

1 **Title: How clonal are bacteria over time?**

2

3 **Author: B. Jesse Shapiro<sup>1</sup>**

4

5 **Affiliation**

6 1. Département de sciences biologiques, Université de Montréal, Montréal, QC H3C 3J7,  
7 Canada

8 \* jesse.shapiro@umontreal.ca

9

10

11

12 **Abstract**

13 Bacteria and archaea reproduce clonally (vertical descent), but exchange genes by recombination  
14 (horizontal transfer). Recombination allows adaptive mutations or genes to spread rapidly within  
15 (or even between) species, and reduces the burden of deleterious mutations. Clonality – defined  
16 here as the balance between vertical and horizontal inheritance – is therefore a key microbial trait,  
17 determining how quickly a population can adapt and the size of its gene pool. Here, I discuss  
18 whether clonality varies over time and if it can be considered a stable trait of a given population. I  
19 show that, in some cases, clonality is clearly not static. For example, non-clonal (highly  
20 recombining) populations can give rise to clonal expansions, often of pathogens. However, an  
21 analysis of time-course metagenomic data from a lake suggests that a bacterial population's past  
22 clonality (as measured by its genetic diversity) is a good predictor of its future clonality. Clonality  
23 therefore appears to be relatively – but not completely – stable over evolutionary time.

24

25 **Introduction**

26 Here, I revisit the question posed in the title of a classic paper by John Maynard Smith and  
27 colleagues [1]: How clonal are bacteria, and more specifically how does clonality vary among  
28 different microbial populations and over time? First, what do we mean by clonality? Perfectly  
29 clonal bacteria replicate by cell division (vertical descent) and evolve by random mutations that  
30 occur during DNA replication. In this theoretical population, there is negligible horizontal  
31 transfer of DNA by recombination across the resulting tree of vertical descent. Very few (if any)  
32 natural bacterial populations fit this idealized, theoretical definition of clonality. Or, as discussed  
33 below, they might only fit it for a short amount of time. However, knowing where a bacterial

34 population of interest happens to fall along a spectrum of clonality can help us understand its  
35 biology, and even make predictions about its evolution.

36 The opposite of clonality is panmixis – a situation in which the rate of horizontal transfer  
37 is higher than the rate of vertical cell division, resulting in random association (linkage  
38 equilibrium) among loci in the genome [1,2]. However, rates of horizontal transfer  
39 (recombination) vary widely across the genome, such that a population can be mostly clonal,  
40 except for a few loci in the genome [3]. These loci came to be termed genomic islands – a  
41 metaphor I will build upon below. Some of the first islands identified were called pathogenicity  
42 islands because they contained virulence factors [4]. However, non-pathogenic environmental  
43 bacteria also contain islands, conferring adaptation to different ecological niches. For example,  
44 genes in *Prochlorococcus* genomic islands confer adaptation to light and nutrient conditions  
45 [5,6]. But islands need not confer niche adaptation to their host genome; they can be neutral to  
46 host fitness or even detrimental, selfish parasites. Here, I define genomic islands broadly as any  
47 piece of DNA that is transferred horizontally (by either homologous or nonhomologous  
48 recombination) from cell to cell and therefore evolves independently (*i.e.* is unlinked) from the  
49 rest of the genome.

50 I will begin by extending the use of island analogies to include continents, peninsulas and  
51 archipelagos (**Table 1**). I will then use these analogies to discuss to what extent microbial  
52 populations are clonal or panmictic, and how often they transition between the two regimes.

53

#### 54 **Are some islands really peninsulas?**

55 In the classic analogy, an island is totally disconnected from the mainland, meaning that genes in  
56 the island evolve independently of the genome (**Table 1**). Examples of islands that fit this strict  
57 independence might include integrated phages and other “selfish” elements, or genes that reside  
58 in a particular niche but not in a particular genome (*e.g.* a gene ecology model [7]). Peninsulas  
59 provide an analogy that might better describe how islands are related to microbial genomes. A  
60 peninsula (or “presque-île,” from the French for “almost island”) is a geographic term for a very  
61 narrow strip of land connected to (but distinct from) the mainland. In my analogy, an island is  
62 evolutionarily independent of the mainland genome, but their fates may become linked, forming a  
63 peninsula. For example, a bacterium may acquire a gene from a vast microbial gene pool. This  
64 gene allows the bacterium to invade a new ecological niche, triggering a clonal expansion in  
65 which the fate of the gene and its new host genome are linked, at least for the duration of the  
66 clonal expansion. One such example could be *Yersinia pestis*, which acquired a single gene  
67 allowing flea-borne transmission and triggering a clonal expansion in the form of Plague

68 pandemics [8]. Another peninsula, the prophage-encoded cholera toxin, and its links to the  
69 mainland *Vibrio cholerae* genome [9,10], is discussed below.

70

### 71 **Are some genomes archipelagos?**

72 The very concept of one or a few islands implies a contrast with the large, clonal genomic  
73 mainland or continent. But some microbial genomes may contain so many islands that there is no  
74 mainland, only a vast archipelago (**Table 1**). A striking recent example is a population of  
75 hotspring cyanobacteria in which virtually every gene in the genome evolved independently due  
76 to frequent recombination [11], leading the authors to call the population “quasi-sexual” (in other  
77 words, panmictic). Frequent recombination was confirmed by another group, using different  
78 methods to study the same cyanobacteria [12]. This group also found that despite a history of  
79 panmixis over long time scales, populations are more clonal over shorter time scales. Similarly,  
80 the Asian ocean population of *Vibrio parahaemolyticus* also forms a panmictic gene pool, with  
81 each recombination block of ~1.8 kbp evolving independently [13]. However, the panmictic gene  
82 pool occasionally gives rise to pandemic clones. In another example, we found that almost every  
83 gene in a population of *Vibrio cyclitrophicus* genomes showed signs of recombination over  
84 relatively recent time scales [14]. Such apparently high rates of recombination in natural  
85 populations were mysterious at first, contradicting recombination rates measured in the lab  
86 [15,16] and predicted by theory [7,17]. However, theoretical models (discussed below) suggest  
87 mechanisms capable of explaining how genes can spread through populations more rapidly by  
88 recombination than by clonal expansion [18-20].

89

### 90 **Clonal expansions from panmictic pools**

91 I propose that archipelagos are not necessarily static over time, and that archipelagos can  
92 sometimes coalesce into continents. Given the right ecological opportunity, a genome from a  
93 panmictic gene pool can escape the "gravitational pull" of recombination and take off into a  
94 clonal expansion. An example mentioned earlier is *V. cholerae*, a genetically diverse group of  
95 coastal marine bacteria, some of which cause cholera. Virulence is mainly determined by two loci  
96 in the genome: the cholera toxin and the toxin-coregulated pilus. Both genes are frequently gained  
97 and lost by recombination [21,22], but are always found in one lineage of *V. cholerae* – the  
98 lineage causing severe disease with pandemic potential, known as the phylocore genome (PG)  
99 group [10]. It remains a mystery why the PG lineage evolved once, and only once. If PG *V.*  
100 *cholerae* really did evolve just once, this would be surprising because *V. cholerae* draws on a  
101 diverse, global gene pool and can be considered panmictic [23]. Therefore multiple different

102 lineages would be expected to acquire the two (or perhaps a handful of) genetic elements required  
103 for pandemic disease. This leads to the hypothesis that pandemic cholera emergence is *selection*  
104 *limited* rather than *diversity limited*. In other words, benign *V. cholerae* strains constantly acquire  
105 virulence genes. However, these strains rarely encounter the right ecological niche to flourish,  
106 *e.g.* a human population consuming brackish water. "The right niche" has appeared a few times in  
107 human history: for example in India in the 1800s, when the Classical lineage evolved, and again  
108 in Indonesia in the 1950s, when the El Tor lineage evolved [24]. When the right conditions  
109 appear, the PG lineage, along with its virulence factors, takes off in a clonal expansions which  
110 continue to wreak havoc today (*e.g.* cholera pandemics from the 1800s to today, all caused by the  
111 PG clonal group). The virulence factors, previously islands in an archipelago, became a peninsula  
112 connected to the PG mainland. The linkage between virulence factors and PG remains imperfect  
113 because different variants of the cholera toxin continue to flow in and out of the PG continent  
114 [10,21]; hence the toxin remains a peninsula, not firmly part of the mainland.

115 *V. cholerae* is a particularly well-characterized example of a panmictic gene pool giving  
116 rise to a clonal expansion, but similar evolutionary dynamics are seen in other pathogens as well  
117 (*e.g.* *V. parahaemolyticus* described above [13]). Enterotoxigenic *Escherichia coli* (ETEC) seems  
118 to behave similarly, with deep branches of the phylogeny obscured by frequent recombination  
119 and plasmid exchange, but more recent branches experiencing mostly clonal descent, with tight  
120 linkage between virulence factors and the genomic mainland [25]. These observations are  
121 consistent with an ancient, panmictic gene pool giving rise to clonal expansions, which can last  
122 for decades or centuries.

123

#### 124 **The balance between recombination and selection**

125 Let us consider the evolutionary forces that determine clonality: natural selection and  
126 recombination. The effect of recombination on clonality is straightforward: more recombination  
127 means less clonality. The effect of natural selection is more complex, but is defined here simply  
128 as a force which favors clonal expansions of adaptive mutants within an ecological niche.

129 Selection, as defined here, therefore includes ecological effects. When driven by ecological  
130 selection, clonal expansions are called selective sweeps, in which one clone outcompetes all  
131 others, purging genetic diversity in the population.

132 Recombination and selection interact to determine the clonality of a population.

133 Recombination rates depend both on the ability of DNA to enter a cell and be incorporated into  
134 the genome (the baseline rate) and the ability of that DNA to be retained by a balance of genetic  
135 drift and natural selection (the realized rate). Some bacteria, such as *Helicobacter pylori*, have

136 realized recombination rates that are much higher than point mutation rates, exchanging at least  
137 10% of their genome within a single four-year human infection [26]. Others, such as  
138 *Staphylococcus aureus* [27,28] and *Mycobacterium tuberculosis* are decidedly more clonal [29-  
139 31]. Recombination rates (both realized and baseline) vary widely across the genome. Of 10  
140 pathogenic bacterial species studied, all had identifiable recombination 'hot' regions, although  
141 their length, genomic location and gene content varied [3]. Genes of different functions had  
142 different realized recombination rates, implying a role for natural selection on gene function in  
143 determining whether newly acquired genes are retained.

144

### 145 **Modeling the recombination-selection balance**

146 When rates of recombination are relatively low compared to rates of natural selection on adaptive  
147 genes within niches, entire genomes will sweep to fixation before they can be shuffled by  
148 recombination. Following previous modelling work,  $s$  is defined here as the selective coefficient  
149 of a niche-adaptive allele, and  $r$  is the recombination rate, per locus per generation [7,32,33]. The  
150  $s \gg r$  regime is well described in the Stable Ecotype Model [17], which predicts that most of the  
151 genome will follow a single, clonal phylogeny. Genome-wide sweeps thus increase clonality and  
152 can be considered a hallmark of clonal populations (**Table 1, Figure 1A**). In the  $r \gg s$  regime,  
153 individual genes (rather than entire genomes) will sweep to fixation (*i.e.* reach 100% frequency)  
154 in ecological niches to which they are adapted, without affecting genetic diversity elsewhere in  
155 the genome (**Figure 1B**). The first serious theoretical attempt to reconcile the observations of  
156 gene-specific sweeps with low recombination rates was the "Adapt Globally, Act Locally" model  
157 [18,20], in which globally adaptive genes (adaptive in multiple different niches) trigger genome-  
158 wide sweeps within a niche before being transferred to the next niche. When viewed across  
159 niches, the result is a gene-specific sweep. A recent model suggested another mechanism by  
160 which these gene sweeps can occur at moderate (not unrealistically high) rates of recombination  
161 [20]. In this model, a microbial habitat is bombarded with genetically maladapted migrants,  
162 allowing gene sweeps to occur, although the adaptive allele never reaches 100% frequency due to  
163 the constant input of migrants. In another model, Takeuchi et al. [19] show that gene sweeps can  
164 occur when  $r$  is either very high or – counter-intuitively – when  $r$  is very low, but only when  
165 negative frequency-dependent selection (NFDS) reduces the rate of genome-wide selective  
166 sweeps. NFDS might be commonly imposed on bacteria and archaea by viral (phage) predation,  
167 providing a selective advantage to rare alleles of phage receptor genes [34].

168 One parameter not thoroughly explored in any of these models is the effective population  
169 size ( $N_e$ ). Populations with small  $N_e$  are dominated by drift, and natural selection is inefficient.

170 They may experience genome-wide sweeps independently of  $r$  and  $s$ . As discussed below,  $N_e$  is  
171 probably an important determinant of clonality.

172

### 173 **Genome-wide and gene-specific sweeps in nature**

174 To date, empirical evidence for gene-specific and genome-wide sweeps has come mostly from  
175 cross-sectional studies of a single population of genomes at a single point in time, with  
176 recombination and selection inferred backward in time [11,13,14,35]. Sequencing microbial  
177 genomes or metagenomes sampled over time – already a typical practice in genomic  
178 epidemiology (e.g. [28,36]) – promises to elucidate the rates of gene-specific and genome-wide  
179 sweeps in nature (**Figure 1**).

180 In a pioneering study, Bendall et al. [37] sampled a lake over nine years and followed  
181 single-nucleotide polymorphism (SNP) and gene frequencies in 30 bacterial populations by  
182 metagenomic sequencing. They inferred that one of the populations (*Chlorobium*-111) had  
183 undergone a near-complete genome-wide sweep over the nine-year study, with most SNP  
184 diversity purged genome-wide (**Figure 1A**). By identifying regions of the genome with  
185 unexpectedly low diversity compared to the genome-wide average, they inferred that gene-  
186 specific sweeps had occurred in six of the populations, but these sweeps occurred before the  
187 beginning of the nine-year study. During the nine-year study, they observed examples “where a  
188 few adjacent SNPs trended toward fixation while genome-wide diversity was maintained”  
189 (**Figure 1B**). They took this observation as consistent with gene-specific selective sweeps, but did  
190 not attempt to determine whether the sweeps were due to selection or drift. Similarly, the inferred  
191 genome-wide sweep could have been caused by selection or drift.

192 Is the frequency of genome-wide sweeps controlled mainly by the balance of  
193 recombination and selection, or could drift (controlled by the effective population size,  $N_e$ ) play  
194 an important role as well? In their Figure 2 legend, Bendall et al. observe that “populations with  
195 many SNPs were not necessarily sequenced deeper than those with few SNPs.” This statement  
196 was simply meant to assure the reader that SNP calling was not biased by sequence coverage, but  
197 it also suggests that microbial populations tend to be far from a standard neutral model with  
198 constant population size (*i.e.* the Wright-Fisher coalescent [38]). Taking SNP density as a rough  
199 measure of  $N_e$  and sequence coverage as a rough measure of the census population size, the data  
200 show that census population size is a poor predictor of  $N_e$ , and thus that populations are likely  
201 impacted by population bottlenecks and/or selection – although it is difficult to distinguish  
202 between them.

203           As a whole, the study showed that both genome-wide and gene-specific sweeps can occur  
204 in different microbial populations from the same environment. Whether microbial populations  
205 behaved differently due to differences in their ecology (*i.e.* regime of natural selection) or in their  
206 baseline recombination rates remains a question for future study; Cohan recently suggested that  
207 ecological differences could be more important [39]. Specifically, he suggests that ecological  
208 "generalists" could have more opportunities for diversification, and thus be relatively immune to  
209 genome-wide sweeps, compared to ecological specialists. I suggest that these generalists may  
210 simply have larger  $N_e$ , such that diversity is rarely purged by drift (*e.g.* bottlenecks), and that  
211 diversity is mainly governed by selection and recombination. It is possible that ecological  
212 generalists tend to have large  $N_e$ , or that  $N_e$  and ecology exert independent effects (*e.g.* ecological  
213 generalists can have low  $N_e$  and still resist genome-wide sweeps). The fact that one genome-wide  
214 sweep was observed over a nine year period suggests that such events might be relatively rapid  
215 but rare (only observed in one of 30 populations). Meanwhile, gene-sweeps might be more  
216 common historically (affecting six of 30 populations), but could take longer to proceed to  
217 completion.

218

### 219 **Is clonality a stable trait?**

220 As described in the *V. cholerae* example, some pathogenic bacterial populations can switch  
221 between panmictic and clonal lifestyles [11,13,14,19,25,34,35]. Therefore clonality can vary over  
222 time, but how much and how often? To quantify the stability of clonality over time, not just in  
223 pathogens but in free-living environmental bacteria, I re-analyzed the lake time-course of Bendall  
224 et al. [37]. Because estimates of selection and recombination rates were not readily available for  
225 this dataset, I defined clonality based on the total genetic diversity (SNP density) in a population,  
226 which scales with  $N_e$ , and includes the effects of both drift and of genome-wide selective sweeps.  
227 Frequent genome-wide selective sweeps ( $s \gg r$ ) and/or small  $N_e$  result in clonality. I identified  
228 20 "old, diverse" populations as those with a high density of SNPs (>1500 SNPs/Mbp) at the  
229 beginning of the time-course. These populations are likely "old" because they have gone a  
230 relatively long time since the last genome-wide purge of genetic diversity and are relatively non-  
231 clonal. They include the six populations inferred to have undergone gene-specific sweeps [37].  
232 The remaining ten populations were defined as "young, low-diversity," having more recently  
233 experienced a genome-wide purge of diversity. The "old, diverse" populations have a relatively  
234 low ratio of nonsynonymous (N) to synonymous (S) SNPs, suggesting large  $N_e$  and ample time  
235 for purifying selection to remove (mostly deleterious) nonsynonymous mutations (**Figure 2**). In  
236 contrast, the "young, low-diversity" populations are more likely to have high N:S ratios,

237 suggesting smaller  $N_e$  and less time for purifying selection to act. The cutoff of 1500 SNPs/Mbp  
238 (0.15% SNP density) is somewhat arbitrary, but seems to correspond to an inflection point in  
239 **Figure 2**, and also falls squarely within the 0-0.30% SNP density, above which *E. coli* evolution  
240 appears to transition from clonal to recombining [40].

241 With young (presumably clonal, or low  $N_e$ ) and old (less clonal, higher  $N_e$ ) populations  
242 thus defined, I asked whether “old, diverse” populations tended to maintain their diversity  
243 through the 9-year period of the study. In their paper, Bendall et al. defined two alternative  
244 population types: 1) those that maintained stable SNP diversity over nine years, and 2) those that  
245 experienced significant fluctuations in diversity due to clonal expansions – defined when one, but  
246 not all timepoints are dominated by a single allele ( $\geq 95\%$  frequency) at  $>40\%$  of SNP sites in the  
247 genome [37]. By definition, the 19 populations of the first type did not experience genome-wide  
248 sweeps during the study, while the 11 populations of the second type did experience genome-  
249 wide purges of diversity, which were transient in ten cases and apparently permanent in one case  
250 (*Chlorobium*-111). Strikingly, 17 out of 20 “old, diverse” populations maintained their diversity  
251 over the nine-year study, compared to only two out of ten “young, low-diversity” populations  
252 (Fisher test, Odds Ratio = 19.4,  $P < 0.001$ ). This result suggests that populations with a history of  
253 genome-wide sweeps tend to experience subsequent genome-wide sweeps, and those that have  
254 maintained genetic diversity in the past tend to maintain their diversity into the future. In other  
255 words, clonality can be considered a relatively stable microbial trait. However, we must take care  
256 in taking a genome-wide sweep as evidence for clonality ( $s \gg r$ ). In small effective population  
257 sizes, genome-wide sweeps can occur due to drift and bottlenecks, independently of  $s$  and  $r$ .  
258 Therefore, many of “young, low diversity” populations (**Figure 2**) may have lost diversity over  
259 nine years due to low  $N_e$ , not due to low  $r$ .

260

### 261 **History repeats itself**

262 It appears that pathogens are more likely than free-living bacteria to undergo clonal expansions,  
263 due in part to their ecology and transmission dynamics [41]. Free-living aquatic bacteria, on the  
264 other hand, seem to be more likely to live in large, panmictic populations and behave like  
265 archipelagos [11,13,14,37]. If clonality is indeed a stable trait, this implies that history will repeat  
266 itself, and that the future behavior of microbial populations can be predicted with some  
267 confidence from their past behavior. Diverse populations tend to stay diverse. Clonal populations  
268 (that experience frequent genome-wide sweeps) tend to stay clonal. But history is not doomed to  
269 repeat itself forever. As we have seen, clonal expansions, such as pandemic *V. cholerae*, may  
270 originate when a panmictic gene pool (an archipelago) coalesces into a clonal continent, with



271 virulence factors linked as peninsulas. Many such pathogenic clones have been documented, with  
272 life spans of decades to thousands of years [25,28,42,43]. Other pathogens, such as *Streptococcus*  
273 *pneumoniae*, may retain their panmictic population structure throughout an outbreak [44-46].  
274 Why some pathogens are clonal and others are panmictic is an open question, but surely depends  
275 on the balance between recombination and selection, on the effective population size, and on the  
276 time scales considered.

277

278 **Acknowledgments.** I am grateful to Rex Malmstrom, Yan Boucher, and Salvador Almagro-  
279 Moreno for their thoughtful comments which improved the manuscript. I was supported by a  
280 Canada Research Chair.

281 **Table 1. Extended island metaphors of microbial genome evolution.**  
 282

Geographic metaphor	Genetic unit to which the metaphor applies	Type of selective sweep experience by the unit	Dominant mode of genetic transmission	Example
Island	Gene	Gene-specific	horizontal	genes in the <i>V. cholerae</i> integron [22,23]
Peninsula	Gene	Genome-wide	vertical (clonal)	the cholera toxin gene, acquired horizontally, then linked to a clonal <i>V. cholerae</i> genome [9,21]
Continent	Genome	Genome-wide	vertical (clonal)	clonal expansions of <i>S. aureus</i> [28], <i>M. tuberculosis</i> [31,43]
Archipelago	Genome	Gene-specific	horizontal	hotspring cyanobacteria [11], ocean vibrios [13,14], pneumococcus [44,46]

283  
 284  
 285

## Figures

286

287 **Figure 1. Temporal dynamics of genome-wide and gene-specific selective sweeps inferred**  
 288 **from metagenomic data.** Genetic diversity can be measured by mapping metagenomic reads to a  
 289 reference genome, identifying SNPs, and calculating the allele frequencies at each SNP position  
 290 in the genome over time. The lowest possible genetic diversity occurs when a single allele is  
 291 present in 100% of metagenomic reads. Alternatively, diversity could be defined in terms of gene  
 292 presence/absence, based on relative coverage of a gene in the reference genome by metagenomic  
 293 reads. **(a)** In a hypothetical genome-wide selective sweep, all positions in the genome tend toward  
 294 low diversity. **(b)** In a hypothetical gene-specific selective sweep, only one or a few positions in  
 295 the genome tend toward low diversity, while the rest of the genome maintains high or  
 296 intermediate diversity.

297

298 **Figure 2. Past diversity predicts future diversity.** Based on data from Table 2 of Bendall et al.  
 299 [37], the nonsynonymous to synonymous (N:S) SNP ratio was plotted against the total SNP  
 300 density (SNPs per megabasepair) for each of 30 bacterial populations. A pseudocount of 1 was  
 301 added to both N and S counts. These 30 populations were divided into 20 “old, diverse”  
 302 populations (>1500 SNPs/Mbp) and 10 “young, low-diversity” populations (<1500 SNPs/Mbp).  
 303 Each point represents one of the 30 populations, colored in black if diversity was maintained over  
 304 a 9-year metagenomic time-course, or in red if it was not. Seventeen out of 20 “old, diverse”  
 305 populations maintained their diversity over the 9-year study, compared to only 2 out of 10  
 306 “young, low-diversity” populations (Fisher test, Odds Ratio = 19.4,  $P < 0.001$ ). The same result is  
 307 obtained drawing a cutoff based on N:S rather than SNPs/Mbp: populations with  $N:S < 0.5$  tend  
 308 to maintain their diversity, whereas those with  $N:S > 0.5$  tend to be purged of diversity over the 9  
 309 years (Fisher test, Odds Ratio = 5.94,  $P = 0.042$ ). Consistent with previous observations that N:S  
 310 depends on the evolutionary time available for purifying selection to act [47,48], N:S is  
 311 negatively correlated with SNPs/Mbp, a proxy for evolutionary time or the time since the last  
 312 genome-wide purge of genetic diversity in this dataset (Pearson’s correlation of  $\log_{10}$  transformed  
 313 data,  $r = -0.81$ ,  $P = 5.6e-8$ ).

314 **References**

315

316

317 1. Smith JM, Smith NH, O'Rourke M, Spratt BG: **How clonal are bacteria?** *Proc Natl Acad*  
318 *Sci USA*. 1993, 90: 4384–4388.

319 2. Guttman DS, Dykhuizen DE: **Clonal Divergence in Escherichia coli as a Result of**  
320 **Recombination, Not Mutation.** *Science*. 1994, 266: 1380–1383.

321 3. Yahara K, Didelot X, Jolley KA, Kobayashi I, Maiden MCJ, Sheppard SK, et al.: **The**  
322 **landscape of realized homologous recombination in pathogenic bacteria.** *Mol Biol*  
323 *Evol*. 2016, 33: 456–471. doi:10.1093/molbev/msv237

324 4. Hacker J, Blum-Oehler G, Mühldorfer I, Tschäpe H: **Pathogenicity islands of virulent**  
325 **bacteria: structure, function and impact on microbial evolution.** *Mol Microbiol*. 1997,  
326 23: 1089–1097.

327 5. Coleman ML, Sullivan MB, Martiny AC, Steglich C, Barry K, Delong EF, et al.:  
328 **Genomic islands and the ecology and evolution of Prochlorococcus.** *Science*. 2006,  
329 311: 1768–1770.

330 6. Coleman ML, Chisholm SW: **Ecosystem-specific selection pressures revealed through**  
331 **comparative population genomics.** *Proc Natl Acad Sci USA*. 2010, 107: 18634–18639.  
332 doi:10.1073/pnas.1009480107

333 7. Shapiro BJ, Polz MF: **Ordering microbial diversity into ecologically and genetically**  
334 **cohesive units.** *Trends Microbiol*. 2014, 22: 235–247. doi:10.1016/j.tim.2014.02.006

335 \* A comprehensive discussion of the roles of selection and recombination in structuring microbial  
336 diversity.

337

338 8. Rasmussen S, Allentoft ME, Nielsen K, Orlando L, Sikora M, Sjögren K-G, et al.: **Early**  
339 **Divergent Strains of Yersinia pestis in Eurasia 5,000 Years Ago.** *Cell*. 2015, 163: 571–  
340 582. doi:10.1016/j.cell.2015.10.009

341 9. Waldor MK, Mekalanos JJ: **Lysogenic conversion by a filamentous phage encoding**  
342 **cholera toxin.** *Science*. 1996, 272: 1910–1914.

343 10. Chun J, Grim CJ, Hasan NA, Lee JH, Choi SY, Haley BJ, et al.: **Comparative genomics**  
344 **reveals mechanism for short-term and long-term clonal transitions in pandemic**  
345 **Vibrio cholerae.** *Proc Natl Acad Sci USA*. 2009, 106: 15442–15447.  
346 doi:10.1073/pnas.0907787106

347 11. Rosen MJ, Davison M, Bhaya D, Fisher DS: **Fine-scale diversity and extensive**  
348 **recombination in a quasisexual bacterial population occupying a broad niche.**  
349 *Science*. 2015, 348: 1019–1023. doi:10.1126/science.aaa4456

350 \* Deep sequencing of 90 cyanobacterial marker genes reveals extensive recombination, with each  
351 genome composed of a random mix of alleles from the gene pool.

352

- 353 12. Melendrez MC, Becraft ED, Wood JM, Olsen MT, Bryant DA, Heidelberg JF, et al.:  
354 **Recombination Does Not Hinder Formation or Detection of Ecological Species of**  
355 ***Synechococcus* Inhabiting a Hot Spring Cyanobacterial Mat.** *Front Microbiol.* 2016,  
356 6: 544. doi:10.1128/AEM.72.1.723-732.2006
- 357 13. Cui Y, Yang X, Didelot X, Guo C, Li D, Yan Y, et al.: **Epidemic Clones, Oceanic Gene**  
358 **Pools, and Eco-LD in the Free Living Marine Pathogen *Vibrio parahaemolyticus*.** *Mol*  
359 *Biol Evol.* 2015, 32: 1396–1410. doi:10.1093/molbev/msv009
- 360 \* In-depth population genomic evidence of a panmictic ocean gene pool containing non-clonal  
361 ecological populations.  
362
- 363 14. Shapiro BJ, Friedman J, Cordero OX, Preheim SP, Timberlake SC, Szabo G, et al.:  
364 **Population Genomics of Early Events in the Ecological Differentiation of Bacteria.**  
365 *Science.* 2012, 336: 48–51. doi:10.1126/science.1218198
- 366 15. Vulić M, Dionisio F, Taddei F, Radman M: **Molecular keys to speciation: DNA**  
367 **polymorphism and the control of genetic exchange in enterobacteria.** *Proc Natl Acad*  
368 *Sci USA.* 1997, 94: 9763–9767.
- 369 16. Majewski J, Cohan FM: **DNA sequence similarity requirements for interspecific**  
370 **recombination in *Bacillus*.** *Genetics.* 1999, 153: 1525–1533.
- 371 17. Wiedenbeck J, Cohan FM: **Origins of bacterial diversity through horizontal genetic**  
372 **transfer and adaptation to new ecological niches.** *FEMS Microbiology Reviews.* 2011,  
373 35: 957–976. doi:10.1111/j.1574-6976.2011.00292.x
- 374 18. Majewski J, Cohan FM: **Adapt globally, act locally: the effect of selective sweeps on**  
375 **bacterial sequence diversity.** *Genetics.* 1999, 152: 1459–1474.
- 376 19. Takeuchi N, Cordero OX, Koonin EV, Kaneko K: **Gene-specific selective sweeps in**  
377 **bacteria and archaea caused by negative frequency-dependent selection.** *BMC*  
378 *Biology.* 2015, 13: 20. doi:10.1186/s12915-015-0131-7
- 379 \* Mathematical modeling shows how individual genes can sweep through populations with low  
380 recombination rates in the presence of negative frequency-dependent selection.  
381
- 382 20. Niehus R, Mitri S, Fletcher AG, Foster KR: **Migration and horizontal gene transfer**  
383 **divide microbial genomes into multiple niches.** *Nature Communications.* 2015, 6: 8924.  
384 doi:10.1038/ncomms9924
- 385 21. Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, et al.: **Evidence for**  
386 **several waves of global transmission in the seventh cholera pandemic.** *Nature.* 2011,  
387 477: 462–U111. doi:10.1038/nature10392
- 388 22. Orata FD, Kirchberger PC, Méheust R, Barlow EJ, Tarr CL, Boucher Y: **The Dynamics**  
389 **of Genetic Interactions between *Vibrio metoecus* and *Vibrio cholerae*, Two Close**  
390 **Relatives Co-Occurring in the Environment.** *Genome Biology and Evolution.* 2015, 7:  
391 2941–2954. doi:10.1093/gbe/evv193/-/DC1

- 392 \* A clear demonstration that genes in the integron are recombined between species just as often  
393 as within species, and that 24% of ‘core’ genes have crossed species boundaries.  
394
- 395 23. Boucher Y, Cordero OX, Takemura A, Hunt DE, Schliep K, Baptiste E, et al.: **Local**  
396 **Mobile Gene Pools Rapidly Cross Species Boundaries To Create Endemicity within**  
397 **Global *Vibrio cholerae* Populations.** *mBio*. 2011, 2: e00335–10.  
398 doi:10.1128/mBio.00335-10
- 399 24. Boucher Y, Orata FD, Alam M: **The out-of-the-delta hypothesis: dense human**  
400 **populations in low-lying river deltas served as agents for the evolution of a deadly**  
401 **pathogen.** *Front Microbiol*. 2015, 6: 1120. doi:10.3389/fmicb.2015.01120
- 402 25. Mentzer von A, Connor TR, Wieler LH, Semmler T, Iguchi A, Thomson NR, et al.:  
403 **Identification of enterotoxigenic *Escherichia coli* (ETEC) clades with long-term**  
404 **global distribution.** *Nature Genetics*. 2014, 46: 1321–1326. doi:10.1038/ng.3145
- 405 26. Cao Q, Didelot X, Wu Z, Li Z, He L, Li Y, et al.: **Progressive genomic convergence of**  
406 **two *Helicobacter pylori* strains during mixed infection of a patient with chronic**  
407 **gastritis.** *Gut*. 2015, 64: 554–561. doi:10.1136/gutjnl-2014-307345
- 408 27. Vos M, Didelot X: **A comparison of homologous recombination rates in bacteria and**  
409 **archaea.** *The ISME journal*. 2009, 3: 199–208. doi:10.1038/ismej.2008.93
- 410 28. Baines SL, Holt KE, Schultz MB, Seemann T, Howden BO, Jensen SO, et al.:  
411 **Convergent adaptation in the dominant global hospital clone ST239 of methicillin-**  
412 **resistant *Staphylococcus aureus*.** *mBio*. 2015, 6: e00080. doi:10.1128/mBio.00080-15
- 413 29. Hershberg R, Lipatov M, Small PM, Sheffer H, Niemann S, Homolka S, et al.: **High**  
414 **Functional Diversity in *Mycobacterium tuberculosis* Driven by Genetic Drift and**  
415 **Human Demography.** *PLOS Biol*. 2008, 6: e311.  
416 doi:10.1371/journal.pbio.0060311.st006
- 417 30. Namouchi A, Didelot X, Schock U, Gicquel B, Rocha EPC: **After the bottleneck:**  
418 **Genome-wide diversification of the *Mycobacterium tuberculosis* complex by mutation,**  
419 **recombination, and natural selection.** *Genome Research*. 2012, 22: 721–734.  
420 doi:10.1101/gr.129544.111
- 421 31. Casali N, Nikolayevskyy V, Balabanova Y, Harris SR, Ignatyeva O, Kontsevaya I, et al.:  
422 **Evolution and transmission of drug-resistant tuberculosis in a Russian population.**  
423 *Nature*. 2014, 46: 279–286. doi:10.1038/ng.2878
- 424 32. Shapiro BJ, David LA, Friedman J, Alm EJ: **Looking for Darwin's footprints in the**  
425 **microbial world.** *Trends in Microbiology*. 2009, 17: 196–204.  
426 doi:10.1016/j.tim.2009.02.002
- 427 33. Friedman J, Alm EJ, Shapiro BJ: **Sympatric Speciation: When Is It Possible in**  
428 **Bacteria?** *PLOS ONE*. 2013, 8: e53539. doi:10.1371/journal.pone.0053539.g004
- 429 34. Cordero OX, Polz MF: **Explaining microbial genomic diversity in light of evolutionary**  
430 **ecology.** *Nature Reviews Microbiology*. 2014, 12: 263–273. doi:10.1038/nrmicro3218

- 431 35. Cadillo-Quiroz H, Didelot X, Held NL, Herrera A, Darling A, Reno ML, et al.: **Patterns**  
432 **of Gene Flow Define Species of Thermophilic Archaea.** *PLOS Biol.* 2012, 10:  
433 e1001265. doi:10.1371/journal.pbio.1001265.t001
- 434 36. Croucher NJ, Didelot X: **The application of genomics to tracing bacterial pathogen**  
435 **transmission.** *Current Opinion in Microbiology.* 2015, 23: 62–67.  
436 doi:10.1016/j.mib.2014.11.004
- 437 37. Bendall ML, Stevens SL, Chan L-K, Malfatti S, Schwientek P, Tremblay J, et al.:  
438 **Genome-wide selective sweeps and gene-specific sweeps in natural bacterial**  
439 **populations.** *The ISME journal.* 2016, advance online publication, 8 January 2016;  
440 doi:10.1038/ismej.2015.241
- 441 \*\* A pioneering application of time-course metagenomics to infer genome-wide and gene-  
442 specific selective sweeps in natural lake bacterial populations.  
443
- 444 38. Kingman J: **On the genealogy of large populations.** *J Appl Prob.* 1982, 19: 27–43.
- 445 39. Cohan FM: **Bacterial Speciation: Genetic Sweeps in Bacterial Species.** *Curr Biol.* 2016,  
446 26: R112–R115. doi:10.1016/j.cub.2015.10.022
- 447 40. Dixit PD, Pang TY, Studier FW, Maslov S: **Recombinant transfer in the basic genome**  
448 **of *Escherichia coli*.** *Proc Natl Acad Sci USA.* 2015, 112: 9070–9075.  
449 doi:10.1073/pnas.1510839112
- 450 41. Achtman M, Wagner M: **Microbial diversity and the genetic nature of microbial**  
451 **species.** *Nature Reviews Microbiology.* 2008, 6: 431–440. doi:10.1038/nrmicro1872
- 452 42. Snitkin ES, Zelazny AM, Thomas PJ, Stock F, NISC Comparative Sequencing Program,  
453 Henderson DK, et al.: **Tracking a Hospital Outbreak of Carbapenem-Resistant**  
454 ***Klebsiella pneumoniae* with Whole-Genome Sequencing.** *Science Translational*  
455 *Medicine.* 2012, 4: 148ra116–148ra116. doi:10.1126/scitranslmed.3004129
- 456 43. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, et al.: **Out-of-Africa**  
457 **migration and Neolithic coexpansion of *Mycobacterium tuberculosis* with modern**  
458 **humans.** *Nature Genetics.* 2013, 45: 1176–1182. doi:10.1038/ng.2744
- 459 44. Croucher NJ, Finkelstein JA, Pelton SI, Mitchell PK, Lee GM, Parkhill J, et al.:  
460 **Population genomics of post-vaccine changes in pneumococcal epidemiology.** *Nature*  
461 *Genetics.* 2013, 45: 656–663. doi:10.1038/ng.2625
- 462 45. Chewapreecha C, Harris SR, Croucher NJ, Turner C, Marttinen P, Cheng L, et al.: **Dense**  
463 **genomic sampling identifies highways of pneumococcal recombination.** *Nature*  
464 *Genetics.* 2014, 46: 305–309. doi:10.1038/ng.2895
- 465 46. Marttinen P, Croucher NJ, Gutmann MU, Corander J, Hanage WP: **Recombination**  
466 **produces coherent bacterial species clusters in both core and accessory genomes.**  
467 *Microbial Genomics.* 2015, 1. doi:10.1099/mgen.0.000038.

468 \* The first simulation of bacterial evolution including mutation, allelic exchange (homologous  
469 recombination) and gene gain/loss (non-homologous recombination).  
470

471 47. Rocha EP, Smith JM, Hurst LD, Holden MT, Cooper JE, Smith NH, et al.: **Comparisons**  
472 **of dN/dS are time dependent for closely related bacterial genomes.** *J Theor Biol.* 2006,  
473 239: 226–235.

474 48. Kryazhimskiy S, Plotkin JB: **The Population Genetics of dN/dS.** *PLOS Genet.* 2008, 4:  
475 e1000304.

476

Figure 1

(a) Genome-wide sweep

(b) Gene-specific sweeps

Legend

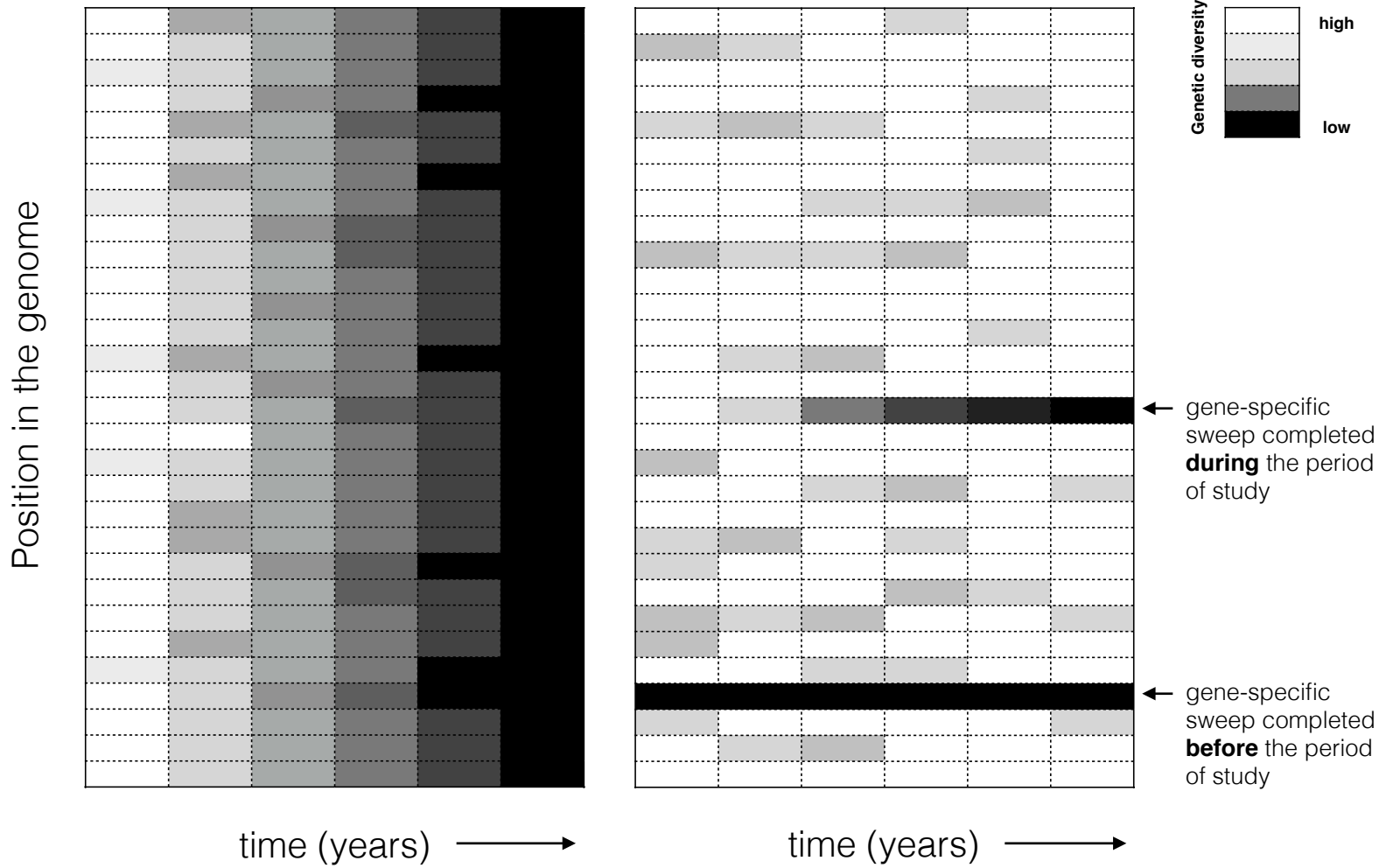




Figure 2

