

1 **Homologous recombination substantially delays sequence but not gene content divergence of**  
2 **prokaryotic populations**

3 Jaime Iranzo, Yuri I. Wolf, Eugene V. Koonin\*, Itamar Sela\*

4 National Center for Biotechnology Information, National Library of Medicine, National Institutes of  
5 Health, Bethesda, MD 20894, USA

6 \*For correspondence: [itamar.sela@nih.gov](mailto:itamar.sela@nih.gov); [koonin@ncbi.nlm.nih.gov](mailto:koonin@ncbi.nlm.nih.gov)

7

8

9 **Abstract**

10 Evolution of bacterial and archaeal genomes is a highly dynamic process that involves extensive gain and  
11 loss of genes. Therefore, phylogenetic trees of prokaryotes can be constructed both by the traditional  
12 sequence-based methods (gene trees) and by comparison of gene compositions (genome trees).  
13 Comparing the branch lengths in gene and genome trees with identical topologies for 34 clusters of  
14 closely related bacterial and archaeal genomes, we found that the terminal branches of gene trees were  
15 systematically compressed compared to those of genome trees. Thus, sequence evolution seems to be  
16 significantly delayed with respect to genome evolution by gene gain and loss. The extent of this delay  
17 widely differs among bacterial and archaeal lineages. We develop and explore mathematical models  
18 demonstrating that the delay of sequence divergence can be explained by sequence homogenization  
19 that is caused by homologous recombination. Once evolving genomes become isolated by barriers that  
20 impede homologous recombination, gene and genome evolution processes settle into parallel  
21 trajectories, and genomes diverge, resulting in speciation. This model of prokaryotic genome evolution  
22 gives a mechanistic explanation of our previous finding that archaeal genomes contain a class of genes  
23 that turn over rapidly, before significant sequence divergence occurs, and provides a framework for  
24 correcting phylogenetic trees, to make them consistent with the dynamics of gene turnover.

25

26

## 27 Introduction

28 Evolution of bacterial and archaeal genomes is a highly dynamic process that involves extensive gain and  
29 loss of genes, with turnover rates comparable to if not exceeding the rate of nucleotide substitution<sup>1-3</sup>.  
30 Gene gain and loss occur through insertion and deletion of genome segments of variable size, including  
31 large genomic islands, via mechanisms of non-homologous recombination, often involving mobile  
32 genetic elements<sup>4,5</sup>. The gene gain and loss events can be used to generate “gene content trees” that  
33 reflect the evolution of microbial pangenomes and complement traditional phylogenetic trees  
34 constructed from sequence alignments of highly conserved marker genes. From such trees, a gene  
35 turnover clock can be defined. The gene turnover clock ticks at a rate that does not necessarily correlate  
36 with the rate of the traditional, sequence-based molecular clock.

37 Evolution of prokaryotic populations is strongly affected by homologous recombination, which is  
38 regarded as a major contributor to maintaining genetic cohesion by preventing sequence divergence via  
39 gene conversion<sup>6,7</sup>. However, because barriers to recombination do not necessarily affect the whole  
40 genome, bacterial strains can diverge at some loci while remaining cohesive at others<sup>8</sup>. The rate of  
41 population divergence and, eventually, speciation thus depends on the dynamics of recombination  
42 barrier emergence across genomes.

43 Efficient homologous recombination between two genomes requires the presence of (nearly) identical  
44 nucleotide sequences flanking the exchanged genomic regions. The minimum length of these flanks  
45 depends on the species, with typical values around 25-100 nucleotides<sup>9,10</sup>. As genomes diverge, the  
46 probability to find fully conserved flanking sequences decreases, and so does the efficiency of  
47 recombination<sup>11,12</sup>. The existence of genetic barriers to homologous recombination was initially  
48 observed in experimental studies which have shown that sequence divergence of over 5% can prevent  
49 most recombination events in some bacteria<sup>13-15</sup>. Subsequently, comparative genomic analyses have  
50 confirmed that barriers to homologous recombination are widespread and can be used to define  
51 biological species in bacteria and archaea<sup>16-18</sup>. Mechanistic modeling of the molecular processes involved  
52 in homologous recombination has shown that barriers to recombination can build up spontaneously if  
53 the balance of mutation and recombination favors the sequence variability in the population<sup>19-21</sup>.  
54 Barriers to recombination can also arise after the acquisition of new genes<sup>8,22</sup>, especially if the newly  
55 acquired genes are involved in niche specialization<sup>23,24</sup>. In this case, the barriers seem to result, in part,  
56 from selection against gene conversion events that would lead to the loss of the recently acquired,  
57 beneficial genes.

58 Here, we investigate how the fraction of genes shared by closely related bacterial and archaeal genomes  
59 decays with the phylogenetic distance and show that molecular and gene turnover clocks are  
60 incongruent at short evolutionary times. To elucidate the origin of this discrepancy, we develop a  
61 mathematical model of genome evolution that describes the dynamics of sequence divergence in the  
62 presence of gene conversion. The model predicts the existence of a recombination-driven delay in the  
63 molecular clock, the magnitude of which corresponds to the time required for barriers to recombination  
64 to spread across the genome. By fitting the model to genomic data, we obtain estimates of such  
65 recombination-driven delays in 34 groups of closely related bacteria and archaea, and show that the  
66 incongruence between the molecular and gene turnover clocks disappears if the former is corrected to  
67 account for the estimated delay. Finally, we investigate the tempo and the factors that contribute to the  
68 establishment of barriers to recombination in the populations of diverse bacteria and archaea.

69

70

## 71 **Results**

### 72 Discrepancy between the molecular clock and the gene turnover clock

73 The evolution of gene content in closely related genomes has been recently investigated by  
74 mathematical modeling and comparative genomics<sup>25-28</sup>. These analyses have shown that the fraction of  
75 genes shared by a pair of genomes decays exponentially with time as the genomes diverge, and that  
76 genes can be classified in two categories based on their turnover rates<sup>26</sup>. Here we extend these  
77 approaches to sets of 3 and more genomes. For illustration, let us consider a simple scenario in which all  
78 genes have the same turnover rate  $\lambda$ . The fraction of genes shared by a pair of genomes is then

$$79 \quad I_2 \sim e^{-\lambda D_2} \quad [1]$$

80 where  $D_2$  is the total evolutionary tree distance, which in the case of 2 genomes is equal to twice the  
81 distance from the last common ancestor. We show in the Methods that a similar formula describes the  
82 divergence in gene content for groups of 3 or more genomes. Specifically, the fraction of genes shared  
83 by  $k$  genomes decays as

$$84 \quad I_k \sim e^{-\lambda D_k} \quad [2]$$

85 where  $D_k$  is the total evolutionary distance spanned by those  $k$  genomes. Given a phylogenetic tree, the  
86 total distance  $D_k$  is the sum of branch lengths of the subtree that includes the  $k$  genomes. The most  
87 notable aspect of this result is that the dynamics of gene content divergence is independent of the  
88 number of genomes considered ( $k$ ). As a result, plots of the fraction of shared genes ( $I_k$ ) as a function of  
89 the total evolutionary distance ( $D_k$ ) for different sample sizes collapse into a single curve (Fig. 1a). This  
90 property remains valid under very general models of gene turnover, including the case where the rates  
91 of gene gain and loss differ across gene families (see Methods), under the condition that the tree  
92 distance is proportional to the rate of the gene turnover clock.

93 To test this theoretical prediction, we analyzed the profiles of gene sharing in 34 groups of closely  
94 related genomes from Bacteria and Archaea. As a proxy for the evolutionary time, we used the branch  
95 lengths of high-resolution sequence similarity phylogenetic trees built from concatenated alignments of  
96 single-copy core genes. Then, we sampled subsets of genomes and represented the fraction of shared  
97 genes as a function of the total tree distance. At odds with the theoretical expectation, we found that  
98 gene-sharing decay curves depend on the number of sampled genomes: as more genomes are added,  
99 the curves shift down and the fraction of shared genes becomes smaller than expected (left panel on  
100 Fig. 1c shows a representative case; see also Supplementary Fig. S1). As illustrated by Fig. 1b, such a  
101 non-overlapping pattern could be easily explained if the lengths of the terminal branches in the  
102 phylogenetic trees were systematically underestimated. In more general terms, the curves in Fig. 1c  
103 involve two different clocks: the molecular clock, that is used to infer the branch lengths in the  
104 phylogenetic tree; and the gene turnover clock, that governs the stochastic process of gene loss and,  
105 with it, the decay in the fraction of shared genes. Thus, the absence of overlap among gene sharing  
106 curves reflects a non-linear relationship between these two clocks, or more specifically, a delay in the  
107 molecular clock relative to the gene turnover clock.

108 To further explore the differences between the molecular clock and the gene turnover clock, we  
109 generated an alternative set of phylogenetic trees by rescaling branch lengths so that they are  
110 proportional to the number of gene gains and losses occurred during the evolution of a lineage (Fig.  
111 2a,b; see Methods for details). Such “gene content trees” are topologically congruent with the gene  
112 turnover clock, as long as the genomes are close enough to assume that gene turnover rates have  
113 remained approximately constant since the last common ancestor (this approach does not require,  
114 however, that gain and loss rates are homogeneous across genes). The non-linear relationship between  
115 the molecular clock and the gene turnover clock becomes evident when comparing pairwise distances  
116 among leaves in both sets of trees (Fig. 2c). Sequence divergence only starts building up after a transient  
117 period during which genes are gained and lost, although the duration of this transient phase varies  
118 across taxa.

119

### 120 A recombination barrier model of genome divergence explains the delay in the molecular clock

121 The observation that early divergence of strains proceeds through a transient stage in which  
122 substitutions (effectively) do not accumulate motivated us to study the dynamics of sequence  
123 divergence in the presence of homologous recombination. Rather than focusing on the mechanistic  
124 details of recombination, we formulated a phenomenological model for the fraction of loci within a  
125 genome that become isolated with respect to another genome from the same original population (see  
126 Methods). Loci that are not affected by barriers to recombination experience periodical gene conversion  
127 that reverts them to the “population average”<sup>21,29</sup>. Our model captures this fact by assuming that  
128 recombining loci provide a negligible contribution to the genome-wide sequence divergence. However,  
129 as soon as barriers to recombination are established, the affected loci start diverging from the ancestral  
130 population at a rate that is proportional to the substitution rate. As a result, the overall sequence  
131 divergence at a given time results from the contribution of all loci that are isolated by barriers to  
132 recombination, weighted by the time elapsed since the respective locus crossed the barrier.

133 We show in the Methods that homologous recombination within a population induces a delay in the  
134 molecular clock, such that the average divergence of a genome with respect to a member of the same  
135 population grows as

$$136 \quad \Delta(t) = 2\mu(t - \tau(t)) \quad [3]$$

137 where  $t$  is the time since the last common ancestor. Mathematically, the delay  $\tau(t)$  is a concave and  
138 saturating function, the detailed form of which depends on the dynamics of the evolution of  
139 recombination barriers (for exact expressions for some simple scenarios, see Supplementary  
140 Information). For sufficiently long times, this term reaches a constant value  $\tau_\infty$ , which is the long-term  
141 evolutionary delay of the molecular clock induced by homologous recombination. The delay in the  
142 molecular clock accounts for the amount of unobserved variation that is erased by gene conversion  
143 during the early phases of divergence from the ancestral population.

144 A notable consequence of the delay in the molecular clock is that the terminal branches of phylogenetic  
145 trees inferred from sequence analysis appear shorter than expected given the actual evolutionary times.  
146 Specifically, terminal branches are shortened by a distance  $\mu\tau(t_A)$ , whereas internal branches are  
147 shortened by a distance  $\mu\tau(t_A) - \mu\tau(t_B)$ , where  $t_A$  and  $t_B$  are consecutive branching times measured

148 from the tips towards the root. Because both  $\tau(t_A)$  and  $\tau(t_B)$  tend to the same value  $\tau_\infty$  as time passes,  
149 the recombination-driven delays cancel out at long evolutionary times and deep internal branches  
150 remain approximately unchanged.

151 To show that recombination-driven delays are the plausible cause for the lack of linearity between the  
152 molecular and gene turnover clocks, we used the recombination barrier model to correct the branches  
153 of sequence-based trees, accounting for the effects of recombination (see Methods). Then, we sampled  
154 subsets of genomes and reassessed the divergence in gene content as a function of the corrected tree  
155 distances. Remarkably, correction of the sequence trees led to gene-sharing decay curves that do not  
156 depend on the sample size, as predicted by the theory (right panel on Fig. 1c shows a representative  
157 case). When extending the same approach to all groups of genomes using taxa-specific delays (see  
158 below) we found that the mean separation among the curves obtained for different sample sizes  
159 decreased by 30% (permutation test  $p = 0.012$ ; Supplementary Fig. S2). These results are compatible  
160 with the existence of a recombination-driven delay in the molecular clock, which causes a systematic  
161 shortening of the terminal branches of phylogenetic trees built from sequence alignments.

162

### 163 The dynamics of escape from recombination

164 We leveraged the differences between the substitution and gene turnover clocks to obtain quantitative  
165 estimates of the recombination-driven delays in different taxa. Starting from (uncorrected) sequence  
166 similarity trees (Fig. 2a) and gene content trees (Fig. 2b), we retrieved all pairwise distances among  
167 leaves and plotted the distances from the sequence similarity tree against those from the gene content  
168 tree (Fig. 2c). The recombination-driven delays were obtained by fitting the recombination barrier  
169 model to such plots. To better understand how barriers to recombination spread along the genome, we  
170 evaluated several scenarios for the temporal establishment of such barriers (see Methods). We found  
171 that the data in 18 out of the 34 studied groups are best explained by an “autocatalytic” scenario, in  
172 which the rate at which barriers spread accelerates as more and more loci become isolated  
173 (Supplementary Table S1). The autocatalytic scenario also provides good fits in 13 more groups,  
174 although the inclusion of an extra parameter required in this scenario is not statistically justified if the  
175 delays are very small. Under the autocatalytic scenario, the fraction of sites susceptible to homologous  
176 recombination follows a sigmoidal curve in time, with a relatively sharp transition (with few exceptions)  
177 from a state of fully recombining genomes to a state in which all sites freely diverge.

178 The estimates of the long-term evolutionary delay  $\tau_\infty$  range from less than 0.01 to more than 0.5  
179 underreported substitutions per site, with broad variation among prokaryotic groups and sometimes  
180 even within the same genus (Fig. 3 and Supplementary Table S1). In approximately half of the groups,  
181 molecular evolution appears to be strongly delayed, possibly by the pull of recombination, as indicated  
182 by the fact that  $\tau_\infty$  is larger than the depth of the sequence similarity tree. In these cases, there are few  
183 pairs of genomes that have reached the regime of free (linear) divergence, and the estimation of the  
184 upper 95% confidence bound for the long-term evolutionary delay becomes unfeasible. The 5 groups of  
185 Firmicutes included in the analysis (covering bacilli, clostridia and streptococci) belong to this category.  
186 In contrast, representatives from the genus *Pseudomonas* are characterized by a linear divergence  
187 regime, with little or no signs of recombination-driven delay.

188 The magnitude of the recombination-driven delay is tightly linked to the time frame over which  
189 sequence divergence takes place. In taxa with little or no delay, variations in the time at which different  
190 genes cross the recombination barrier are negligible; from the perspective of divergence times, all genes  
191 in these taxa start diverging roughly at the same time (Fig. 4a, left). Conversely, in taxa with lengthy  
192 delays, between-gene variations in divergence times can be comparable to the total evolutionary depth  
193 of the taxon (Fig. 4a, right). Accordingly, it can be expected that differences in sequence divergence  
194 across genes will be larger in taxa with long recombination-driven delays. To test whether genomic data  
195 are compatible with this prediction, we first calculated, for a set of 100 nearly universal gene families,  
196 the gene family- and taxa-corrected evolutionary rates (i.e. the gene-specific substitution rates  
197 normalized by the gene family- and taxon-averaged substitution rates). It can be shown that the  
198 standard deviation of evolutionary rates corrected in this way is equal to the coefficient of variation of  
199 the times over which genes have been diverging (see Supplementary Information). As predicted by the  
200 model, there is a significant positive correlation between the recombination-driven delays and the  
201 variance of evolutionary rates ( $R = 0.64$ ,  $p < 0.001$ ). Moreover, when comparing taxa with long and short  
202 delays (relative to the evolutionary depth), we found that variance of evolutionary rates is significantly  
203 greater in the taxa with long delays (Fig. 4b;  $p = 0.002$ , Student's T-test).

204 To elucidate potential causes for the large variation in the recombination-driven delays found in the  
205 data, we searched for associations with genomic and ecological features, such as genome size, number  
206 of mobile genetic elements, gene turnover rate, lifestyle (free-living or host-associated), natural  
207 competence for transformation, and effective population size (Supplementary Fig. S3). Among those, we  
208 only found a strong negative correlation between the recombination-driven delay and the relative rate  
209 of gene turnover with respect to substitutions (Fig. 4c, Spearman's  $\rho = -0.86$ ,  $p < 0.001$ ), which suggests  
210 that fast gene turnover facilitates the spread of barriers to recombination. The statistical power for the  
211 analysis of ecological traits (lifestyle and effective population size) was low, and therefore, it cannot be  
212 ruled out that these factors contribute to the spread of recombination barriers as well.

213

## 214 Discussion

215 Evolution of prokaryotes occurs at two levels: at the sequence level, through substitutions and small  
216 indels, and at the genome level, through the transfer and loss of genes and groups of genes<sup>4,30-32</sup>.  
217 Whereas evolution at the sequence level is traditionally used to determine phylogenetic relationships  
218 among prokaryotes and to assign new genomes to taxonomic groups, it is the gene repertoire (and  
219 therefore evolution at the genome level) which determines metabolic capacities, ecological properties  
220 and pathogenicity of bacterial strains<sup>33-35</sup>. By studying sequence and gene content divergence in closely  
221 related groups of prokaryotes, we found that sequence evolution is often delayed with respect to  
222 genome evolution, although the magnitude of the delay broadly varies across bacterial and archaeal  
223 lineages. We show that the delay in sequence evolution is likely to result from gene conversion that  
224 homogenizes the core genome while, at the same time, accessory genes can be gained and lost. These  
225 results are fully compatible with our previous findings indicating that archaeal genomes contain a subset  
226 of genes that turn over extremely rapidly, before detectable sequence divergence occurs<sup>26</sup>. The delay  
227 on sequence divergence caused by homologous recombination provides a mechanistic explanation for  
228 the previously observed "instantaneous" gene turnover in prokaryotes<sup>36</sup>. These results are also

229 compatible with the observation that genes are gained and lost at higher rates on the tips of  
230 phylogenetic trees<sup>2</sup>.

231 Notwithstanding the long-standing debate on the applicability of the species concept in prokaryotes,  
232 several recent studies strongly suggest that speciation does occur in bacteria and is a crucial factor  
233 shaping the earth microbiome<sup>16,37</sup>. The establishment of barriers to recombination is a pivotal step in the  
234 early divergence of closely related prokaryotic strains that eventually leads to speciation<sup>19,23</sup>. In the  
235 absence of such barriers, sequence divergence is prevented by the cohesive effect of intra-strain  
236 homologous recombination; only after the barriers arise and recombination ceases, substitutions start  
237 to accumulate at an appreciable rate. By modeling sequence evolution in the presence of gene  
238 conversion, we show here that the temporal dynamics for the spread of barriers to recombination  
239 directly affects the molecular clock. Specifically, homologous recombination sets back the molecular  
240 clock by an amount of time that, in the long-term, equals the average waiting time for the establishment  
241 of recombination barriers. Homologous recombination during early divergence also leads to the  
242 compression of the tips of phylogenetic trees. Our model provides a framework for correcting such  
243 trees, to make them consistent with the dynamics of gene turnover.

244 The causes that underlie the spread of barriers to recombination are complex and likely involve genomic  
245 and ecological factors. Theoretical models show that barriers can simply emerge if mutations generate  
246 diversity faster than recombination erodes it. However, it is a matter of debate whether mutation and  
247 recombination alone can explain the formation of non-recombining species in nature<sup>20,21</sup>. Our finding  
248 that recombination-driven delays negatively correlate with the rate of gene turnover supports the  
249 hypothesis that gene gain and loss facilitates the establishment of barriers to homologous  
250 recombination, by promoting niche differentiation<sup>18</sup> and/or by interfering with gene conversion at  
251 flanking loci<sup>22-24</sup>. Moreover, the autocatalytic (or sigmoidal) dynamics that best describes the fraction of  
252 recombination-free loci in our model is consistent with the previous findings indicating that barriers to  
253 recombination are initiated by the acquisition of lineage-specific genes and subsequently spread from  
254 the vicinity of those genes towards the rest of the genome<sup>8</sup>.

255 Recombination-driven delays are indicative of the time frame over which different genes start to diverge  
256 during the split of prokaryotic lineages. Lineages with short delays are characterized by low variance in  
257 the gene-specific (gene family- and taxon-corrected) evolutionary rates, which implies that most of their  
258 genes began diverging over a brief period of time. In contrast, the higher variances observed in lineages  
259 with long delays imply that, in those lineages, divergence of genes occurred in an asynchronous way and  
260 spanned longer periods of time. To illustrate this point, it has been estimated that the split between  
261 *Escherichia* and *Salmonella* took place around 140 million years ago and spanned over 70 million years<sup>8</sup>,  
262 which is in close agreement with our finding that the delay in the *Escherichia/Salmonella* group is  
263 approximately half of its evolutionary depth. The high variability of the sequence divergence delays  
264 across taxa implies that the time required for lineages to split and form new species is strongly taxon-  
265 dependent. For taxa with long delays, it appears natural to think of speciation events being driven by the  
266 emergence of strong recombination barriers, whereas taxa with short delays would follow a more  
267 continuous speciation dynamics, with new species forming as the genomes gradually diverge<sup>20,21</sup>.  
268 Although our analysis focused on prokaryotic genomes, the conclusions appear to be general and can be  
269 extended to the evolution of other genomes that are substantially affected by horizontal gene transfer  
270 and homologous recombination, for example, viruses.



271

## 272 **Methods**

### 273 Recombination-driven delay in sequence divergence

274 We adopted a phenomenological approach to modeling the effect of intra-population homologous  
275 recombination on the genetic diversity of a population. Under this approach, each genomic region is  
276 either subject to periodical recombination with other members of the population or has reached a point  
277 where homologous recombination is not possible anymore. In the latter case, the respective genomic  
278 region is assumed to have crossed (diverged beyond) the recombination barrier. There was no attempt  
279 to explicitly model the molecular mechanisms that underlie barriers to recombination, and therefore,  
280 the present model is, in principle, applicable to scenarios in which ecological factors, rather than  
281 accumulated sequence divergence, drive the isolation of evolving populations. Instead, given a pair of  
282 genomes, we used a general function  $f(t)$  to describe the fraction of each genome that is subject to  
283 recombination. The probability that a region of the genome crosses the recombination barrier exactly at  
284 time  $t$  (where the time starts counting at the last common ancestor of the two genomes) is

$$285 \quad P(t) = -df/dt \quad [4]$$

286 Considering that only the regions that have crossed the recombination barrier make a long-term  
287 contribution to sequence divergence, and that the number of substitutions in those regions is  
288 proportional to the time elapsed since the recombination barrier was crossed, the overall sequence  
289 divergence between a pair of genomes becomes

$$290 \quad \Delta_2(t) = 2\mu \int_0^t (t-u)P(u)du \quad [5]$$

291 where  $\mu$  is the average substitution rate. Integration by parts leads to the final result

$$292 \quad \Delta_2(t) = 2\mu(t - \tau(t)) \quad [6]$$

293 where  $\tau(t) = \int_0^t f(u) du$  represents a recombination-driven delay in the molecular clock.

294 To contrast the model with genomic data, we explored three more specific scenarios by assigning  
295 explicit functional forms to the rate at which regions of the genome cross the recombination barrier.  
296 Thus, we considered a power law time-dependency of the rate  $R(t) = \lambda t^\gamma$  (which includes the case of a  
297 constant rate); a linear time-dependency plus a constant term  $R(t) = \lambda_0 + \lambda_1 t$ ; and an “autocatalytic”  
298 scenario  $R(f) = \lambda_0 - \lambda_1 f$  in which the rate increases as more and more regions of the genome cross  
299 the recombination barrier. Given a functional expression for  $R(t, f)$ , the fraction of the genome subject  
300 to recombination is obtained by solving the differential equation  $df/dt = R(f, t)f$  (see Supplementary  
301 Information for further details).

### 302 Model of genome content divergence

303 The number of genes,  $x$ , in a prokaryotic genome was modeled as a stochastic birth-death process, with  
304 genome-wide gain rate  $P^+$  and loss rate  $P^-$ <sup>25</sup>. Under this model, the number of genes in a genome is  
305 described by the differential equation

$$306 \quad dx/dt = P^+ - P^- \quad [7]$$

307 To facilitate comparison of gene content across genomes, each genome was represented by a vector  $\mathbf{X}$   
 308 with elements that assume values of 1 or 0. Each entry represents a gene, *i.e.* an ATGC-COG, where 1 or  
 309 0 indicate the presence or absence, respectively, of that ATGC-COG in the genome. Genome size  $x$  is  
 310 then given by the sum of all elements in  $\mathbf{X}$ .

311 The number of shared genes in a pair of genomes, or pairwise intersection,  $I_2$  is defined as

$$312 \quad I_2(t) = \langle \mathbf{X}_1 \cdot \mathbf{X}_2 \rangle \quad [8]$$

313 where  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are vectors that represent the two genomes, the angled brackets indicate averaging  
 314 over multiple realizations of the stochastic process, and the dot operation stands for a scalar product.  
 315 The dynamics of pairwise intersections is given by

$$316 \quad dI_2/dt = 2\langle (d\mathbf{X}_1/dt) \cdot \mathbf{X}_2 \rangle \quad [9]$$

317 where we used the fact that both averages are equal  $\langle (d\mathbf{X}_1/dt) \cdot \mathbf{X}_2 \rangle = \langle \mathbf{X}_1 \cdot (d\mathbf{X}_2/dt) \rangle$ . Assuming that  
 318 genes are acquired from an (effectively) infinite gene pool<sup>38</sup>, we have

$$319 \quad \langle (d\mathbf{X}_1/dt) \cdot \mathbf{X}_2 \rangle = -P^- \cdot I_2(t)/x \quad [10]$$

320 and substituting this relation into the differential equation for pairwise intersections of Eq.(9), we  
 321 obtain

$$322 \quad dI_2/dt = -2P^- \cdot I_2(t)/x \quad [11]$$

323 The solution of this equation shows that pairwise genome intersections decay exponentially as

$$324 \quad I_2(t) = I_2(0) e^{-vt} \quad [12]$$

325 with decay constant  $v = 2P^-/x$ . Assuming a molecular clock, the time  $t$  can be translated into tree  
 326 pairwise distance as  $D_2 = 2t/t_0$  and the pairwise similarity decays exponentially with the tree distance  
 327  $D_2$  as

$$328 \quad I_2(d) = x e^{-\lambda D_2} \quad [13]$$

329 with decay constant  $\lambda = t_0 P^-/x$ . Note that the ratio  $P^-/x$  gives the per-gene loss rate.

330 The intersection of  $k$  genomes,  $I_k$ , is the number of orthologous genes that are shared by all  $k$  genomes.  
 331 It is formally defined as

$$332 \quad I_k = \langle \text{intersect}(\mathbf{X}_1 \dots \mathbf{X}_k) \rangle \quad [14]$$

333 Similar to pairwise genome intersections, the time derivative of  $k$ -intersections is given by

$$334 \quad dI_k/dt = k\langle \text{intersect}(d\mathbf{X}_1/dt \dots \mathbf{X}_k) \rangle = -k I_k(t)P^-/x \quad [15]$$

335 Solving the differential equation, we obtain an exponential decay for the  $k$ -intersections

$$336 \quad I_k(t) = I_k(0) e^{-kP^-t/x} \quad [16]$$

337 If time is inferred from a tree

$$338 \quad I_k(D_k) = x e^{-\lambda \cdot D_k} \quad [17]$$

339 where  $\lambda$  is proportional to  $P_-/x$  and  $D_k$  is the sum of branch lengths in the tree that encompasses the  $k$   
340 genomes (see section 2 of the Supplementary Information for formal derivation). This expression can be  
341 extended to genomes composed of fast and slow evolving genes, and becomes

$$342 \quad I_k(D_k) = x_1 \cdot e^{-\lambda_1 \cdot D_k} + x_2 \cdot e^{-\lambda_2 \cdot D_k} \quad [18]$$

343 where  $x_1$  and  $x_2$  are the average numbers of genes of each class.

#### 344 Genomic data

345 We used the Alignable Tight Genomic Clusters (ATGC) database<sup>39</sup> to define groups of closely related  
346 bacterial and archaeal genomes. By construction, ATGCs are independent of taxonomic affiliation and  
347 meet the objective criteria of high synteny and low divergence (synonymous substitution rate  $dS < 1.5$  in  
348 protein-coding genes). We selected 36 ATGCs that matched the following criteria: i) maximum pairwise  
349 tree distance is at least 0.1 substitutions per site, and ii) the phylogenetic tree contains more than two  
350 clades, such that pairwise tree distances are centered around more than two typical values. Two of the  
351 36 genome clusters were identified as outliers and were excluded from the dataset. The 34 genome  
352 clusters analyzed in this study are listed in Supplementary Table S1. To facilitate computational analysis,  
353 we subsampled large ATGCs to keep at most 20 representative genomes per ATGC. ATGC-specific  
354 Clusters of Orthologous Genes (ATGC-COGs) were downloaded from the ATGC database and  
355 postprocessed to obtain finer grain gene families by reconstructing approximate phylogenetic trees  
356 from original ATGC COG alignments and splitting them into subtrees with minimum paralogy. ATGC-  
357 specific phyletic profiles were built by registering, as a binary matrix, the presence or absence of each  
358 ATGC-COG in each genome within the ATGC (multiple genes from a single genome that belong to the  
359 same ATGC-COG were counted once).

#### 360 Tree construction

361 High-resolution phylogenetic trees based on the concatenated alignments of single-copy core genes  
362 were downloaded from the ATGC database<sup>39</sup>. We refer to these trees as sequence similarity-based  
363 trees. The phylogenomic reconstruction software Gloome<sup>40</sup> was used to obtain trees based on the gene  
364 content similarity among the members of each ATGC. As the input for Gloome, we used the phyletic  
365 profiles for the presence or absence of each ATGC-COG, and the sequence similarity-based trees from  
366 the ATGC database; options were set to optimize the tree branch lengths under a genome evolution  
367 model with 4 categories of gamma-distributed gain and loss rates. This procedure resulted in 2 trees per  
368 ATGC, both with the same topology but with different branch lengths (one based on sequence  
369 divergence, the other on gene content divergence). All trees were inspected for extremely long and  
370 short branches, and clades responsible for such branches (typically 1 or 2 genomes in 5 out of the 34  
371 ATGCs) were manually removed to avoid possible artifacts in the following steps.

#### 372 Tree comparison and model fitting

373 For each ATGC, we computed all pairwise distances among tree leaves in the sequence similarity- and  
374 gene content-based trees. Then, we compared the observed relationship between both sets of distances  
375 with the expectations of the recombination barrier model under four scenarios for the recombination  
376 barrier crossing rate (constant, linearly increasing with time, linearly increasing plus a constant term,  
377 and proportional to the fraction of the genome that has already crossed, which leads to an  
378 “autocatalytically” accelerated crossing rate). For each scenario, we fitted the model parameters using

379 non-linear least-squares optimization (implemented in Matlab R2018b), with sequence similarity-based  
 380 tree distances as independent variable and gene content-tree distances as dependent variable. The  
 381 choice of sequence similarity-based tree distances as the independent variable was motivated by the  
 382 need to fulfil the assumption of homoscedasticity in the non-linear regression model. Additionally, we  
 383 studied the fit of a heuristic power law model  $y = bx^\alpha$ . To compare the goodness of fit provided by  
 384 different models, we calculated the Akaike Information Criterion (AIC) as  $AIC = 2k +$   
 385  $n(\ln(2\pi RSS/n) + 1)$ , where  $k$  is the number of parameters,  $n$  is the number of observations, and  $RSS$   
 386 is the residual sum of squares<sup>41</sup>. The 95% confidence interval for the delay parameter  $\tau_\infty$  was obtained  
 387 by finding the values of  $\tau_\infty$  such that the residual sum of squares becomes  $RSS = RSS^* (1 +$   
 388  $1.96^2/(n - 1))$ , where  $RSS^*$  is the residual sum of squares produced by the best fit of the model<sup>42</sup>.

### 389 Correction of tree branch lengths

390 To account for unobserved variation, tree branch lengths were corrected by using the autocatalytic  
 391 barrier spread model with the parameters inferred in the previous step. In that model, the observable  
 392 divergence increases with time according to the function

$$393 \quad f(t) = t - \tau_\infty(1 - \ln(1 + e^{\phi(1-\xi t/\tau_\infty)})/\ln(1 + e^\phi)) \quad [19]$$

394 with

$$395 \quad \xi = (1 + e^{-\phi}) \ln(1 + e^\phi) / \phi \quad [20]$$

396 In the simplest case of an ultrametric tree, the corrected height (the distance from the tip) of a node  $i$ ,  
 397  $\tilde{h}_i$ , can be calculated by applying the inverse function to the original height  $h_i$ , that is

$$398 \quad \tilde{h}_i = f^{-1}(h_i) \quad [21]$$

399 Branch lengths would then be obtained by subtracting the depths of the parent and child nodes.  
 400 Extending this idea to non-ultrametric trees, we first defined the parental height of a node,  $H_i$ , as the  
 401 distance between the node's parent and the tip, that is,  $H_i = h_i + b_i$ , where  $b_i$  is the length of the  
 402 branch that connects node  $i$  to its parent node. Parental heights were calculated as weighted averages  
 403 from the tips to the root, such that  $H_i = b_i$  for the leaves and

$$404 \quad H_i = \frac{w_{i1}H_{i1} + w_{i2}H_{i2}}{w_{i1} + w_{i2}} + b_i \quad [22]$$

405 for internal nodes. Subindices  $i1$  and  $i2$  refer to the child nodes of node  $i$ ; weights were computed as  
 406  $w_i = b_i$  for leaves and  $w_i = w_{i1} + w_{i2} + b_i$  for internal nodes (thus, the weight of a node is equal to the  
 407 total length of the subtree that contains that node as its root plus the length of the branch that connects  
 408 it to its parent node). Next, corrected values of the parental heights were obtained by computing,  
 409 numerically (MATLAB R2018b), the inverse function  $\tilde{H}_i = f^{-1}(H_i)$ . The corrected branch lengths for the  
 410 tips are simply  $\tilde{b}_i = \tilde{H}_i$ . Finally, we proceeded from tips to root and obtained the corrected branch  
 411 lengths associated with internal nodes as

$$412 \quad \tilde{b}_i = \tilde{H}_i - \frac{\tilde{w}_{i1}\tilde{H}_{i1} + \tilde{w}_{i2}\tilde{H}_{i2}}{\tilde{w}_{i1} + \tilde{w}_{i2}} \quad [23]$$

413 where the new weights  $\tilde{w}_{i1}$  and  $\tilde{w}_{i2}$  were recalculated at each step using the corrected branch lengths.

### 414 Analysis of gene content

415 For a set of  $k$  genomes that belong to the same ATGC, the gene content overlap was calculated as the  
416 number of ATGC-COGs shared by all the genomes divided by the mean number of ATGC-COGs per  
417 genome. The total sequence divergence of a set of  $k$  genomes was calculated as the sum of all branch  
418 lengths in the sequence similarity subtree that results from selecting the corresponding leaves in the  
419 whole-ATGC tree. Curves for the temporal decay of the fraction of shared genes were obtained by  
420 plotting the gene content overlap against the total sequence divergence for all possible combinations of  
421  $k$  genomes within the ATGC. Smooth curves were obtained by fitting a cubic spline model with 5 knots  
422 (placed in both extremes and in the 25-, 50-, and 75-percentiles of the data x-values) using the SLM tool  
423 (D'Errico, August 10, 2017(<http://www.mathworks.com/matlabcentral/fileexchange/24443>)) in MATLAB  
424 R2018b. The mean separation among the curves obtained for different values of  $k$  was calculated as

$$425 \quad S = \left( \frac{\sum_{k=2}^{n-2} \sum_{k'=k+1}^{n-1} \int_{a_{kk'}}^{b_{kk'}} (f_k(x) - f_{k'}(x))^2 dx}{\sum_{k=2}^{n-2} \sum_{k'=k+1}^{n-1} (b_{kk'} - a_{kk'})} \right)^{1/2} \quad [24]$$

426 where  $n$  is the number of genomes in the ATGC,  $f_k$  and  $f_{k'}$  are the curves that result from fitting the  
427 spline model to sets of  $k$  and  $k'$  genomes, respectively, and  $a_{kk'}$  and  $b_{kk'}$  are the bounds of the x-axis  
428 interval in which both  $f_k$  and  $f_{k'}$  are defined<sup>43</sup>. A value of  $S = 0$  corresponds to the exact coincidence of  
429 the curves for all values of  $k$ , which is expected in the absence of evolutionary delays. To assess whether  
430 correction of branch lengths for unobserved variation reduces the separation among curves, we  
431 calculated the relative change in separation as  $(S_{original} - S_{corrected}) / \max(S_{original}, S_{corrected})$ ,  
432 which takes the value of 1 when correction leads to complete collapse, positive values  $\leq 1$  when  
433 correction reduces separation, and negative values  $\geq -1$  when correction increases separation among  
434 curves. The statistical significance of the relative change in separation was assessed using a permutation  
435 test that involved calculating the median relative change across ATGCs and comparing it with  $10^6$   
436 randomized datasets in which the “original” and “corrected” labels were randomly reassigned in each  
437 ATGC.

#### 438 Quantification of evolutionary rates

439 For the analysis of evolutionary rates, we focused on a previously published list of 100 nearly universal  
440 gene families<sup>44</sup>, defined as clusters of orthologous genes or COGs<sup>45</sup>. To minimize possible confounding  
441 effects due to paralogy, we identified all the ATGC-COGs that match any of the universal COGs and  
442 restricted the analysis to those COGs that are represented by a single ATGC-COG in at least 30 of the 34  
443 analyzed ATGCs. Multiple sequence alignments for the selected ATGC-COGs were downloaded from the  
444 ATGC database and processed to extract all pairwise distances (Nei-Tamura method). Only index  
445 orthologs from the ATGC database, i.e. a single sequence per ATGC-COG per genome, were included. For  
446 each ATGC-COG, pairwise distances between sequences were plotted against pairwise distances in the  
447 phylogenetic tree, and a linear regression model with zero intercept was applied to obtain the relative  
448 evolutionary rate of the ATGC-COG with respect to the ATGC average. To minimize the impact of rare  
449 instances of gene replacement, which manifest as a non-linear relationship between sequence and tree  
450 pairwise distances, we discarded the ATGC-COGs with the fit to the regression model  $R^2 < 0.9$ . To  
451 account for COG-specific evolutionary rates, the relative evolutionary rate of each ATGC-COG was  
452 divided by the mean of all ATGC-COGs that match the same COG. The result is the ATGC-COG residual  
453 evolutionary rate, that is, the ATGC-COG evolutionary rate corrected by COG- and ATGC-wise averages.  
454 For each ATGC, the dispersion of evolutionary rates was quantified as the standard deviation of the  
455 residual evolutionary rates of its ATGC-COGs.

## 456 Author contributions

457 JI and IS performed research; JI, YIW, EVK and IS analyzed the data; JI, EVK and IS wrote the manuscript  
458 that was edited and approved by all authors.

459

## 460 Acknowledgements

461 The authors thank Koonin group members for helpful discussions. The authors' research is supported by  
462 intramural program funds of the National Institutes of Health.

463

## 464 References

- 465 1 Puigbò, P., Lobkovsky, A. E., Kristensen, D. M., Wolf, Y. I. & Koonin, E. V. Genomes in turmoil:  
466 quantification of genome dynamics in prokaryote supergenomes. *BMC biology* **12**, 66,  
467 doi:10.1186/s12915-014-0066-4 (2014).
- 468 2 Hao, W. & Golding, G. B. The fate of laterally transferred genes: life in the fast lane to adaptation  
469 or death. *Genome research* **16**, 636-643, doi:10.1101/gr.4746406 (2006).
- 470 3 Vos, M., Hesselman, M. C., Te Beek, T. A., van Passel, M. W. J. & Eyre-Walker, A. Rates of Lateral  
471 Gene Transfer in Prokaryotes: High but Why? *Trends in microbiology* **23**, 598-605,  
472 doi:10.1016/j.tim.2015.07.006 (2015).
- 473 4 Koonin, E. V. & Wolf, Y. I. Genomics of bacteria and archaea: the emerging dynamic view of the  
474 prokaryotic world. *Nucleic Acids Res* **36**, 6688-6719, doi:10.1093/nar/gkn668 (2008).
- 475 5 Soucy, S. M., Huang, J. & Gogarten, J. P. Horizontal gene transfer: building the web of life.  
476 *Nature reviews. Genetics* **16**, 472-482, doi:10.1038/nrg3962 (2015).
- 477 6 Hoetzinger, M. & Hahn, M. W. Genomic divergence and cohesion in a species of pelagic  
478 freshwater bacteria. *BMC genomics* **18**, 794, doi:10.1186/s12864-017-4199-z (2017).
- 479 7 Hanage, W. P., Spratt, B. G., Turner, K. M. & Fraser, C. Modelling bacterial speciation. *Philos*  
480 *Trans R Soc Lond B Biol Sci* **361**, 2039-2044, doi:10.1098/rstb.2006.1926 (2006).
- 481 8 Retchless, A. C. & Lawrence, J. G. Temporal fragmentation of speciation in bacteria. *Science* **317**,  
482 1093-1096, doi:10.1126/science.1144876 (2007).
- 483 9 Shen, P. & Huang, H. V. Homologous recombination in *Escherichia coli*: dependence on substrate  
484 length and homology. *Genetics* **112**, 441-457 (1986).
- 485 10 Kung, S. H., Retchless, A. C., Kwan, J. Y. & Almeida, R. P. Effects of DNA size on transformation  
486 and recombination efficiencies in *Xylella fastidiosa*. *Applied and environmental microbiology* **79**,  
487 1712-1717, doi:10.1128/AEM.03525-12 (2013).
- 488 11 Majewski, J. Sexual isolation in bacteria. *FEMS microbiology letters* **199**, 161-169,  
489 doi:10.1111/j.1574-6968.2001.tb10668.x (2001).
- 490 12 Vulic, M., Dionisio, F., Taddei, F. & Radman, M. Molecular keys to speciation: DNA polymorphism  
491 and the control of genetic exchange in enterobacteria. *Proceedings of the National Academy of*  
492 *Sciences of the United States of America* **94**, 9763-9767 (1997).
- 493 13 Matic, I., Rayssiguier, C. & Radman, M. Interspecies gene exchange in bacteria: the role of SOS  
494 and mismatch repair systems in evolution of species. *Cell* **80**, 507-515 (1995).
- 495 14 Majewski, J. & Cohan, F. M. DNA sequence similarity requirements for interspecific  
496 recombination in *Bacillus*. *Genetics* **153**, 1525-1533 (1999).

- 497 15 Majewski, J., Zawadzki, P., Pickerill, P., Cohan, F. M. & Dowson, C. G. Barriers to genetic  
498 exchange between bacterial species: *Streptococcus pneumoniae* transformation. *Journal of*  
499 *bacteriology* **182**, 1016-1023 (2000).
- 500 16 Bobay, L. M. & Ochman, H. Biological species are universal across Life's domains. *Genome*  
501 *biology and evolution*, doi:10.1093/gbe/evx026 (2017).
- 502 17 Dykhuizen, D. E. & Green, L. Recombination in *Escherichia coli* and the definition of biological  
503 species. *Journal of bacteriology* **173**, 7257-7268 (1991).
- 504 18 Cadillo-Quiroz, H. *et al.* Patterns of gene flow define species of thermophilic Archaea. *PLoS Biol*  
505 **10**, e1001265, doi:10.1371/journal.pbio.1001265 (2012).
- 506 19 Falush, D. *et al.* Mismatch induced speciation in *Salmonella*: model and data. *Philos Trans R Soc*  
507 *Lond B Biol Sci* **361**, 2045-2053, doi:10.1098/rstb.2006.1925 (2006).
- 508 20 Fraser, C., Hanage, W. P. & Spratt, B. G. Recombination and the nature of bacterial speciation.  
509 *Science* **315**, 476-480, doi:10.1126/science.1127573 (2007).
- 510 21 Dixit, P. D., Pang, T. Y. & Maslov, S. Recombination-Driven Genome Evolution and Stability of  
511 Bacterial Species. *Genetics* **207**, 281-295, doi:10.1534/genetics.117.300061 (2017).
- 512 22 Lawrence, J. G. & Retchless, A. C. The interplay of homologous recombination and horizontal  
513 gene transfer in bacterial speciation. *Methods Mol Biol* **532**, 29-53, doi:10.1007/978-1-60327-  
514 853-9\_3 (2009).
- 515 23 Shapiro, B. J. *et al.* Population genomics of early events in the ecological differentiation of  
516 bacteria. *Science* **336**, 48-51, doi:10.1126/science.1218198 (2012).
- 517 24 Marttinen, P. & Hanage, W. P. Speciation trajectories in recombining bacterial species. *PLoS*  
518 *computational biology* **13**, e1005640, doi:10.1371/journal.pcbi.1005640 (2017).
- 519 25 Sela, I., Wolf, Y. I. & Koonin, E. V. Theory of prokaryotic genome evolution. *Proceedings of the*  
520 *National Academy of Sciences of the United States of America* **113**, 11399-11407,  
521 doi:10.1073/pnas.1614083113 (2016).
- 522 26 Wolf, Y. I., Makarova, K. S., Lobkovsky, A. E. & Koonin, E. V. Two fundamentally different classes  
523 of microbial genes. *Nature microbiology* **2**, 16208, doi:10.1038/nmicrobiol.2016.208 (2016).
- 524 27 Sela, I., Wolf, Y. I. & Koonin, E. V. Genome plasticity, a key factor of evolution in prokaryotes.  
525 *bioRxiv* doi: 10.1101/357400, doi:doi: 10.1101/357400 (2018).
- 526 28 Iranzo, J., Cuesta, J. A., Manrubia, S., Katsnelson, M. I. & Koonin, E. V. Disentangling the effects  
527 of selection and loss bias on gene dynamics. *Proceedings of the National Academy of Sciences of*  
528 *the United States of America* **114**, E5616-E5624, doi:10.1073/pnas.1704925114 (2017).
- 529 29 Koonin, E. V. & Wolf, Y. I. The fundamental units, processes and patterns of evolution, and the  
530 tree of life conundrum. *Biology direct* **4**, 33, doi:10.1186/1745-6150-4-33 (2009).
- 531 30 Koonin, E. V. The Turbulent Network Dynamics of Microbial Evolution and the Statistical Tree of  
532 Life. *J Mol Evol* **80**, 244-250, doi:10.1007/s00239-015-9679-7 (2015).
- 533 31 Polz, M. F., Alm, E. J. & Hanage, W. P. Horizontal gene transfer and the evolution of bacterial and  
534 archaeal population structure. *Trends in genetics : TIG* **29**, 170-175,  
535 doi:10.1016/j.tig.2012.12.006 (2013).
- 536 32 Booth, A., Mariscal, C. & Doolittle, W. F. The Modern Synthesis in the Light of Microbial  
537 Genomics. *Annual review of microbiology* **70**, 279-297, doi:10.1146/annurev-micro-102215-  
538 095456 (2016).
- 539 33 Ochman, H., Lawrence, J. G. & Groisman, E. A. Lateral gene transfer and the nature of bacterial  
540 innovation. *Nature* **405**, 299-304, doi:10.1038/35012500 (2000).
- 541 34 Pal, C., Papp, B. & Lercher, M. J. Adaptive evolution of bacterial metabolic networks by  
542 horizontal gene transfer. *Nat Genet* **37**, 1372-1375, doi:10.1038/ng1686 (2005).

- 543 35 Rodriguez-Valera, F., Martin-Cuadrado, A. B. & Lopez-Perez, M. Flexible genomic islands as  
544 drivers of genome evolution. *Curr Opin Microbiol* **31**, 154-160, doi:10.1016/j.mib.2016.03.014  
545 (2016).
- 546 36 Hoetzinger, M., Schmidt, J., Jezberova, J., Koll, U. & Hahn, M. W. Microdiversification of a Pelagic  
547 Polynucleobacter Species Is Mainly Driven by Acquisition of Genomic Islands from a Partially  
548 Interspecific Gene Pool. *Applied and environmental microbiology* **83**, doi:10.1128/AEM.02266-16  
549 (2017).
- 550 37 Hanage, W. P., Fraser, C. & Spratt, B. G. Sequences, sequence clusters and bacterial species.  
551 *Philos Trans R Soc Lond B Biol Sci* **361**, 1917-1927, doi:10.1098/rstb.2006.1917 (2006).
- 552 38 Baumdicker, F., Hess, W. R. & Pfaffelhuber, P. The infinitely many genes model for the  
553 distributed genome of bacteria. *Genome biology and evolution* **4**, 443-456,  
554 doi:10.1093/gbe/evs016 (2012).
- 555 39 Kristensen, D. M., Wolf, Y. I. & Koonin, E. V. ATGC database and ATGC-COGs: an updated  
556 resource for micro- and macro-evolutionary studies of prokaryotic genomes and protein family  
557 annotation. *Nucleic Acids Res.* **45**, D210-D218, doi:10.1093/nar/gkw934 (2017).
- 558 40 Cohen, O. & Pupko, T. Inference of gain and loss events from phyletic patterns using stochastic  
559 mapping and maximum parsimony--a simulation study. *Genome biology and evolution* **3**, 1265-  
560 1275, doi:10.1093/gbe/evr101 (2011).
- 561 41 Burnham, K. P., Anderson, D. R. & Burnham, K. P. *Model selection and multimodel inference : a*  
562 *practical information-theoretic approach*. 2nd edn, (Springer, 2002).
- 563 42 Cowan, G. *Statistical data analysis*. (Clarendon Press, 1998).
- 564 43 Bhattacharjee, S. M. & Seno, F. A measure of data collapse for scaling. *J Phys a-Math Gen* **34**,  
565 6375-6380 (2001).
- 566 44 Wolf, Y. I., Snir, S. & Koonin, E. V. Stability along with extreme variability in core genome  
567 evolution. *Genome biology and evolution* **5**, 1393-1402, doi:10.1093/gbe/evt098 (2013).
- 568 45 Galperin, M. Y., Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Expanded microbial genome coverage  
569 and improved protein family annotation in the COG database. *Nucleic Acids Res* **43**, D261-269,  
570 doi:10.1093/nar/gku1223 (2015).

571

572

573

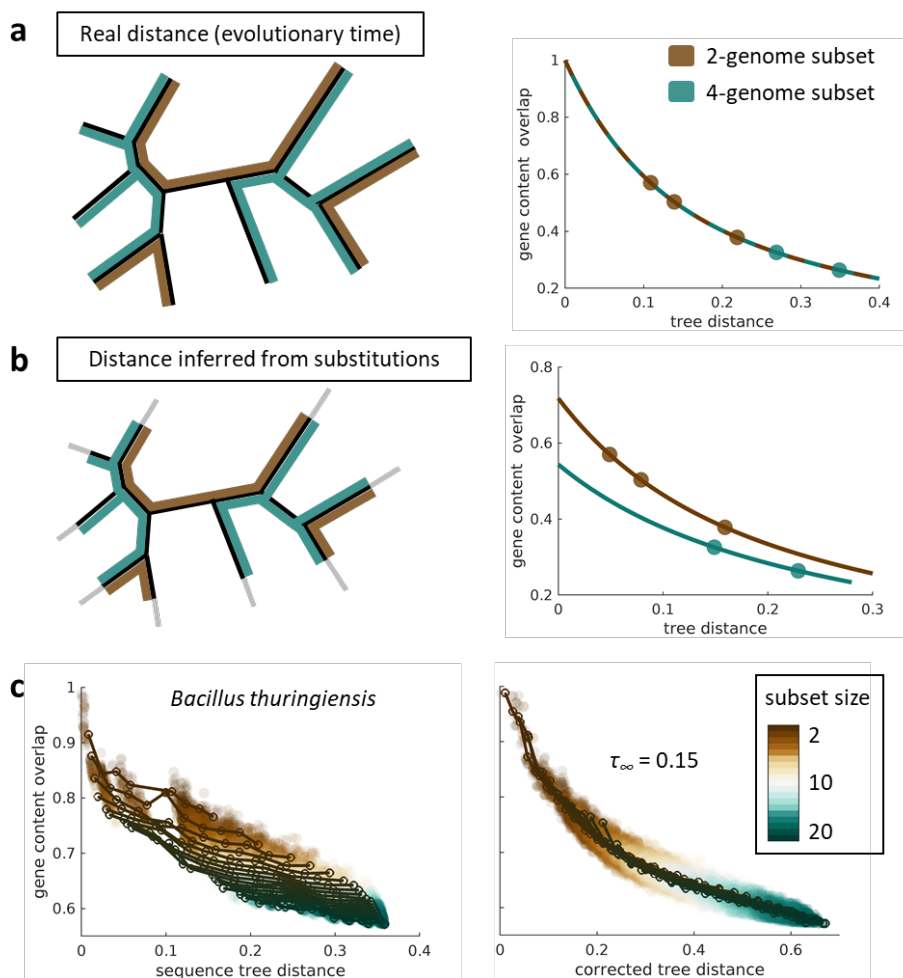


574 **Figure 1. Homologous recombination leads to compression of terminal branches in sequence-based**  
575 **phylogenetic trees that can be detected through the analysis of gene content decay curves.**

576 a) If tree distances are proportional to the true evolutionary time, the fraction of genes shared by a  
577 subset of genomes will decay with the total length of the subtree, and the decay curves will be the same  
578 regardless of the number of genomes in the subset. For illustration purposes, three subsets of 2  
579 genomes are highlighted in brown, and two subsets of 4 genomes are highlighted in green.

580 b) Homologous recombination between pairs of closely related genomes erases recent sequence  
581 divergence which results in an underestimation of the evolutionary times associated with terminal tree  
582 branches. Such underestimation leads to gene content decay curves that depend on the number of  
583 genomes included in the subset. Accordingly, the decay curve of subsets of 4 genomes is different from  
584 the decay curve of subsets of 2 genomes.

585 c) The gene content decay curves of the *Bacillus thuringiensis/cereus/anthracis* group are compatible  
586 with a scenario of recombination-driven shortening of the terminal tree branches (left plot, based on  
587 the tree from Fig. 2a). On the right, if the recombination model is used to correct for unobserved  
588 variation (fit in Fig. 2c, left panel), overlapping decay curves are obtained.



589

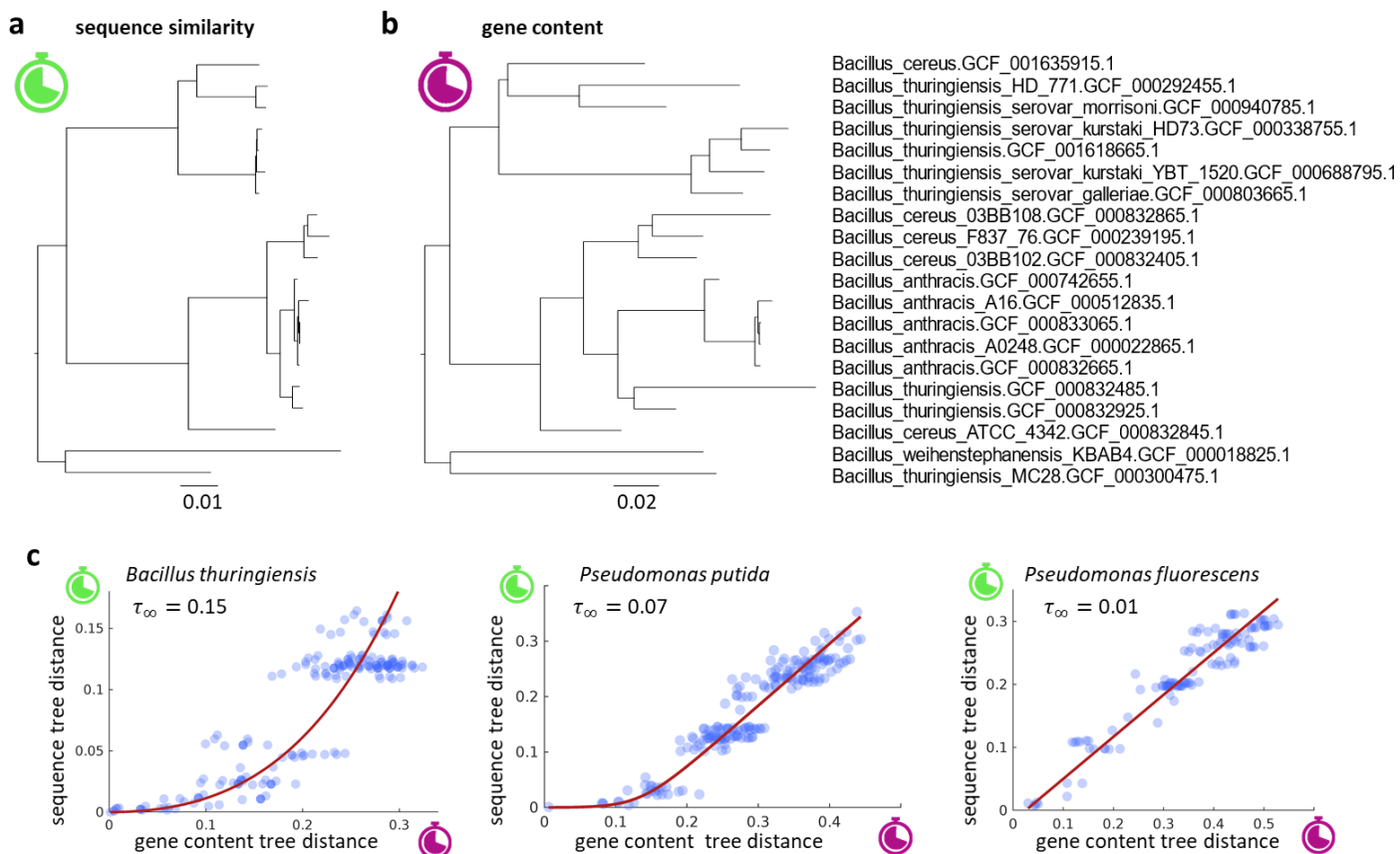
590

591 **Figure 2. Lack of collinearity between sequence and gene content distances supports a recombination-**  
592 **drive delay in the molecular clock.**

593 a) Phylogenetic tree of the *Bacillus thuringiensis/cereus/anthracis* group based on the concatenated  
594 alignment of genes shared by all members of the group.

595 b) The same tree, with branch lengths proportional to the number of gene gain and loss events  
596 estimated by phylogenomic analysis.

597 c) Comparison of pairwise distances between leaves in the sequence similarity and gene content trees,  
598 for three representative groups (the left-most plot corresponds to the trees in a and b). The red line is  
599 the fit of the recombination-driven model of sequence divergence with a long-term delay in the  
600 molecular clock equal to  $\tau_{\infty}$ .

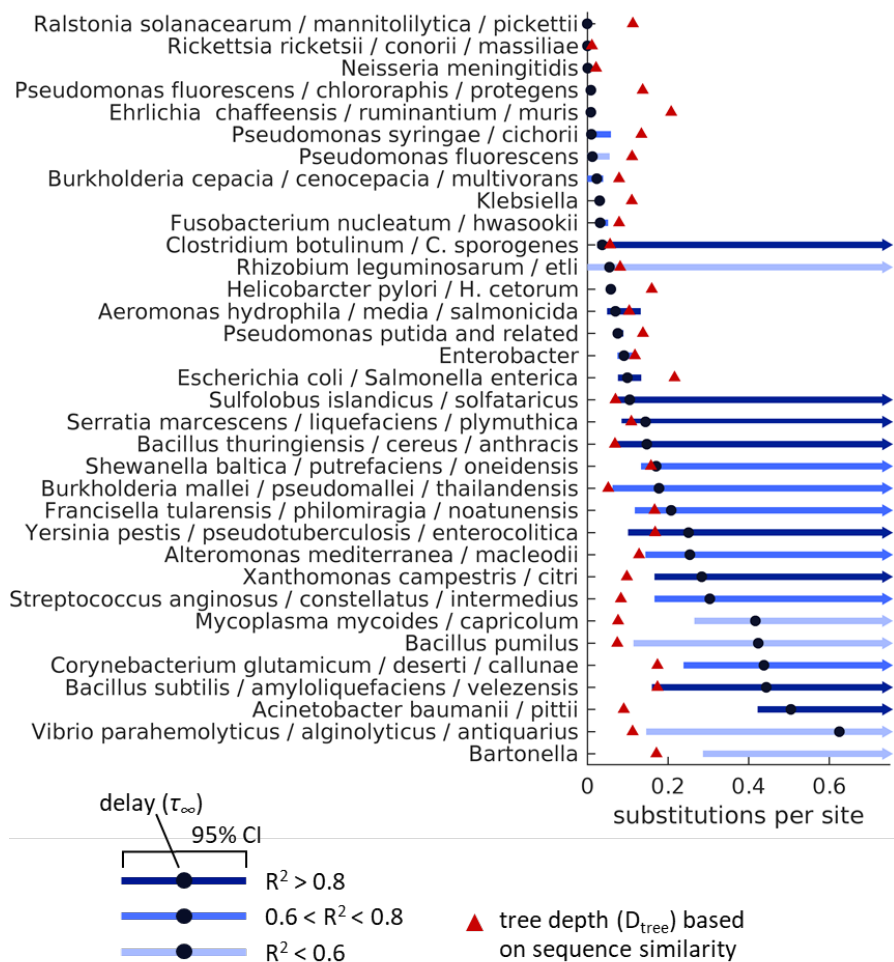


601

602

603 **Figure 3. Recombination-driven delay ( $\tau_\infty$ ) in the molecular clock in different groups of bacteria.**

604 Black circles indicate the best fit of the long-term delay parameter  $\tau_\infty$  based on the comparison of  
 605 sequence similarity and gene content trees. Blue lines show the 95% confidence intervals. Red triangles  
 606 indicate the total depth of the sequence similarity tree. Values of the delay above or very close to the  
 607 total tree depth imply that most genomes in a group are strongly bound by homologous recombination;  
 608 an upper 95% confidence bound cannot be calculated in those cases.



609

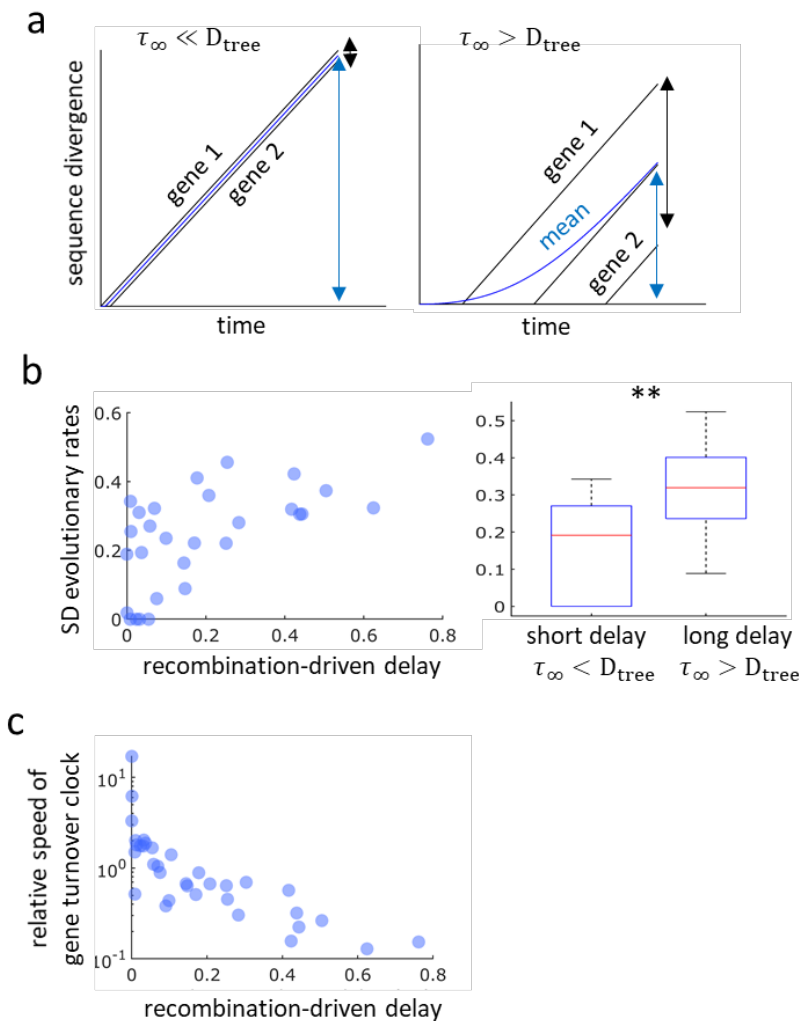
610

611 **Figure 4: Recombination-driven delay leads to over-dispersion of evolutionary rates within a taxon (a,**  
612 **b) and negatively correlates with the gene turnover rate (c).**

613 (a) The divergence of individual genes (black lines) grows linearly after the establishment of barriers to  
614 gene conversion, whereas the overall divergence (blue lines) grows non-linearly as the barriers spread  
615 across the genome. The mean (blue arrows) and standard deviation (black arrows) of the gene-specific  
616 divergences are determined by the values of the recombination-driven delay ( $\tau_{\infty}$ ) and the tree depth  
617 ( $D_{\text{tree}}$ ).

618 (b) Standard deviations of the residual evolutionary rates (corrected by gene- and taxon-wise rates) in  
619 linearly diverging ( $\tau_{\infty} > D_{\text{tree}}$ ) and strongly delayed taxa ( $\tau_{\infty} \ll D_{\text{tree}}$ ); Pearson's correlation coefficient  
620  $R = 0.64$  ( $p < 0.001$ ). \*\* Statistically significant with  $p = 0.002$  (Student's T test).

621 (c) Negative association between the recombination-driven delays and the relative rates of gene  
622 turnover with respect to substitutions; Spearman's correlation coefficient  $\rho = -0.86$  ( $p < 0.001$ ).



623

