-synonymous mutations<sup>14,15</sup>. More generally, interpreting dS as  $N_e$  is questionable if there is

widespread population substructure and uneven sampling of sequenced strains (although this has
been debated: see Discussion). For this reason, we focused on dN/dS as an inversely related
proxy for Ne, and we considered dS as a measure of the divergence between the subset of
analysed genomes per species, but not necessarily as representative of species-wide Ne. Under
this interpretation, the observed positive correlation between dS and both si and sp (Figure 3c
and 3d) is expected simply because these are all measures of genome divergence.

We next recapitulated the previously observed association between dN/dS and standard 224 measures of pangenome diversity<sup>2,3</sup>, and then explored whether dN/dS is also associated with 225 si/sp. We found that the mean number of genes per species was not significantly associated with 226 dN/dS (Figure 4a), but both genomic fluidity, and si were significantly negatively correlated 227 (Figure 4b and c; Spearman correlations, P < 0.05). Although the mean number of genes has 228 been considered in this context previously, it is a measure of overall pangenome size, and is not a 229 direct measure of gene content diversity, which could explain why we did not observe a 230 significant relationship. In contrast, the latter two observations agree with past work<sup>2,3</sup>, but, as 231 discussed above, the biological interpretation of these associations is unclear. Notably, we also 232 found  $s_i/s_p$  to be negatively associated with dN/dS (Spearman's  $\rho$ =-0.22; P < 0.001; Figure 4d), 233 although less strongly than s<sub>i</sub> alone (Spearman's  $\rho$ =-0.54; P < 0.001). These results were 234 qualitatively robust to the number of genomes subset when computing si and sp (Extended Data 235 Figure 7). In addition, although dS was significantly associated with s<sub>i</sub> (Figure 3c), it was not 236 significantly associated with  $s_i/s_p$  (Spearman  $\rho=0.07$ ; P=0.08). Taken together, these results 237 highlight that si remains associated with dN/dS even after normalization by sp, but that the 238 association with dS is lost after this normalization. If pseudogene presence/absence diversity is 239 assumed to be a proxy for neutral gene content diversity, this finding suggests that intact 240 singleton gene prevalence is particularly associated with selection efficacy (dN/dS), and not 241 simply with strain divergence. These results are consistent with si/sp behaving somewhat 242 analogously to dN/dS as a measure of the efficacy of selection. As a higher fraction of rare genes 243 (relative to pseudogenes) are retained when selection is more effective, this is consistent with 244 many singleton genes conferring adaptive benefits, and/or some singleton pseudogenes being 245 slightly deleterious. 246

There are certain species included in our analyses that are outliers with low values of dS and/or high values of dN/dS (**Figure 3b**), which could potentially impact these correlations. In

addition, within-species dN/dS systematically varies with dS, which could also impact our
 conclusions. To address these issues, we re-ran our analysis with outliers removed and using
 partial Spearman correlations that control for dS (Figure 5). With minor exceptions, the results
 remained qualitatively unchanged, which demonstrates that these factors are not driving the
 signal.

Another potential issue is that the species included in our analyses, although they span 255 substantial prokaryotic diversity, were biased towards specific groups, particularly 256 Gammaproteobacteria (286 species) and Bacilli (161 species). As there is substantial variation in 257 pangenome diversity and evolutionary metrics at the class level for the species we considered 258 (Extended Data Figure 8), taxonomic biases could potentially be driving the correlations with 259  $s_i/s_p$  we observed. To account for this, we conducted a linear modelling analysis, where a 260 separate model was generated with each of the four pangenome diversity measures as the 261 response, and dS, dN/dS, and taxonomic class as predictors. All models were highly significant 262 (P<0.001; Extended Data Figure 9) and ranged in adjusted R<sup>2</sup> values from 0.197 to 0.420 for 263 the s<sub>i</sub>/s<sub>p</sub> and s<sub>i</sub> models, respectively. All but one class (Bacilli) were significant predictors in at 264 least one model, and the classes Clostridia, Bacteroidia, and Chlamydiia were significant 265 predictors across all four models. Similarly, dS was a significant predictor of all pangenome 266 diversity metrics except for si/sp. In contrast, dN/dS was a significant predictor for all pangenome 267 diversity metrics except for the mean number of genes. This analysis demonstrated that, despite 268 class-specific differences in pangenome diversity, our overall inferences are robust to taxonomic 269 class as a confounder. 270

271

254

A final caveat of these analyses is that higher si/sp values could be explained by selection 272 to preserve rare intact genes or to purge rare pseudogenes. If there were selection for pseudogene 273 loss, then the pseudogene content per genome would be expected to be lower in species with 274 higher selection efficacy. Contrary to this prediction, the mean percent of species' genomes 275 covered by pseudogenes was not significantly associated with dN/dS (Spearman's  $\rho = 0.063$ ; P =276 0.104; Figure 6a), which is inconsistent with a model of widespread slightly deleterious 277 pseudogenes that are purged only in species with sufficiently high Ne. However, pseudogene 278 coverage is negatively but weakly associated with dS (Spearman's  $\rho = -0.090$ ; P = 0.020; Figure 279

**6b**), which highlights that this interpretation is dependent on the assumption that dN/dS is a more appropriate proxy than dS for N<sub>e</sub> across these genomes (see Discussion). Together, the lack of association between pseudogene content and dN/dS, and the weak association with dS, argue against adaptive purging of rare pseudogenes as a major driver of variation in s<sub>i</sub>/s<sub>p</sub>.

284

#### 285 Discussion

The ability to distinguish neutral and adaptive models of pangenome evolution has been hindered 286 by a lack of tools to test for selection acting on gene content. This contrasts with an established 287 toolkit of tests for selection at the nucleotide and protein levels, including dN/dS and its 288 extensions. Here we propose pseudogenes as a reference for distinguishing neutral and adaptive 289 forces acting on pangenomes – particularly rare genes. We show that the association between 290 pangenome diversity and synonymous-site variation disappears after correcting for pseudogene 291 diversity with the s<sub>i</sub>/s<sub>p</sub> metric, while the association with dN/dS is maintained. This indicates that 292 a higher proportion of intact singleton genes (relative to singleton pseudogenes) are present when 293 selection is more effective. This is consistent with many rare intact genes, but not all, conferring 294 host-adaptive functions. These genes are more likely to be retained when selection is efficient<sup>3</sup> 295 (such as in *E. coli*), and more likely to degenerate neutrally and become pseudogenes in species 296 with lower Ne (such as obligate intracellular bacteria). Our results could also be explained by 297 widespread slightly deleterious rare pseudogenes, which can be purged only in species with high 298 Ne, but we did not detect a significant association between dN/dS and pseudogene content (and 299 only a weak association with dS), making this less likely. 300

A common explanation for widespread selection on rare accessory genes is adaptation to highly specialized niches<sup>16–18</sup>. While genes recently acquired through horizontal gene transfer are often hypothesised to confer niche-specific adaptations<sup>4</sup>, it is challenging to make highconfidence inferences without knowing the background of all recently transferred genes that were not retained – and are thus, by definition, unobservable. By focusing on pseudogenes, which are observable but likely to evolve mostly by drift, we can establish a (nearly) neutral background against which to discern potentially niche-specific adaptations.

We relied on the assumption that any selection pressures acting upon pseudogenes tend to be of much lower magnitude compared to intact genes. In other words, we assumed that, overall, the pseudogene instances we identified do not reflect adaptive gene loss<sup>19</sup> (which is unlikely to

substantially increase with selection efficacy, as described above), nor do they contain beneficial
regulatory sequences for modulating gene expression levels<sup>20</sup>. This second possibility would be
inconsistent with the positive association we observed between si/sp and selection efficacy.
Instead, our results are consistent with rare pseudogenes evolving under a regime closer to
neutrality relative to rare intact genes.

Our enrichment test results highlight that a significant proportion of rare accessory genes 316 are under selection to be retained. Notably, 19% of ultra-rare intact genes are in COG categories 317 significantly depleted for pseudogenes. We stress that this is a rough approximation and does not 318 imply that precisely 19% of ultra-rare intact genes have adaptive value. We hypothesise that 319 many such genes are under effective purifying selection, while relaxed purifying selection could 320 account for the observed enrichment of transposons among pseudogenes. Similarly, the 321 enrichment of selenocysteine-containing dehydrogenases among pseudogenes could similarly 322 reflect relaxed or sporadic purifying selection on these elements, which is interesting as 323 selenium, selenocysteine's defining component, is sporadically used across the prokaryotic 324 tree<sup>21</sup>. 325

Gene-level selection could also account for certain observations. For instance, the DNA 326 partitioning protein highly enriched in intact ultra-rare genes, COG1192, is a known plasmid-327 encoded element predicted to be involved in plasmid partitioning<sup>22</sup>. There could be an 328 ascertainment bias toward identifying such genes as intact, because were they pseudogenized or 329 lost the entire plasmid might not be transferred to daughter cells. Similar biases could also 330 account for why prophage and plasmid-associated elements in the mobilome more generally are 331 depleted among pseudogenes, although these elements are also more likely to be adaptive to the 332 host genome<sup>23,24</sup>. 333

Pseudogene diversity can be influenced by many factors, including life history. For instance, obligate intracellular bacteria are characterized by widespread degeneration of their genome, followed by streamlining<sup>25</sup>. Depending on a species' stage in this evolutionary process, its genome could be enriched or depleted for pseudogenes relative to other bacteria. This likely accounts for certain s<sub>i</sub>/s<sub>p</sub> outliers, such as the obligate intracellular bacteria *Rickettsia prowazekii*, which had the lowest s<sub>i</sub>/s<sub>p</sub> ratio. Accordingly, our framework could be improved by incorporating per-species parameters of pseudogene gain and loss dynamics.

Last, an important assumption underlying these results is that dN/dS is an accurate proxy 341 for selection efficacy, and thus Ne, while dS does not represent species Ne. This is 342 counterintuitive, as nucleotide diversity at neutral sites is the standard metric for calculating Ne. 343 However, the validity of this metric in prokaryotes has been debated<sup>26,27</sup>. The estimated dS 344 values would only be generalizable to the overall species (i.e., including all unsampled genomes) 345 if this sampling is representative of the actual diversity in nature. For instance, if sequenced 346 genomes are more likely to represent strains adapted to distinct local environments, and be 347 depleted for near-identical strains from the same environment, then this would not be 348 representative. It is highly unlikely, particularly given the low numbers of genomes considered, 349 for the strains considered in these analyses to accurately reflect the strain diversity and 350 population substructure across species. Instead, we believe dS is more appropriately considered a 351 measure of divergence time among the subset of genomes analysed, but not generalizable across 352 the entire species. In contrast, dN/dS across core genes is more appropriate to generalize across a 353 species when calculated based on a subset of genomes, as this is expected to be similar for all 354 pairwise strain comparisons (albeit with variation depending on strain divergence time), 355 regardless of whether the genomes are representative of the overall species' strain diversity. 356 However, dN/dS is also an imperfect measure, particularly because synonymous sites do not 357 evolve completely neutrally<sup>28</sup>. Regardless, it is uncontroversial to assume that selection acting on 358 synonymous sites is, on average, much weaker compared to on non-synonymous sites. 359 Accordingly, although dN/dS may be inappropriate to use for explicitly calculating Ne, the 360 species' relative ranks for this measure will correspond inversely to their relative Ne rank, all else 361 being equal. 362

Despite these caveats, our work highlights the value of using pseudogene diversity as a 363 neutral null<sup>29</sup> for evaluating the evolutionary forces acting upon intact accessory genes. 364 Establishing true neutrality in microbial genomes is challenging<sup>30</sup>, but the clear association we 365 identified between dN/dS and  $s_i/s_p$  suggests that pseudogene presence/absence diversity can 366 provide insight into how rare accessory genes evolve. Crucially, rare genes in nearly all 367 functional categories are less likely to be pseudogenes when there is no redundant gene copy in 368 the same genome, indicating that even very rare accessory genes are commonly under selection 369 to maintain an intact copy in the genome. Using this pseudogene-based comparative approach, 370 we show that a neutral pangenome model can be rejected and identify which types of rare genes, 371

based on their functional annotation and which species encode them, are more likely to beretained.

374

#### 375 Methods

#### 376 Dataset processing – In-depth pangenome analysis

We conducted an analysis of 10 bacterial species with a relatively high number of genomes 377 (ranging from 135-6,916). We selected these species from the set identified for the broad 378 379 pangenome analysis (see below), but that were also represented by > 100 genomes that were not phylogenetically redundant. For these data, we clustered both intact genes and pseudogenes with 380 cd-hit<sup>31</sup> version 4.8.1 with an identity cut-off of 95% over at least 90% of both compared 381 sequences. This clustering was performed on all genes and pseudogenes across all ten species. 382 We assigned clusters to pangenome partitions as described in the main text. Note that we defined 383 the ultra-rare partition to include doubletons, and not only singletons, to account for cases where 384 two highly similar strains are present with the same ultra-rare gene. As there were > 100385 genomes considered for each species within this analysis, doubletons also correspond to highly 386 rare genes. 387

We functionally annotated each resulting cluster with COG IDs and categories<sup>11</sup> using 388 eggNOG-mapper<sup>32</sup> version 2.1.6 (based on eggNOG orthology data<sup>33</sup> version 5.0.2) with 389 DIAMOND<sup>34</sup> version 2.0.14 and these parameter options: --score 60, --pident 40, --query cover 390 20, --subject\_cover 20, --tax\_scope auto, and --target\_orthologs all. This was performed for 391 individual elements separately (i.e. the original sequences rather than the cluster representatives), 392 and for database sequence matches to pseudogene hits. We focused on the database sequence 393 matches for pseudogene hits, as eggNOG-mapper annotates protein sequences, which is 394 problematic for most pseudogenes as the protein-coding information is generally lost. 395 Accordingly, annotating the corresponding database hits per pseudogene is a more reliable way 396 of assigning putative function. We used majority rule of all member sequences per cluster to 397 assign individual COG IDs and categories, and the same approach for assigning functions to 398 individual pseudogene sequences based on database sequence annotations. We assigned COG 399 categories based on a mapping of COG IDs from the COG 2020 database release. This was 400 performed as the raw output COG categories were based on an earlier version of the database 401 that did not include mobilome (category X) annotations. 402

#### Generalized linear mixed models 404 Generalized linear mixed models were fit in R using the glmmTMB<sup>35</sup> package v1.1.5, one for the 405 ultra-rare, other-rare, and shell pangenome partitions, respectively. Only COG-annotated 406 elements were included in these models, excluding those annotated by the (rare) A, B, Y, and Z 407 COG categories only. We used the binomial family and nlminb optimization algorithm with 408 1000 set for both iter.max and eval.max. The full R-style formula for each model was: 409 410 pseudogene ~ COG-category + non-redundant-status + COG-category: non-redundant-status + (1 411 | species) + (1 | COG-category: species) + (1 | non-redundant-status: species) 412 413 In this formula, random effects are specified as those in parentheses including "1|" and 414 interaction terms are indicated with ":". The response was a Boolean variable indicating whether 415 each element is a pseudogene. The COG-category variable is categorical indicating the one-letter 416 COG category code that each element belongs to. In cases where elements were members of 417 multiple categories, duplicate rows were created for each category. The Transcription category 418 (K) was selected as the first level, to be used for the intercept, as it was the most consistently 419 abundant COG category across all three partitions (third in the other-rare and shell, and fourth in 420 ultra-rare). The non-redundant-status variable was a Boolean variable indicating whether each 421 element was not redundant with another intact element of the same COG ID (gene family, not 422 category) in the same genome. This negative formulation of redundancy (i.e. whether an element 423 is not redundant, rather than whether it is redundant) was chosen as most elements were 424 redundant, and so we decided to set the default level in each model (False) to be more 425 representative. The species variable corresponded to the name of the species encoding each 426 element. 427 We also fit simpler models with subsets of these variables and computed Akaike 428 Information Criterion (AIC) values for each model, that allowed us to compare across models 429 and investigate whether more complex models provide significantly more information. We 430 visualized the AICs per model based on normalized scores that transformed the minimum model 431 AIC per partition to be 0 and the maximum model AIC per partition to be 1. 432

To estimate the rough percentage of intact genes in pseudogene-depleted vs. pseudogeneenriched categories, we classified all intact genes by whether they were included in a tested category (such as genes within a certain COG category that were non-redundant) that was significant in each GLMM. We then tallied the numbers of genes categorized as significantly pseudogene-depleted vs. pseudogene-enriched relative to the total number of genes tested. Note that genes found across multiple COG categories were duplicated in the input table for each COG category, and so would contribute multiple times to the tally of total genes.

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

446

## 447 Dataset processing – broad pangenome analysis

We downloaded all genomes used in this study from the Genome Taxonomy Database<sup>12</sup> release 448 202. We identified all species in this database with at least ten high quality genomes, based on 449 these criteria: (1) marked as passing the minimum information about a metagenome-assembled 450 genome<sup>36</sup> check; (2) Check $M^{37}$  completeness > 98% and contamination < 1%; (3) fewer than 451 1000 contigs; (4) contig N50 > 5000; (6) fewer than 100,000 ambiguous bases. We also 452 restricted our analyses to genomes in RefSeq (rather than those in GenBank only), except for 453 Wolbachia pipientis genomes, which were numerous but primarily limited to GenBank. For 454 species with more than twenty genomes, we randomly sampled down to twenty genomes. We 455 identified 670 species that fit these criteria and downloaded the corresponding genomes. Certain 456 genomes had been relabelled or removed from NCBI since the release of Genome Taxonomy 457 Database release 202, which resulted in a minimum of nine genomes per species (we eliminated 458 two species with fewer than nine genomes). We annotated all genomes with Prokka<sup>38</sup> version 459 1.14.5 with the -kingdom, --compliant, and -rfam options. We also specified the --metagenome 460 flag for all genomes with 50 or more contigs. We ran Panaroo<sup>39</sup> version 1.3.0 on all output GFFs, 461 with the -remove-invalid-genes and --clean-mode strict options. We then ran Pseudofinder<sup>40</sup> 462 v1.1.0 on the Prokka-output GenBank files to identify all putative pseudogenes, using protein 463

sequences from the UniRef90 database<sup>41</sup> (UniProt KB release 2022 01) as a reference database. 464 We restricted the output to intergenic pseudogenes specifically, as the other pseudogene types 465 identified by Pseudofinder correspond to divergent intact coding sequences (in length or 466 modularity), which are difficult to interpret as truly degenerating sequences, and could simply 467 represent functionally divergent proteins. We performed three filtering steps on the output 468 intergenic pseudogenes. Specifically, we excluded all (1) pseudogene calls within 500 bp of 469 contig ends, (2) pseudogenes of called length < 100 bp or > 5000 bp, and (3) pseudogenes that 470 471 substantially differed from the mean size of all matching database hits (mean database size – pseudogene size was inclusively required to be between -500 bp and 2000 bp). Pseudogenes 472 were clustered with cd-hit using the same settings as described above. Where possible, these 473 commands were parallelized with GNU Parallel<sup>42</sup> version 20161222. 474

475

## 476 Pangenome metric computation

The mean numbers of singletons (whether of intact genes or pseudogenes) per species were identified after repeated subsampling to nine strains per species and then comparing the overlapping genes/pseudogenes. This procedure was repeated for up to 100 replicates (or until the maximum number of strain combinations was reached) and the number of singletons per genome was computed across all replicates. Note that for a supplementary analysis this subsampling was also conducted for subsamples of three and 20 genomes.

the estimated number of singleton genes for a given species, given a number of

subsampled genomes k,  $U_k$ , is defined as:  $U_k = \frac{\sum_{j=1}^r \left(\frac{\sum_{m=1}^k u_{jm}}{k}\right)}{r}$ , where r is the total number of 484 subsampled replicates and  $u_{im}$  is the number of genome-specific genes found in genome m 485 (based on genomes subsampled in replicate *j*). If there are N genomes in total for a given species, 486 then  $r = \min(100, \binom{N}{2})$ . The mean percentage of intact gene singletons per species can then be 487 calculated as:  $s_i = 100 * \frac{U_k}{C}$ , where G is the mean number of genes per genome (across all N 488 genomes). This same procedure was repeated for pseudogenes, except that the numbers of 489 singleton pseudogenes were computed per subsample replicate, and the mean percentage of 490 pseudogenes per species  $(s_p)$  was calculated based on the average number of pseudogenes per 491 genome. To be clear, this formulation means that the s<sub>i</sub>/s<sub>p</sub> metric corresponds to a comparison of 492

the percentage of singleton intact and pseudogene calls overall per species, rather than of calls
within each individual genome.

495

#### 496 *Evolutionary metric computation*

We performed codon-aware multiple-sequence alignment of all ubiquitous and single-copy genes 497 sequences per-species with muscle<sup>43</sup> version 3.8.1551, based on the HyPhy<sup>44</sup> version 2.5.36 498 codon-aware workflow (https://github.com/veg/hyphy-analyses/tree/master/codon-msa). We then 499 concatenated the core gene alignments per species with a Python script 500 (cat core genome msa.py) and computed pairwise dN/dS and dS for each combination of strain 501 pairs per species with an additional script (mean\_pairwise\_dnds.py). Both scripts, and the bash 502 commands for running the codon-aware alignments, are available in v1.1.0 of this repository: 503 https://github.com/gavinmdouglas/handy\_pop\_gen. The latter script identifies potential non-504 synonymous and synonymous mutation sites between each sequence pair using the NG86 505 approach<sup>45</sup>. We computed the mean values across all pairwise strain comparisons, resulting in a 506 single measure of dN/dS and dS per species. 507

508

#### 509 Linear models

We built linear models using the lm function in R to predict pangenome diversity, based on (per 510 species) either the mean number of genes, the genomic fluidity, s<sub>i</sub>, or s<sub>i</sub>/s<sub>p</sub>. The predictors 511 included dS, dN/dS, and taxonomic class. Classes with  $\leq 5$  member species were collapsed into 512 the "Other" category, which was set as the intercept for the models. One species, *Rickettsia* 513 prowazekii, was excluded from this analysis due to values of zero for si and si/sp. We transformed 514 all continuous variables to be normally distributed, except for the mean number of genes, which 515 was already normally distributed. We performed a square-root transformation of the genomic 516 fluidity, si, si/sp, and dS values. The dN/dS values were especially right skewed and required a 517 negative inverse transformation (-1 \* 1/x), where x is each dN/dS value) to be normalized. We 518 then converted each continuous variable to standardized units, by mean-centring and dividing by 519 the standard deviation. This step means that the model outputs refer to units of standard deviation 520 per variable, which makes it possible to compare the magnitude of coefficients across models 521 with different response variables. 522

#### 524 General analyses

- <sup>525</sup> No tests for statistical power were conducted to determine the sample sizes required for this
- study, but we used genomes from all available species in the Genome Taxonomy Database of
- <sup>527</sup> sufficient quality. All statistical analyses were conducted in R v4.2.2. Figures were generated
- with  $ggplot2^{46}$  v3.4.0, with the exception of the heatmaps, which were created with the
- <sup>529</sup> ComplexHeatmap<sup>47</sup> package v2.14.0.
- 530

## 531 **Data Availability**

- Key data files are openly available on Zenodo<sup>48</sup> (<u>https://doi.org/10.5281/zenodo.7942836</u>). All
- analysed genomes are publicly available as part of NCBI RefSeq/GenBank (with accession IDs
- listed in the Zenodo repository). Additional databases used in this study include the eggNOG 5
- database for eggNOG-mapper (<u>http://eggnog5.embl.de</u>) and UniProt KB release 2022\_01
- 536 (https://www.uniprot.org/release-notes/2022-02-23-release).
- 537

## 538 Code Availability

- 539 The code used for the analyses in this manuscript is openly available GitHub at
- 540 <u>https://github.com/gavinmdouglas/pangenome\_pseudogene\_null.</u>
- 541

## 542 Acknowledgements

- <sup>543</sup> We would like to thank Ford Doolittle for providing motivating ideas, and for advice and
- feedback throughout this project. We would also like to thank Louis-Marie Bobay for reading a
- <sup>545</sup> draft of this manuscript and providing feedback, and Adam Eyre-Walker for providing
- constructive comments. GMD was supported by a Natural Sciences and Engineering Research
- 547 Council of Canada (NSERC) Postdoctoral Fellowship and BJS is supported by an NSERC
- 548 Discovery Grant.
- 549

#### 550 Author Contributions Statement

<sup>551</sup> Both GMD and BJS designed the study and wrote the manuscript. GMD conducted all analyses.

552

## 553 Competing Interests Statement

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

556	<u>Tables</u>
557	Not applicable.
558	
559	<b>Figure Legend</b>

# ends/Captions

Figure 1: Distributions of gene or pseudogene sequence clusters by species and frequency in the 560 pangenome, restricted to clusters that could be COG-annotated. Mixed elements are sequence 561 clusters that include both pseudogenes and intact genes in the same cluster. Percentages 562 correspond to the breakdown per species within a given element type (i.e. intact, mixed, or 563 pseudogene) and raw counts are shown in parentheses. 564

565

Figure 2: Summary of significant coefficients (P < 0.05) in generalized linear mixed model with 566 singleton and doubleton (ultra-rare) element state (intact or pseudogene) as the response. This 567 model was based on 213,912 separate elements. The predictors were each element's annotated 568 COG category, whether the element is redundant with an intact gene of the same COG ID (i.e. 569 gene family, not COG category) in the same genome, and the interaction between these 570 variables. The non-redundant coefficients represent the sum of the overall non-redundant 571 coefficient and the interaction of non-redundancy and each COG category. Bars represent the 572 estimated logit (log-odds) coefficient values: estimates > 0 indicate an increased probability of 573 an element being classified as a pseudogene. Error bars represent one standard error, which is a 574 point estimate per coefficient (rather than reflecting a distribution of coefficients). 575

576

Figure 3: Distributions of singleton-based pangenome diversity and molecular evolution metrics. 577 (a) Mean percentage of intact genes and pseudogenes that are singletons (i.e. genome-specific) 578 per species. The mean percent singletons (for both intact genes and pseudogenes) per species 579 was based on repeated subsampling to nine genomes (for up to 100 replicates). Possible (but 580 non-exhaustive) drivers of higher or lower  $s_i/s_p$  ratios are indicated alongside coloured arrows. 581 Species mentioned in the main text are indicated. (b) Relationship between synonymous 582 substitution rates (dS), a measure of strain divergence, and the ratio of the non-synonymous to 583 synonymous substitution rates (dN/dS), coloured by s<sub>i</sub>/s<sub>p</sub>. Relationship between dS and (c) the 584 mean percent intact singletons and (d) the mean percent pseudogene singletons, shaded by 585

<sup>586</sup> dN/dS. Across all panels, each point represents one of 668 prokaryotic species (>= 9 genomes <sup>587</sup> each). Two-tailed Spearman correlation coefficients and *P*-values are indicated on panels b-d, <sup>588</sup> and correspond to the comparison of the variables shown on the x and y axes. The more exact *P*-<sup>589</sup> values for panels b-d are all  $P < 2.2 \times 10^{-16}$ .

590

Figure 4: Associations between pangenome diversity metrics and estimated efficacy of selection 591 (dN/dS). Each panel presents the association between the ratio of non-synonymous to 592 synonymous substitution rates (dN/dS; across each species' core genome), plotted on a  $log_{10}$ 593 scale, and one of the following measures: (a) the mean number of genes per genome, (b) 594 genomic fluidity, (c) the mean percent of intact singletons, and the percentage of singleton intact 595 genes normalized by the percentage of singleton pseudogenes per species. Each point is one of 596 668 prokaryotic species. The two-tailed Spearman correlation coefficients and P-values are 597 indicated. The more exact *P*-values output for panels b-d are  $P < 2.2 \times 10^{-16}$ ,  $P < 2.2 \times 10^{-16}$ , and 598  $P=5.484 \text{ x } 10^{-9}$ , respectively. 599

600

Figure 5: Spearman's correlations between molecular evolution and pangenome diversity metrics. Each cell represents the correlation coefficient for a pairwise comparison of variables. Coloured cells are significant (P < 0.05), and non-significant cells are dark grey. The left plot includes all species while the right plot provides the results based on a subset of species, with outliers for dS and dN/dS removed. The column dN/dS (dS partial) corresponds to two-tailed partial Spearman correlations between dN/dS and each variable, controlling for dS.

607

**Figure 6**: Weak relationships between the mean percent of each species' genome covered by pseudogenes and the (a) within-species ratio of the non-synonymous to synonymous substitution rates (dN/dS) and (b) the within-species synonymous substitution rate (dS). Each point corresponds to one of 668 species. Values of  $s_i/s_p$  are overlaid on a  $log_{10}$  scale on both panels. The result of two-tailed Spearman correlation tests between the variables plotted on the x and y are indicated (including the two-tailed partial Spearman correlation controlling for dS in panel a).

- 614
- 615

### 616 **References**

617	1.	Innamorati, K. A., Earl, J. P., Aggarwal, S. D., Ehrlich, G. D. & Hiller, N. L. The Bacterial
618		Guide to Designing a Diversified Gene Portfolio. in The Pangenome: Diversity, Dynamics
619		and Evolution of Genomes (eds. Tettelin, H. & Medini, D.) 51-87 (Springer, 2020).
620		doi:10.1007/978-3-030-38281-0_3.
621	2.	Sela, I., Wolf, Y. I. & Koonin, E. V. Theory of prokaryotic genome evolution. Proc. Natl.
622		Acad. Sci. U. S. A. 113, 11399–11407 (2016).
623	3.	Bobay, L. M. & Ochman, H. Factors driving effective population size and pan-genome
624		evolution in bacteria. BMC Evol. Biol. 18, 153 (2018).
625	4.	McInerney, J. O., McNally, A. & O'Connell, M. J. Why prokaryotes have pangenomes. Nat.
626		<i>Microbiol.</i> <b>2</b> , 170402 (2017).
627	5.	Kimura, M. & Crow, J. F. The number of alleles that can be maintained in a finite
628		population. Genetics 49, 725–738 (1964).
629	6.	Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent on effective
630		population size. ISME J. 11, 1719–1721 (2017).
631	7.	Vos, M. & Eyre-Walker, A. Are pangenomes adaptive or not? Nat. Microbiol. 2, 1576
632		(2017).
633	8.	Danneels, B., Pinto-Carbó, M. & Carlier, A. Patterns of nucleotide deletion and insertion
634		inferred from bacterial pseudogenes. Genome Biol. Evol. 10, 1792–1802 (2018).
635	9.	Kuo, C. H. & Ochman, H. The extinction dynamics of bacterial pseudogenes. <i>PLoS Genet.</i> 6,
636		e1001050 (2010).
637	10.	Wolf, Y. I., Makarova, K. S., Lobkovsky, A. E. & Koonin, E. V. Two fundamentally
638		different classes of microbial genes. Nat. Microbiol. 2, 1-6 (2016).

639	11.	Galperin, M. Y. et al. COG database update: focus on microbial diversity, model organisms,
640		and widespread pathogens. Nucleic Acids Res. 49, D274–D281 (2021).
641	12.	Parks, D. H. et al. A standardized bacterial taxonomy based on genome phylogeny
642		substantially revises the tree of life. Nat. Biotechnol. 36, 996–1004 (2018).
643	13.	Kislyuk, A. O., Haegeman, B., Bergman, N. H. & Weitz, J. S. Genomic fluidity: An
644		integrative view of gene diversity within microbial populations. BMC Genomics 12, (2011).
645	14.	Rocha, E. P. C. et al. Comparisons of dN/dS are time dependent for closely related bacterial
646		genomes. J. Theor. Biol. 239, 226-235 (2006).
647	15.	Kryazhimskiy, S. & Plotkin, J. B. The Population Genetics of dN/dS. PLoS Genet. 4,
648		e1000304 (2008).
649	16.	Boucher, Y. et al. Local Mobile Gene Pools Rapidly Cross Species Boundaries To Create
650		Endemicity within Global Vibrio cholerae Populations. mBio 2, e00335-10 (2011).
651	17.	Niehus, R., Mitri, S., Fletcher, A. G. & Foster, K. R. Migration and horizontal gene transfer
652		divide microbial genomes into multiple niches. Nat. Commun. 6, 8924 (2015).
653	18.	Smillie, C. S. et al. Ecology drives a global network of gene exchange connecting the human
654		microbiome. Nature 480, 241–244 (2011).
655	19.	Hottes, A. K. et al. Bacterial Adaptation through Loss of Function. PLoS Genet. 9, e1003617
656		(2013).
657	20.	Oren, Y. et al. Transfer of noncoding DNA drives regulatory rewiring in bacteria. Proc. Natl.
658		Acad. Sci. U. S. A. 111, 16112–16117 (2014).
659	21.	Peng, T., Lin, J., Xu, YZ. & Zhang, Y. Comparative genomics reveals new evolutionary
660		and ecological patterns of selenium utilization in bacteria. ISME J. 10, 2048–2059 (2016).

661	22.	A Schlüter <i>et al.</i> Erythromycin Resistance-Conferring Plasmid pRSB105, Isolated from a
662		Sewage Treatment Plant, Harbors a New Macrolide Resistance Determinant, an Integron-
663		Containing Tn402-Like Element, and a Large Region of Unknown Function. Appl. Environ.
664		<i>Microbiol.</i> <b>73</b> , (2007).
665	23.	Bobay, L. M., Rocha, E. P. C. & Touchon, M. The adaptation of temperate bacteriophages to
666		their host genomes. Mol. Biol. Evol. 30, 737-751 (2013).
667	24.	McKerral, J. C. et al. The Promise and Pitfalls of Prophages. bioRxiv 2023.04.20.537752
668		(2023) doi:10.1101/2023.04.20.537752.
669	25.	Giovannoni, S. J., Cameron Thrash, J. & Temperton, B. Implications of streamlining theory
670		for microbial ecology. ISME J. 8, 1553–1565 (2014).
671	26.	Daubin, V. & Moran, N. A. Comment on 'The Origins of Genome Complexity'. Science
672		<b>306</b> , 978 (2004).
673	27.	Lynch, M. & Conery, J. S. Response to Comment on 'The Origins of Genome Complexity'.
674		Science <b>306</b> , 978 (2004).
675	28.	Sharp, P. M., Bailes, E., Grocock, R. J., Peden, J. F. & Sockett, R. E. Variation in the
676		strength of selected codon usage bias among bacteria. Nucleic Acids Res. 33, 1141-1153
677		(2005).
678	29.	Koonin, E. V. Splendor and misery of adaptation, or the importance of neutral null for
679		understanding evolution. BMC Biol. 14, 114 (2016).
680	30.	Rocha, E. P. C. Neutral Theory, Microbial Practice: Challenges in Bacterial Population
681		Genetics. Mol. Biol. Evol. 35, 1338–1347 (2018).
682	31.	Li, W. & Godzik, A. CD-HIT: A fast program for clustering and comparing large sets of
683		protein or nucleotide sequences. Bioinformatics 22, 1658-1659 (2006).

684	32.	Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J. eggNOG-
685		mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the
686		Metagenomic Scale. Mol. Biol. Evol. 38, 5825–5829 (2021).
687	33.	Huerta-Cepas, J. et al. eggNOG 5.0: A hierarchical, functionally and phylogenetically
688		annotated orthology resource based on 5090 organisms and 2502 viruses. Nucleic Acids Res.
689		<b>47</b> , D309–D314 (2019).
690	34.	Buchfink, B., Reuter, K. & Drost, HG. Sensitive protein alignments at tree-of-life scale
691		using DIAMOND. Nat. Methods 18, 366-368 (2021).
692	35.	Brooks, M., E. et al. glmmTMB Balances Speed and Flexibility Among Packages for Zero-
693		inflated Generalized Linear Mixed Modeling. R J. 9, 378 (2017).
694	36.	Bowers, R. M. et al. Minimum information about a single amplified genome (MISAG) and a
695		metagenome-assembled genome (MIMAG) of bacteria and archaea. Nat. Biotechnol. 35,
696		725–731 (2017).
697	37.	Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
698		assessing the quality of microbial genomes recovered from isolates, single cells, and
699		metagenomes. Genome Res. 25, 1043-55 (2015).
700	38.	Seemann, T. Prokka: rapid prokaryotic genome annotation. <i>Bioinformatics</i> <b>30</b> , 2068–2069
701		(2014).
702	39.	Tonkin-Hill, G. et al. Producing polished prokaryotic pangenomes with the Panaroo pipeline.
703		<i>Genome Biol.</i> <b>21</b> , 180 (2020).
704	40.	Syberg-Olsen, M. J., Garber, A. I., Keeling, P. J., McCutcheon, J. P. & Husnik, F.
705		Pseudofinder: Detection of Pseudogenes in Prokaryotic Genomes. Mol. Biol. Evol. 39,
704		msac153(2022)
/06		liiste 135 (2022).

- 41. The UniProt Consortium. The Universal Protein Resource. *Nucleic Acids Res.* 36, D190–
  D195 (2008).
- 42. Tange, O. GNU Parallel: the command-line power tool. *Login USENIX Mag.* 36, 42–47
  (2011).
- 43. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
  throughput. *Nucleic Acids Res.* 32, 1792–1797 (2004).
- 44. Kosakovsky Pond, S. L. *et al.* HyPhy 2.5—A Customizable Platform for Evolutionary
  Hypothesis Testing Using Phylogenies. *Mol. Biol. Evol.* 37, 295–299 (2020).
- 45. Nei, M. & Gojobori, T. Simple methods for estimating the numbers of synonymous and
- nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.* **3**, 418–426 (1986).
- 46. Wickham, H. ggplot2: Elegant Graphics for Data Analysis. (Springer-Verlag New York,
  2016).
- 47. Gu, Z. Complex heatmap visualization. *iMeta* **1**, e43 (2022).
- 48. Douglas, G. M. & Shapiro, B. J. Data and code for 'Pseudogenes act as a neutral reference
- for detecting selection in prokaryotic pangenomes'. Zenodo repository.
- 722 <u>https://doi.org/10.5281/zenodo.7942836</u>. (2023).